

9. Transaction Processing

Transaction management:

- Transaction
- Concurrency control
- Recovery system

Transaction:

- A *transaction* is an execution of a user program, and is seen by the DBMS as a series or list of actions.
- The actions that can be executed by a transaction include reads and writes of database objects, whereas actions in an ordinary program could involve user input, access to network devices, user interface drawing, etc.(Or)
- It is defined as a collection of operations that form a single logical unit of work is called transaction.

It is a foundation for concurrent execution and recovery from system failure in a DBMS

Eg: Select, commit, create etc.

The ACID Properties: A DBMS must ensure four important properties of transaction to maintain data in the face of concurrent access and System failures.

- ACID is used to refer to these 4 properties.
- A**tomicity: All actions in the transaction happen, or none happen.
- C**onsistency: If each transaction is consistent, and the DB starts consistent, it ends up consistent.
- I**solation: Execution of one transaction is isolated from that of all others.
- D**urability: If a transaction commits, its effects persist, if system failure occurs.

1
[9.1]

Very valuable properties of DBMSs

- without these properties, DBM's would be much less useful

Transaction access data using two operations.

- **Read(x):** which transfers data item x from DB to a logical buffer belonging to transaction that executed read operation.
- **Write(x):** which transfers data item x to DE from a logical buffer belonging to transaction that executed write operation.
- Eg: Transfer of 100/-Rs from account A to account B.

T1: Read (A)

A: =A-100;

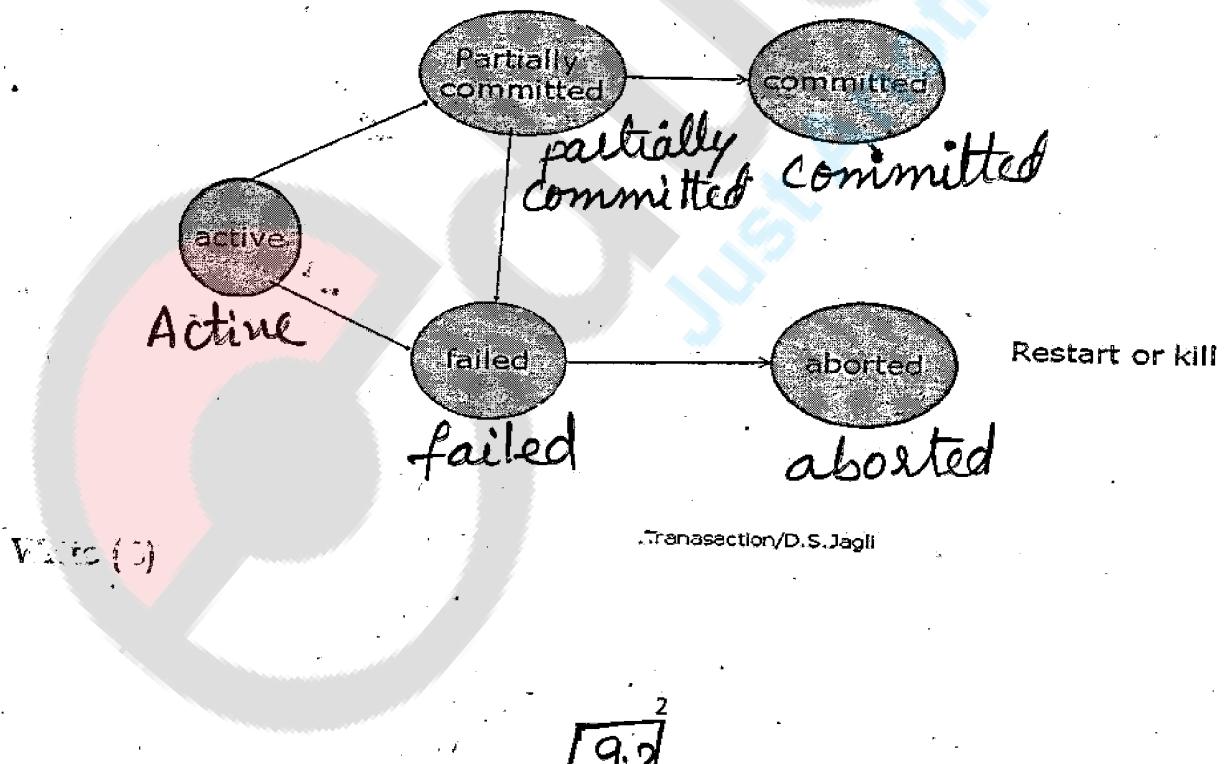
Write (A);

