# Concurrency Control

# DBMS Concurrency Control

Concurrency Control is the management procedure that is required for controlling concurrent execution of the operations that take place on a database.

But before knowing about concurrency control, we should know about concurrent

execution. # Concurrent Execution in DBMS

- o In a multi-user system, multiple users can access and use the same database at one time, which is known as the concurrent execution of the database. It means that the same database is executed simultaneously on a multi-user system by different users.

- o While working on the database transactions, there occurs the requirement of using the database by multiple users for performing different operations, and in that case, concurrent execution of the database is performed.

- o The thing is that the simultaneous execution that is performed should be done in an interleaved manner, and no operation should affect the other executing operations, thus maintaining the consistency of the database. Thus, on making the concurrent execution of the transaction operations, there occur several challenging problems that need to be solved.

# Problems with Concurrent Execution

In a database transaction, the two main operations are **READ** and **WRITE** operations. So, there is a need to manage these two operations in the concurrent execution of the transactions as if these operations are not performed in an interleaved manner, and the data may become inconsistent. So, the following problems occur with the Concurrent Execution of the operations:

## Problem 1: Lost Update Problems (W - W Conflict)

The problem occurs *when two different database transactions perform the read/write operations on the same database items in an interleaved manner (i.e., concurrent execution) that makes the values of the items incorrect hence making the database inconsistent*.

**For example:**

**Consider the below diagram where two transactions $T_X$ and $T_Y$, are performed on the same account A where the balance of account A is $300.**

| Time | Tx | Ty |
|------|------|------|
| $t_1$ | READ (A) | — |
| $t_2$ | A = A - 50 | |
| $t_3$ | — | READ (A) |
| $t_4$ | — | A = A + 100 |
| $t_5$ | — | — |
| $t_6$ | WRITE (A) | — |
| $t_7$ | | WRITE (A) |

LOST UPDATE PROBLEM

- o At time t1, transaction $T_X$ reads the value of account A, i.e., $300 (only read). o At time t2, transaction $T_X$ deducts $50 from account A that becomes $250 (only deducted and not updated/write).

- o Alternately, at time t3, transaction $T_Y$ reads the value of account A that will be $300 only because $T_X$ didn't update the value yet.

- o At time t4, transaction $T_Y$ adds $100 to account A that becomes $400 (only added but not updated/write).

- o At time t6, transaction $T_X$ writes the value of account A that will be updated as $250 only, as $T_Y$ didn't update the value yet.

- o Similarly, at time t7, transaction $T_Y$ writes the values of account A, so it will write as done at time t4 that will be $400. It means the value written by $T_X$ is lost, i.e., $250 is lost.

Hence data becomes incorrect, and database sets to inconsistent.

## Dirty Read Problems (W-R Conflict)

The dirty read problem occurs *when one transaction updates an item of the database, and somehow the transaction fails, and before the data gets rollback, the updated database item is accessed by another transaction. There comes the Read-Write Conflict between both transactions.*

**For example:**
**Consider two transactions $T_X$ and $T_Y$ in the below diagram performing read/write operations on account A where the available balance in account A is $300:**

| Time | Tx | Ty |
|:---:|:---:|:---:|
| $t_1$ | READ (A) | — |
| $t_2$ | A = A + 50 | — |
| $t_3$ | WRITE (A) | — |
| $t_4$ | — | READ (A) |
| $t_5$ | SERVER DOWN ROLLBACK | — |

**DIRTY READ PROBLEM**

o  At time t1, transaction $T_X$ reads the value of account A, i.e., $300.

o  At time t2, transaction $T_X$ adds $50 to account A that becomes $350.

o  At time t3, transaction $T_X$ writes the updated value in account A, i.e., $350.

o  Then at time t4, transaction $T_Y$ reads account A that will be read as $350.

o  Then at time t5, transaction $T_X$ rollbacks due to server problem, and the value changes back to $300 (as initially).

o  But the value for account A remains $350 for transaction $T_Y$ as committed, which is the dirty read and therefore known as the Dirty Read Problem.

# Unrepeatable Read Problem (W-R Conflict)

*Also known as Inconsistent Retrievals Problem that occurs when in a transaction, two different values are read for the same database item.*

**For example:**

**Consider two transactions, $T_X$ and $T_Y$, performing the read/write operations on account A, having an available balance = $300. The diagram is shown below:**

| Time | Tx | Ty |
|------|------|------|
| $t_1$ | READ (A) | — |
| $t_2$ | — | READ (A) |
| $t_3$ | — | A = A + 100 |
| $t_4$ | — | WRITE (A) |
| $t_5$ | READ (A) | — |

**UNREPEATABLE READ PROBLEM**

- At time t1, transaction $T_X$ reads the value from account A, i.e., $300.

- At time t2, transaction $T_Y$ reads the value from account A, i.e., $300.

- At time t3, transaction $T_Y$ updates the value of account A by adding $100 to the available balance, and then it becomes $400.

- At time t4, transaction $T_Y$ writes the updated value, i.e., $400.

- After that, at time t5, transaction $T_X$ reads the available value of account A, and that will be read as $400.

- It means that within the same transaction $T_X$, it reads two different values of account A, i.e., $ 300 initially, and after updation made by transaction $T_Y$, it reads $400. It is an unrepeatable read and is therefore known as the Unrepeatable read problem.

