

UNIT-5

* Deadlock:

" is a condition where two or more transactions are waiting indefinitely for one another to give up locks.

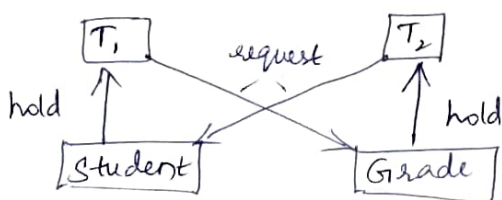
- For example, In the student table, transaction T_1 holds lock on some rows & in grade table.

↓
needs to update some rows

Simultaneously T_2 transaction holds lock on some rows in grade table & needs to update the rows in the student table held by T_1 .

- So, here main problem is, T_1 is waiting for T_2 to release its lock & similarly T_2 is waiting for T_1 to release its lock.

Here DBMS detects the deadlock & aborts one of the transaction.

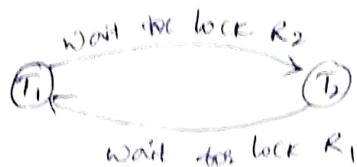


* Deadlock detection:

In a DB, when a transaction waits indefinitely to obtain a lock then the DBMS should detect whether the transaction is involved in deadlock or not. The lock manager maintains a WFG (Wait For Graph) to detect the deadlock cycle in the DB.

Wait for Graph:

- A graph is created based on the transaction & their lock.
- If the created graph has a cycle or closed loop, then there is a deadlock.
- WFG maintained by the system for every transaction which is waiting for some data held by the others. The system keeps checking the graph if there is any cycle in the graph.



* Deadlock - Avoidance:

- When a DB is stuck in a deadlock state, then it is better to avoid the DB rather than aborting or restarting the DB. It will waste of time & resources.
- Deadlock avoidance mechanism is used to detect any deadlock situation in advance.
- wait for graph is used for detecting the deadlock situation but this method is suitable only for smaller DB. For larger DB, deadlock prevention method can be used.

* Deadlock prevention:

- " " is suitable for large DB.
- If the resources are allocated in such a way that deadlock never occurs, then the deadlock can be prevented.
- DBMS analyzes the operations of the transaction whether they can create a deadlock situation or not. If it forms a deadlock, then DBMS never allows that transaction to be executed.

Failure classification:

②

To find that where the problem has occurred, we generalize a failure into the following categories:

(i) Transaction failure

(ii) System crash

(iii) Disk failure

(i) Transaction failure:

It occurs when it fails to execute or when it reaches a point from where it can't go any further. If a transaction or process is failed then that is called as transaction failure.

Reasons for transaction failure are:

Logical error: If a transaction can't complete its execution due to code error or internal logic error.

Syntax error: If DBMS itself terminates an active transaction because the DB system is not able to execute it.

Eg: System aborts its active transaction in case of deadlock or resource unavailability.

(ii) System crash:

System failure can occur due to power failure or other h/w or s/w failure.

Fail stop assumption: In system crash, non-volatile storage is assumed not to be corrupted.

(iii) Disk failure:

- It occurs where hard disk drives or storage devices fail frequently.

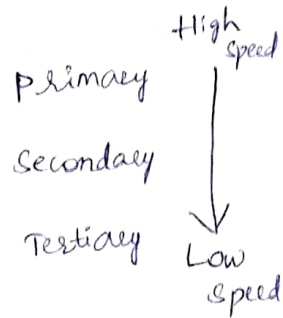
It was the common problem in early days of technology evolution.

- Disk failure occurs due to the formation of bad sectors, disk head crash & unreachability to the disk or any other failure, which destroys the disk storage.

