

UNIT-V

**(Transaction Processing, Concurrency
Control, and Recovery)**

Concurrency Control

- How to design those schedules which ensure conflict serializability and other properties?
- **Timestamp Ordering Protocols**
- Basic idea of time stamping is to decide the order between the transaction before they enter in the system using a stamp (time stamp)
- With each transaction t_i , in the system, we associate a unique fixed timestamp, denoted by $TS(t_i)$. If a new transaction t_j , enters into the system with a timestamp $TS(t_j)$, then always $TS(t_i) < TS(t_j)$.

- Let $W\text{-timestamp}(Q)$ is the largest time-stamp of any transaction that executed $\text{write}(Q)$ successfully.
- Let $R\text{-timestamp}(Q)$ is the largest time-stamp of any transaction that executed $\text{read}(Q)$ successfully.
- **Read Operation Request**
- If $TS(T_i) < W\text{-timestamp}(Q)$, then T_i needs to read a value of Q that was already overwritten. Hence, the read operation is rejected, and T_i is rolled back.
- If $TS(T_i) \geq W\text{-timestamp}(Q)$, then the read operation is executed, and $R\text{-timestamp}(Q)$ is set to the maximum of $R\text{-timestamp}(Q)$ and $TS(T_i)$.

- **Write Operation Request**

- If $TS(T_i) < R\text{-timestamp}(Q)$, This means that the value of Q is already read. Hence, the write operation is rejected, and T_i is rolled back.
- If $TS(T_i) < W\text{-timestamp}(Q)$, then T_i is attempting to write an obsolete value of Q . Hence, this write operation is rejected, and T_i is rolled back.
- If $TS(T_i) \geq R\text{-timestamp}(Q)$, then the write operation is executed, and $W\text{-timestamp}(Q)$ is set to $\max(W\text{-timestamp}(Q), TS(T_i))$.
- If $TS(T_i) \geq W\text{-timestamp}(Q)$, then the write operation is executed, and $W\text{-timestamp}(Q)$ is set to $\max(W\text{-timestamp}(Q), TS(T_i))$.

- Problem: Dirty reads allowed; Order of commit not discussed; Hence Cascadelessness and recoverability are questionable.
- Advantage: Deadlock free, conflict and view serializability ensured
- Same for Thomas write rule except for conflict serializability.

- **THOMAS WRITE RULE**

- It is an improvement in time stamping protocol, which makes some modification and may generate those protocols that are even view serializable.
- It is a Modified version of the timestamp-ordering protocol in which Blind write operations may be ignored under certain circumstances.
- The protocol rules for read operations remain unchanged.

- Write Request Rules:
- if $TS(T_i) < W\text{-timestamp}(Q)$, then T_i is attempting to write an obsolete value of $\{Q\}$. Rather than rolling back T_i as the timestamp ordering protocol would have done, this $\{\text{write}\}$ operation can be ignored.
- This modification is valid because with $TS(T_i) < W\text{-timestamp}(Q)$, the value written by this transaction will never be read by any other transaction performing $\text{Read}(Q)$ ignoring such obsolete write operation is considerable.
- Thomas' Write Rule allows greater potential concurrency. Allows some view-serializable schedules that are not conflict serializable.

T_3	T_4	T_6
read(Q)	write(Q)	
write(Q)		
		write(Q)

Lock Based Protocols

- To ensure isolation is required that data items be accessed in a mutually exclusive manner.
- Locking is the most fundamental approach to ensure this.
- Idea is first obtain a lock on the desired data item then if lock is granted then perform the operation and then unlock it.

- Two modes of locking
- **Shared mode:** If transaction T_i has obtained a shared-mode lock (denoted by S) on any data item Q, then T_i can read, but cannot write Q.
- Any other transaction can also acquire a shared mode lock on the same data item.
- **Exclusive mode:** If transaction T_i has obtained an exclusive-mode lock (denoted by X) on any data item Q, then T_i can both read and write Q,
- Any other transaction cannot acquire either a shared or exclusive mode lock on the same data item.

- Compatibility Matrix: Helps Concurrency Control Manager to grant the locks as per rules.

Current State of lock of data items				
Requested Lock		Exclusive	Shared	Unlocked
	Exclusive	N	N	Y
	Shared	N	Y	Y
	Unlock	Y	Y	-

- Lock based protocol **do not ensure serializability** as granting and releasing of lock do not follow any order and any transaction any time may go for lock and unlock.

