

→ crashes can occur while

- the T is executing the logtail update, &
- while recovery action is being taken.

$$A = 1000\$, \quad B = 2000\$, \quad r = 200\$$$

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

$\langle T_0, A, 950 \rangle$

$\langle T_0, B, 2050 \rangle$

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

$\langle T_0, A, 950 \rangle$

$\langle T_0, B, 2050 \rangle$

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

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

$\langle T_1, C, 600 \rangle$

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

$\langle T_0, A, 950 \rangle$

$\langle T_0, B, 2050 \rangle$

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

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

$\langle T_1, C, 600 \rangle$

$\langle T_1, \text{commit} \rangle$

(a) no redo actions need to be taken.

(b) since, no commit record appears in log.

The log records of the incomplete transaction T_0 can be deleted from the log.

(c) redo(T_0) must be performed since $\langle T_0, \text{commit} \rangle$ is present.

also the log records of the incomplete transaction T_1 can be deleted from the log.

(d) if the crash occurs just after the log record $\langle T_1, \text{commit} \rangle$ is written to stable storage,

when the system comes back up, two commit records are in the log:

\therefore so system must perform redo(T_0) & redo(T_1) in the order in which their commit records appear in the log.

Immediate Database Modification:

It allows database update of an uncommitted transaction.

To be made as the writes are issued.

- In the event of a transaction failure, the system must use the old-value field of the log record.
~~No~~ undo operation required.
- Before T_i starts its execution, the system writes $\langle T_i \text{ start} \rangle$ to the log.
- During its execution, any write(x) operation by T_i is preceded by the writing of the new update record to the log.
- When T_i partially commits, the system writes the record $\langle T_i \text{ commit} \rangle$ to the log.
- Update log record must be written before database item is written
 - we assume that the log record is output directly to stable storage.
 - before execution of an output(B) operator, the log records corresponding to B be written onto 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.

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

$\langle T_0, A, 1000, 950 \rangle$

$\langle T_0, B, 2000, 2050 \rangle$

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

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

$\langle T_1, C, 700, 600 \rangle$

log

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

$\langle T_0, A, 1000, 950 \rangle$

$\langle T_0, B, 2000, 2050 \rangle$

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

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

$$\begin{aligned} A &= 950 \\ B &= 2050 \end{aligned}$$

- Recovery actions are
- (a) $\text{undo}(T_0)$: B is restored to 2000 And A to 1000.
 - (b) $\text{undo}(T_1)$ and $\text{redo}(T_0)$; C is set to 700, & A & B are set to 950 and 2050.
 - (c) $\text{redo}(T_0)$ and $\text{redo}(T_1)$.

Checkpoints:

When a system failure occurs, we must consult the log to determine those transactions that need to be redone and those that need to be undone.

→ We need to search the entire log to determine this info.

problems in recovery procedure ~~are~~ are

- ① Searching the entire log is time-consuming.
- ② We might unnecessarily redo transactions which have already ~~done~~ output their updates to the database.
- ③

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.

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 .

→ Using the log, the system can handle any failure that does not result in the loss of info. in nonvolatile storage.

Recovery procedure has two parts

- $\text{undo}(T_i)$ - restores the value of all data items updated by transaction T_i to the old values, going backwards
- $\text{redo}(T_i)$ - sets the value of all data items updated by transaction T_i to the new values, going forward
- Both operations must be idempotent.

even if the operation is executed multiple times, the effect is the same as if it is executed once

After a failure has occurred, the recovery scheme consults the log to determine which transactions need to be redone and which need to be undone,

- 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 $\langle T_i \text{ commit} \rangle$.
- 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

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

$\langle T_0, A, 1000, 950 \rangle$

$\langle T_0, B, 2000, 2050 \rangle$

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

$\langle T_0, A, 1000, 950 \rangle$

$\langle T_0, B, 2000, 2050 \rangle$

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

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

$\langle T_1, C, 700, 600 \rangle$

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

$\langle T_0, A, 1000, 950 \rangle$

$\langle T_0, B, 2000, 2050 \rangle$

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

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

$\langle T_1, C, 700, 600 \rangle$

$\langle T_1, \text{commit} \rangle$

(a)

1.7

→ Scan

ARIES: Crash Recovery algorithm

ARIES

- ① It uses a log sequence number (LSN) to identify log records, & the use of LSNs in database pages to identify which operations have been applied to a database page.