Thus, in order to maintain consistency in the database and avoid such problems that take place in concurrent execution, management is needed, and that is where the concept of Concurrency Control comes into role.

# Concurrency Control

Concurrency Control is the working concept that is required for controlling and managing the concurrent execution of database operations and thus avoiding the inconsistencies in the database. Thus, for maintaining the concurrency of the database, we have the concurrency control protocols.

## Concurrency Control Protocols

The concurrency control protocols ensure the *atomicity, consistency, isolation, durability* and *serializability* of the concurrent execution of the database transactions. Therefore, these protocols are categorized as:

- o Lock Based Concurrency Control Protocol

- o Time Stamp Concurrency Control Protocol

- o Validation Based Concurrency Control Protocol

# 1. Lock-Based Protocol

In this type of protocol, any transaction cannot read or write data until it acquires an appropriate lock on it. There are two types of lock:

**1. Shared lock:**

- o It is also known as a Read-only lock. In a shared lock, the data item can only read by the transaction. o It can be shared between the transactions because when the transaction holds a lock, then it can't update the data on the data item.

**2. Exclusive lock:**

- o In the exclusive lock, the data item can be both reads as well as written by the transaction. o This lock is exclusive, and in this lock, multiple transactions do not modify the same data simultaneously.
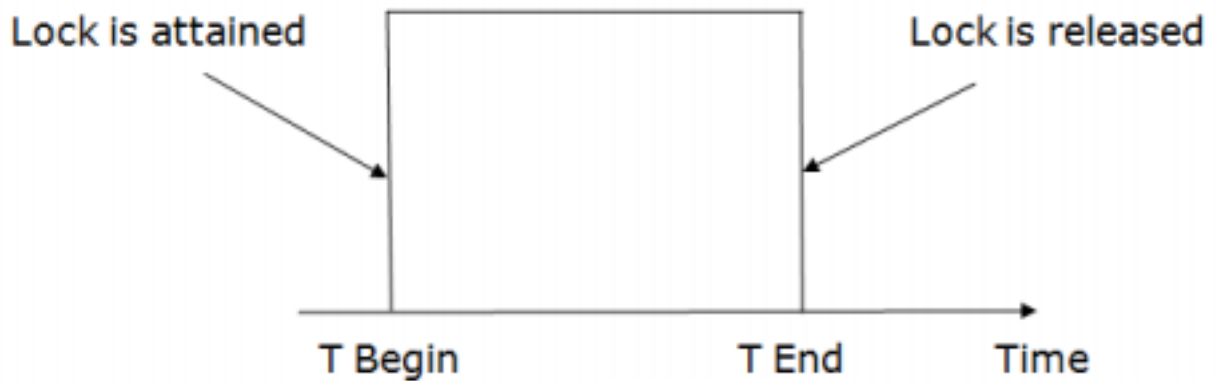
# There are four types of lock protocols available:
## 1. Simplistic lock protocol

It is the simplest way of locking the data while transaction. Simplistic lock-based protocols allow all the transactions to get the lock on the data before insert or delete or update on it. It will unlock the data item after completing the transaction.
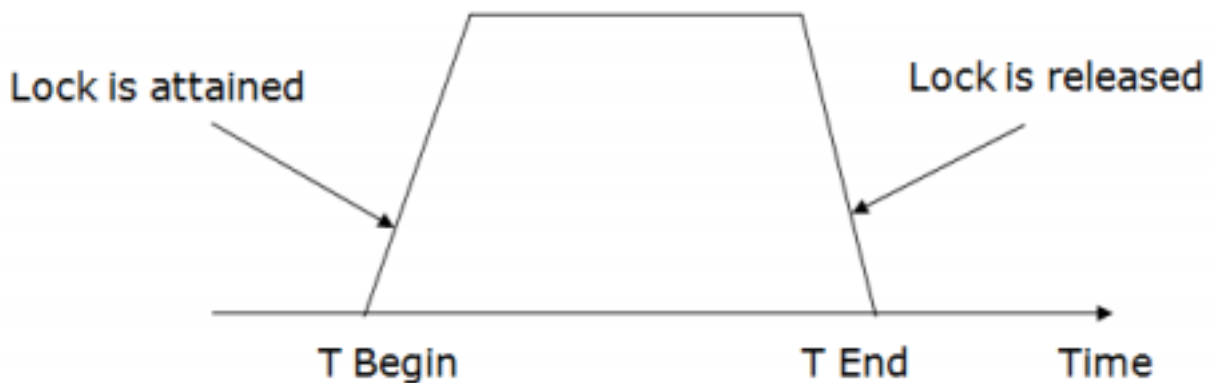
## 2. Pre-claiming Lock Protocol

- o Pre-claiming Lock Protocols evaluate the transaction to list all the data items on which they need locks.

- o Before initiating an execution of the transaction, it requests DBMS for all the lock on all those data items.

- o If all the locks are granted then this protocol allows the transaction to begin. When the transaction is completed then it releases all the lock.

- o If all the locks are not granted then this protocol allows the transaction to rolls back and waits until all the locks are granted.

Lock is attained          Lock is released

T Begin          T End          Time

## 3. Two-phase locking (2PL)

- o  The two-phase locking protocol divides the execution phase of the transaction into three parts. o  In the first part, when the execution of the transaction starts, it seeks permission for the lock it requires.