T_1	T_2
LOCK-X(A)	
READ(A)	
WRITE(A)	
UNLOCK(A)	
	LOCK-S(B)
	READ(B)
	UNLOCK(B)
LOCK-X(B)	
READ(B)	
WRITE(B)	
UNLOCK(B)	
	LOCK-S(A)
	READ(A)
	Unlock (A)

- Problem is that transactions release locks too soon, resulting in loss of total isolation and atomicity.
- To guarantee serializability, need an additional protocol concerning the positioning of lock and unlock operations in every transaction.
- Transaction follows 2PL protocol if all locking operations precede first unlock operation in the transaction
- Deadlock
- More duration of locked variable-> concurrency is compromised

- **Two phase locking protocol (Basic 2PL)**
- The protocol ensures that each transaction issue lock and unlock requests in two phases, note that each transaction will be 2 phased schedule.
- Growing phase- A transaction may obtain locks, but not release any locks.
- Shrinking phase- A transaction may release locks, but may not obtain any new locks.
- Initially a transaction is in growing phase and acquires lock as needed and in between can perform operation to reach lock point and once a transaction releases a lock, it can issue no more lock requests i.e. it enters the shrinking phase.
- The order of serializability is the order of T_i reaching at the lock point
- Problem: Deadlock; Non-recoverability, cascading rollback

T_1	T_2
LOCK-X(A)	
READ(A)	
WRITE(A)	
	LOCK-S(B)
	READ(B)
LOCK-X(B)	
READ(B)	
WRITE(B)	
	LOCK-S(A)
	READ(A)
	UNLOCK(B)
UNLOCK(A)	
UNLOCK(B)	
	UNLOCK(A)

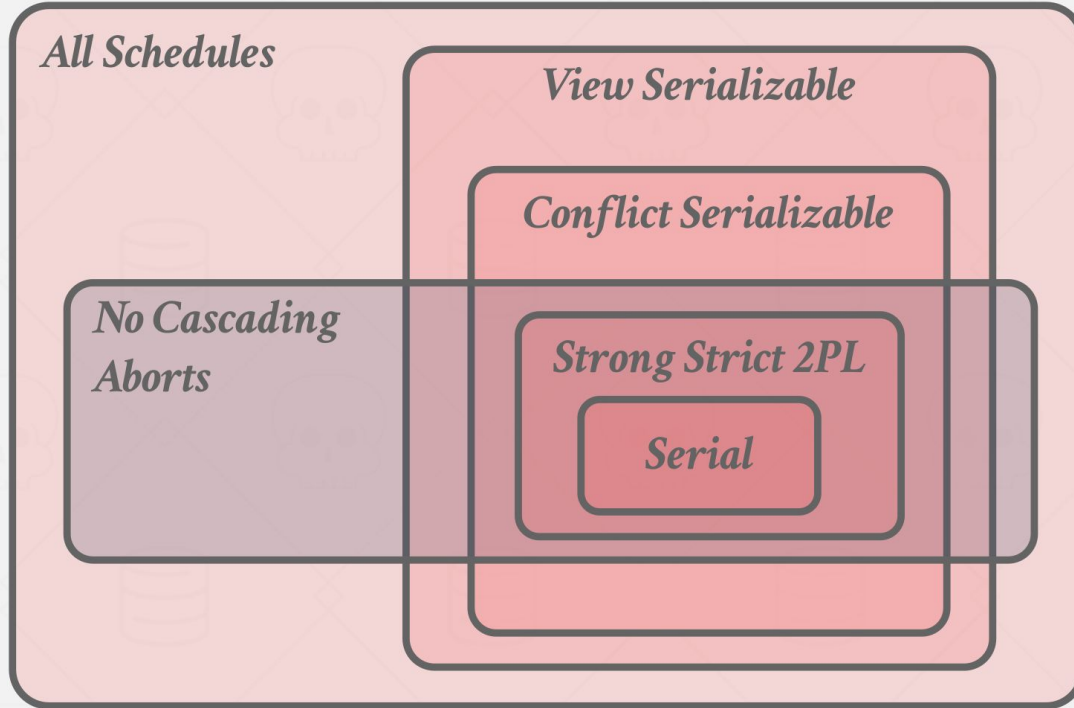
- **Conservative 2PL**

- The idea is there is no growing phase
- Transaction start directly from lock point,
- i.e. transaction must first acquire all the required locks then only it can start execution.
- If all the locks are not available then transaction must release the acquired locks and must wait.(Deadlock free)
- Shrinking phase will work as usual, and transaction can unlock any data item at anytime.
- Problem: we must have a knowledge in future to understand what is data required so that we can use it + recoverability and cascadelessness is compromised+ concurrency compromised because all locks held.