② LSN = logical address of record in the log

Page LSN : stored in page

→ LSN of most recent update to page

RevolSN : stored in log record

dirty page: buffer version has updated but yet reflected on disk

→ T's add log records.

→ checkpoints are performed periodically.

contains → active transaction list,

→ LSN of most recent log records of T, and

→ list of dirty pages in the buffer.

↳ to determine where redo should start
check

Recovery phases:

- Analysis phase → forward from last checkpoint ↓
- Redo pass → forward from RedoLSN, which is determined in analysis pass.
- Undo pass → backwards from end of log, undoing incomplete transactions.

Analysis pass:

RedoLSN = min(LSNs of dirty pages recorded in checkpoint)

→ if no dirty pages, RedoLSN = LSN of checkpoint

→ Scan log forwards from last checkpoint.

→ find transaction to be rolled back.

ARIES: A Transaction Recovery Method

Supporting fine granularity Locking and partial Rollback using write-ahead Logging

ARIES: (Algorithm for Recovery and Isolation Exploiting Semantics)

Write-ahead logging.

① T_i enters the commit stage after the $\langle T_i \text{ commit} \rangle$ log record has been output to stable storage.

② Before the $\langle T_i \text{ commit} \rangle$ log record can be output to stable storage, all log records pertaining to T_i must have been output to stable storage.

③ Before a block of data in memory can be op'ed to db, all log records pertaining to data in that block must have been op'ed to stable storage.

AR
Recovery with concurrent transactions.

i) Interaction with Concurrency Control.

ii) Transaction Rollback - using log. (immediate ~~log~~ db modifications)

iii) checkpoints → we require that the checkpoint log record be of the form $\langle \text{checkpoint- } L \rangle$, L is a list of transaction active at the time of the checkpoint.

→ no updates are performed while the checkpoint is in progress.

iv) fuzzy checkpoint: is a check-point allowed to perform updates even while buffer blocks are being written out.

v) Restart Recovery:

Buffer management:

Log-Record Buffering:

- Every log record is o/p to stable storage at the time it is created : It imposes a high overhead on system execution.
- O/p to stable storage is in units of blocks.
i.e. a log record is much smaller than a block.
- write log records temporarily in the log buffer in ^{main} memory.
- Multiple o/p to stable storage in a single op operation.
- we must impose additional requirements on the recovery technique to ensure transaction atomicity.

- ① Transact. T_i enters the commit record has been output to stable storage' after the $\langle T_i \text{ commits} \rangle$ log.
- ② Before the $\langle T_i \text{ Commit} \rangle$ log record can be o/p to stable storage, all log records pertaining to Transact. T_i must have been o/p to stable storage.
- ③ Before a block of data in main memory can be o/p to the db, all log records pertaining to data in that block must have been o/p to stable storage.

This rule is called

Write-ahead Logging (WAL)

- Writing the partially filled buffer block log to disk is referred to as a log file.

for each transaction T_i in L , if T_i is not in redo-list then
it adds T_i to the undo-list.

The recovery proceeds as follows

- ① The system scans the log from the most recent backward,
and performs an "undo" for each log record that belongs to T_i
on the undo-list, redo-list is ignored in this phase.
The scan stops when $\langle T_i \text{ start} \rangle$ records have been found for
every transaction T_i in the undo-list.

- ② The system locates the most recent <checkpoint $L \rightarrow$ record on log
- ③ The system scans the log forward from the most recent
<checkpoint $L \rightarrow$ record and performs redo for each log record
that belongs to T_i that is on the redo-list.

→ After the system has undone all transactions on the undo-list,
it redoes those transaction in the redo-list.

→ It is impl. to undo the transaction in the undo-list before
redoing transaction in the redo-list, otherwise a problem may occur
Suppose $A = 10$.

T_1 update A to 20 and aborted. Transl. roll back would set A to 10.

