

122. second normal form

entity C (p#, pname, empno, ename, job, chg/day, nday)

p# → pname - I

p#, empNo → pname, ename, Job, chg/day, nday - BCNF

empNo → ename, Job, chg/day - F

Job → chg/day - II

p#	pname	empNo	ename	Job	chg/day	nday
1	I	101	R	Acc	1000	30
1	I	103	A	Ma	4000	15
1	I	105	B	cl	500	30
1	I	106	C	cl	500	25
2	H	101	R	Acc	1000	30
2	H	103	A	Ma	4000	10
2	H	107	B	cl	500	30
2	H	109	S	cl	500	30
3	E	105	D	cl	500	30
3	E	103	A	Ma	4000	30
3	E	108	E	cl	500	15
4	P	104	X	client	1000	30

p#	pname
1	I
2	H
3	E
4	P

EmpNo	Ename	Job	Chg/Day
101	R	Acc	1000
103	A	Ma	4000
105	B	cl	500
106	C	cl	500
107	B	cl	500
109	S	cl	500
108	E	cl	500
104	X	client	1000

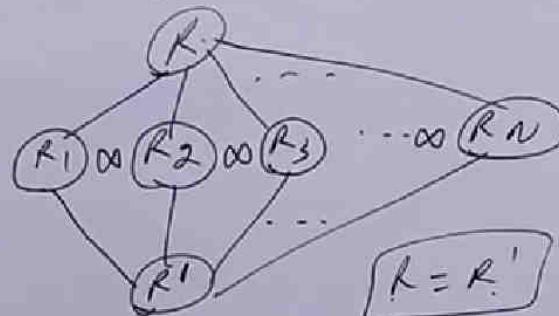
p#	empNo	nday
1	101	30
1	103	15
:	:	:

RA → ✓
IA → ✓
EA → ✓
YA → ✓
DP → ✓

138. 5th normal form

5th Normal Form (PJNF)

→ 5th Normal form is based on concept of Join Dependencies.



- I NF → Atomic values
- II NF → Partial dependencies
- III NF → Transitive dependencies
- IV NF → M.V.P
- V NF → Join dependencies

$$\begin{array}{c} R(A,B,C,D) \\ \downarrow \\ R_1(A,B,C,D) \end{array}$$

Join Dependency - ✓

→ A table is said to have Join dependency if it is possible to decompose it into 2 or more subtables & by performing natural join b/w all these tables we get original table

→ Join dependency is trivial if one of the subtable is original table itself
5th N.F. ✓

→ For each Nontrivial Join Dependency each Subtable must be Superkey of R.

IS 3NF Perfect?

Faculty:-

Fid	Office	Dept	Rank	DOH
f101	2	1	3	4
f101	6	5	7	8
f102	2	1	3	6
f105	7	5	6	4
f103	2	1	7	8
f104	6	5	4	9

faculty(fid, office, dept, rank, doh)

fid, dept → office, rank, doh $\xrightarrow{\text{BCNF}}$ Office → Dept $\xrightarrow{\text{II}}$ CK $\xrightarrow{\text{I}}$ fid, dept
CK $\xrightarrow{\text{II}}$ fid, office

A II

fid, office, Dept, Rank, DOH

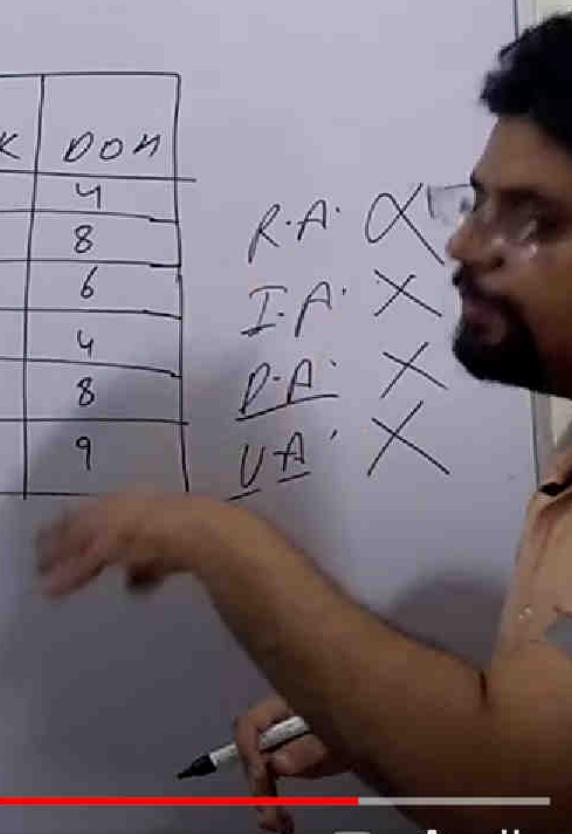
(Office, Dept)

fid office Rank DOH

Fid	Office	Rank	DOH
f101	2	3	4
f101	6	7	8
f102	2	3	6
f105	7	6	4
f103	2	7	8
f104	6	4	9

R.A. X
I.A. X
D.A. X
U.A. X

Office	Dept
2	1
6	5
7	5



139. introduction to transactions

TRANSACTIONS

- Transaction is a part of HLL program which is either completely executed or not executed at all
- We can identify this part by 2 System calls i.e "begin transaction" & "end transaction"
- whatever we write in between begin transaction & end transaction constitute a single transaction
- A HLL program can have more than one transaction also.
- We can write this program in any HLL like COBOL, C, C++, Java, SQL etc



- Transaction accesses Data base "item".
- Database item can be cell of table or single tuple or collection of tuples or entire table or entire database(collection of tables)
- Size of database item is known as granularity
- For simplifying we assume granularity to be cell of table.
- In simplest scenario transaction will read some items & write some items.
- If transaction only read items but does not write any item in database then it is called as readonly transaction

139. introduction to transactions

TRANSACTIONSRead item(n) :- ??

① Get disk address of block that contain the item n.

② Read on block (Transfer block from H.D. to RAM)

③ copy the value of item n from RAM into program variable n.

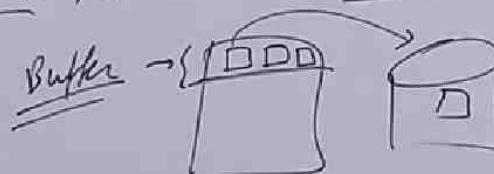
Write item(n) :-

① same ✓

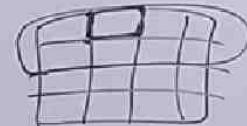
② same ✓

③ Copy value of program variable n into block in RAM containing item 'n'

④ Write this block on H.D.

Buffer:-

int n;
n ← item n
n + 50
n - 100

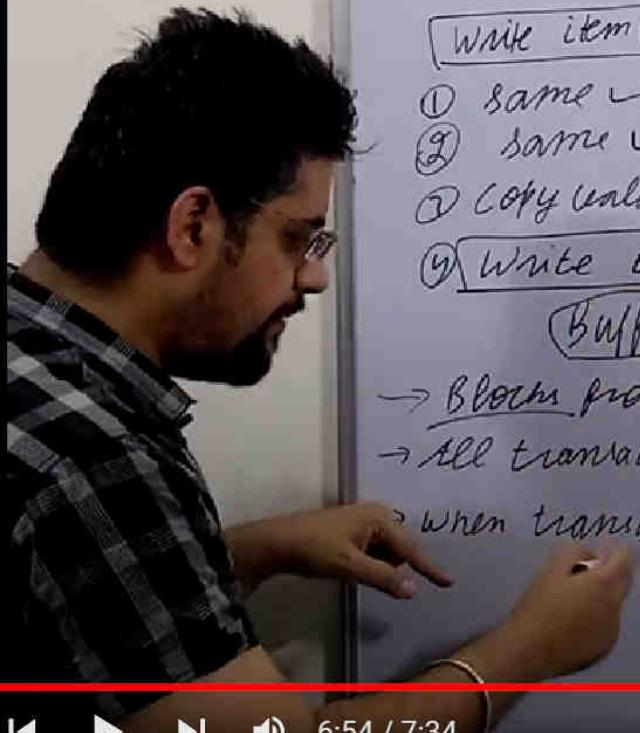


item n ← n

→ Blocks from HD are stored at a common place in RAM called buffer.

→ All transactions will read & write data items in buffer only.

→ When transaction commits then respective block from buffer is stored in database.



140. concept of read and write item

Read item(n) \rightarrow Name of item (DB is a collection of named items)

- ① Get disk address of block that contain item n. (How?? implementation dependent)
- ② Read in block (i.e transfer block from H.D. to RAM) [Why entire block??]
- ③ Copy value of item n from RAM into program variable n

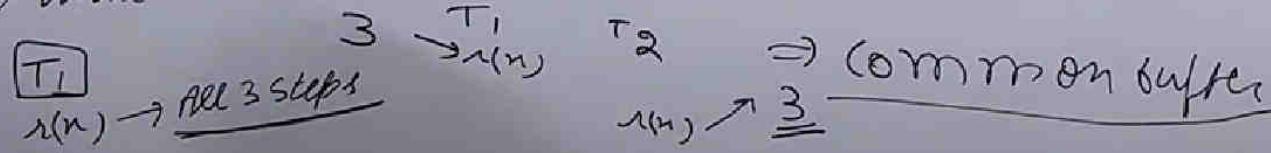
Write item(n) :-

① same

② same

③ copy value of program variable n into block in RAM containing item n.

④ Write block in H.D. [This step make changes in Database]



$r(n)$ \rightarrow may or may not need all 3 steps

If we assume block containing item n is still in RAM then step 1 & 2 not needed
block not in RAM : it has been replaced by some other block.

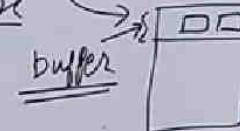
if T_2 do $r(n)$ immediately after $r(n)$ of T_1 ??

int n;

n \leftarrow item n

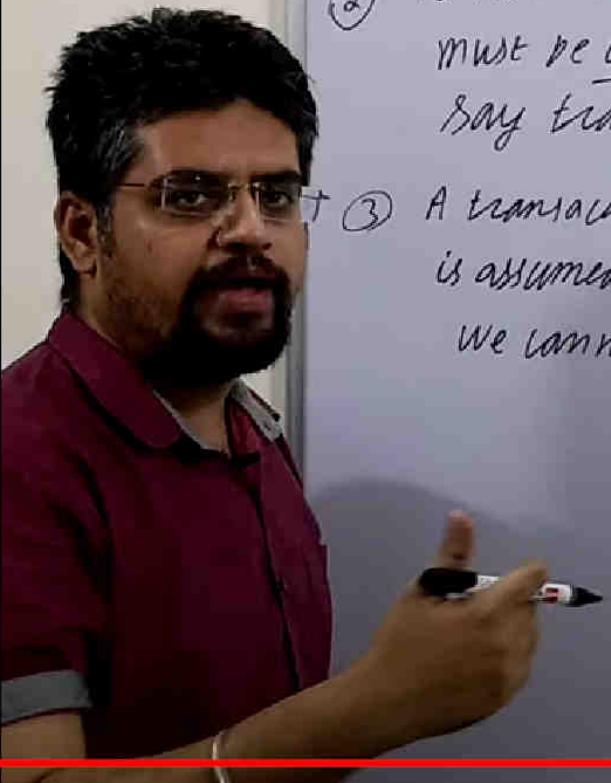
n = 500;

n + 200;



141. few terms in transactions

Few terms in transactions.

- 
- ① If any kind of failure occurs & transaction is not completed successfully such a transaction is termed as ABORTED.
 - ② To maintain atomicity property, changes made by aborted transaction on database must be undone. Once the changes made by aborted transaction is undone we say transaction has been rolled back. [Responsibility of recovery scheme to manage transaction aborts]
 - ③ A transaction that completes its execution successfully is said to be committed & its effect is assumed to be permanently recorded in the database.
We cannot undo the effect of committed transaction by aborting it

142. acid properties

TRANSACTIONS4 basic properties of transactions

① Atomicity → Transaction will be either completely executed or not executed at all.

② Consistency

③ Isolation

④ Durability

"ACID" properties

→ If before executing transaction database is in consistent state

then after execution of transaction it must remain in consistent state only.

→ Consistent state means all constraints are fulfilled.

$$\text{Eg. } A + B = 100$$

$$\begin{array}{l} T_2 \\ \text{r(A)} \\ \text{A} = \text{A} + 10 \\ \text{w(A)} \end{array}$$

$$\begin{array}{l} \times \\ \text{A} + \text{B} = 110 \end{array}$$

Commit

Transfer 500Rs from

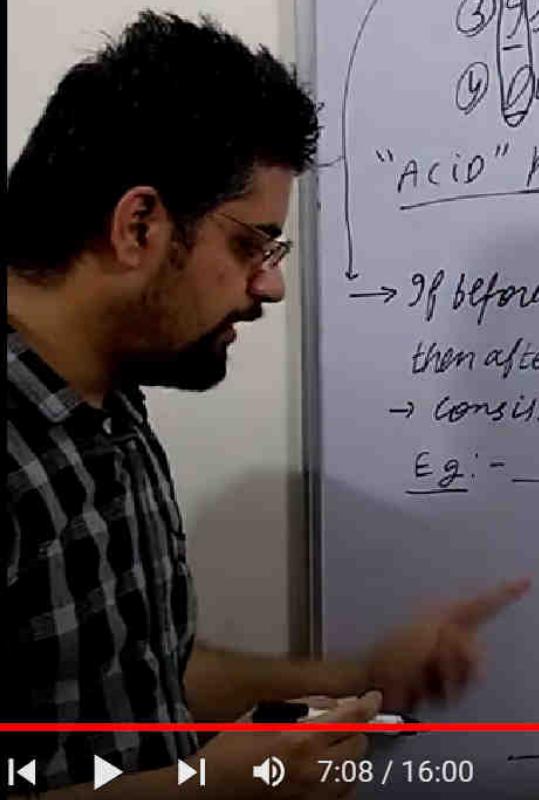
amount A to B

T₁
r(A)

$A = A - 500$
w(A)

r(B)
 $B = B + 500$

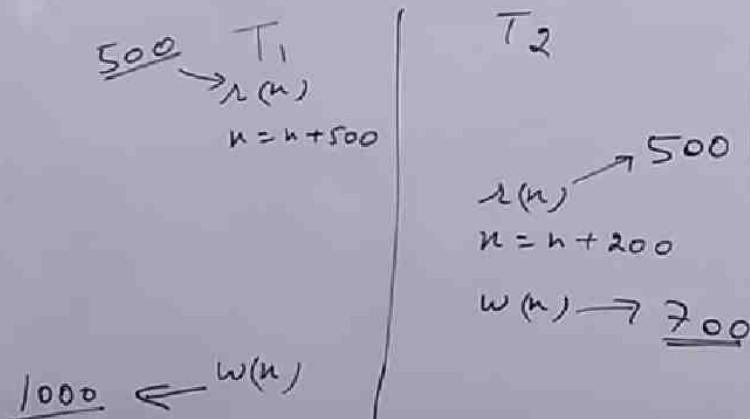
w(B)
Commit



142. acid properties

TRANSACTIONSIsolation:-

For every pair of transactions t_1 & t_2 it appears to t_1 that either t_1 has completed its execution before t_2 starts or t_1 will start its execution after t_2 completes its execution.