Read (B);

B: =B+100;

Write (B);

State diagram of Transaction



2
9.2

Atomicity of Transactions

- A transaction might commit after completing all its actions, or it could abort (or be aborted by the DBMS) after executing some actions.
- Atomicity* means a transaction executes when all actions of the transaction are completed fully, or none are. This means there are no partial transactions.
- Atomic Transactions: a user can think of a transaction as always either executing all its actions, or not executing any actions at all.
 1. One approach: DBMS logs all actions so that it can undo the actions of aborted transactions.
 2. Another approach: Shadow Pages

Transaction Consistency

- "Consistency" - data in DBMS is accurate in modeling real world, follows integrity constraints.
- Consistency* involves beginning a transaction with a 'consistent' database, and finishing with a 'consistent' database.
- For example, in a bank database, money should never be "created" or "deleted" without an appropriate deposit or withdrawal. Every transaction should see a consistent database.
- System checks ICs and if they fail, the transaction rolls back (i.e., is aborted).
 1. DBMS enforces some ICs, depending on the ICs declared in CREATE TABLE statements.
 2. Beyond this, DBMS does not understand the semantics of the data. (e.g., it does not understand how the interest on a bank account is computed).

Isolation

- Multiple users can submit transactions, Each transaction executes as if it was running by itself.

- Concurrency is achieved by DBMS, which interleaves actions (reads/writes of DB objects) of various transactions.
- *Isolation* ensures that a transaction can run independently, without considering any side effects that other concurrently running transactions might have.
- Many techniques have been developed. Fall into two basic categories:
 - Pessimistic – don't let problems arise in the first place
 - Optimistic – assume conflicts are rare, deal with them *after* they happen.

Durability - Recovering From a Crash

- *Durability* defines the persistence of committed data: once a transaction commits, the data should persist in the database even if the system crashes before the data is written to non-volatile storage.
- System Crash - short-term memory lost
- Disk Crash - "stable" data lost
 - Need back ups; raid-techniques can help avoid this.
- At the end all committed updates and only those updates are reflected in the database.
 - Some care must be taken to handle the case of a crash occurring during the recovery process!

Concurrency Control:

Isolation and Serializability

Scheduling Transactions: A *schedule* is a series of (possibly overlapping) transactions.

- *Serial schedule:* A schedule that does not interleave the actions of different transactions.
 - i.e., you run the transactions serially (one at a time)

- Equivalent schedules:** For any database state, the effect and output of executing the first schedule is identical to the effect of executing the second schedule.
- Serializable schedule:** A schedule that is equivalent to some serial execution of the transactions.

With a serializable schedule you only see things that could happen in situations where you were running transactions one-at-a-time.

Anomalies with Interleaved Execution

Conflict Serializable Schedules

- Two operations conflict if they are by different transactions, they are on the same object, and at least one of them is a write.
- Two schedules are conflict equivalent if they involve the same actions of the same transactions, and every pair of conflicting actions is ordered the same way. Schedule S is conflict serializable if S is conflict equivalent to some serial schedule.

Some Examples of Conflicts

A conflict exists when two transactions access the same item, and at least one of the accesses is a write.

1. lost update problem

T: transfer ₹100 from A to C: R(A) W(A) R(C) W(C)

S: transfer ₹100 from B to C: R(B) W(B) R(C) W(C)

2. inconsistent retrievals problem (dirty reads violate consistency)

T: transfer ₹100 from A to C: R(A) W(A) R(C) W(C)

S: compute total balance for A and C: R(A) R(C)

3. nonrepeatable reads

T: transfer ₹100 from A to C: R(A) W(A) R(C) W(C)

S: check balance and withdraw ₹100 from A: R(A) R(A) W(A)

Legal Interleaved Schedules:

Examples

$T < S$

1. Avoid *lost update* problem

T: transfer \$100 from A to C: $R(A)$ $W(A)$ $R(C)$ $W(C)$
S: transfer \$100 from B to C: $R(B)$ $W(B)$ $R(C)$ $W(C)$