T_2 update A to 30 and committed, the system crashed
The state of the log at the time of the crash is

$\langle T_1, A, 10, 20 \rangle$

$\langle T_2, A, 10, 30 \rangle$

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

redo pass is performed first A is set to 30,
in the undo pass, A is set to 10, which is wrong

Advanced Recovery techniques

Restart Recovery:

Recovery actions:

- 1) In the redo phase, the system replays updates of all transactions by scanning the log forward from the last checkpoint. The records are of the form $\langle T_i, X_j, v_1, v_2 \rangle$ & $\langle T_i, X_j, v_1 \rangle$
- 2) In the undo phase, the system rolls back all transactions in the undo-list. It performs rollback by scanning the log backward from the end.

ARIES: (Algorithm for Recovery and Isolation Exploiting Semantics)

It is used to reduce the time taken for recovery, and to reduce the overheads of checkpointing.

ARIES algo:

- It uses a log sequence number (LSN) to identify log records.
- It supports physiological redo operations.
- It uses a dirty page table to minimize unnecessary redos during recovery. Dirty pages are those that have been updated in memory, and the disk is not up-to-date.
- Each log record in ARIES has a log sequence number (LSN).
- Each page maintains an identifier called the page LSN, whenever an operation occurs on a page, the page LSN field stores the LSN of its log record.
- During the redo phase of recovery, any log records with LSN less than or equal to the page LSN of a page should not be executed on the page. Since their actions are already reflected on the page. So ARIES can avoid reading many pages. So recovery time is reduced significantly.

→ each log record also contains the LSN of the previous log record of the same transaction in prevLSN field, it permits log records of trans. to be fetched backward.

Compensation log records (CLR's) - have a field, called the undoNextLsn, those log records that needs to be undone next,

→ the Dirty page table - contains the list of pages that have been updated in the buffer. whenever a page is flushed to disk, the page is removed from the Dirty page table.

Recovery algorithm: ARIES recovery from a system crash in 3 phases.

- 1) analysis phase: - This phase determines which transactions to be undo, which pages were dirty at the time of crash, and the LSN from which the redo pass should start.
- 2) redo pass: It performs redo a transaction to bring the db

Database buffering:

(18)

- main memory is much smaller than the entire db.
- It may be necessary to overwrite a block B_1 in main memory, when another block B_2 needs to be brought into memory.
So B_1 must be output prior to the input of B_2 .
- So all log records pertaining to data in B_1 must be output to stable storage before B_1 is o/p. Thus, the sequence of actions by the system would be:
 - output log records to stable storage until all log records pertaining to block B_1 have been output.
 - output block B_1 to disk.
 - Input block B_2 from disk to main memory.
- no writes to the block B_1 be in progress while the system carries out this sequence of actions.
- So acquire an exclusive lock on the block; so that no one is updating the block. locks are released once the block o/p has completed. These locks are called as "Latches".
 - <To Start>
 - <To, A, 1000, 950>.

To issues a read(B), Assume that the block on which B resides is not in main memory, and that main memory is full.
so A's records block is chosen to be o/p to disk.
Before output that block to disk, the log record must be output to stable storage. The system can use the log record during recovery to bring the db back to a consistent state.

failure with loss of nonvolatile storage.

Basic scheme is to dump the entire contents of the db to stable storage periodically.

The system uses the most recent dump in restoring the db to previous consistent state.

→ no transaction may be active during the dump procedure, and a procedure must take place.

1) output all log records currently residing in main memory

2) O/P all buffer blocks onto the disk.

3) Copy the contents of the db to stable storage.

4) output a log record (dump) onto the stable storage.

→ no undo operations need to be executed.

→ a dump of the db contents is also referred to as an archival dump.

dump procedure is very costly, because.

→ entire db must be copied to stable storage, resulting in considerable data transfer.

→ transaction processing is halted during the dump procedure so CPU cycles are wasted.

Storage Structure:

- Storage types:
- Volatile storage: Info. residing in volatile storage does not usually survive system crashes. Ex. main memory & cache memory.
→ fast accessing.
- Nonvolatile storage: → Info. residing in nonvolatile storage survives system crashes. Ex. disk and magnetic tapes.
- Stable storage: Info. residing in stable storage is never lost.

Stable storage implementation:

We need to replicate the needed info. in several non-volatile storage media with independent failure modes, and to update the info. in a controlled manner to ensure that failure during data transfer does not damage the needed info.

Storage media can be protected from failure during data transfer. Block transfer b/w memory and disk storage can result in.

- Successful completion: The transferred info. arrived safely at its destination.

- Partial failure: A failure occurred in the midst of transfer, and the dest. block has incorrect info.