Durability:-

Once transaction commits (i.e. completed successfully) changes it has made to the database must persist even if there are some system failures.

Property

Atomicity

Responsibility

Transaction management component of DBMS

Consistency

→ Application programmer

Isolation

→ Concurrency control component

Durability

→ Recovery component.

143. why concurrency control is needed

why concurrency control is needed?

→ We want to execute transactions in interleaved fashion so that their execution becomes faster

T_1	$r(n) \rightarrow I/O$
	$n = n + 50 \rightarrow CPU$
	$w(n) \rightarrow I/O$

T_2	$r(y) \rightarrow I/O$
	$y = y + 100 \rightarrow CPU$
	$w(y) \rightarrow I/O$

→ So transactions mainly do I/O operations
2 hence executing them in interleaved fashion will give more efficiency.

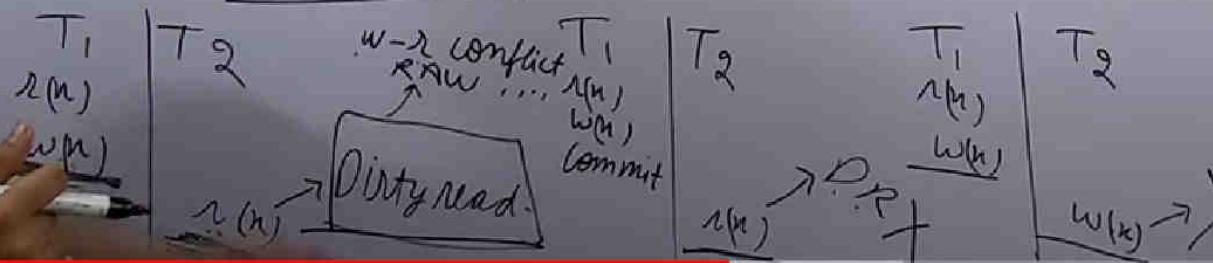
Few problems :-

→ Temporary update problem :-
Dirty Read problem :-

→ If Transaction reads item written by uncommitted transaction then that read operation is called dirty read.

T_1	$r(n)$	T_2	$r(y)$	T_1	$r(n)$	T_2	$r(y)$
					$\xrightarrow{500} w(n)$		$\xrightarrow{700} w(y)$
							$\xrightarrow{1000} commit$

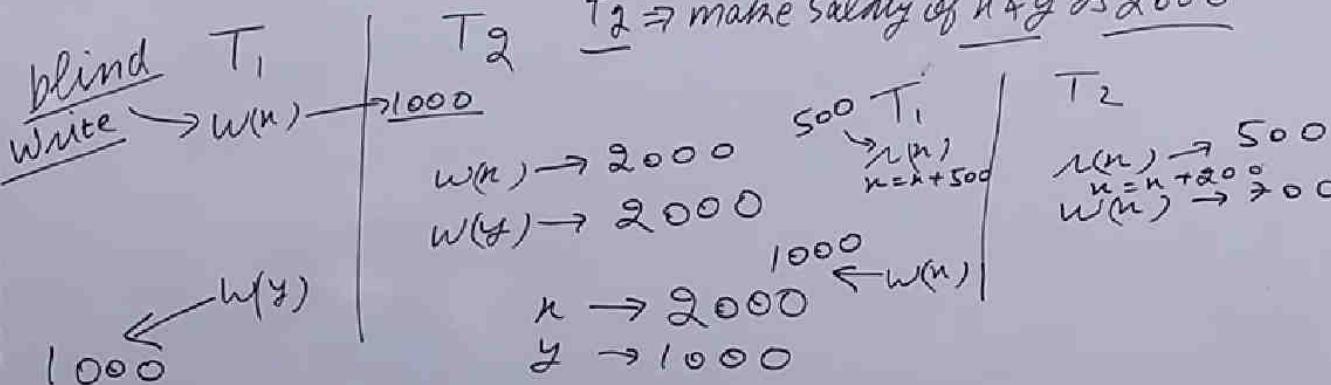
X → r(y) Abort



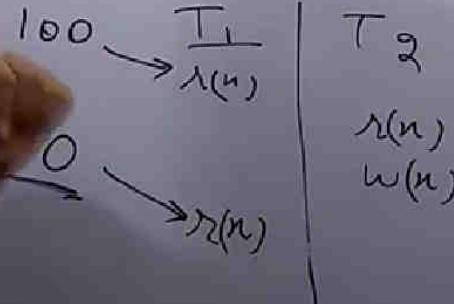
143. why concurrency control is needed

why concurrency control is needed?

② Lost update problem:- - $T_1 \rightarrow$ make salary of $n & y$ as 1000



③ Unrepeatable read problem:-



(g) Incorrect Scansory problem.

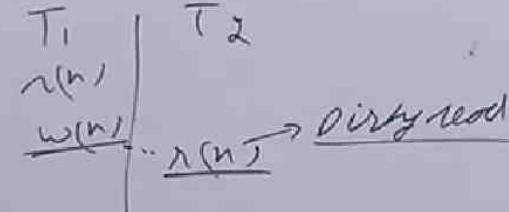
T_1	T_2	T_3	T_4	T_5
→ h	+	+	+	+
0	4	2(4)	1(6)	0(6)

Scansory

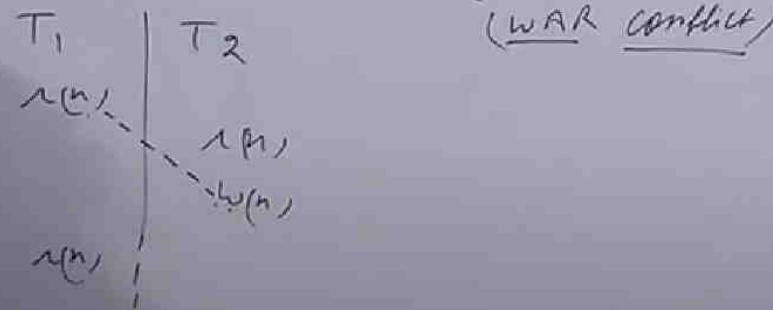
143. why concurrency control is needed

Summary

- ① dirty read problem or Temporary update problem (w-r conflict)
(RAW conflict)



- ② Unrepeatable read problem (r-w conflict)
(WAR conflict)



- ③ Lost update problem or OVERRWITING UNCOMMITTED DATA (w-w conflict)
(WW conflict)



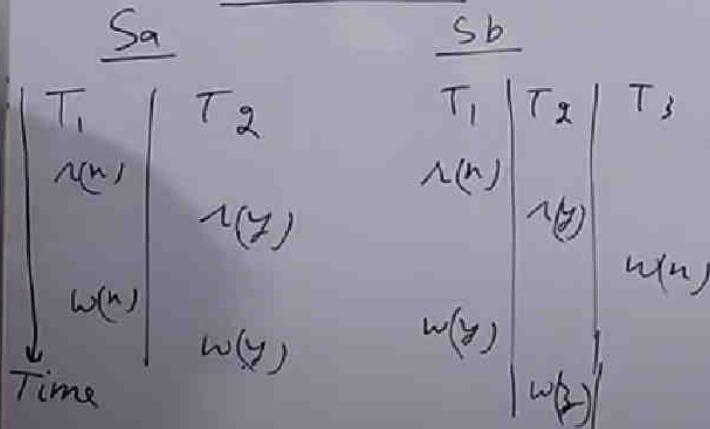
144. schedule

Concurrency control techniques:-

→ CCT are used to ensure consistency of database when transactions are executed concurrently or in interleaved fashion [final answer must be correct].

Few terms:-① Schedule :-

Set of transactions that are executed together in this system.



Sa: r₁(x); r₂(y); w₁(x); w₂(y)

Sb: r₁(x); r₂(y); w₃(z); w₁(y); w₂(z)

145. Types of schedules and intro to serial schedule

- Serial Schedule
- Serializable ...
- Conflict equivalent schedule
- Conflict Serializable
- View Equivalent
- View Serializable
- Recoverable schedule -
- Cascadable " "
- Strict

① Serial schedule :-

T ₁	X	T ₁	X	T ₁	✓
T ₂		T ₂		T ₂	T ₂
r(n)		r(n)		r(n)	r(n)
w(n)		w(n)		w(n)	w(n)
Commit		IL		abort	
r(n)	T ₁ → T ₂	T ₂ → T ₁	w(n)		
w(n)					
Commit					

$n \times n - 1 \times n - 2 \dots 1$ commit

- A schedule S is said to be serial if no transaction starts till running transaction ends [commit or abort]
- So only one transaction is active at a time & commit or abort of active transaction initiates execution of next transaction

→ If we consider each transaction to be correct (i.e consistency property holds) then all possible serial schedules ($n!$) are correct.

→ Main disadvantage of serial schedules is that if transaction waits for an I/O then CPU cannot be switched to another transaction. Hence lot of CPU cycles are wasted.

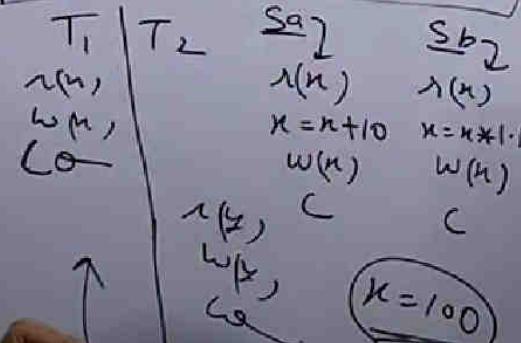


$$\frac{3 \times 2 \times 1}{3 \times 2 \times 1} = \underline{\underline{6}} \quad \frac{n \times n - 1 \times n - 2 \dots 1}{n \times n - 1 \times n - 2 \dots 1} = \underline{\underline{n!}}$$

146. serializable schedule

Types of schedules

- Serial Schedule ✓
- Serializable ... ✓
- Conflict equivalent schedule
- Conflict serializable
- View equivalent
- View serializable
- Recoverable schedule -
- Cascadable .. . -
- Strict .. . -

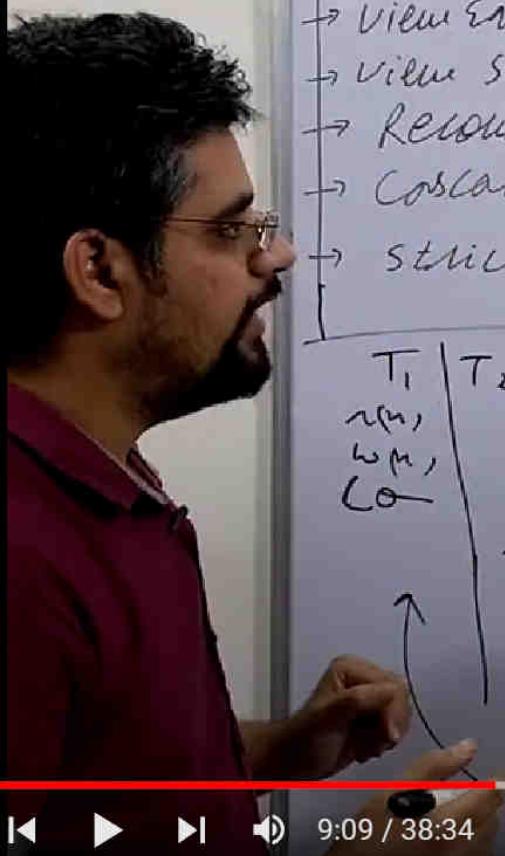
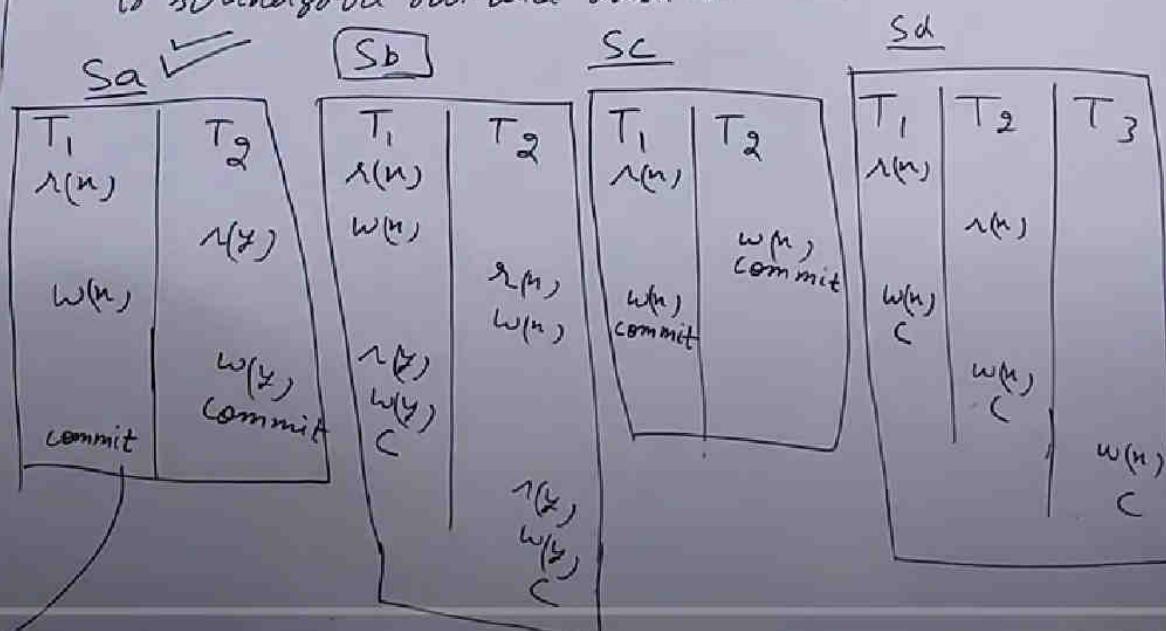


$$x = n + 10$$

$$x = n + 11$$

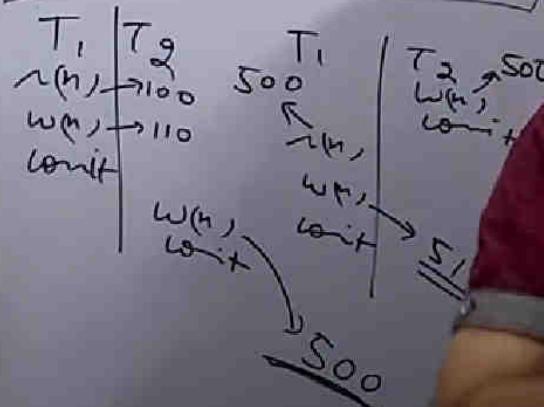
Serializable schedule

- A schedule S is said to be serializable if its effect on database is guaranteed to be same as some serial schedule.
- This means if we have a schedule of ' n ' transactions then its result must be same as one of the $n!$ possible serial schedules.
- From defn we can easily say that every serial schedule is serializable but vice versa not true.