- o  In the second part, the transaction acquires all the locks. The third phase is started as soon as the transaction releases its first lock.

- o  In the third phase, the transaction cannot demand any new locks. It only releases the acquired locks.



Lock is attained          Lock is released

T Begin          T End          Time

There are two phases of 2PL:

**Growing phase:** In the growing phase, a new lock on the data item may be acquired by the transaction, but none can be released.
**Shrinking phase:** In the shrinking phase, existing lock held by the transaction may be released, but no new  locks can be acquired.

In the below example, if lock conversion is allowed then the following phase can happen:

1. Upgrading of lock (from S(a) to X (a)) is allowed in growing phase.

2. Downgrading of lock (from X(a) to S(a)) must be done in shrinking phase.

**Example:**

| | T1 | T2 |
|---|---|---|
| 0 | LOCK-S(A) | |
| 1 | | LOCK-S(A) |
| 2 | LOCK-X(B) | |
| 3 | —— | —— |
| 4 | UNLOCK(A) | |
| 5 | | LOCK-X(C) |
| 6 | UNLOCK(B) | |
| 7 | | UNLOCK(A) |
| 8 | | UNLOCK(C) |
| 9 | —— | —— |

The following way shows how unlocking and locking work with 2-PL.
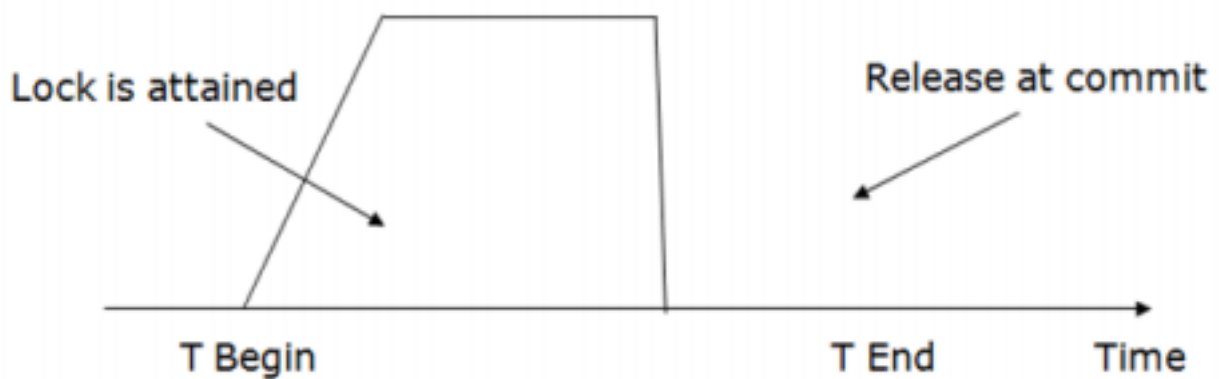
**Transaction T1:**

o **Growing phase:** from step 1-3

o **Shrinking phase:** from step 5-7

o **Lock point:** at 3

**Transaction T2:**
o **Growing phase:** from step 2-6

o **Shrinking phase:** from step 8-9

o **Lock point:** at 6

# 4. Strict Two-phase locking (Strict-2PL)

- o The first phase of Strict-2PL is similar to 2PL. In the first phase, after acquiring all the locks, the transaction continues to execute normally.

- o The only difference between 2PL and strict 2PL is that Strict-2PL does not release a lock after using it.

- o Strict-2PL waits until the whole transaction to commit, and then it releases all the locks at a time.

- o Strict-2PL protocol does not have shrinking phase of lock release.



It does not have cascading abort as 2PL does.

# 2. Timestamp Ordering Protocol

The timestamp-ordering protocol ensures serializability among transactions in their conflicting read and write operations. This is the responsibility of the protocol system that the  conflicting pair of tasks should be executed according to the timestamp values of the transactions.

- · The timestamp of transaction $T_i$ is denoted as $TS(T_i)$.
- · Read time-stamp of data-item X is denoted by R-timestamp(X).
- · Write time-stamp of data-item X is denoted by W-timestamp(X).

Timestamp ordering protocol works as follows −

- · **If a transaction Ti issues a read(X) operation −**

    - o If $TS(Ti) < W\text{-timestamp}(X)$
        - ▪ Operation rejected.
    - o If $TS(Ti) >= W\text{-timestamp}(X)$
        - ▪ Operation executed.
    - o All data-item timestamps updated.

**· If a transaction Ti issues a write(X) operation −**

- o If $TS(Ti) < R\text{-timestamp}(X)$
  - ▪ Operation rejected.
- o If $TS(Ti) < W\text{-timestamp}(X)$
  - ▪ Operation rejected and Ti rolled back.
- o Otherwise, operation executed.

## Thomas' Write Rule

This rule states if $TS(Ti) < W\text{-timestamp}(X)$, then the operation is rejected and $T_i$ is rolled back.

Time-stamp ordering rules can be modified to make the schedule view serializable.

Instead of making $T_i$ rolled back, the 'write' operation itself is ignored.

# 3. Optimistic Methods of Concurrency Control:

The optimistic method of concurrency control is based on the assumption that conflicts of database operations are rare and that it is better to let transactions run to completion and only check for conflicts before they commit.
An optimistic concurrency control method is also known as validation or certification methods. No checking  is done while the transaction is executing. The optimistic method does not require locking or timestamping  techniques. Instead, a transaction is executed without restrictions until it is committed. In optimistic methods,  each transaction moves through the following phases:

   a. Read phase.
   b. Validation or certification phase.
   c. Write phase.