- **Rigorous 2PL**
- Requires that all locks be held until the transaction commits.
- This protocol requires that locking be two phase and also all the locks taken be held by transaction until that transaction commit.
- Hence there is no shrinking phase in the system.
- Recoverable and cascadeless schedule created (Hint: Dirty read avoided)
- Problem: Deadlock because growing phase exists + No unlock hence concurrency is compromised

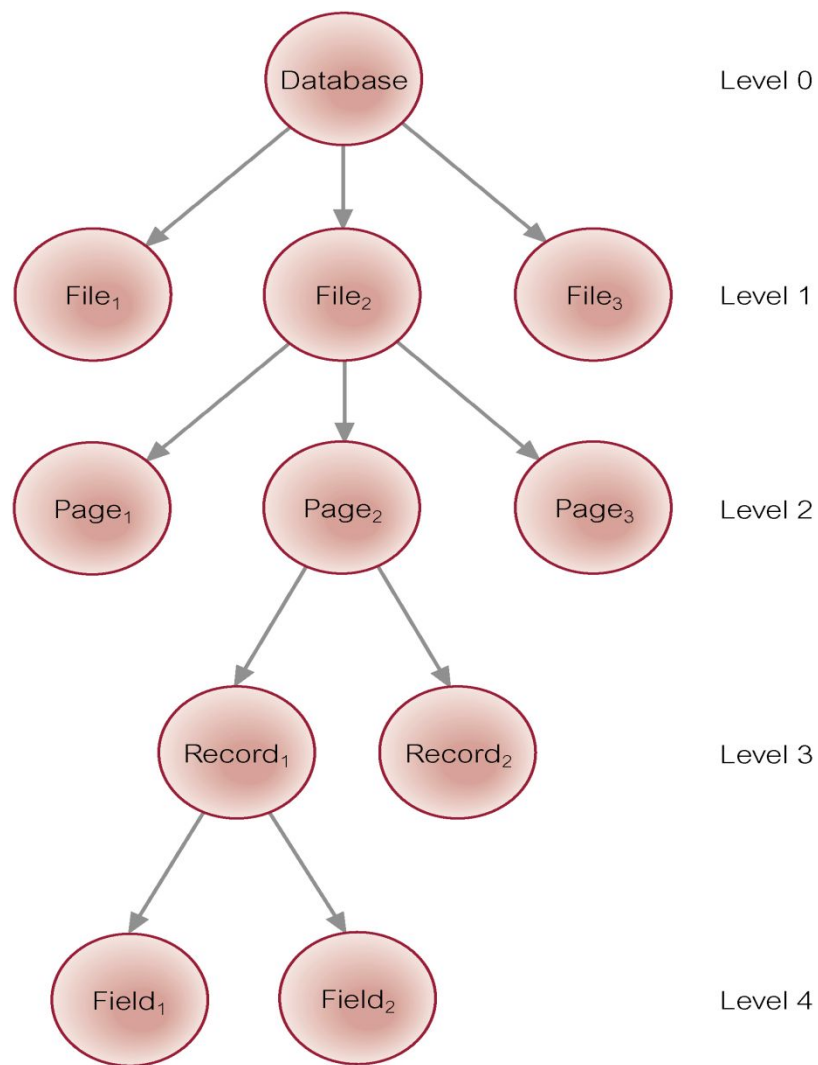
- **Strict 2PL**
- All exclusive-mode locks needed by transaction are held until that transaction commits.
- This requirement ensures that any data written by an uncommitted transaction are locked in exclusive mode.
- This prevents any other transaction from reading the data.
- Shared locks are released in partial shrinking phase
- So it is simplified form of rigorous 2PL
- Properties similar to Rigorous 2PL
- Problem: Deadlock because growing phase exists