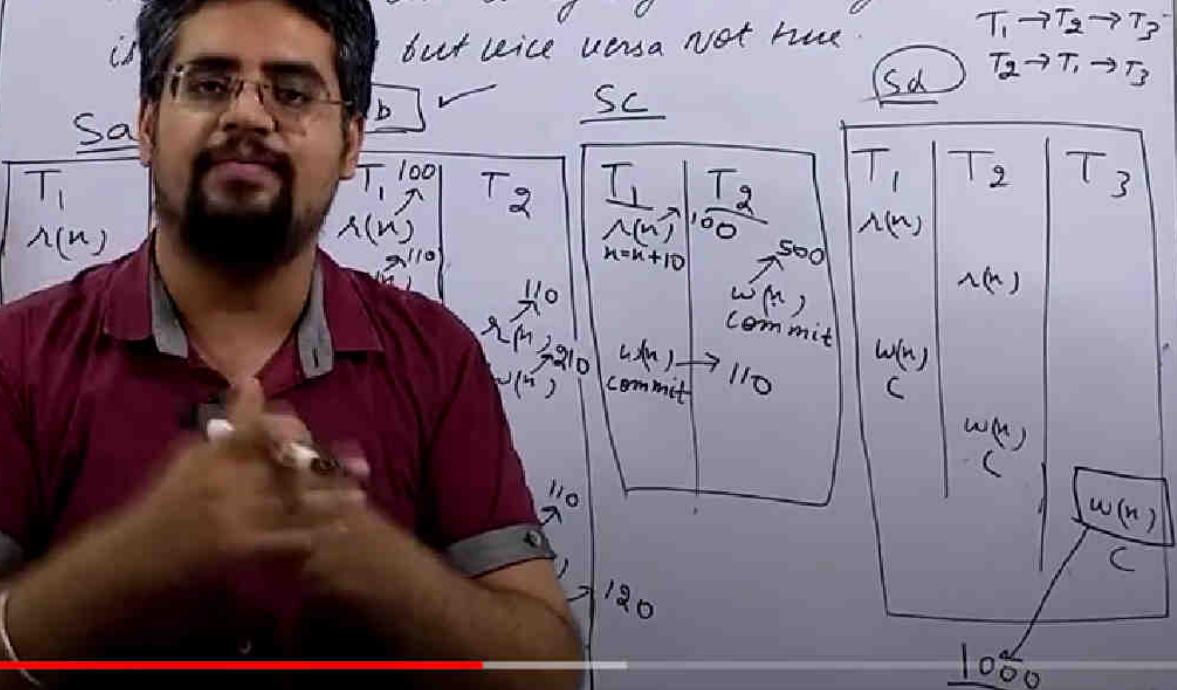
146. serializable schedule

- Serial Schedule ✓
- Serializable ... ✓
- Conflict equivalent schedule
- Conflict serializable
- View equivalent
- View serializable
- Recoverable schedule
- Cascadable " "
- Strict



Serializable schedule

- A schedule S is said to be serializable if its effect on database is guaranteed to be same as some serial schedule.
- This means if we have a schedule of ' n ' transactions then its result must be same as one of the $n!$ possible serial schedules.
- From this we can easily say that every serial schedule is serializable but vice versa not true.



146. serializable schedule

Types of schedules

- Serial Schedule ✓
- Serializable ... ✓
- Conflict equivalent schedule
- Conflict serializable
- View equivalent
- View serializable
- Recoverable schedule
- Cascadeness " ..
- Strict

$T_1 \begin{array}{|c} \text{r}(n) \\ \text{w}(m) \end{array}$ | $T_2 \begin{array}{|c} \text{r}(k) \\ \text{w}(l) \end{array}$

$\Rightarrow \frac{2!}{2!}$

Serializable schedule

→ Checking schedule is serializable or not is an undecidable problem i.e. no one in this world will ever be able to write algorithm for this problem.

Non serial schedules ($>> n!$)

Correct/Serializable

$$\boxed{A, B, C, D} \quad \frac{4!}{2!2!}$$

In correct/Non serializable

$$\boxed{A, B, C, D} \quad \frac{4!}{3!1!} = 1$$

Serial \Rightarrow n Transactions $\sim \frac{n!}{n!}$



$$\frac{n!}{n!} \cdot \frac{n-1!}{n-1!} \cdot \frac{n-2!}{n-2!} \cdots \approx \frac{1}{n!}$$

$$\frac{T_1 \rightarrow n}{T_2 \rightarrow m} \quad \frac{(n+m)!}{n! * m!}$$

$$\frac{4!}{2!2!} \cdot \frac{3!}{2!1!} \cdot \frac{2!}{1!1!} = \frac{4 \times 3 \times 2 \times 1}{2 \times 1 \times 2 \times 1} = 6$$

Total schedule: $\frac{n_1! n_2! n_3! \cdots n_n!}{n_1! n_2! n_3! \cdots n_n!}$

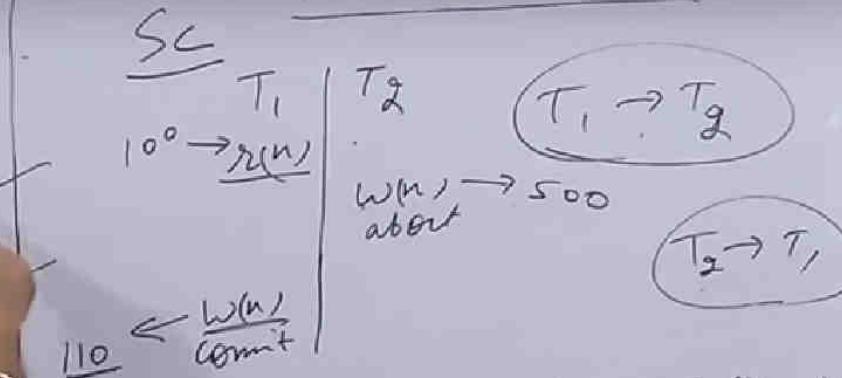
146. serializable schedule

Types of schedules

- Serial Schedule
- Serializable ...
- Conflict Equivalent Schedule
- Conflict Serializable

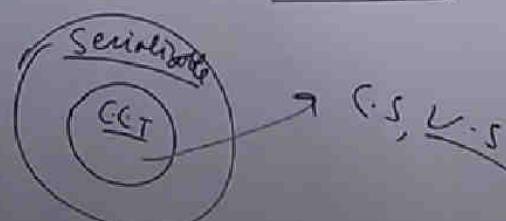
- View Equivalent
- Non Serializable
- Derivable schedules
- ...

Serializable schedule



→ C.C.T must guarantee serializable schedules.

→ But checking, schedule is serializable or not is an undecidable problem. So, we have other types of schedules which are subset of serializable schedules that can be checked by C.C.T.



C.S., V.S.

→ C.C.T allows only correct schedules to be executed in the system. But all correct schedules are not C.C.T.



37:57 / 38:34



146.2 gate 2012 question on serializability

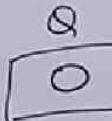
$T_1: R(P)$

$R(Q)$

if $P=0$ then $Q=Q+1$

$W(Q)$

T_1 $Q=1$



$T_2: R(Q)$

$R(P)$

if $Q=0$ then $P=P+1$

$W(P)$

Any non serial interleaving of

T_1 & T_2 leads to

$T_1 \rightarrow T_2 \Rightarrow Q=1, P=0$

$T_2 \rightarrow T_1 \Rightarrow P=1, Q=0$

② Before T_2 last operation
1st operation of T_1 is done \Rightarrow

GATE 2014

T_1

T_2

$R(u)$

$W(u)$
commit

$W(u)$
commit

$W(y)$

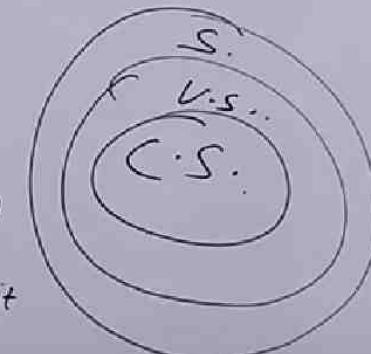
$R(z)$

Commit

T_3

T_4

$R(u)$
 $R(y)$
commit



Conflict Serializable?
Recoverable?

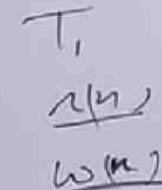
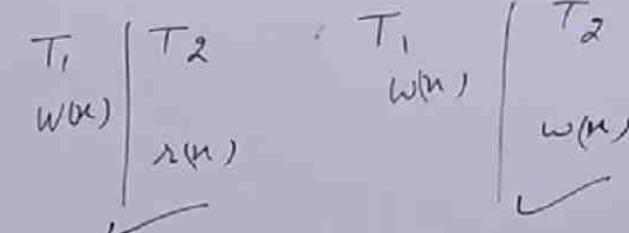
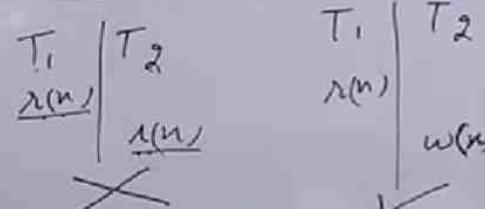
$P=1, Q=1$

147. conflict equivalent schedule

Conflict equivalent Schedules

minimum 2 schedules

Conflicting operations



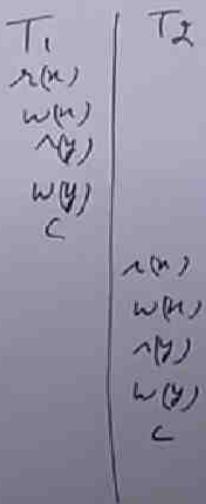
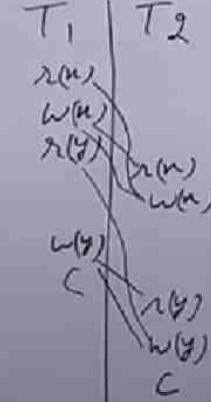
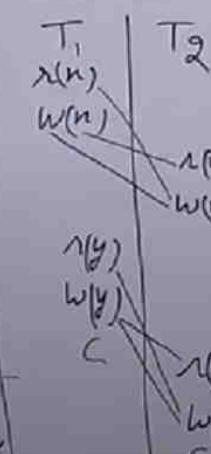
2 operations are said to be conflicting if following 3 cond'n are satisfied :-

- ① They belong to different transactions
- ② They access same database item
- ③ At least one of them must be write operation

Defn.: 2 schedules are said to be conflict equivalent

if order of every conflicting operation is same in both schedules

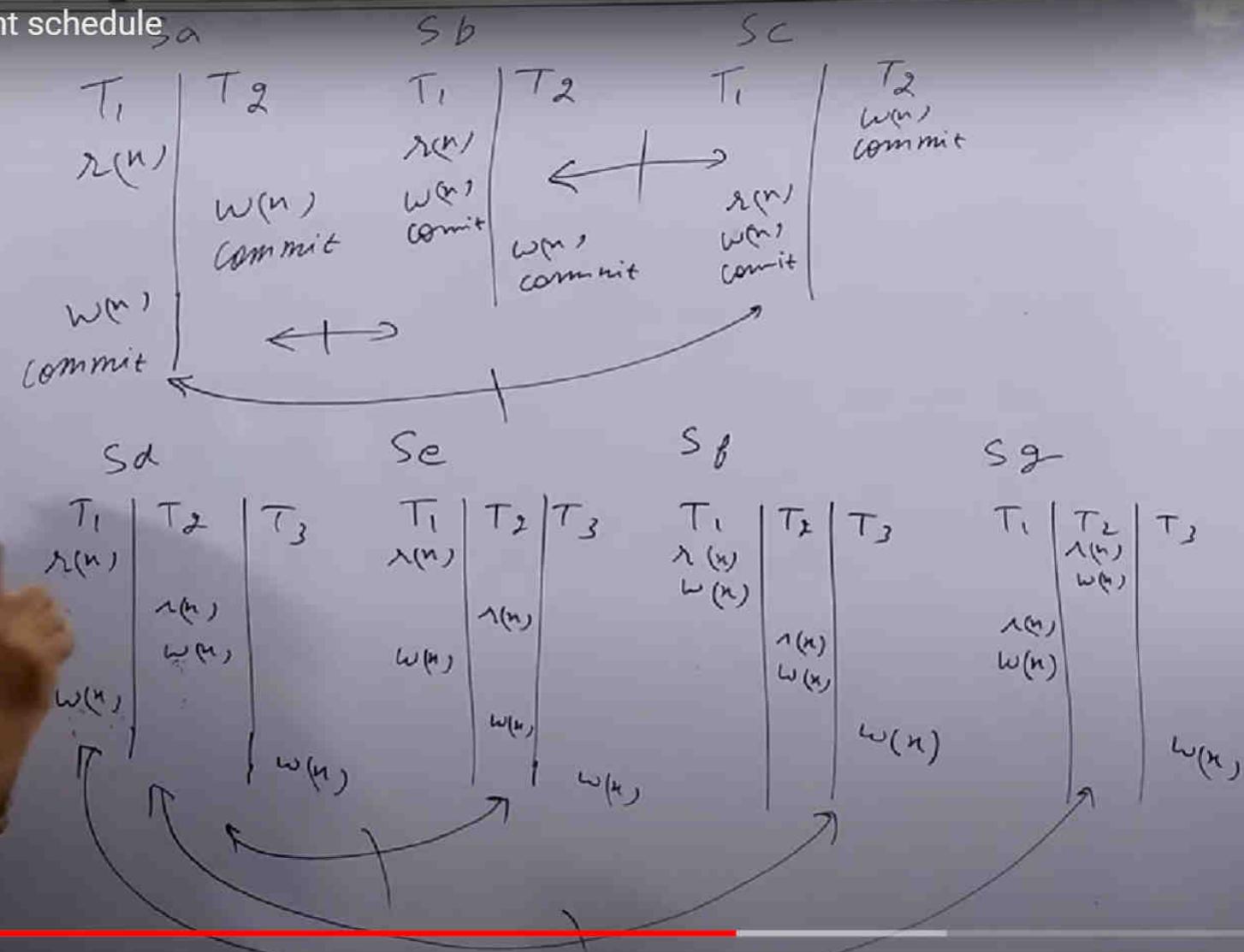
Defn.: 2 schedules are said to be conflict equivalent if we are able to construct one schedule from another by series of non conflicting swaps.