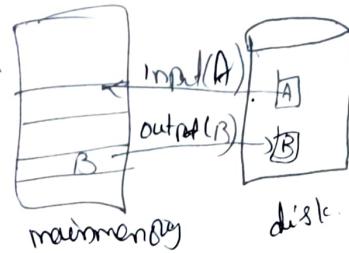
- Total failure: The failure occurred early during the transfer, so that the dest. block remains intact.

If a data-transfer-failure occurs, the system detects it and invokes a recovery procedure to restore the block to a consistent state.

~~Total~~

Data Access

- The db system resides permanently on nonvolatile storage and is partitioned into fixed-length storage units called 'blocks'.
Blocks are the units of data transfer to and from disk, and may contain several data items.
- Transaction i/p intel. from the disk to main memory, & then o/p the intel. back onto the disk.
- the i/p and o/p operations are done in block units.
- the blocks residing on the disk are referred to as physical block.
- " " " main memory " " buffer block
- i) input(B) transfers the physical block B to main memory
ii) output(B) " " buffer block B to the disk, and replaces the appropriate physical block there.
- Each transaction T_i has a private work area in which copies of all the data items accessed and updated by T_i are kept.
Each data item X kept in the work area of T_i is denoted by x_i .
 T_i interacts with the db system by transferring data to and from its work area to the system buffer.
We transfer data by these two operations:
- i) $\text{read}(x)$ assigns the value of data item X to the local variable x_i
 - (a) If block B_x on which X resides is not in main memory, it issues $\text{input}(B_x)$,
 - (b) It assigns to x_i the value of X from buffer block.
 - ii) $\text{write}(x)$ assigns the value of local variable x_i to data item x in the buffer block



Recovery with Concurrent Transactions:

→ we roll back a failed transaction, T_i by using the log,

The system scans the log backward,

for every log record of the form $\langle T_i, x_j, v_1, v_2 \rangle$ found in the log,
the system restores the data item x_j to its old value v_1 .

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

$\langle T_i, A, 20, 30 \rangle$

Checkpoints:

→ checkpoints are used to reduce the no. of log records
that the system must scan when it recovers from a crash.

Consider the following trans./s during recovery

→ those trans./s that started after the most recent checkpoint.

→ the one trans./s, if any, that was active at the time of the
most recent checkpoint.

$\langle \text{checkpoint } L \rangle$, L - is list of trans./s active at the time
of checkpoint.

In fuzzy checkpoint, trans./s are allowed to perform updates
even while buffer blocks are being written out.

Restart Recovery:

→ undo list - consists of trans./s to be undone

→ redo list - " " " redone.

the system scans the log backward, examining each record,
until it finds the first $\langle \text{checkpoint} \rangle$ record.

→ for each record found of the form $\langle T_i \text{ commit} \rangle$, it adds T_i to redo-list
" " " " " " " " $\langle T_i \text{ start} \rangle$, if T_i is not in redo-list
then it adds T_i to undo-list.

- If block B_x on which 'x' resides is not in main memory,
- it issues input (B_x) .
 - (b) it assigns the value of x_i to X in Buffer B_x .

Recovery and Atomicity:

Recovery with concurrent Transactions:

check points:

transl.s are not allowed to perform any update action, such as writing to a buffer block & writing a log record, while a checkpoint is in progress.

- The $\langle \text{checkpoint} \rangle$ record in the log allows the system to streamline its recovery procedure.
- for a transl., the $\langle T_i \text{ commit} \rangle$ record appears in the log below the $\langle \text{checkpoint} \rangle$ record. Thus, at recovery time, there is no need to perform a redo operation on T_i .
- After a failure has occurred, the recovery scheme examines the log to determine the most recent transaction T_i that started executing before the most recent checkpoint took place, by searching the backward, from the end of the log, until it finds the first $\langle \text{checkpoint} \rangle$ record, then it continues the search backward until it finds the next $\langle T_i \text{ start} \rangle$ record.

once the system has identified transl. T_i ,

- for all transactions T_k in T that have no $\langle T_k \text{ commit} \rangle$ record in the log, execute $\text{Undo}(T_k)$
- " .. " have $\langle T_k \text{ commit} \rangle$ record appears in the log, execute $\text{Redo}(T_k)$.