## *a. Read phase:*

In a Read phase, the updates are prepared using private (or local) copies (or versions) of the granule. In this phase, the transaction reads values of committed data from the database, executes the needed computations,  and makes the updates to a private copy of the database values. All update operations of the transaction are  recorded in a temporary update file, which is not accessed by the remaining transactions.
It is conventional to allocate a timestamp to each transaction at the end of its Read to determine the set of transactions that must be examined by the validation procedure. These set of transactions are those who have  finished their Read phases since the start of the transaction being verified

## *b. Validation or certification phase :*

In a validation (or certification) phase, the transaction is validated to assure that the changes made will not affect the integrity and consistency of the database.

If the validation test is positive, the transaction goes to the write phase. If the validation test is negative, the transaction is restarted, and the changes are discarded. Thus, in this phase the list of granules is checked for conflicts. If conflicts are detected in this phase, the transaction is aborted and restarted. The validation algorithm must check that the transaction has :
   · Seen all modifications of transactions committed after it starts.
   · Not read granules updated by a transaction committed after its start.

## c.Write phase :

In a Write phase, the changes are permanently applied to the database and the updated granules are made public. Otherwise, the updates are discarded and the transaction is restarted. This phase is only for the Read Write transactions and not for Read-only transactions.

## *Advantages of Optimistic Methods for Concurrency Control :*

i. This technique is very efficient when conflicts are rare. The occasional conflicts result in the transaction roll back.
ii. The rollback involves only the local copy of data, the database is not involved and thus there will  not be any cascading rollbacks.

## *Problems of Optimistic Methods for Concurrency Control :*

i. Conflicts are expensive to deal with, since the conflicting transaction must be rolled back. ii. Longer transactions are more likely to have conflicts and may be repeatedly rolled back because of  conflicts with short transactions.

## *Applications of Optimistic Methods for Concurrency Control :*

i. Only suitable for environments where there are few conflicts and no long transactions. ii. Acceptable for mostly Read or Query database systems that require very few update transactions

# DBMS | Concurrency Control Protocol | Multiple Granularity Locking

In the various Concurrency Control schemes have used different methods and every individual Data item as the unit on which synchronization is performed. A certain drawback of this technique is if a transaction $T_i$ needs to access the entire database, and a locking protocol is used, then $T_i$ must lock each item in the database. It is less efficient, it would be more simpler if $T_i$ could use a single lock to lock the entire database. But, if it consider the second proposal, this should not in fact overlook the certain flaw in the proposed method. Suppose another transaction just needs to access a few data items from a database, so locking the entire database  seems to be unnecessary moreover it may cost us loss of Concurrency, which was our primary  goal in the first place. To bargain between Efficiency and Concurrency. Use Granularity.
Let's start by understanding what is meant by Granularity.

**Granularity –** It is the size of data item allowed to lock. Now *Multiple Granularity* means hierarchically breaking up the database into blocks which can be locked and can be track what need to lock and in what fashion. Such a hierarchy can be represented graphically as a tree. For example, consider the tree, which consists of four levels of nodes. The highest level represents the entire database. Below it are nodes of type **area**; the database consists of exactly  these areas. Area has children nodes which are called files. Every area has those files that are  its child nodes. No file can span more than one area.
Finally, each file has child nodes called records. As before, the file consists of exactly those

records that are its child nodes, and no record can be present in more than one file. Hence, the levels starting from the top level are:

· database
· area
· file
· record

**Figure –** Multi Granularity tree Hierarchy

Consider the above diagram for the example given, each node in the tree can be locked individually. As in the 2-phase locking protocol, it shall use shared and exclusive lock modes. When a transaction locks a node, in either shared or exclusive mode, the transaction also implicitly locks all the descendants of that node in the same lock mode. For example, if transaction $T_i$ gets an explicit lock on file $F_c$ in exclusive mode, then it has an implicit lock in exclusive mode on all the records belonging to that file. It does not need to lock the individual records of $F_c$ explicitly. this is the main difference between Tree Based Locking and Hierarchical locking for multiple granularity.

Now, with locks on files and records made simple, how does the system determine if the root node can be locked? One possibility is for it to search the entire tree but the solution nullifies

the whole purpose of the multiple-granularity locking scheme. A more efficient way to gain this knowledge is to introduce a new lock mode, called *Intention lock mode*.

**Intention Mode Lock –**

In addition to **S** and **X** lock modes, there are three additional lock modes with multiple granularity:

· **Intention-Shared (IS):** explicit locking at a lower level of the tree but only with shared locks. · **Intention-Exclusive (IX):** explicit locking at a lower level with exclusive or shared locks. · **Shared & Intention-Exclusive (SIX):** the sub-tree rooted by that node is locked explicitly in shared mode and explicit locking is being done at a lower level with exclusive mode locks.