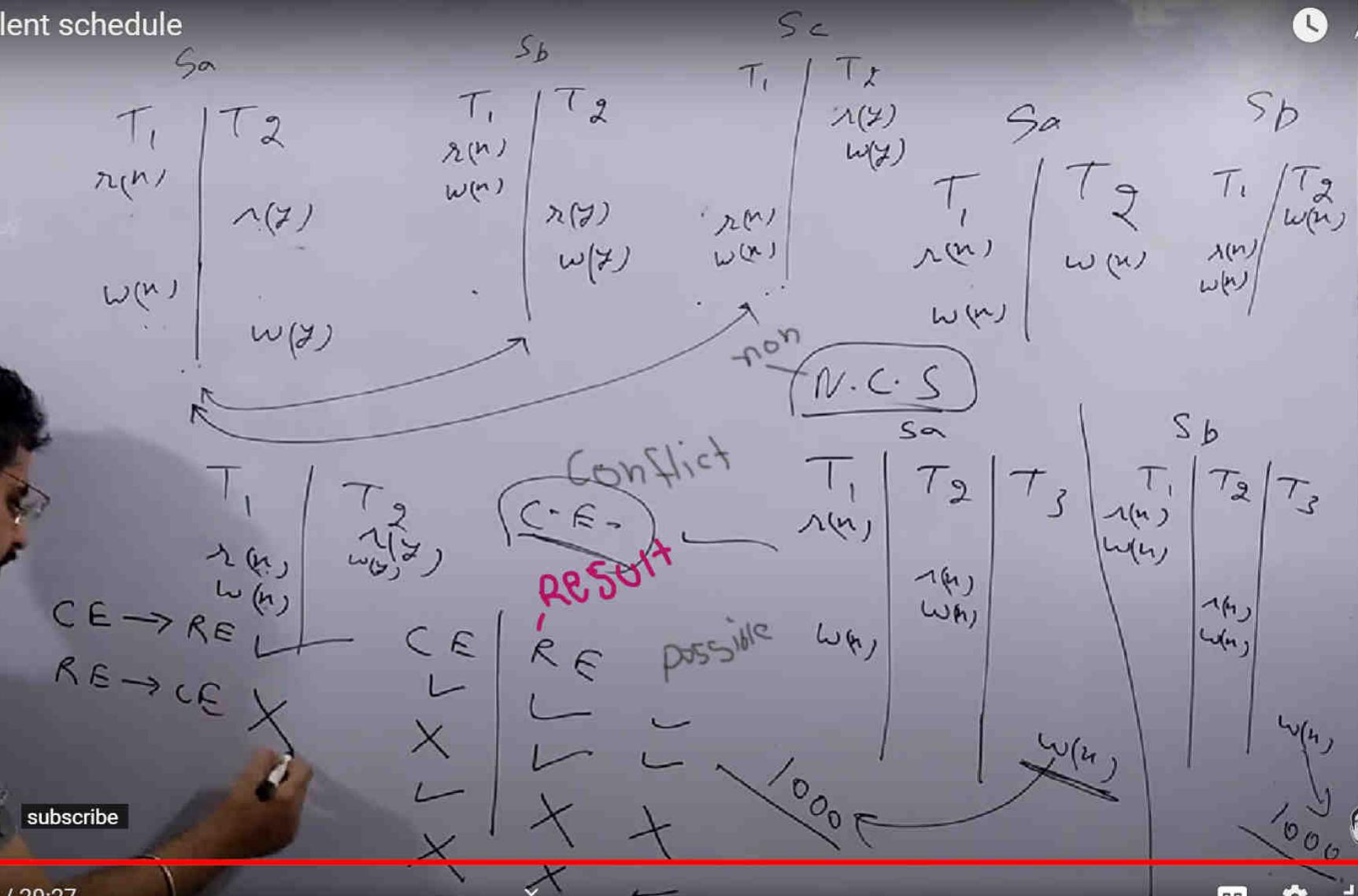
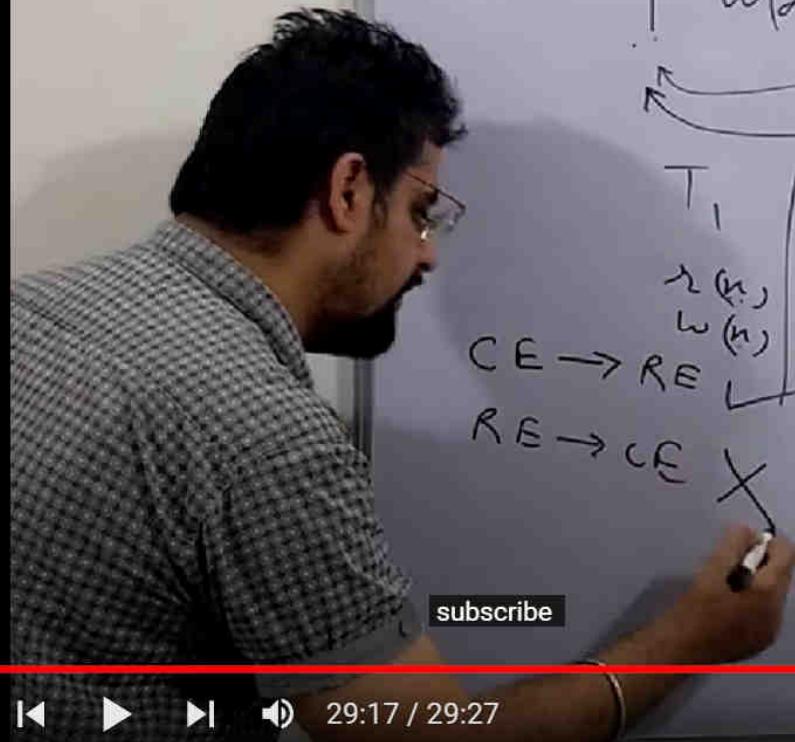
S_b

S_a

147. conflict equivalent schedule



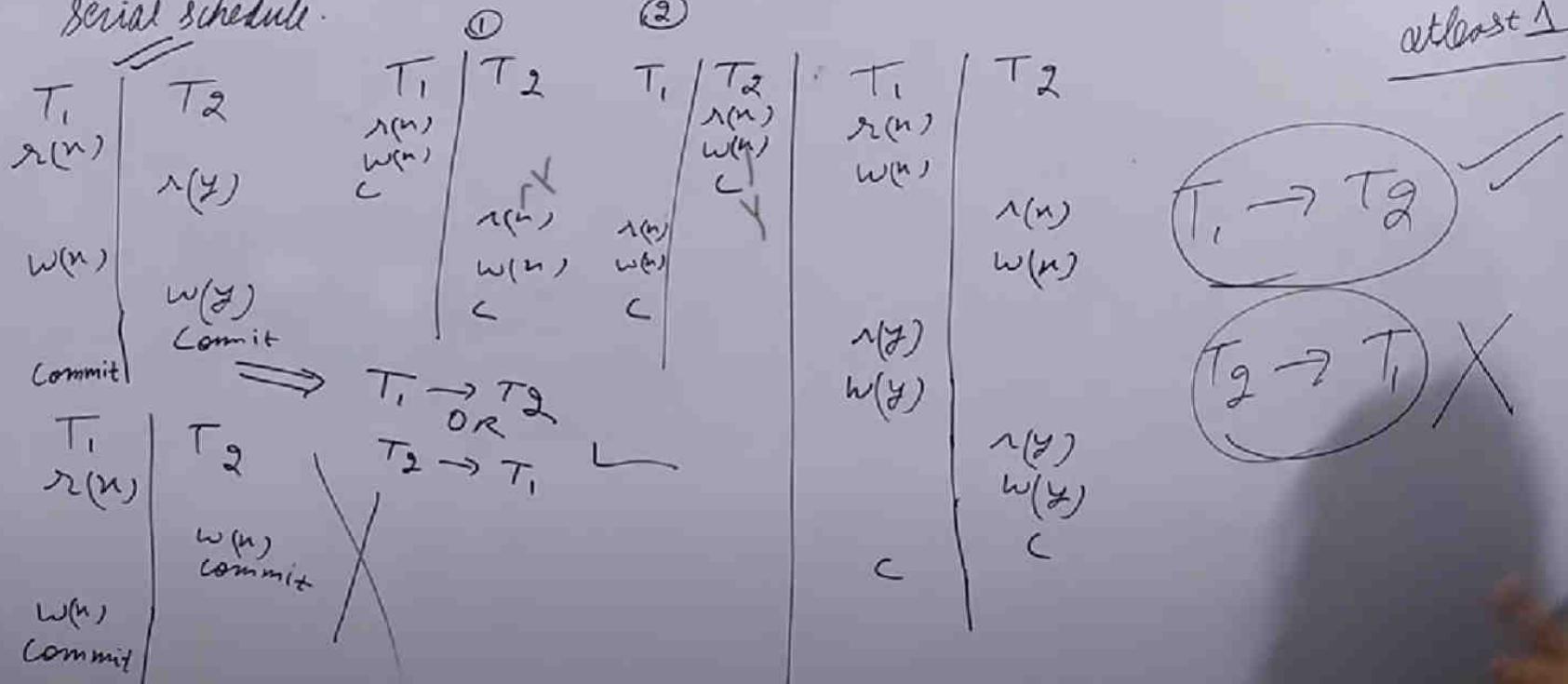
147. conflict equivalent schedule



148. conflict serializability

Conflict Serializable schedule

→ A schedule S is said to be conflict serializable if it is conflict equivalent to some serial schedule.



148. conflict serializability

Conflict Serializable Schedule

① A schedule S is said to be conflict serializable if it is conflict equivalent to some serial schedule ✓

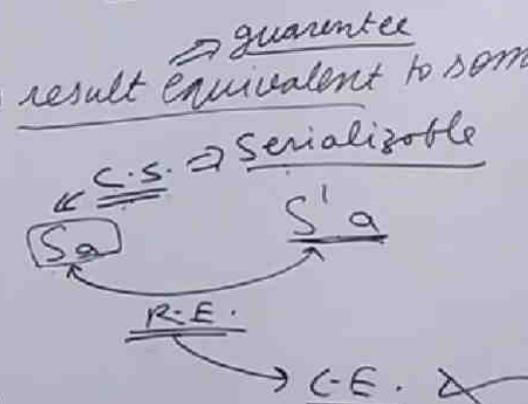
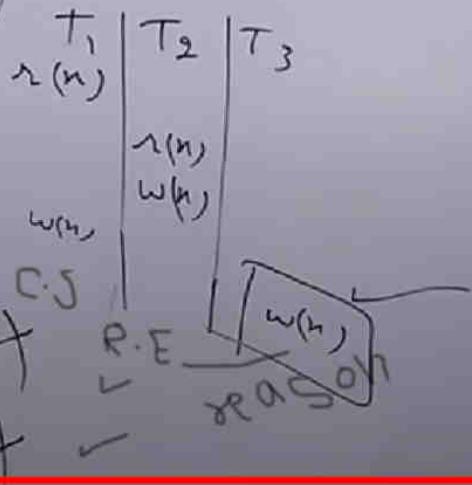
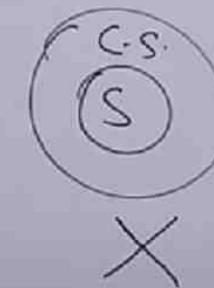
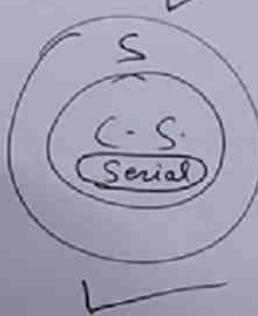
② A schedule S is said to be serializable if it is result equivalent to some serial schedule. $\xrightarrow{\text{C.S.} \Rightarrow \text{Serializable}}$ guarantee

C.E. \rightarrow R.E. ✓
R.E. \rightarrow C.E. ✗

Conclusion ??

$T_1 \mid T_2$
 $r(n), w(n), \text{Commit}$

$r(n)$
 $w(n), \text{Commit}$



148.1 number of conflict serializable schedules (GATE 2017)

GATE 2017 $T_1 \rightarrow r(n), w(n), r(y), w(y)$ $T_2 \rightarrow r(y), w(y), r(z), w(z)$ $\overline{T_2 \rightarrow T_1}$ $\overline{T_1 \rightarrow T_2}$ T_2 T_1 $r(n)$ Total No. of C.S. Schedules ?? $6(4!) \text{ w/ } 1 \text{ more}$ 5y

$\Leftrightarrow \left\{ \begin{array}{l} 5C_2 + 5C_1 \\ 4C_2 + 4C_1 \\ 3C_2 + 3C_1 \\ 4C_2 + 4C_1 \\ 3C_2 + 3C_1 \\ 3C_2 + 3C_1 \\ + \\ 1 \end{array} \right.$

 $w(n)$ $r(y)$ $w(y)$ $\frac{1(4)}{w(y)}$

149. how to check conflict serializability

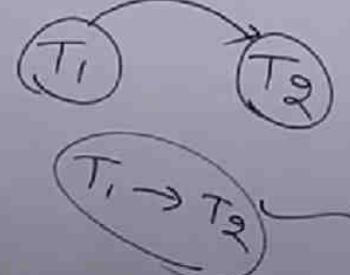
Ques How to check schedule is conflict serializable or not?

Algorithm:-

- Step 1 Each transaction is represented as node of graph (Precidence graph)
- Step 2 If we have $r(n)$ operation in T_i then try to find $w(n)$ operation in T_j below $r(n)$
If we find the same then draw directed edge from $T_i \rightarrow T_j$ only once for all
Conflicting operations.
- Step 3 If we have cycle in this graph then schedule is not c.s. else it is c.s.
& serializability order is given by topological ordering of nodes of graph.