* Storage Structure:

(i). Primary Storage:

- " " is used to access the data quickly.
- Volatile storage (ie, It will not store the data permanently). In any system failure or power cut, the data will be lost.
- Eg: main memory & cache memory

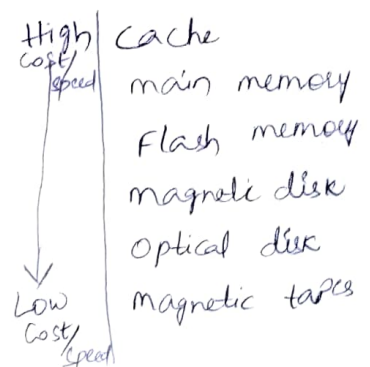


(ii). Secondary Storage:

- " " also called as online storage. It allows the user to store the data permanently.
- Non volatile memory (In power failure or system crash) data will be there. There won't be any loss of data.
- Eg: Flash memory, magnetic disk storage (USB)

(iii). Tertiary Storage:

- " " is external from the computer system. It has slowest speed. But it is capable of storing large amount of data. It is known as offline storage.
- It is used for data backup
- Eg: Optical storage, Tape storage (CD, DVD)



* Recovery and Atomicity:

(2)

When a system crashes, it may have several transactions being executed & various files opened for them to modify the data items.

- Transactions are made of various operations, which are atomic in nature.
- When DBMS recovers from a crash, it should maintain the following.
 - (i). It should check the states of all the transactions, which were being executed.
 - (ii). A transaction may be in the middle of some operation, which were being executed. DBMS must ensure the atomicity of transaction in this case.
 - (iii). It should check whether the transaction can be completed now or it needs to be rolled back.
 - (iv). No transactions would be allowed to leave DBMS in an inconsistent state.

There are two types of techniques, which can help DBMS in recovering as well as maintaining atomicity of a transaction.

- (i). Maintaining the logs of each transaction & writing them onto some stable storage before actually modifying the DB.
- (ii). Maintaining shadow paging, where changes are done on a volatile memory & later, actual DB is updated.

- The process of restoring the DB to a correct state in the event of a failure is known as DB recovery.
- The DB has to be restored back to consistent state from its inconsistent state that has been caused due to failure.
- Two techniques are used for recovery from transaction failure
 - (i). Deferred update (No UNDO/REDO)
 - (ii). Immediate update (UNDO/REDO)

* Log based Recovery:

- Log is a sequence of records. Log of each transaction is maintained in some stable storage so that if any failure occurs, then it can be recovered.
- If any operation performed on the DB, then it will be recorded in the log. But the process of storing the logs should be done before the actual transaction is applied in the DB.
- Let's assume there is a transaction to modify the city of a student.
 - (i). When the transaction is initiated, then it writes start log
 $\langle T_n, \text{start} \rangle$
 - (ii). When transaction modifies the city from 'Hyd' to 'Delhi'
 $\langle T_n, \text{City}, 'Hyd', 'Delhi' \rangle$
 - (iii). When transaction is finished, then it writes another log to indicate the end of the transaction.
 $\langle T_n, \text{commit} \rangle$

Deferred DB modification:

- " " technique occurs if transaction does not modify the DB until it has committed.
- All logs are created & stored in stable storage & DB is updated when a transaction commits.
- It modifies the DB after completion of transaction.

Recover system uses following operation:

$\text{Redo}(T_i)$: All data items updated by T_i are set to new value

Immediate DB modification:

(4)

- occurs if DB modification occurs while the transaction is still active
- DB is modified immediately after every operation it follows an actual DB modification.
- Undo(T_i): All data items updated by T_i are set to old value.
- Redo(T_i): " " " " " new value.

* When the system is crashed, then the system consults the log to find which transaction needs to be undone & which needs to be redone.

(i). If the log contains the record $\langle T_i, \text{start} \rangle$ & $\langle T_i, \text{commit} \rangle$ (or)

$\langle T_i, \text{commit} \rangle$ then T_i needs to be redone.

(ii). If the log contains the record $\langle T_i, \text{start} \rangle$ but doesn't contain the

record either $\langle T_i, \text{commit} \rangle$ or $\langle T_i, \text{abort} \rangle$ then T_i needs to be undone.

* Recovery with Concurrent transactions:

- Concurrency control means multiple transactions can be executed at the same time & then interleaved logs occur. But, there may be changes in transaction results so maintain the order of execution of those transactions.
- During recovery, it would be very difficult for the recovery system to rollback all the logs & then start recovering.
- Recovery with concurrent transactions can be done in following four ways.
 - (i) Interaction with Concurrency control
 - (ii) Transaction rollback
 - (iii) Checkpoints
 - (iv) Restart recovery

(i). Interaction with Concurrency Control:

- The recovery depends on concurrency control scheme that is used.
- So, to rollback a failed transaction, we must undo the updates performed by the transaction.