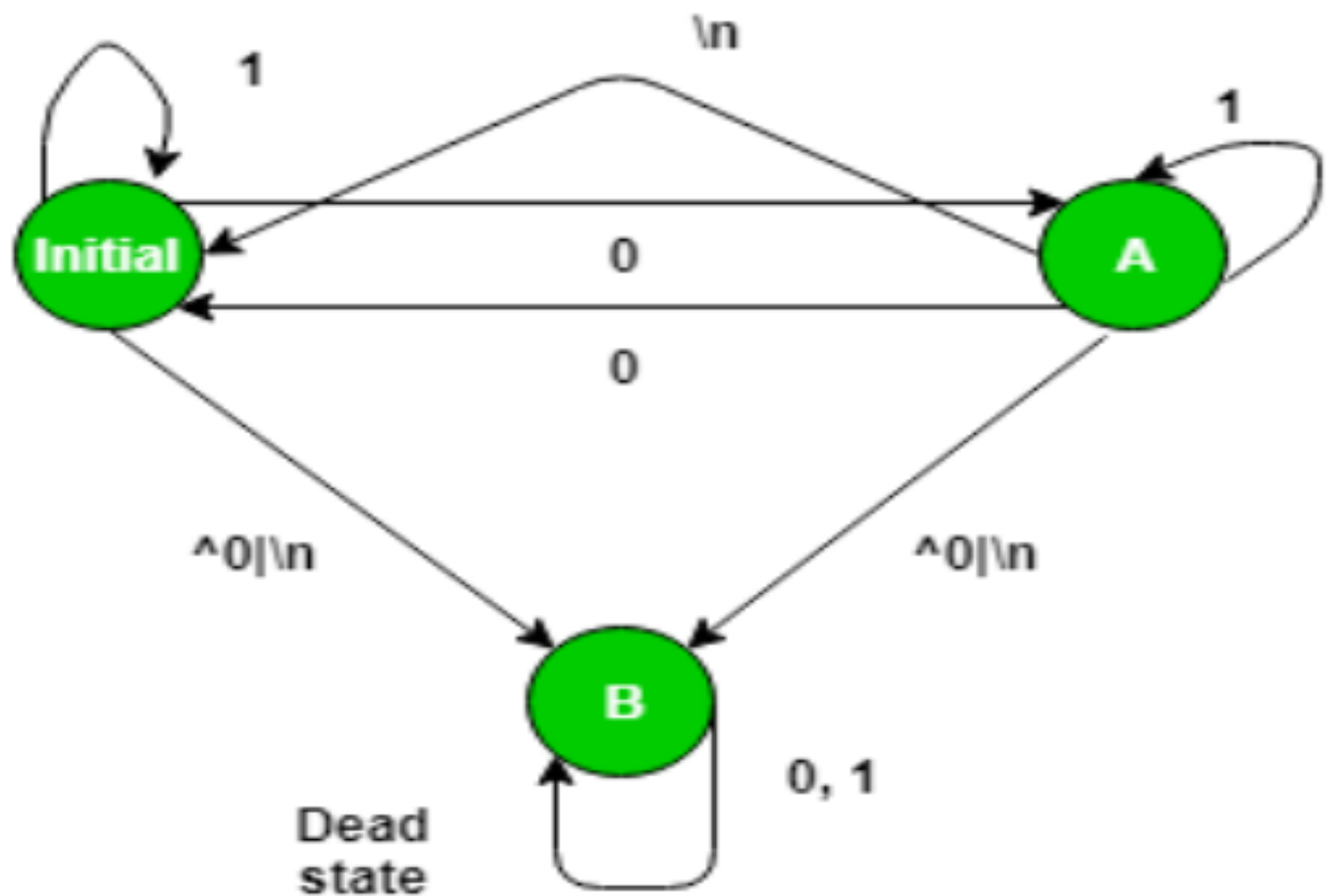The compatibility matrix for these lock modes are described below:



**Figure –** Multi Granularity tree Hierarchy

The multiple-granularity locking protocol uses the intention lock modes to ensure serializability. It requires that a transaction $T_i$ that attempts to lock a node must follow these protocols:

1. Transaction $T_i$ must follow the lock-compatibility matrix.
2. Transaction $T_i$ must lock the root of the tree first, and it can lock it in any mode.
3. Transaction $T_i$ can lock a node in S or IS mode only if $T_i$ currently has the parent of the node locked in either IX or IS mode.
4. Transaction $T_i$ can lock a node in X, SIX, or IX mode only if $T_i$ currently has the parent of the node

locked in either IX or SIX mode.

5. Transaction $T_i$ can lock a node only if $T_i$ has not previously unlocked any node (i.e., $T_i$ is two phase). 6. Transaction $T_i$ can unlock a node only if $T_i$ currently has none of the children of the node locked. Observe that the multiple-granularity protocol requires that locks be acquired in top-down (root-to-leaf) order, whereas locks must be released in bottom-up (leaf to-root) order. As an illustration of the protocol, consider the tree given above and the transactions:

· Say transaction $T_1$ reads record $R_{a2}$ in file $F_a$. Then, $T_2$ needs to lock the database, area $A_1$, and $F_a$ in IS mode (and in that order), and finally to lock $R_{a2}$ in S mode.

· Say transaction $T_2$ modifies record $R_{a9}$ in file $F_a$. Then, $T_2$ needs to lock the database, area $A_1$, and file $F_a$ (and in that order) in IX mode, and at last to lock $R_{a9}$ in X mode.

· Say transaction $T_3$ reads all the records in file $F_a$. Then, $T_3$ needs to lock the database and area $A_1$ (and in that order) in IS mode, and at last to lock $F_a$ in S mode.