Q1 T_1
 $r(n)$
 $w(n)$
 $r(y)$
 $w(y)$

T_2
 $r(x)$
 $w(x)$



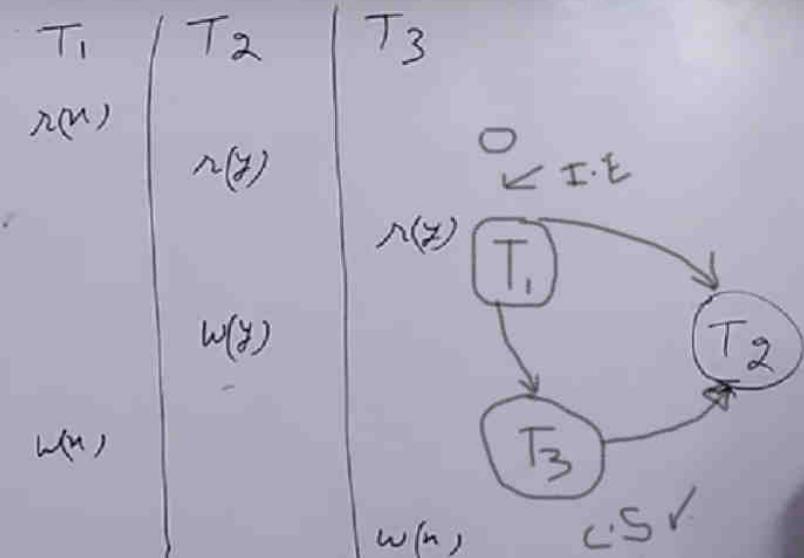
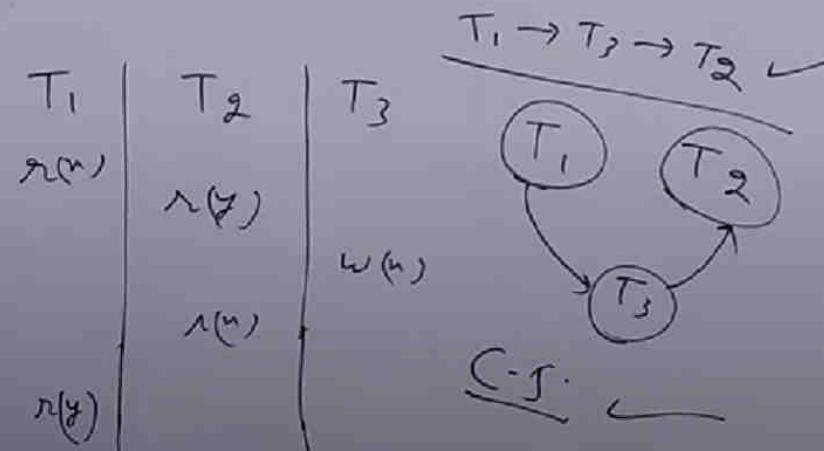
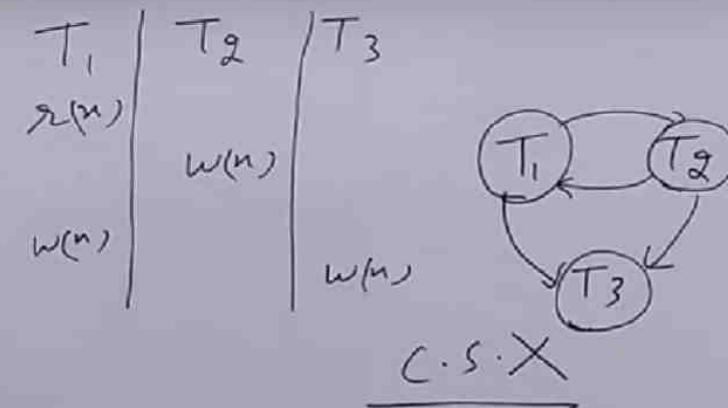
Q2 T_1
 $r(n)$
 $w(n)$

T_2
 $w(n)$



c.s X

149. how to check conflict serializability

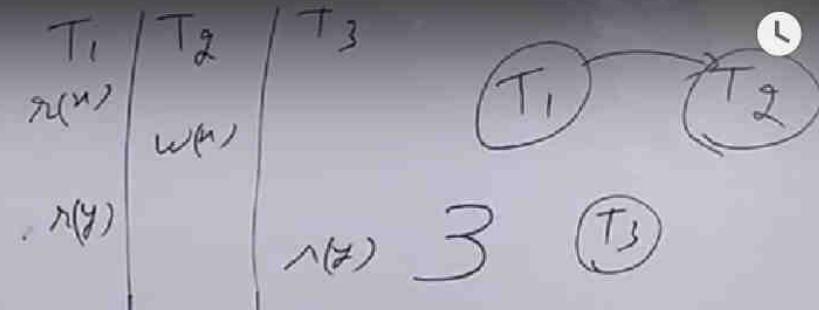
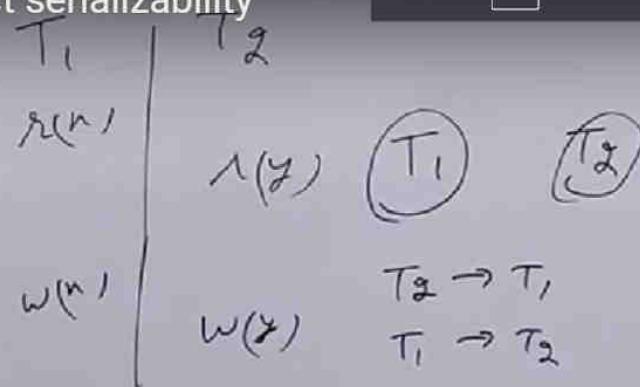


T₁ → T₃ → T₂

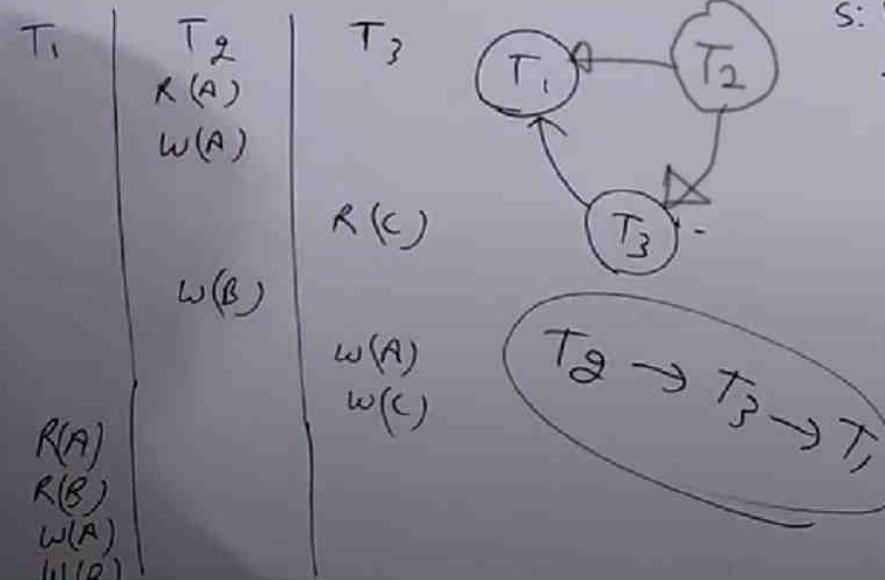
q/n coming
edge

149. how to check conflict serializability

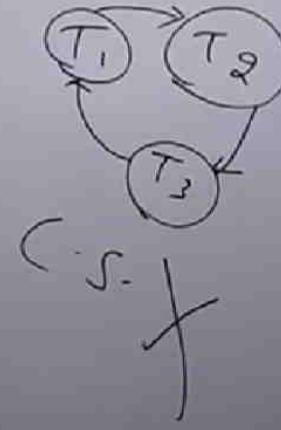
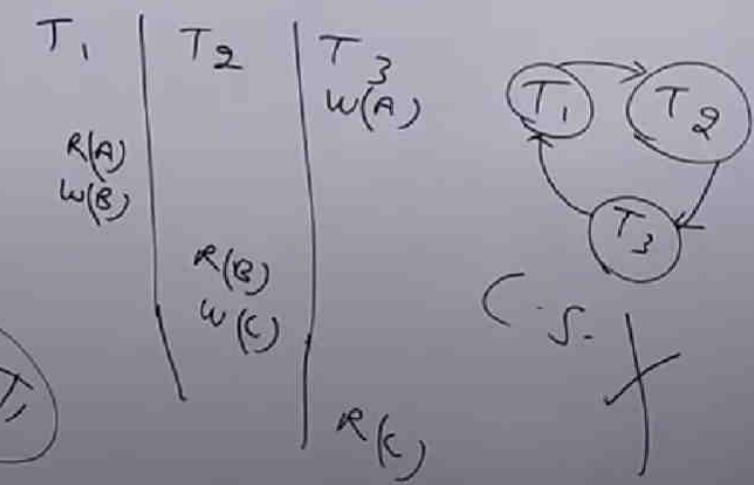
Press Esc to exit full screen



& Find all possible serializable orders -



S: $w_3(A); R_1(A); w_1(B); R_2(B); w_2(C); R_3(C)$

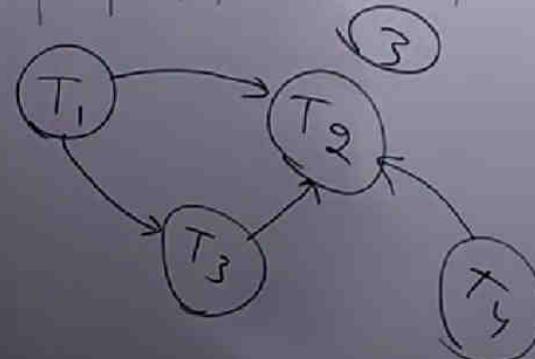


149. how to check conflict serializability

T ₁	T ₂	T ₃	T ₄
r(A)			r(A)
	r(A)		
w(B)		n(A)	
	w(A)		
		R(B)	
w(B)			

⇒ To detect cycle in graph having n vertices we need $O(n^2)$ time.

- Hence checking schedule is C.S. or not is a P class problem i.e fast algo exists for it.
- Most of G.C.T are based on concept of conflict serializability.
- But every correct schedule (serializable) is not C.S.
- G.C.T based on C.S. will reject those schedules which are not C.S.
- Hence lots & lots of correct schedules will not be executed in our system :- they are not C.S.



T ₁ r(n)	T ₂	T ₃	T ₄ → T ₁ → T ₃ → T ₂
	w(n)		T ₁ → T ₄ → T ₃ → T ₂
		w(n)	T ₁ → T ₃ → T ₄ → T ₂
			T ₁ → T ₃ → T ₄ → T ₂

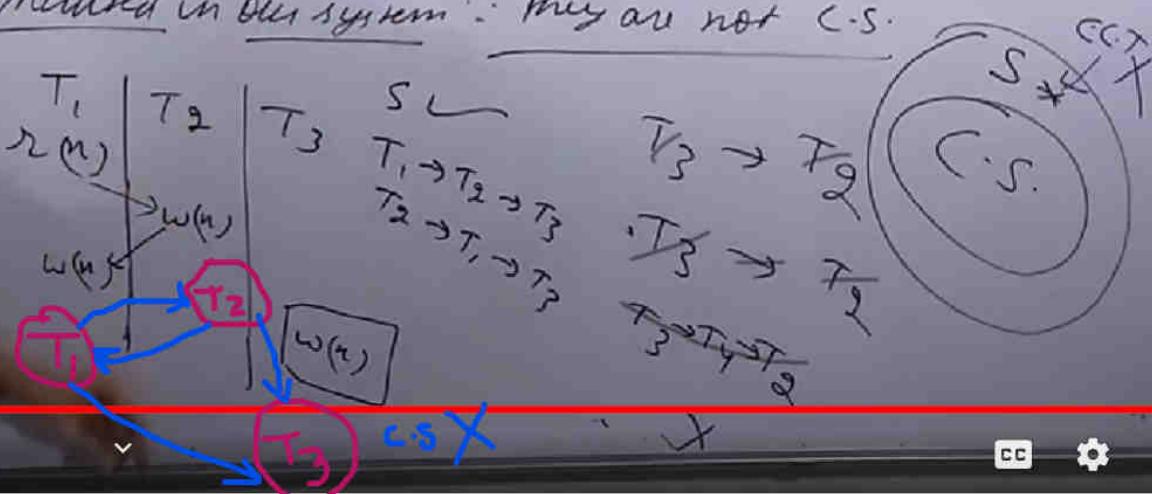
150. why concurrency control techniques are based on conflict serializability

T ₁	T ₂	T ₃	T ₄
r(A)			
	w(B)		
		(A)	

→ To detect cycle in graph having n vertices we need $O(n^2)$ time.

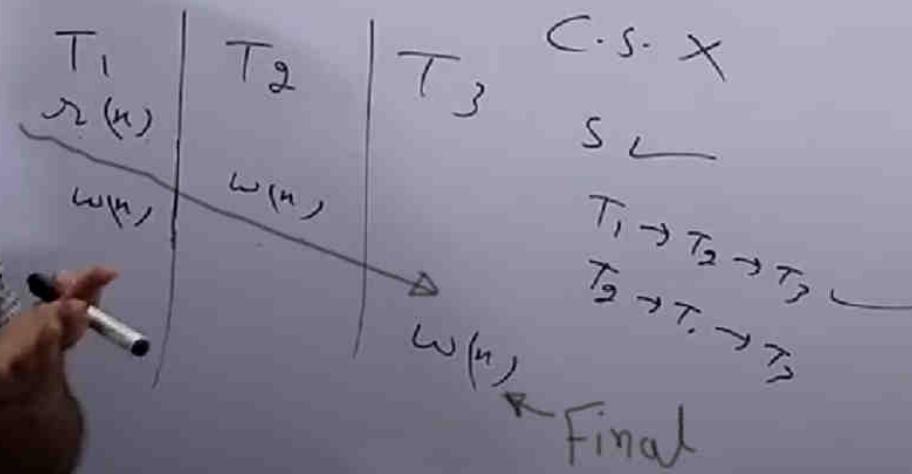
$$2^n \cdot n^2 \cdot n^3$$

- Hence checking schedule is C.S. or not is a plan problem i.e. fast algo exists for it
- Most of CCT are based on concept of Conflict serializability
- But every correct schedule (serializable) is not C.S.
- CCT based on C.S. will reject those schedules which are not C.S.
- Hence lots & lots of correct schedules will not be encountered in our system : they are not C.S.



151. why concept of view serializability

- Definition of C.S. is very restrictive : here $w-w$ is always a conflict.
- But $w-w$ is a conflicting operation only if one of w is final write.
- Hence we have other type of schedule called view serializable.
- According to def'n of View equivalence we have only 2 conflicts i.e. $r-w$ & $w-r$ ($w-w$ is conflict only if one of w is final write).