(ii). Transaction rollback:

- Here, we will rollback a failed transaction by using the log.
- System scans the log backward a failed transaction, for every log record found in the log the system restores the data item.

(iii). Checkpoints:

- Checkpoints is a process of saving a snapshot of the applications state so that it can restart from that point in case of failure.
- checkpoint is a point of time at which a record is written onto the DB from the buffers.
- Checkpoint shortens the recovery process.
- when it reaches the checkpoint, then the transaction will be updated into the DB. & till that point, the entire log file will be removed from the file. Then the log file is updated with the new step of transaction till the next checkpoint & so on.
- checkpoint is used to declare the point before which DBMS was in the consistent state & all the transactions were committed.
- we can use checkpoints to reduce the no. of log records that the system must scan when it recovers from a crash.

(iv) Restart recovery: When system recovers from crash, it consists of 2 lists

- Undo list contains transactions to be undone.
- Redo " " " " redone.
- Initially, both are empty. The system scans the log backward & examines each record

until it finds the first checkpoint record.

(5)

- Keeping & maintaining logs in real-time is difficult. As time passes, the log file will grow too big. ~~to~~ Checkpoint is a mechanism where all previous logs are removed from system & stored permanently in a storage disk.
- checkpoint declares a point before which DBMS was in consistent state & all transactions were committed.
- If the recovery system sees a log with $\langle T_n, \text{start} \rangle$ & $\langle T_n, \text{Commit} \rangle$ or just $\langle T_n, \text{Commit} \rangle$, it puts the transaction in the redo-list.
- If log with $\langle T_n, \text{start} \rangle$ but no commit, it will be going into undo-list.
- All transactions in undo-list are then undone & their logs are removed.
- " " redo-list & their previous logs are removed & then redone before saving their logs.

* Buffer management:

- DB buffer is a temporary storage area in the main memory. It allows storing the data temporarily when moving them from one place to another.
- DB buffer stores a copy of disk blocks.
- Buffer manager is responsible for allocating space to the buffer in order to store data into the buffer.
- If a user request a particular block & the block is available in the buffer, the buffer mgr provides block address in the main memory.
- If the block is not available in the buffer, buffer mgr allocates the block in the buffer.
- If free space is not available, it throws out some existing blocks from

the buffer to allocate the required space for the new block.

- The blocks which are thrown are written back to the disk only if they are recently modified when writing on the disk.
- If the user requests such thrown out blocks, buffer mgr reads the requested block from disk to the buffer & then passes the address of the " " to the user in the main memory.
- The internal actions of the buffer mgr are not visible to the program that may create any problem in disk block requests. The buffer mgr is just like a virtual machine.

(i). Buffer replacement strategy:

If no space is left in the buffer, it is required to remove an existing block from the buffer before allocating the new one.

The various OS uses the LRU (Least Recently Used) scheme.

- In LRU, the block that was least recently used is removed from the buffer & written back to the disk. This type of replacement strategy is known as buffer replacement strategy.

(ii). Pinned blocks:

- If user wants to recover any DB system from the crashes, it is essential to restrict the time when a block is written back to the disk.
- Most recovery systems do not allow the blocks to be written on the disk if the block updation is in progress. Such type of blocks that are not allowed to be written on the disk are known as pinned blocks. Many OS don't support pinned blocks.

(iii) Forced op of blocks:

(6)

- In some cases, it becomes necessary to write the block back to disk even though the space occupied by block in the buffer is not required. When such type of write is required. It is known as forced op of a block.
- It is because sometimes the data stored on the buffer may get lost in some system crashes, but the data stored on disk usually doesn't get affected due to any disk crash.