· Say transaction $T_4$ reads the entire database. It can do so after locking the database in S mode. Note that transactions $T_1$, $T_3$ and $T_4$ *can access the database concurrently*. Transaction $T_2$ can execute concurrently with $T_1$, but not with either $T_3$ or $T_4$. This protocol enhances *concurrency and reduces lock overhead*.Deadlock are still possible in the multiple-granularity protocol, as it is in the two-phase locking protocol. These can be eliminated by using certain deadlock elimination techniques.

# DBMS - Deadlock

In a multi-process system, deadlock is an unwanted situation that arises in a shared resource environment, where a process indefinitely waits for a resource that is held by another process.

For example, assume a set of transactions $\{T_0, T_1, T_2, ...,T_n\}$. $T_0$ needs a resource X to complete its task. Resource X is held by $T_1$, and $T_1$ is waiting for a resource Y, which is held by $T_2$. $T_2$ is waiting for resource Z, which is held by $T_0$. Thus, all the processes wait for each other to release resources. In this situation, none of the processes can finish their task. This situation is known as a deadlock.

Deadlocks are not healthy for a system. In case a system is stuck in a deadlock, the transactions involved in the deadlock are either rolled back or restarted.

## Deadlock Prevention

To prevent any deadlock situation in the system, the DBMS aggressively inspects all the operations, where transactions are about to execute. The DBMS inspects the operations and

analyzes if they can create a deadlock situation. If it finds that a deadlock situation might occur, then that transaction is never allowed to be executed.

There are deadlock prevention schemes that use timestamp ordering mechanism of transactions in order to predetermine a deadlock situation.

### Wait-Die Scheme

In this scheme, if a transaction requests to lock a resource (data item), which is already held with a conflicting lock by another transaction, then one of the two possibilities may occur −

- If $TS(T_i) < TS(T_j)$ − that is $T_i$, which is requesting a conflicting lock, is older than $T_j$ − then $T_i$ is allowed to wait until the data-item is available.

- If $TS(T_i) > TS(t_j)$ − that is $T_i$ is younger than $T_j$ − then $T_i$ dies. $T_i$ is restarted later with a random delay but with the same timestamp.

This scheme allows the older transaction to wait but kills the younger one.

### Wound-Wait Scheme

In this scheme, if a transaction requests to lock a resource (data item), which is already held with conflicting lock by some another transaction, one of the two possibilities may occur −

- If $TS(T_i) < TS(T_j)$, then $T_i$ forces $T_j$ to be rolled back − that is $T_i$ wounds $T_j$. $T_j$ is restarted later with a random delay but with the same timestamp.

- If $TS(T_i) > TS(T_j)$, then $T_i$ is forced to wait until the resource is available.

This scheme, allows the younger transaction to wait; but when an older transaction requests an item held by a younger one, the older transaction forces the younger one to abort and release the item.

In both the cases, the transaction that enters the system at a later stage is aborted.
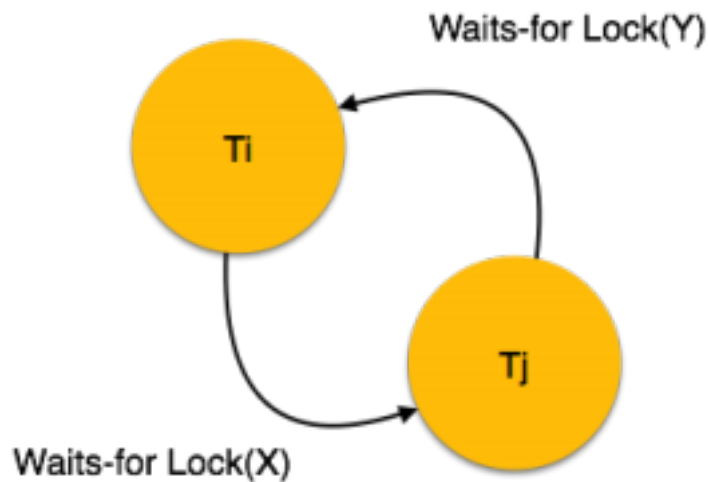
# Deadlock Avoidance

Aborting a transaction is not always a practical approach. Instead, deadlock avoidance mechanisms can be used to detect any deadlock situation in advance. Methods like "wait-for graph" are available but they are suitable for only those systems where transactions are lightweight having fewer instances of resource. In a bulky system, deadlock prevention techniques may work well.

### Wait-for Graph

This is a simple method available to track if any deadlock situation may arise. For each transaction entering into the system, a node is created. When a transaction $T_i$ requests for a lock on an item, say X, which is held by some other transaction $T_j$, a directed edge is created from $T_i$ to $T_j$. If $T_j$ releases item X, the edge between them is dropped and $T_i$ locks the data item.

The system maintains this wait-for graph for every transaction waiting for some data items held by others. The system keeps checking if there's any cycle in the graph.

Waits-for Lock(Y)

Ti

Tj

Waits-for Lock(X)

Here, we can use any of the two following approaches −

- First, do not allow any request for an item, which is already locked by another transaction. This is not always feasible and may cause starvation, where a transaction indefinitely waits for a data item and can never acquire it.

- The second option is to roll back one of the transactions. It is not always feasible to roll back the younger transaction, as it may be important than the older one. With the help of some relative algorithm, a transaction is chosen, which is to be aborted. This transaction is known as the **victim** and the process is known as **victim selection**.