152. view equivalent schedules

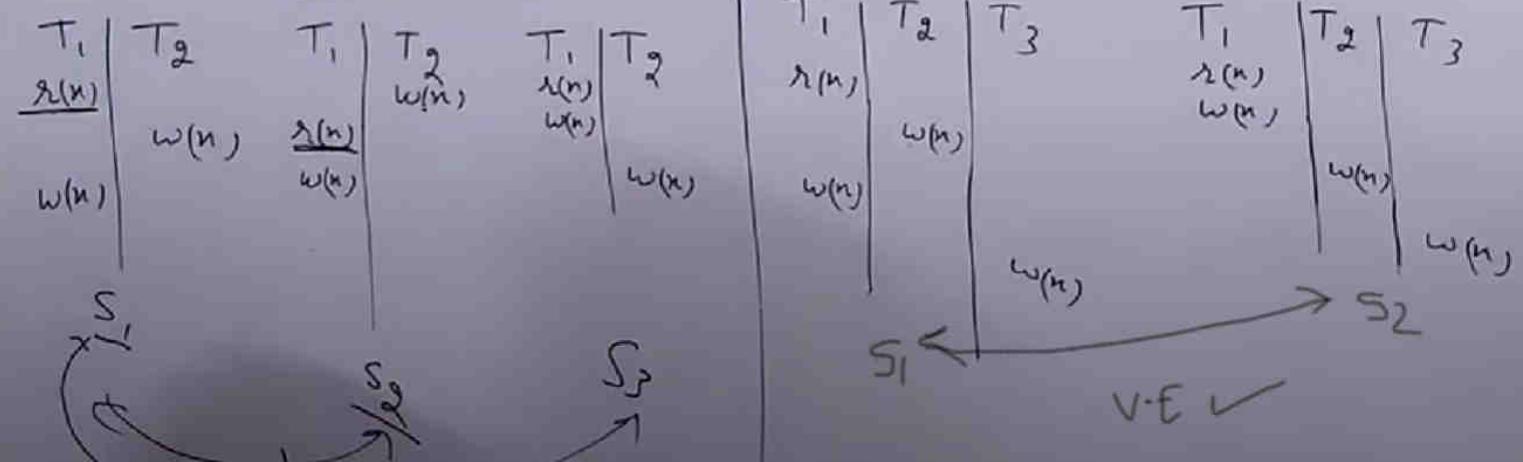
View Equivalent schedules

→ 2 schedules S_1 & S_2 over same set of transactions are view equivalent if following 3 conditions are satisfied.

① If T_1 reads initial value of item n in S_1 , then it must also read initial value of item n in S_2 .

② For each data item n the transaction that performs final write on n in S_1 must also perform the final write on n in S_2 .

③ If transaction T_1 reads value of item n written by T_2 in S_1 then it must also do the same in S_2 [Sequence of R-W & W-R conflicts must be maintained]



152. view equivalent schedules

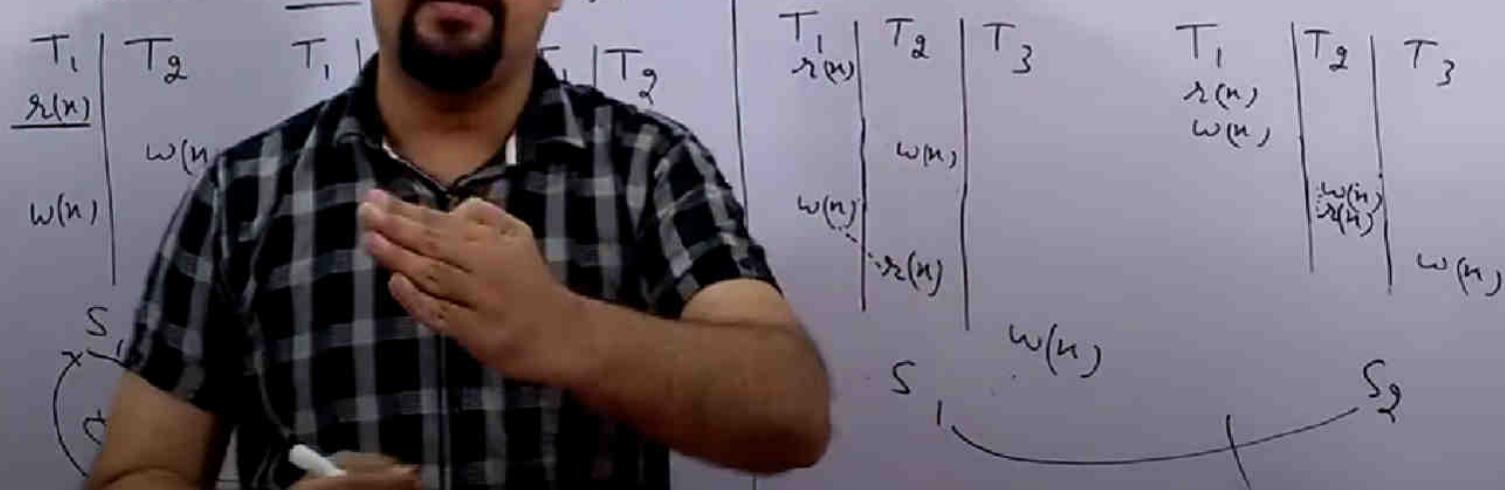
View Equivalent schedules

→ 2 schedules S_1 & S_2 over same set of transactions are view equivalent if following 3 conditions are satisfied.

① If T_1 reads initial value of item n in S_1 , then it must also read initial value of item n in S_2 .

② For each data item n , transaction that performs final write on n in S_1 must also perform the final write on n in S_2 .

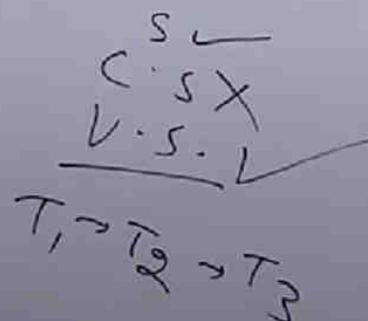
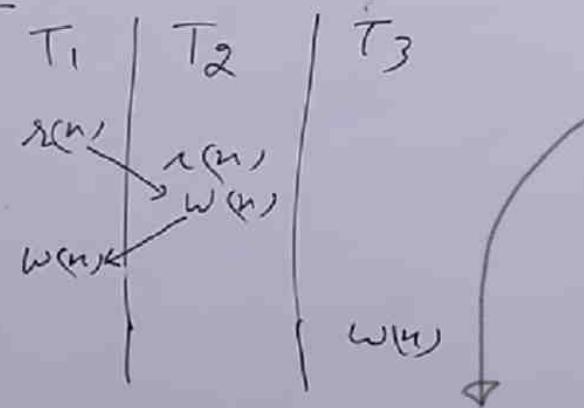
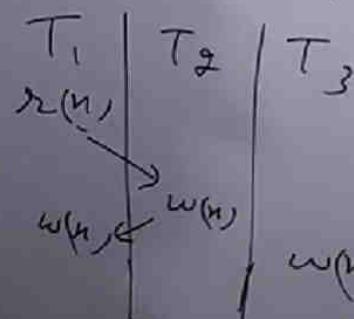
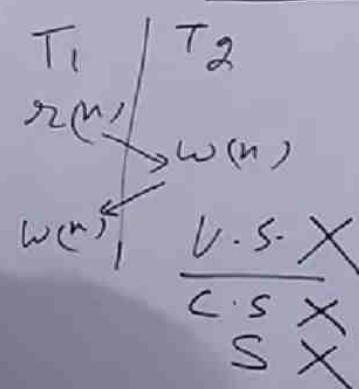
③ If transaction T_2 writes value of item n written by T_2 in S_1 , then it must also do the same in S_2 . [Maintenance of R-W & W-R conflicts must be maintained]



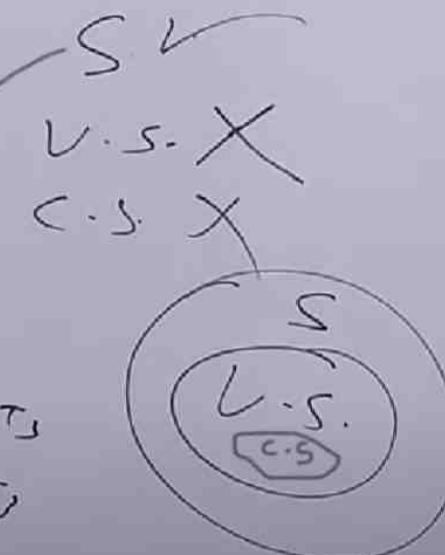
153. view serializability

View Serializable Schedule :-

A schedule S is said to be view serializable if it is view equivalent to some serial schedule.

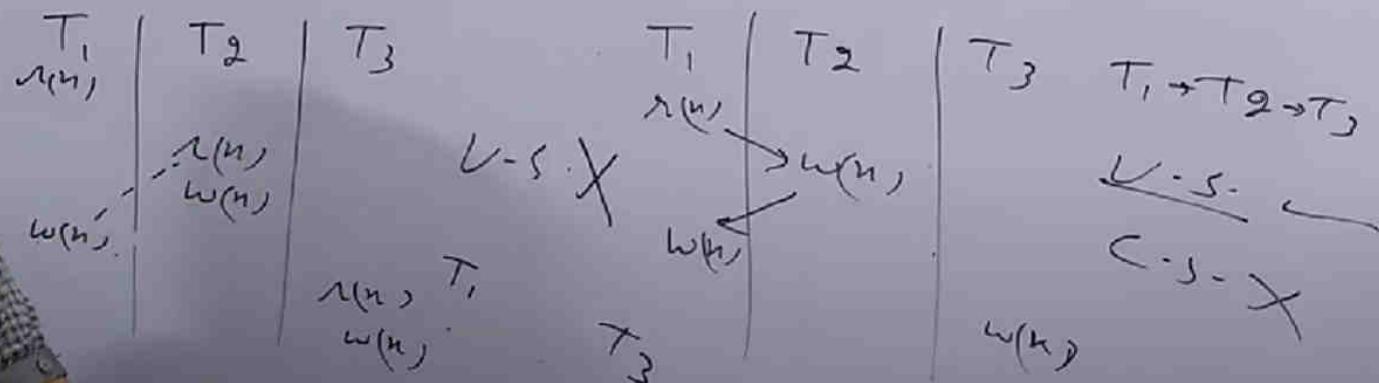


$T_1 \rightarrow T_2 \rightarrow T_3$
 $T_2 \rightarrow T_1 \rightarrow T_3$



153. view serializability

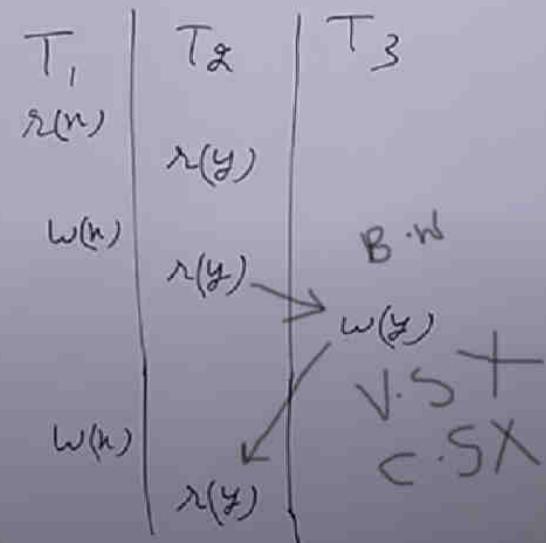
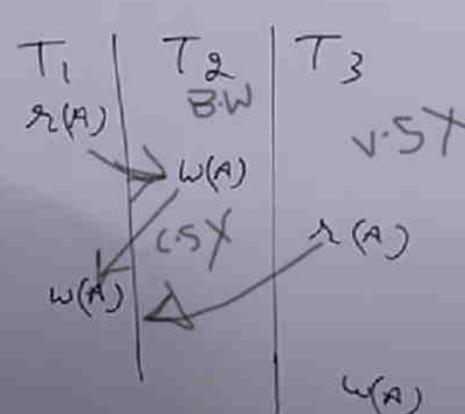
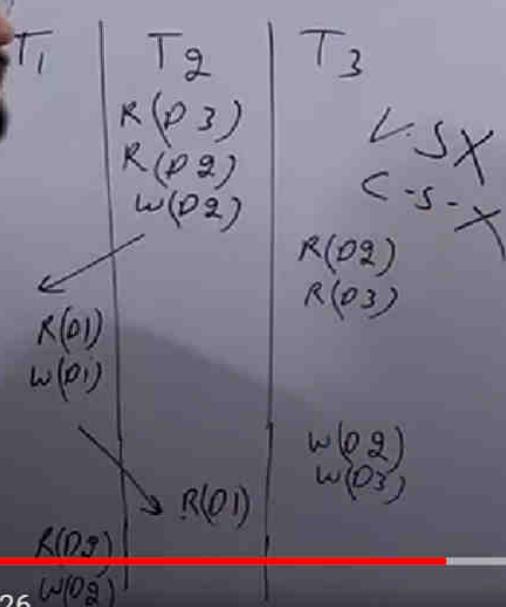
- If we have constraint write assumption then blind writes are not allowed.
i.e w/o reading transaction cannot write any item.
- Under such assumption definition of C.S. & V.S. are exactly same.
- But if we have unconstrained write assumption i.e blind writes are allowed then definition of C.S & V.S. are different.



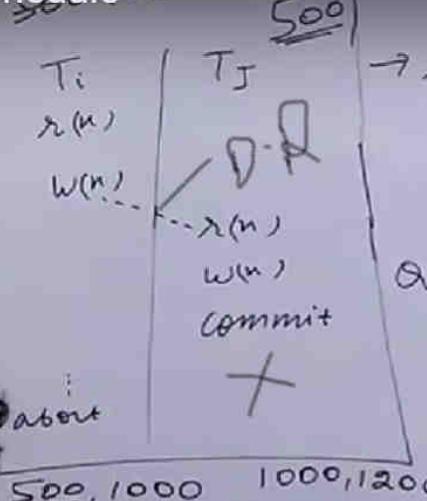
154. how to check view serializability

How to check Schedule is View Serializable or Not?

- ① Check for C.S. If C.S. then V.S.
- ② If schedule is not C.S. then check for blind writes if No blind writes then it will not be V.S. [\because def^h of C.S. & V.S. is same in constraint write assumption]
- ③ If schedule is not C.S. & it has blind writes then apply def^h of V.S.