2. Avoid *inconsistent retrievals* problem

T: transfer \$100 from A to C: $R(A)$ $W(A)$ $R(C)$ $W(C)$
S: compute total balance for A and C: $R(A)$ $R(C)$

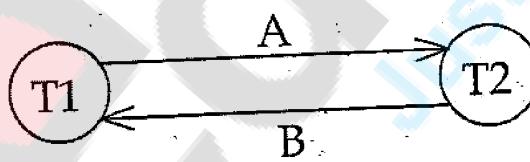
3. Avoid *nonrepeatable reads*

T: transfer \$100 from A to C $R(A)$ $W(A)$ $R(C)$ $W(C)$
S: check balance and withdraw \$100 from A: $R(A)$ $R(A)$ $W(A)$

❖ A schedule that is **not conflict serializable**:

❖ The cycle in the graph reveals the problem. The output of T1 depends on T2, and vice-versa.

T1:	$R(A), W(A),$	$R(B), W(B)$
T2:	$R(A), W(A), R(B), W(B)$	



Dependency graph

Concurrent Executions: Two reasons for concurrency.

- ① Improved throughput and resource utilization.
- ② Reduced waiting time.

Serial Schedule: T_1 is followed by T_2 .

T_1

read(A)

$A := A - 50$

write(A)

read(B)

$B := B + 50$

write(B)

T_2

read(A)

$temp = A * 0.1$

$A = A - temp$

write(A)

read(B)

$B := B + temp$

write(B)

concurrent schedule:

T_1

read(A)

$A := A - 50$

write(A)

T_2

read(A)

$temp = A * 0.1$

write(A)

9.7

Serializability:

① Data base system must control concurrent execution of transactions, to ensure that the db state remains consistent.

<u>T₁</u>	<u>T₂</u>
read(A)	read(A)
write(A)	write(A)

Transactions may perform an arbitrary sequence of operations on a data item i.e residing in the local buffer of the transaction.

The different forms of schedule equivalence

- ① conflict serializability.
- ② view serializability.

19.8

conflict:

if I_1, I_2 are operations by different transactions T_1, T_2 on the same data item 'x', and at least one of these instructions is write operation then I_1, I_2 are conflict.

conflict equivalent:

if schedule S_1 can be transformed into a Schedule S_2 by series of swaps of non conflicting instructions, then S_1 and S_2 are conflict equivalent.

Eg:

		T_1	T_2		T_1	T_2
read(A)				=	read(A)	
write(A)		read(A)			write(A)	
read(B)			write(A)		read(B)	
write(B)					write(B)	
			read(B)			
			write(B)			

$S_1 =$

① Conflict Serializability:

S - Schedule

I_1, I_2 are instructions

T_1, T_2 are transactions

→ I_1, I_2 are using different data items then no problems or conflicts can arise.

→ I_1, I_2 are using same data item then order of I_1, I_2 is matter

→ Four cases need to be consider.

① $I_1 = \text{read}(x)$ $I_2 = \text{read}(x)$

Order of I_1, I_2 doesn't matter for T_1, T_2

② $I_1 = \text{read}(x)$, $I_2 = \text{write}(x)$

Order of I_1, I_2 matters for T_1, T_2 execution

③ $I_1 = \text{write}(x)$, $I_2 = \text{read}(x)$

Order of I_1, I_2 matters for T_1, T_2 execution.

④ $I_1 = \text{write}(x)$, $I_2 = \text{write}(x)$

Order of I_1, I_2 directly affecting the final value of data item - x in schedule.

9.10

conflict serializable:-

The schedule S is conflict serializable if it is conflict equivalent to serial schedule.

view serializability:-

Two schedules S_1, S_2 , where the same set of transactions participates in both schedules.

The schedule S_1, S_2 are said to be view equivalent if following conditions are satisfied

① for each data item X

if T_1 reads initial value of X in Schedule S_1
then T_1 must in Schedule S_2 & read the initial
value of X.

② for each data item X

if T_1 executes read(X) in Schedule S_1 , the
value of X is produced by write(X) in T_2 .
then read(X) of T_1 must be in Schedule S_2 .
value of X is same by write(X) in T_2 .

③ for each data item X

any transaction write final value of X in S_1
must write final value of X in S_2 .