UNIVERSE OF SCHEDULES

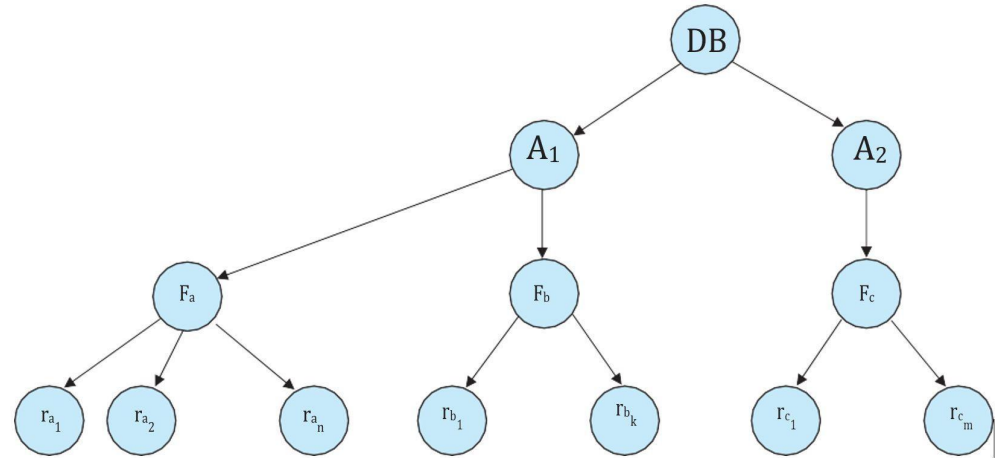


- **Granularity of Data Items (Multiple)**
- Size of data items chosen as unit of protection by concurrency control protocol.
- Ranging from coarse to fine:
 - The entire database.
 - A file.
 - A page (or area or database spaced).
 - A record.
 - A field value of a record.

- Tradeoff:
 - coarser, the lower the degree of concurrency;
 - finer, more locking information that is needed to be stored.
 - Best item size depends on the types of transactions.
- **Hierarchy of Granularity**
- Could represent granularity of locks in a hierarchical structure.
- Root node represents entire database, level 1s represent files, etc.
- When node is locked, all its descendants are also locked.
- DBMS should check hierarchical path before granting lock.



	IS	IX	S	SIX	X
IS	True	True	True	True	False
IX	True	True	False	False	False
S	True	False	True	False	False
SIX	True	False	False	False	False
X	False	False	False	False	False



- **Deadlock**

- An impasse that may result when two (or more) transactions are each waiting for locks held by the other to be released.

- .

Time	T ₁₇	T ₁₈
t ₁	begin_transaction	
t ₂	write_lock(bal_x)	begin_transaction
t ₃	read(bal_x)	write_lock(bal_y)
t ₄	bal_x = bal_x - 10	read(bal_y)
t ₅	write(bal_x)	bal_y = bal_y + 100
t ₆	write_lock(bal_y)	write(bal_y)
t ₇	WAIT	write_lock(bal_x)
t ₈	WAIT	WAIT
t ₉	WAIT	WAIT
t ₁₀	:	WAIT
t ₁₁	:	:

- Three general techniques for handling deadlock:
 - Deadlock prevention.
 - Timeouts.
 - Deadlock detection and recovery.
- **Deadlock Prevention**
- DBMS looks ahead to see if transaction would cause deadlock and never allows deadlock to occur.
- Could order transactions using transaction timestamps:
- Wait-Die - older transaction waits for younger one; older transaction is aborted if necessary. Allows recent jobs (short ones) to complete
- Wound-Wait - younger transaction waits for an older one. If older transaction requests lock held by younger one, younger one is aborted (wounded). Allows long jobs to complete.

- **Timeouts**

- Transaction that requests lock will only wait for a system-defined period of time.
- If lock has not been granted within this period, lock request timeout.
- Problem: DBMS assumes transaction may be deadlocked, even though it may not be, and it aborts and automatically restarts the transaction.

- **Deadlock Detection and Recovery**
- DBMS allows deadlock to occur but recognizes it and breaks it.
- Usually handled by construction of wait-for graph (WFG) showing transaction dependencies:
 - Create a node for each transaction.
 - Create edge $T_i \rightarrow T_j$, if T_i waiting to lock item locked by T_j .
 - Deadlock exists if and only if WFG contains cycle.
 - WFG is created at regular intervals.