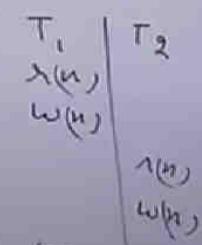
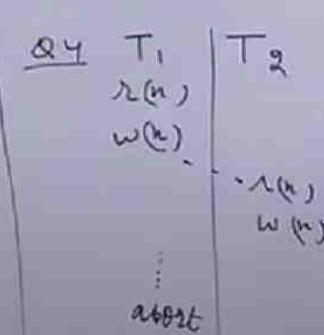
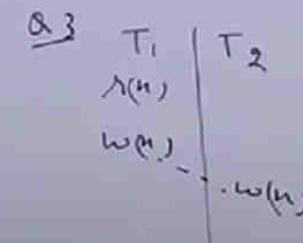
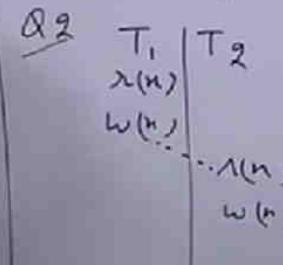
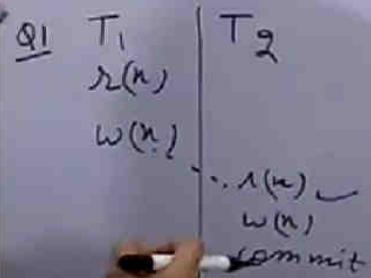
155. recoverable schedule

RECOVERABLE SCHEDULE

→ A schedule is said to be recoverable if transaction T_j reads an item written by uncommitted transaction T_i then T_j must not commit before T_i (commit operation of T_i must appear before commit operation of T_j)

Q How to check schedule is recoverable or not?

- ① Check presence of dirty read if dirty read is absent then schedule is always recoverable.
- ② If dirty read is present then transaction doing dirty read must not commit first.



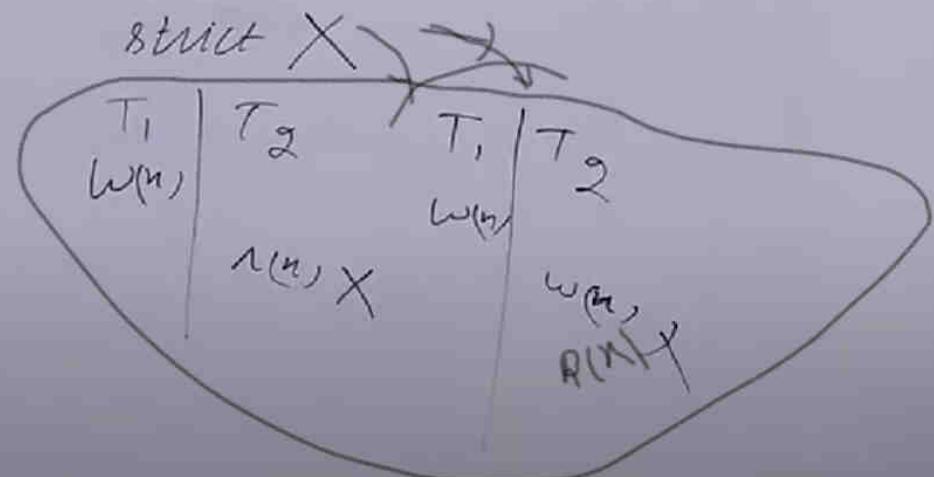
155.1 gate 2016 question on recoverability

GATE 2016
 $S = r_2(n); r_1(n); r_2(y); w_1(n); r_1(y); w_2(n); a_1; a_2$

Recoverable $\leftarrow D.R \rightarrow$ ^{No commit}
{Y{sf}}

Cascading (No cascading abort) \leftarrow D.R X

T_1	r_2
$r(n)$	$r(n)$
$w(n)$	$r(y)$
<u>$r(y)$</u>	<u>$w(n)$</u>
abort	<u>$w(n)$</u> abort



2:49 / 2:50



156. cascadeless schedule



T ₁	T ₂	T ₃	T ₄	T ₅	T ₆
r _(y)					
w _(y)	-	r _(y)			
	r _(y)	-			
		w _(y)	r _(x)	-	
			w _(x)	r _(y)	-
				w _(x)	r _(y)
					f

Cascading abort:-Cascadeless Schedule:-

→ Failure of one transaction leads to failure of series of transactions.

For e.g. if T₁ fails then T₂, T₃, T₄, T₅, T₆ must be aborted

→ A schedule is said to be cascadeless if it avoids cascading aborts.

→ Main reason of cascading abort is dirty read so cascadeless schedules must not allow Dirty read

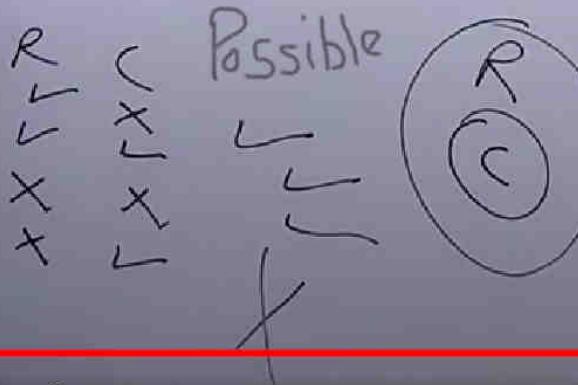
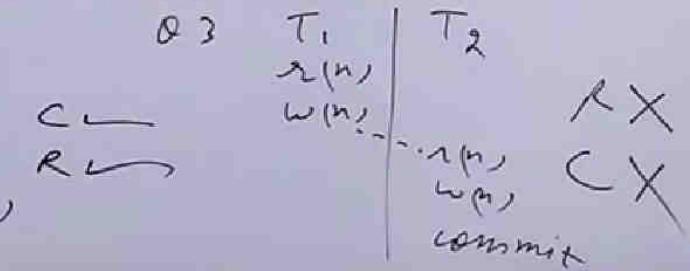
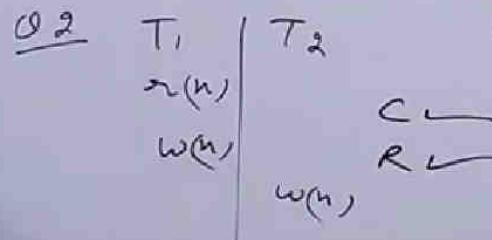
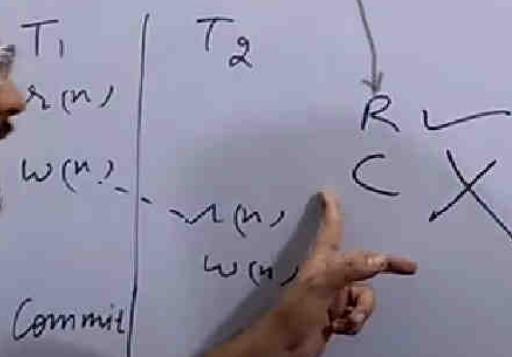
→ In other words, A schedule is said to be cascadeless if transaction reads data written by committed transaction only.

156. cascadeless schedule

\downarrow
(C)

How to check schedule is cascadeless or Not?

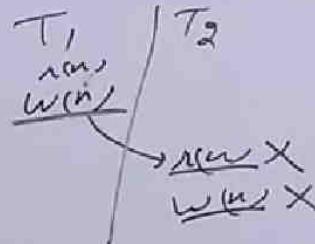
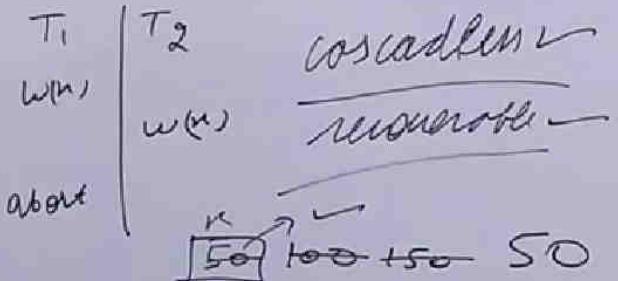
- If dirty present then schedule is not cascadeless else it is cascadeless
- "No Dirty read" is necessary as well as sufficient cond" for cascadeless schedule
- "No dirty read" is only sufficient but not necessary cond" for schedule to be recoverable
 $\text{H} \rightarrow \text{recoverable} \Rightarrow ??$
 $\text{C} \rightarrow \text{"no dirty read"} \Rightarrow \text{F}$



157. strict schedule

Strict Schedule:-

Q Every cascaded schedule must be allowed in database??

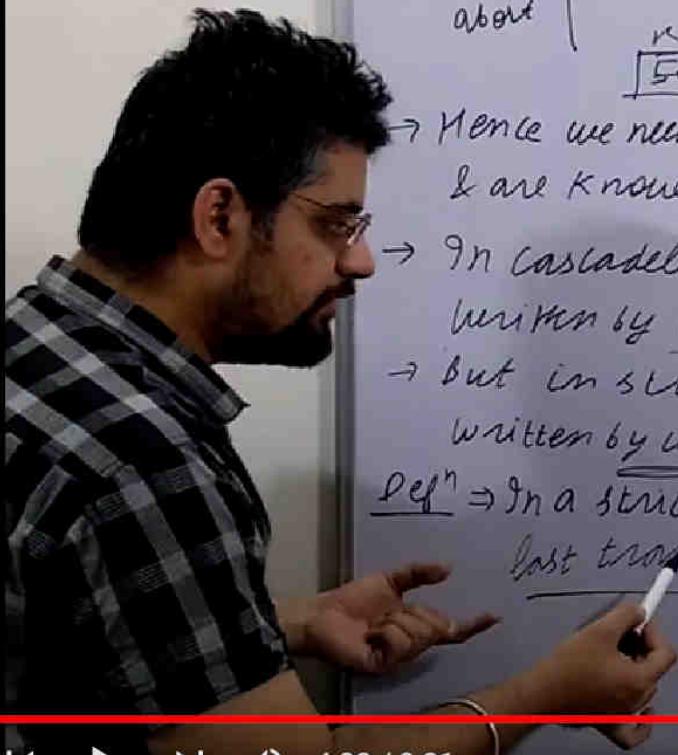


→ Hence we need another type of schedule which are strict than cascaded schedules & are known as strict schedules.

→ In cascaded schedules transactions were not allowed to read any item written by uncommitted transactions (dirty read).

→ But in strict schedules transactions are not allowed to even write any item written by uncommitted transaction.

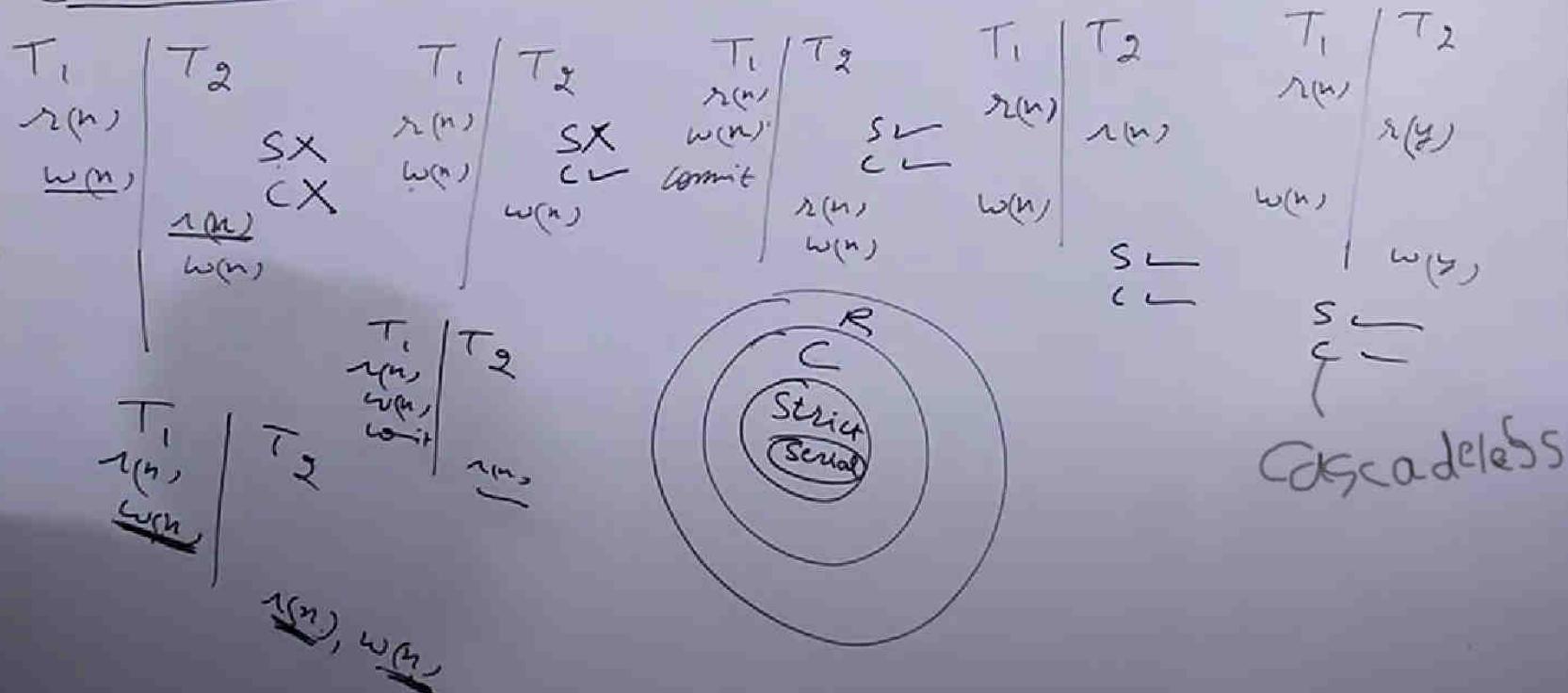
Defn → In a strict schedule, transactions can neither read nor write an item n until the last transaction that wrote X has committed or aborted.





How to check Schedule is strict or not?