View Serializable :-

Schedule S is view serializable if it is view equivalent to a serial schedule.

eg:	T3	T4	T5	T3	T4	T5
	<u>read(x)</u>			<u>read(x)</u>		
	<u>writel(x)</u>	<u>writel(x)</u>		<u>writel(x)</u>	<u>writel(x)</u>	
	<u>waitfor</u>	<u>writel(x)</u>	<u>writel(x)</u>	<u>writel(x)</u>	<u>writel(x)</u>	<u>read(x)</u>

S₁ S₂

View Serializable Schedule

Testing of Serializability :- (Conflict Graph)

designing concurrency control schemes, must show that schedules generated by scheme are serializable

- Determining conflict serializability by precedence graph.

$$G = (V, E)$$

V = all transactions in Schedule
E = set of edges, all edges $T_i \rightarrow T_2$

- ① $T_1 - \text{writel}(x)$ before $T_2 - \text{read}(x)$ must hold one of three conditions
- ② $T_1 - \text{read}(x)$ before $T_2 - \text{executes writel}(x)$
- ③ $T_1 - \text{writel}(x)$ before $T_2 - \text{writel}(x)$.

9.12



<u>T1</u>	<u>T2</u>
read(x)	
write(x)	read(x)

S1

<u>T1</u>	<u>T2</u>
	read(x)
	write(x)

S2

precedence graph



<u>T1</u>	<u>T2</u>
read(A)	
	read(A)
	write(A)
	read(B)
write(A)	
read(B)	
write(B)	write(B)

- ① No - cycle - conflict → serializable
- ② cycle - Not conflict → serializable

Serializability order of transactions can be obtained through topological sorting - which determines a linear order.

Testing view serializability is rather complicated.

9.13

Concurrency control Techniques: concurrency control techniques that are used to ensure isolation property of concurrently executing transactions.

- ① Locking Techniques
- ② Time stamp techniques

① Locking Techniques:

A lock is a variable associated with a data item describes the status of item with respect to possible operations that can be applied.

Lock protocol: set of rules to be followed by each transaction.

Binary Locks: It can have two states $\begin{cases} \text{lock} \\ \text{unlock} \end{cases}$

Types of locks $\begin{cases} \text{Shared} - \text{read only} - S \\ \text{exclusive-write + read} - X \end{cases}$

Lock conversions:

① upgrade lock S to X

② down grade lock X to S.

using binary locks or read/write locks in transactions, doesn't guarantee serializability.

19.14

TWO-PHASE LOCKING (2PL) :- (Basic)

It has two rules (2 phases)

① Expanding or growing phase: - In this phase new locks on items can be acquired but none can be released.

② Shrinking phase: all existing locks can be released but no new locks can be acquired.

Eg:

T1

S-lock(B)
read(B)
unlock(B)
X-lock(A)
read(A)
 $A := A + 100$
write(A)
unlock(A)

T2

S-lock(A)
read(A)
unlock(A)
X-lock(B)
read(B)
 $B := B + 200$
write(B)
unlock(B)

T3

S-locks(B)
read(B)
~~X-lock(A)~~
= unlock(B)
read(A)
 $A = A + B$
write(A)
unlock(A)

T4

~~read(A)~~
S-lock(A)
read(A)
X-lock(B)
unlock(A)
read(B)
 $B = B + A$
write(B)
unlock(B)

S1

① Schedule S1 is non serializable schedule
That uses locking, NO 2PL.

② Schedule S2 is guaranteeing serializability
by 2PL but leads to deadlock.

9:15

2

There are a number of variations of 2PL.

- ① Basic 2PL
- ② Conservative 2PL (≈) static 2PL
- ③ Strict 2PL
- ④ Rigorous 2PL

② static 2PL: - It requires a transaction to lock all the items it access before the transaction begins execution, by predeclaring its read-set and write-set.

- if any of the predeclared items needed cannot be locked, the transaction doesn't lock any item, it waits until all the items are available for locking.

Advantage: It is deadlock free
Disadvantage: Difficult to use in practice.

③ strict 2PL:

- The most popular 2PL
- It guarantees strict schedule.
- Transaction doesn't release any of its X-locks until after it commits or aborts.
- Leads to deadlock
- strict schedule for recoverability.