# DBMS - Data Backup

## Loss of Volatile Storage

A volatile storage like RAM stores all the active logs, disk buffers, and related data. In addition, it stores all the transactions that are being currently executed. What happens if such a volatile storage crashes abruptly? It would obviously take away all the logs and active copies of the database. It makes recovery almost impossible, as everything that is required to recover the data is lost.

Following techniques may be adopted in case of loss of volatile storage −

- We can have **checkpoints** at multiple stages so as to save the contents of the database periodically.

- A state of active database in the volatile memory can be periodically **dumped** onto a stable storage, which may also contain logs and active transactions and buffer blocks.

- <dump> can be marked on a log file, whenever the database contents are dumped from a non-volatile memory to a stable one.

### Recovery

- When the system recovers from a failure, it can restore the latest dump. ·
It can maintain a redo-list and an undo-list as checkpoints.

- It can recover the system by consulting undo-redo lists to restore the state of all transactions up to the last checkpoint.

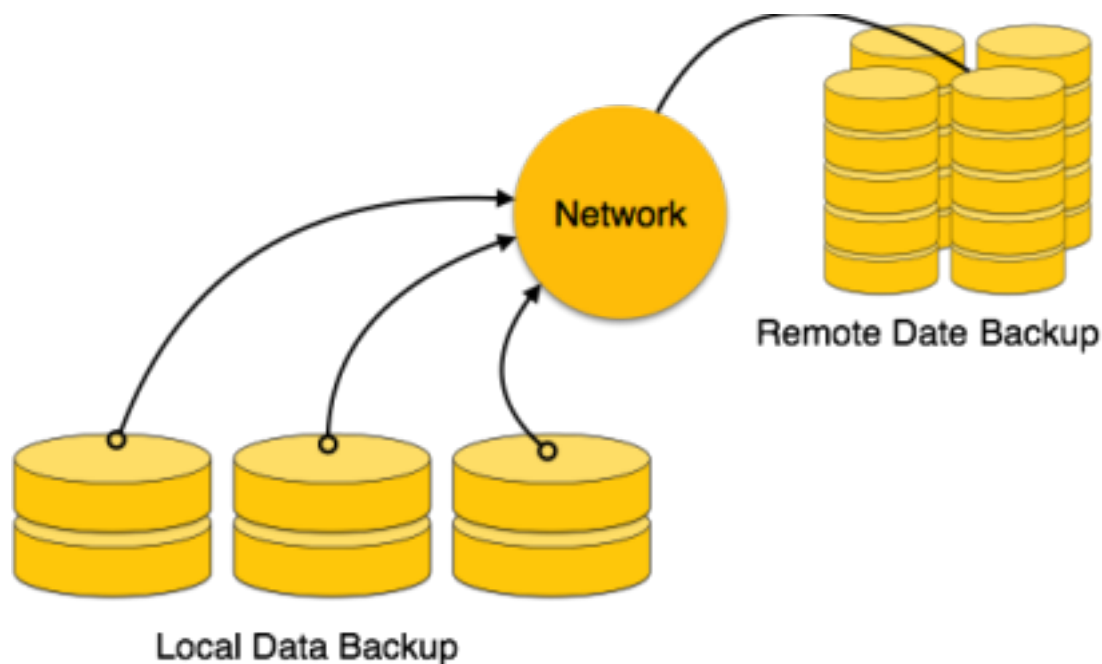# Database Backup & Recovery from Catastrophic Failure

A catastrophic failure is one where a stable, secondary storage device gets corrupt. With the storage device, all the valuable data that is stored inside is lost. We have two different strategies to recover data from such a catastrophic failure −

- Remote backup &minu; Here a backup copy of the database is stored at a remote location from where it can be restored in case of a catastrophe.

- Alternatively, database backups can be taken on magnetic tapes and stored at a safer place. This backup can later be transferred onto a freshly installed database to bring it to the point of backup.

Grown-up databases are too bulky to be frequently backed up. In such cases, we have techniques where we can restore a database just by looking at its logs. So, all that we need to do here is to take a backup of all the logs at frequent intervals of time. The database can be backed up once a week, and the logs being very small can be backed up every day or as frequently as possible.

## Remote Backup

Remote backup provides a sense of security in case the primary location where the database is located gets destroyed. Remote backup can be offline or real-time or online. In case it is offline, it is maintained manually.



Online backup systems are more real-time and lifesavers for database administrators and investors. An online backup system is a mechanism where every bit of the real-time data is backed up simultaneously at two distant places. One of them is directly connected to the system and the other one is kept at a remote place as backup.

As soon as the primary database storage fails, the backup system senses the failure and switches the user system to the remote storage. Sometimes this is so instant that the users can't even realize a failure.

# DBMS - Data Recovery

## Crash Recovery

DBMS is a highly complex system with hundreds of transactions being executed every second. The durability and robustness of a DBMS depends on its complex architecture and its underlying hardware and system software. If it fails or crashes amid transactions, it is expected that the system would follow some sort of algorithm or techniques to recover lost data.

## Failure Classification

To see where the problem has occurred, we generalize a failure into various categories, as follows −

### Transaction failure

A transaction has to abort when it fails to execute or when it reaches a point from where it can't go any further. This is called transaction failure where only a few transactions or processes are hurt.

Reasons for a transaction failure could be −

- **Logical errors** − Where a transaction cannot complete because it has some code error or any internal error condition.