- **Recovery from Deadlock Detection**
- Several issues:
- choice of deadlock victim;
- how far to roll a transaction back;
- avoiding starvation
- Where one transaction (say the longest one) is not explicitly prevented from progressing, but actually never does progress
- Wait-Die: long transactions repeatedly aborted to give priority to shorter transactions
- Wound-Wait: short transactions never started because longer transactions are always in process

- **Optimistic Techniques**

- Locking and Timestamping techniques focus on preventing conflicts
- In some systems, conflict is rare and more efficient to let transactions proceed without delays, so no guarantee of serializability.
- At commit, check is made to determine whether or not a conflict has occurred.
- If there is a conflict, transaction must be rolled back and restarted.
- Potentially allows greater concurrency than traditional protocols.
- Best for systems having more read requests rather than write requests

- **Database Recovery**

- Process of restoring database to a correct state in the event of a failure.
- Need for Recovery Control
- Two types of storage: volatile (main memory) and nonvolatile.
- Volatile storage does not survive system crashes.
- Stable storage represents information that has been replicated in several nonvolatile storage media with independent failure modes.

- **Transactions and Recovery**

- Transactions represent basic unit of recovery.
- Recovery manager responsible for atomicity and durability.
- If failure occurs between commit and database buffers being flushed to secondary storage then, to ensure durability, recovery manager has to redo (rollforward) transaction's updates.
- If transaction had not committed at failure time, recovery manager has to undo (rollback) any effects of that transaction for atomicity.
- Partial undo - only one transaction has to be undone.
- Global undo - all transactions have to be undone.

- **DBMS should provide following facilities to assist with recovery:**
- Backup mechanism, which makes periodic backup copies of database.
- Logging facilities, which keep track of current state of transactions and database changes.
- Checkpoint facility, which enables updates to database in progress to be made permanent.
- Recovery manager, which allows DBMS to restore database to consistent state following a failure.

- **Log File**

- Contains information about all updates to database:
- Transaction records.
- Checkpoint records.
- Often used for other purposes (for example, auditing).
- Log file may be duplexed or triplexed.
- Log file sometimes split into two separate random-access files.
- Potential bottleneck; critical in determining overall performance.

- **Checkpoint**

- Point of synchronization between database and log file. All buffers are force-written to secondary storage.
- Checkpoint record is created containing identifiers of all active transactions.
- When failure occurs, redo all transactions that committed since the checkpoint and undo all transactions active at time of crash.

- The system log, or transaction log, records all the changes or activities happening in the database, ensuring that transactions are durable and can be recovered in the event of a failure.
- Different types of log records represent various stages or actions in a transaction.
- $\langle T_i \text{ start} \rangle$: This log entry indicates that the transaction T_i has started its execution.
- $\langle T_i, X_j, V_1, V_2 \rangle$: This log entry documents a write operation. It states that transaction T_i has changed the value of data item X_j from V_1 to V_2 .
- $\langle T_i \text{ commit} \rangle$: This log entry marks the successful completion of transaction T_i .
- $\langle T_i \text{ abort} \rangle$: This log entry denotes that the transaction T_i has been aborted, either due to an error or a rollback operation.
- Before any write operation modifies the database, a log record of that operation needs to be created. This is to ensure that in case of a failure, the system can restore the database to a consistent state using the log records.

- **Three Main Recovery Techniques**
- Deferred Update
- Immediate Update
- Shadow Paging

- **Deferred Update**

- Updates are not written to the database until after a transaction has reached its commit point. (Delayed Commit)
- If transaction fails before commit, it will not have modified database and so no undoing of changes required. (Less complex recovery)
- May be necessary to redo updates of committed transactions as their effect may not have reached database.
- Reduced I/O operations, as changes are batched and written at once during the commit, saving on intermediate I/O operations.

- **Immediate Update**

- Updates are applied to database as they occur (don't wait for commit).
- As before, to redo updates of committed transactions following a failure (may not have reached DB).
- May need to undo effects of transactions that had not committed at time of failure.
- Essential that log records are written before write to database. Write-ahead log protocol.
- Increased I/O operations, as each change triggers immediate write operations, leading to more frequent I/O operations.

- **Shadow Paging**
- Maintain two page tables during life of a transaction: current page and shadow page table.
- When transaction starts, two pages are the same.
- Shadow page table is never changed thereafter and is used to restore database in event of failure.
- During transaction, current page table records all updates to database.
- When transaction completes, current page table becomes shadow page table.
- Four phases: Initialization->modification->commit->recovery