Diligent 2PL:

- Transaction doesn't release any of its locks (stx)
- guarantees strict schedule
- easy to implement
- leads to deadlock.

Deadlock Handling:

deadlock occurs when each transaction is waiting for some item that is locked by some other transaction

Eg:	T ₁	T ₂
	s-lock(A) read(A)	
	wait for lock(B)	s-lock(B) read(B)
	x-lock(B)	x-lock(A)



wait for graph

① Deadlock prevention:

② Deadlock detection:

① Deadlock prevention :- The protocol which is used in conservative 2PL, requires lock all the items it needs in advance, if any of the items cannot be obtained, none of the items are locked. - The transaction waits then tries again to lock all the items.

9.17

If limits concurrency.

deadlock prevention protocol with time stamp $TS(T)$

The timestamps are typically based on the order in which transactions are started.

The rules

① wait-die ② wound-wait

① wait-die:

if $TS(T_1) < TS(T_2)$ T_1 starts before T_2

T_1 is older T_2 is younger

TS is less for T_1 & more for T_2 .

- if $TS(T_1) < TS(T_2)$

$TS(T_1) > TS(T_2)$

then T_1 is allowed to wait ↑ otherwise ~~abort~~ ^{T_1} (die) and restart it later with same TS .

② wound-wait:

if $TS(T_1) < TS(T_2)$

then abort ~~T_2~~ (wounds) and restart with same TS ; otherwise ~~T_2~~ T_1 is allowed to wait.

- both the cases younger is aborting

- deadlock free techniques.

19.18

5

Deadlock detection & Recovery

Deadlock can be detected by creating a directed graph called wait-for graph.

A deadlock exists in the system if and only if the wait-for graph contains cycle.

Deadlock Recovery:

needs to recover the system from deadlock by performing the following actions:

- ① Select a victim: To determine transaction that will be rolled back to break the deadlock.
The one to get selected should acquire minimum cost factors influence minimum cost are
 - How long the Transaction has compute
 - How many data Items the transaction has used.
 - How many data Items needed.
 - How many transaction will be in rollback.

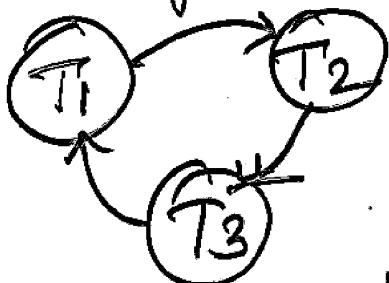
- ② Roll back: to determine how much should the victim transaction get rolled back.

- completely or partially.

- ③ Starvation: If ensure that same victim not to selected repeatedly.

Deadlock Detection:

The system checks if a state of deadlock actually exists, by wait-for graph it can be detected.



wait-for graph

③ Consistency control based on timestamp ordering

Timestamps: assigned TS in order in which the transactions are submitted to system.

- achieve deadlock-free without locks.
- each data item x can have two timestamp values.

④ Read-TS(x):

⑤ Write-TS(x):

⑥ read timestamp: This is the largest timestamp among all timestamps of transactions that have read data item x successfully.

⑦ write timestamp: This is the largest TS among all TS of transactions that have written x .

9. 20

The time stamp ordering protocol:

This is ensuring that any conflicting read + write operations are executed.

This operates as follows

① Suppose Transaction T_1 issues read(x)

ⓐ if $TS(T_1) < \text{write-TS}(x)$

then T_1 is rolled back since value of x is overwritten by a younger transaction. T_1 restarts. otherwise T_1 performs $\text{read}(x)$ & updates the Read-time stamp, to $TS(T_1)$.

② Suppose Transaction T_1 issues write(x)

if $TS(T_1) < \text{Read-TS}(x)$ or $\text{write-TS}(x)$

then T_1 is rolled back.

otherwise T_1 performs $\text{write}(x)$ and updates $\text{write-TS}(x)$ to $TS(T_1)$.

Advantages

- It ensures serializability

- It also ensures freedom from deadlock.

- No waiting for any transaction.

9.21

Thomas's write rule:

The modification to Time stamp ordering protocol, called Thomas's write rule.

① Read Rule is same as basic TO.

② Write Rule i.e T_i issues write(x)

- if $TS(T_i) < \text{Read-TS}(x)$, then T_i is rolled back.