* Failure with loss of non volatile storage:

- The basic measure is to dump the entire contents of the DB to stable storage periodically.
- One approach to dump the DB requires that "no transaction is active during the dumping procedure" & uses a procedure similar to checkpointing.
 - (i). O/p of all the log records currently present in the main memory into the stable storage.
 - (ii). O/p all the buffer blocks into the disk.
 - (iii). Copy all the data present in the DB to the stable storage.
 - (iv). O/p a log record <dump> into the stable storage.
- To recover from the loss of non volatile memory, we restore the DB from the archive & all the transactions that have been committed since the most recent dump are redone. This is known as archival dump.

↓
(Dump of DB contents)
- Dumps of DB & checkpointing are very similar

* ARIES Recovery method: (Algorithm for Recovery & Isolation Exploiting Semantics)

- ARIES is based on the write Ahead log (WAL) protocol.
- Every update operation writes a log record which is one of the following:
 - (i) Undo-only log record:
only the before image is logged. So, undo operation can be done to retrieve the old data.
 - (ii) Redo-only log record:
only after image is logged. So, redo can be attempted.
 - (iii) Undo-redo log record:
Both before & after images are logged.
- Every log record is assigned a unique log sequence number (LSN).
- Every data page has a page LSN field that is set to LSN of the log record corresponding to the last update on the page.
- WAL requires that the log record corresponding to an update make it to stable storage before the data page corresponding to that update is written to disk.
- For performance reasons, each log write is not immediately forced to disk.
- A log tail is maintained in main memory to buffer log writes. The log tail is flushed to disk when it gets full.
- A transaction can't be declared committed until the commit log record makes it to disk.
- Once in a while the recovery subsystem writes a checkpoint record to log.
- The checkpoint record contains the transaction table & dirty page table.
- A master log record is maintained separately. In stable storage, to store the LSN of latest checkpoint record that made it to disk.
- On restart, the recovery subsystem reads the master log record to find checkpoints LSN, starts recovery.

The recovery process consists of 3 phases:

(7)

i) Analysis.

The recovery subsystem determines the earliest log record from which the next pass must start. It also scans the log forward from the checkpoint record to construct a snapshot of what the system looked like at the instant of the crash.

(ii) Redo:

Starting at the earliest LSN, the log is read forward & each update redone.

(iii) Undo:

The log is scanned backward & updates corresponding to lost transactions are undone.

- ARIES is a no-lose type of backup approach i.e., if transaction completes updated values are noted in the disk at that time, it takes some time.

- Recovery mgr called at time of crash.

- In analysis phase identify the dirty pages (not updated or committed) in the buffers & set active transaction at the time of failure.

(P_1, P_3, P_5) are active so needs to restart
 $T_1, T_4 \& T_3$

T_1 - write P_5

T_2 - write P_3

T_2 - commit

T_3 - write P_1

T_4 - write P_3

crash

- In Redo, find all operations done by DBMS before crash & restore the system back to the same state before crash

It aborts all active transactions that are still running. T_1 & T_3 write Operat. written again in order.

- In undo, scans the log backward of active transaction in reverse

T_u with write P_3 undo

T_3 " " P_1 "

T_1 " " P_5 "

Advantages:

- Simple & flexible
- Supports concurrency control protocols
- we can recover the pages independently.

* Remote backup Systems:

" " " provide a wide range of availability, allowing the transaction processing to continue even if the primary site is destroyed by fire, flood or earthquake.

- Data & log records from primary site are continuously backed up into a remote backup site.

- The remote site is also called as "secondary site"

- Remote backup provides security in case primary gets destroyed.

- " " can be offline or online or real-time. If it is offline,

It is maintained manually.

- online backup systems are more real-time & lifesavers for DB admins. Here, every bit of real-time data is backed up simultaneously at two distant places. One is directly connected to system. Other one kept as backup at remote place.

- while designing a remote backup system, following are important

(i) Detection of failure: It is important to detect when primary failed.

(ii). Transfer of control: when primary site fails, backup site takes over the processing & becomes the new primary system. ⑧

(iii). Time to recover:

If the log at remote backup becomes large, recovery will take long time.

(iv). Time to commit:

To ensure that the updates of committed transaction are durable, a transaction should not be announced committed until its log records have reached the backup site.