- **System errors** − Where the database system itself terminates an active transaction because the DBMS is not able to execute it, or it has to stop because of some system condition. For example, in case of deadlock or resource unavailability, the system aborts an active transaction.

### System Crash

There are problems − external to the system − that may cause the system to stop abruptly and cause the system to crash. For example, interruptions in power supply may cause the failure of underlying hardware or software failure.

Examples may include operating system errors.

### Disk Failure

In early days of technology evolution, it was a common problem where hard-disk drives or storage drives used to fail frequently.

Disk failures include formation of bad sectors, unreachability to the disk, disk head crash or any other failure, which destroys all or a part of disk storage.

## Storage Structure

We have already described the storage system. In brief, the storage structure can be divided into two categories −

- **Volatile storage** − As the name suggests, a volatile storage cannot survive system crashes. Volatile storage devices are placed very close to the CPU; normally they are embedded onto the chipset itself. For example, main memory and cache memory are

examples of volatile storage. They are fast but can store only a small amount of information.

· **Non-volatile storage** − These memories are made to survive system crashes. They are huge in data storage capacity, but slower in accessibility. Examples may include hard disks, magnetic tapes, flash memory, and non-volatile (battery backed up) RAM.

# Recovery and Atomicity

When a system crashes, it may have several transactions being executed and various files opened for them to modify the data items. Transactions are made of various operations, which are atomic in nature. But according to ACID properties of DBMS, atomicity of transactions as a whole must be maintained, that is, either all the operations are executed or none.

When a DBMS recovers from a crash, it should maintain the following − · It should check the states of all the transactions, which were being executed.

· A transaction may be in the middle of some operation; the DBMS must ensure the atomicity of the transaction in this case.

· It should check whether the transaction can be completed now or it needs to be rolled back.

· No transactions would be allowed to leave the DBMS in an inconsistent state.

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

· Maintaining the logs of each transaction, and writing them onto some stable storage before actually modifying the database.

· Maintaining shadow paging, where the changes are done on a volatile memory, and later, the actual database is updated.

# Log-based Recovery

Log is a sequence of records, which maintains the records of actions performed by a transaction. It is important that the logs are written prior to the actual modification and stored on a stable storage media, which is failsafe.

Log-based recovery works as follows −

· The log file is kept on a stable storage media.

· When a transaction enters the system and starts execution, it writes a log about it.

$<T_n, Start>$

· When the transaction modifies an item X, it write logs as follows −

$<T_n, X, V_1, V_2>$

It reads $T_n$ has changed the value of X, from $V_1$ to $V_2$.

· When the transaction finishes, it logs −

<T$_n$, commit>

The database can be modified using two approaches −

- **Deferred database modification** − All logs are written on to the stable storage and the database is updated when a transaction commits.
    - **Immediate database modification** − Each log follows an actual database modification. That is, the database is modified immediately after every operation.

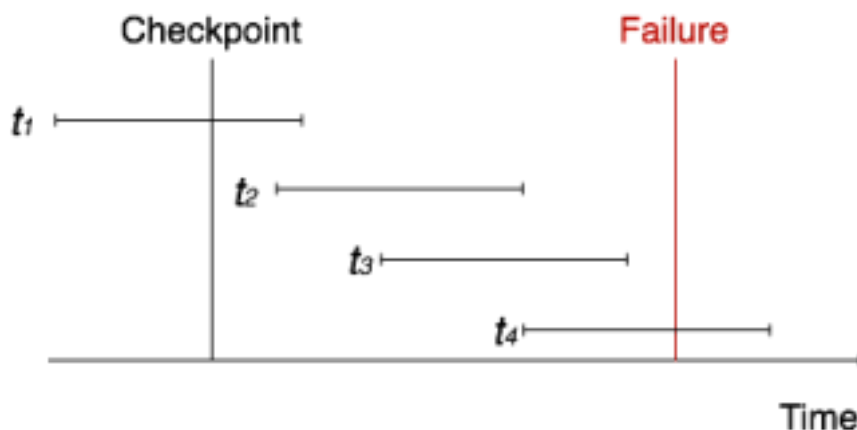# Recovery with Concurrent Transactions

When more than one transaction are being executed in parallel, the logs are interleaved. At the time of recovery, it would become hard for the recovery system to backtrack all logs, and then start recovering. To ease this situation, most modern DBMS use the concept of 'checkpoints'.

## Checkpoint

Keeping and maintaining logs in real time and in real environment may fill out all the memory space available in the system. As time passes, the log file may grow too big to be handled at all. Checkpoint is a mechanism where all the previous logs are removed from the system and stored permanently in a storage disk. Checkpoint declares a point before which the DBMS was in consistent state, and all the transactions were committed.

## Recovery

When a system with concurrent transactions crashes and recovers, it behaves in the following manner −



- The recovery system reads the logs backwards from the end to the last checkpoint. ·

It maintains two lists, an undo-list and a redo-list.

- If the recovery system sees a log with <T$_n$, Start> and <T$_n$, Commit> or just <T$_n$, Commit>, it puts the transaction in the redo-list.

- If the recovery system sees a log with <T$_n$, Start> but no commit or abort log found, it puts the transaction in undo-list.

All the transactions in the undo-list are then undone and their logs are removed. All the transactions in the redo-list and their previous logs are removed and then redone before

saving  their logs.