~~otherwise~~ checks condition

if $TS(T_i) < \text{write-TS}(x)$ then T_i neither performs write operation nor is rolled back. T_i writes an absolute value of x . (outdated value).

otherwise system executes write operation by T_i and sets $\text{write-TS}(x)$ to $TS(T_i)$.

eg

<u>T₂</u>	<u>T₃</u>
read(x)	
writelx	writelnx

9.22

Granularity of Data Item:

A database item could be chosen to be one of the following:

- ① A database record
- ② A field value of a database record
- ③ A disk block
- ④ A whole file
- ⑤ The whole database.

The Granularity can affect the performance of concurrency control and recovery.

Granularity: The size of data item.

Granularity → fine granularity → small size
Granularity → coarse granularity → large size

Coarse granularity:-

The larger the data item size, the lower the degree of concurrency.

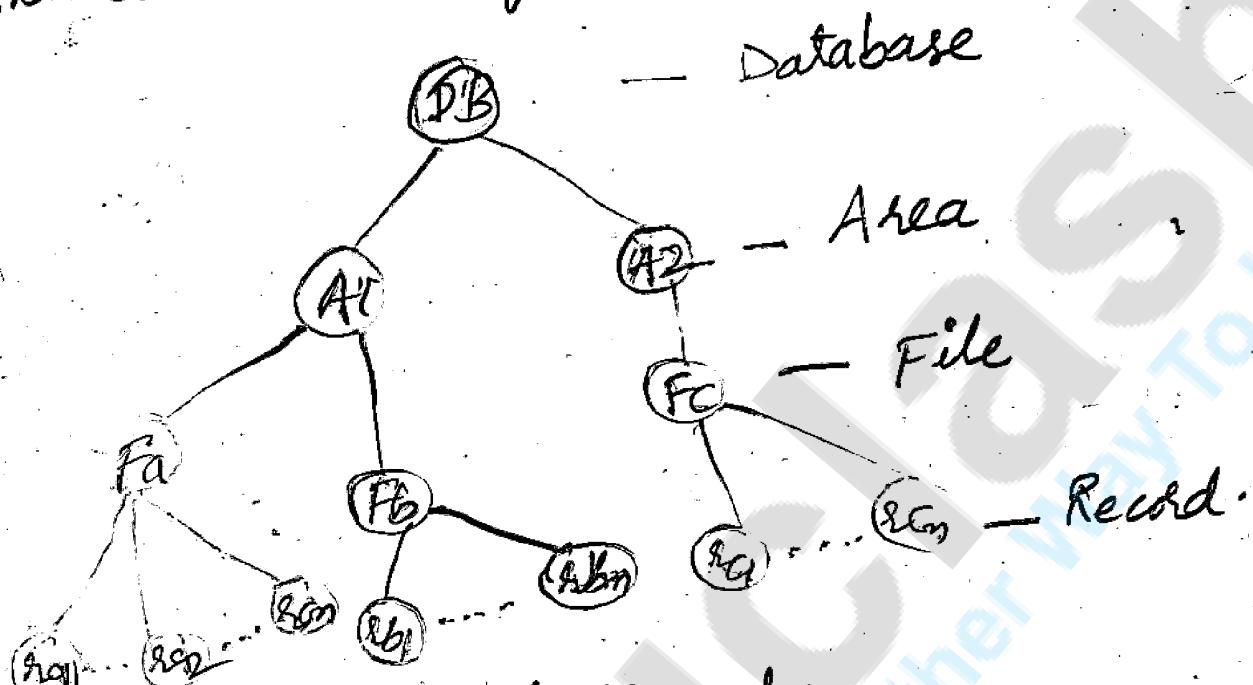
Fine granularity:-

The smaller the data item size, the higher the degree of concurrency. but more lock + unlock operations will be performed, causing a higher overhead and more storage space is required for lock table.

19.23

Multiple Granularity:

A another specialised locking strategy, called multiple granularity locking allows us set locks on objects that contains other objects.



Granularity hierarchy

To make multiple granularity level locking practical, additional types of locks, called intention locks are needed.

The idea behind intention locks is for a transaction to indicate, along the path from root to the desired node, what type of locks it will require from one of the node's descendants.

Types of intention locks: ① IS, ② IX ③ SIX
- Intention shared, Intention exclusive & shared
Intention exclusive.

Compatibility Matrix:

	IS	IX	S	SIX	X
IS	True	T	T	T	F
IX	T	T	F	F	F
S	T	F	T	F	F
SIX	T	F	F	F	F
X	F	F	F	F	F

→ To lock a node in IS (respectively SIX) mode a transaction must first lock all its ancestors in IS (respectively IX) mode.

The multiple-granularity locking protocol :-

~~It~~ ensures serializability.

The MGL protocol consists of the following rules:

- ① The lock compatibility must be adhered to
- ② The root of the tree must be locked first, in any mode.
- ③ A Node N can be locked by transaction T in S or IS mode only if the parent node N is already locked by transaction T in either IS or IX mode.
- ④ A Node N can be locked by transaction T in X, IX or SIX mode only if the parent in either IX or SIX.
- ⑤ It can lock a node only if it has not unlocked any node. (2PL)
- ⑥ Transaction T unlock nodes if child nodes are unlocked.

Recovery system of Transaction failure

Failure classification:

① Transaction failure:

- logical error
- system error

Log-Based Recovery:-

The most widely used structure for recording database modifications is the log. The log is a sequence of log records, recording all the update activities in the database.

There are several types of log records.

- insert log record, start, commit, abort.
- update log record.
- delete log record.

An update log record describes single database write
It has fields

Tid	data-item-id	old value	new value
-----	--------------	-----------	-----------

Represent

$\langle T_i \text{ start} \rangle; \langle T_i, x, v_1, v_2 \rangle$

$\langle T_i \text{ commit} \rangle; \langle T_i \text{ abort} \rangle$

There 3 techniques

- ① Deferred Data base modifications
- ② Immediate Data base modifications
- ③ Check points

whenever a transaction performs a write, it is essential that log record for that write be created before the database is modified.

① Deferred Database modification:

This Technique ensures transaction atomicity by recording all database modifications in the log, but deferring the execution of all write operations of transaction until the transaction partially commits.

Since failure may occur while the updating is taking place, all log records must written to the stable storage.

Eg. T₁

read(A)

$$A = A - 50$$

write(A)

read(B)

$$B = B + 50$$

write(B)

T₂

read(C)

$$C = C - 100$$

write C

<T₁, start>

<T₁, A, 950>

<T₁, B, 1050>

<T₁, commit>

~~<T₁, commit>~~

<T₂, start>

<T₂, C, 600>

<T₂, commit>

Redo(T₁)

sets new values

9.27

② Immediate Database modification

- This Technique allows database modifications to be output to the database while the transaction is still in the active state.
- Data modifications written by active transaction are called uncommitted modifications.
- In the event of a crash or a transaction failure, the system must use the old-value fields of the log records.

Eg: $\langle T_0 \text{ start} \rangle$ ②
 $\langle T_0, A, 1000, 950 \rangle$
 $\langle T_1, B, 2000, \underline{2050} \rangle$ ①
 $\langle T_1, \text{Commit} \rangle$
 $\langle T_2 \text{ start} \rangle$
 $\langle T_2, C, 700, \underline{600} \rangle$ ②
 $\langle T_2, \text{commit} \rangle$ — ③

using the log, the system can handle any failure that does not result in the loss of information in stable storage.

The recovery scheme uses two procedures -

- ① Undo (T_i) restores the value of all data items updated by T_i to the old value.
- ② Redo (T_i) sets of the value of all data items updated by T_i to the new values.

9.28

3. Check points:-

When a system failure occurs, we must refer the log to determine those transactions that need to be undo & redo.

There are 2 difficulties with searching log.

- ① Time consuming
 - ② more redo operations (longer time for recovery).
- To reduce overhead, check point is introduced.

During the execution, the system maintains the log, in addition periodically performs checkpoints.

Sequence of actions in checkpoint Technique

- ① output onto stable storage all log records currently residing in main memory
- ② output the disk all modified buffer blocks
- ③ output onto stable storage a log record <checkpoints>

Eg: <T₁, start>
<T₁, A, 1000, 950>
<T₁, B, 2000, 2050> → ①
<T₁ commit> <check point> → ②
<T₂, start>
<T₂, C, 700, 600> → ③
<T₂ commit> → ④

9.29

Advanced Recovery Techniques:

If DBMS gets crash, we can use ARIES recovery algorithm, which is simple, works well with a wide range of concurrency control mechanism & is used in an increasing no. of database systems.

ARIES:

It is designed to work with a steal, no-force approach.

After crash, restart proceeds in 3 phases.

- ① Analysis: Identify dirty pages in the buffer pool and active transactions at the time of crash.
- ② Redo: Repeat all actions, starting from an appropriate point in the log, restore the database.
- ③ Undo: undoes the actions of transactions that did not commit.

Execution History
Eg:

10 - update: T1 writes p5

20 - update: T2 writes p3

30 - T2 commit

40 - T2 end

50 - update T3 writes p1

60 - update T3 writes p3

* crash, Restart

Analysis:

Redo :

undo :

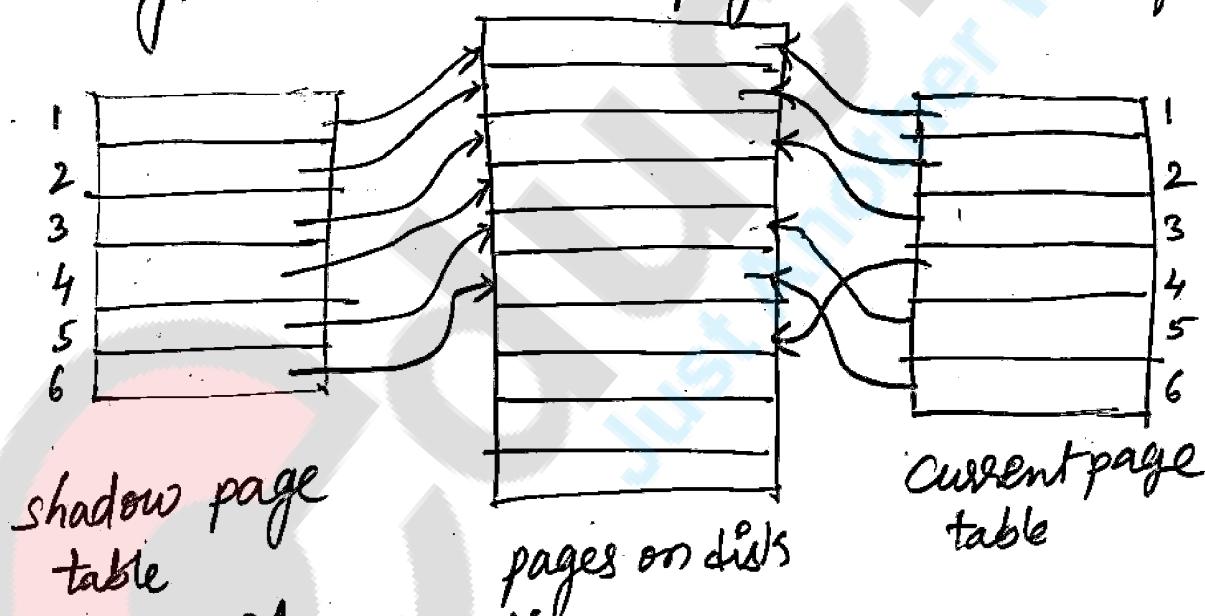
9.30

shadow paging

- shadow paging is an alternative to log-based recovery.
- this is useful if transactions executes serially.

Idea:

- Maintain two pages tables during the lifetime of a transaction - The current page table & shadow page table.
- store the shadow page table in nonvolatile storage i.e. the state of the database prior to transaction execution may be recovered.
- during execution shadow page is not modified.



→ To commit Transaction

- ① flush all modified pages to disk
- ② output current page table to disk
- ③ make the current page table, the new shadow page table
- ④ once pointer to shadow page table has been written
transaction is committed!

Advantages :-

- ① No overhead of writing log records.
- ② Recovery is trivial.

Disadvantages:

- ① Copying the entire page table is very expensive.
- ② Data gets fragmented.
- ③ After every transaction completion, the database pages containing old versions of modified data need to be garbage collected.
- ④ Hard to extend to allow transactions to run concurrently.

WAL : write-ahead log protocol :-

Before writing a page to disk, every update log record that describes a change to this page must be forced to stable storage.

19.32