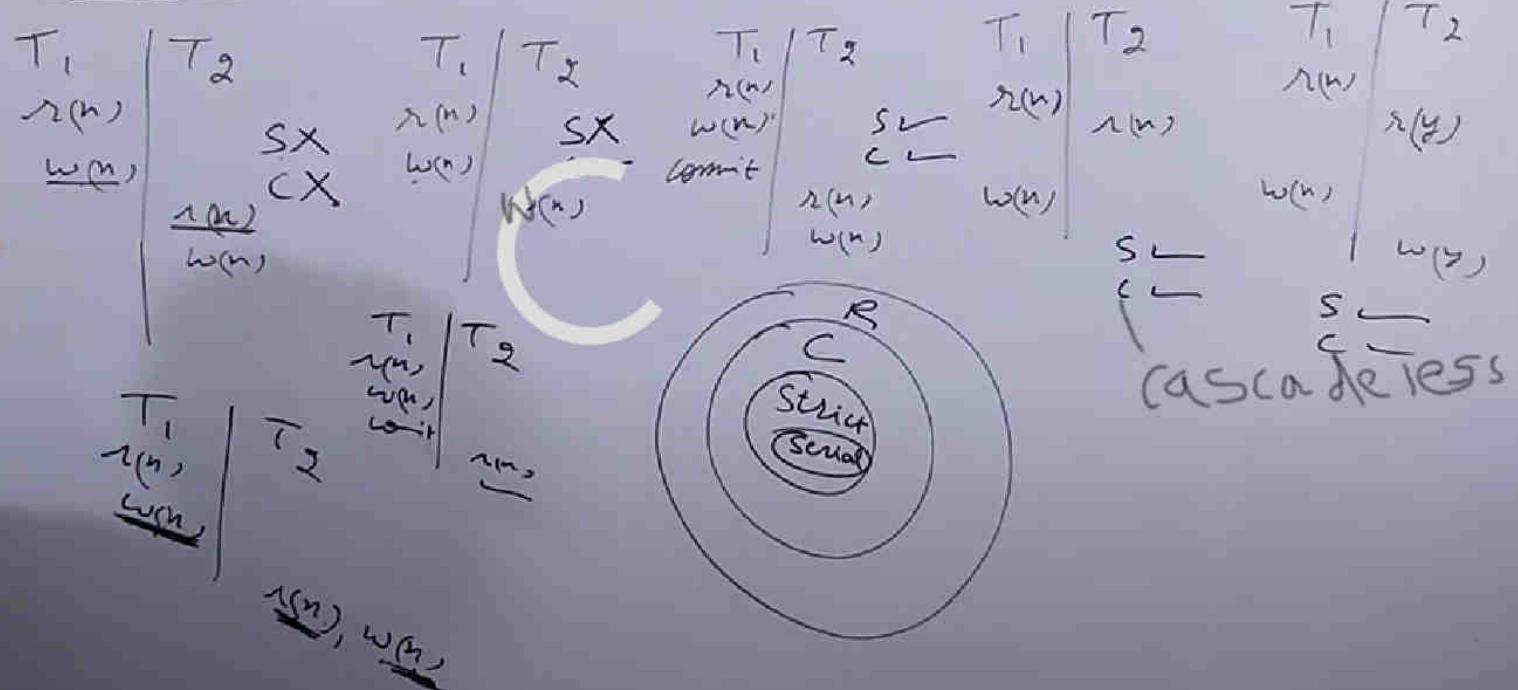
→ If write operation is preceded by any conflicting operation (ie read or write) then that conflicting operation must be done after transaction doing write operation (commit or abort).



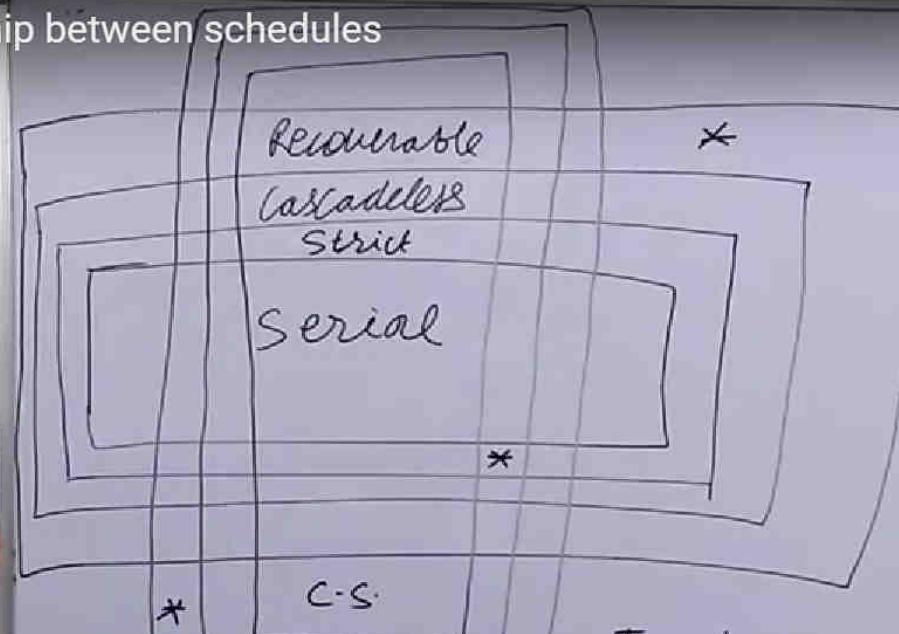
157. strict schedule

How to check Schedule is strict or not?

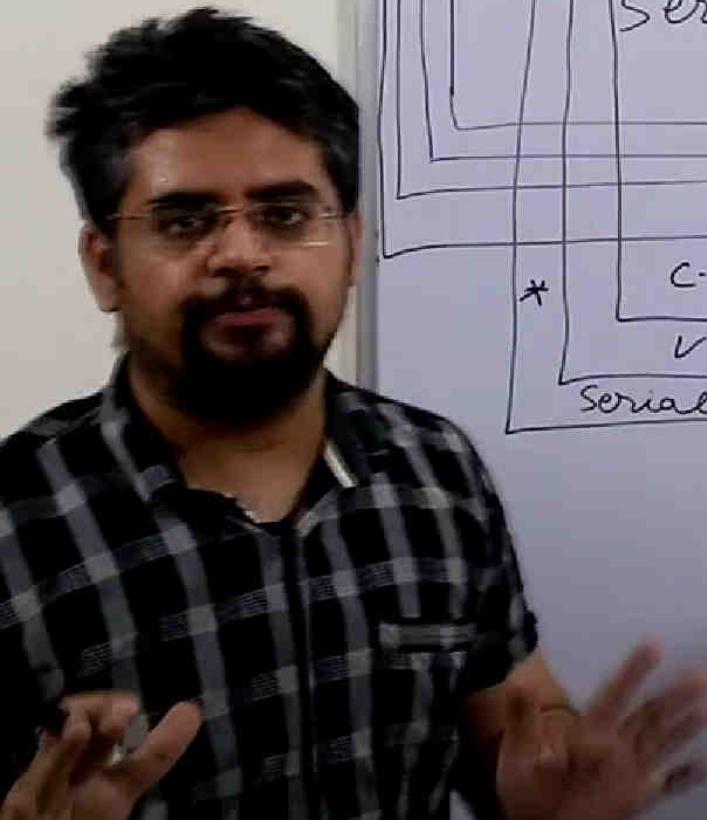
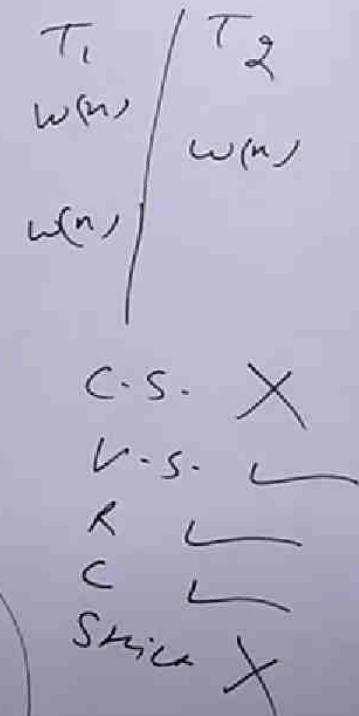
→ If write operation is preceded by any conflicting operation (ie read or write) then that conflicting operation must be done after transaction doing write operation (commit or abort).

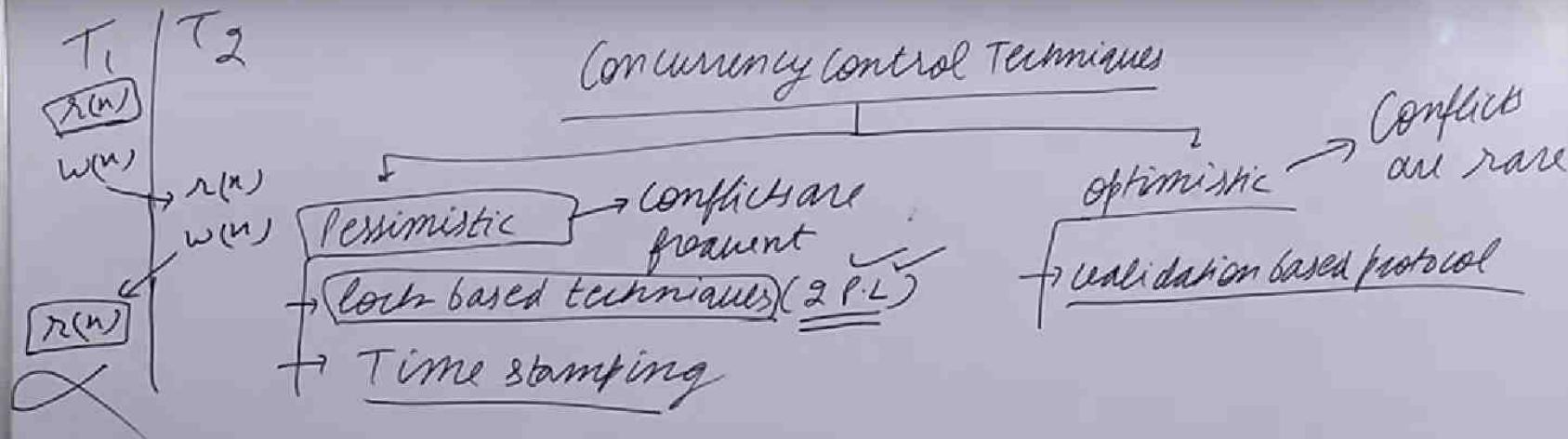


158. relationship between schedules



Relationship w/o all schedules





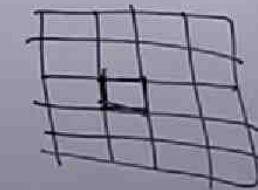
LOCK based techniques:-

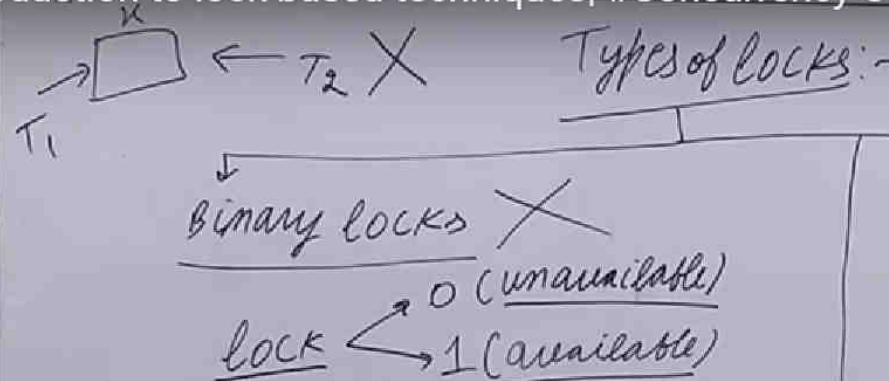
→ LOCK is a variable associated with each database item & it describes the operation which can be performed on item.

→ If item is a cell then lock is associated with each cell

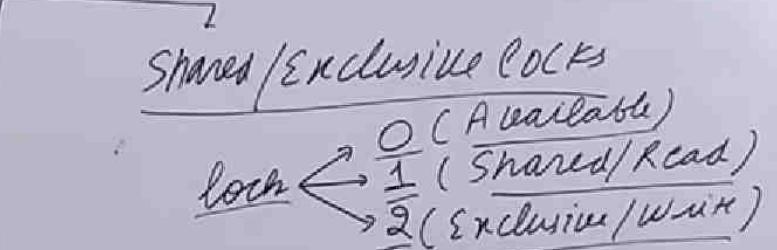
→ tuple tuple

→ table table





- If $lock(n)=0$ then any transaction say T_i has already applied lock on n so any other transaction say T_j can't access item n [neither read nor write]
 - If $lock(n)=1$ then any transaction T_i can apply lock on n & access it (i.e. read or write n)
 - W/o applying lock transaction can't access any item
- Disadvantage is that 2 transactions can't even read same item n .



→ more than one transaction can have shared lock on same item. but atmost one transaction can have write lock on item.

→ If one transaction is having write lock then no other transaction can read or write the same item.

→ If one transaction is having read lock then no other transaction can write the same item.

	$T_1 \rightarrow R$	R	W
$T_2 \downarrow$	R	✓	✗
W_1		✗	✗

- If transaction wants to access any item then it will 1st ask for lock from lock manager. If lock is available then it is granted else transaction will have to wait.
- Locks are granted in FIFO order hence lock based techniques are starvation free.

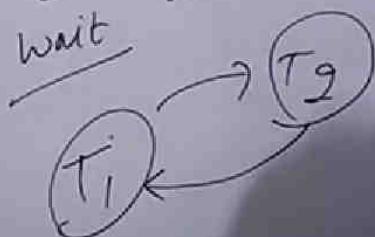
WHY LOCKING PROTOCOLS?

T_1, T_2, T_3 shared

LOCK n

$T_4 \rightarrow$ Exclusive lock n

wait

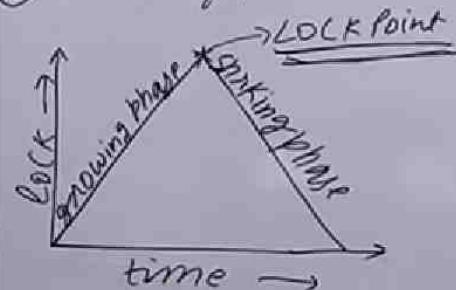


T_1	T_2	T_1	T_2
$r(n)$ $x = n - 500$ $w(n) \dots$ <u>unlock(y)</u>	$n \rightarrow 500 -$ $y \rightarrow 500 +$ $r(n)$ $r(y)$ $display(n+y)$ $unlock(n)$ $unlock(y)$	$r(n)$ $w(n)$ $w(y)$	$r(y) \downarrow$ $r(n) \downarrow$
$r(y)$ $y = y + 500$ $w(y)$ <u>unlock(y)</u>			

- Transaction must not release lock too early or too late
- Hence we need protocol (set of rules) which must be followed by each transaction.

2 Phase locking

- As the name suggest here we have 2 phases i.e. growing phase & shrinking phase
- ① Growing phase → Transaction can only acquire lock & it will not release any lock.
 - ② Shrinking phase → Transaction can only release lock & it cannot acquire lock.



T ₁	T ₂
LOCK-X(A) read(A) LOCK-S(B) read(B) write(A) unlock(A)	
lock-S(B) r(B)	LOCK-X(A) read(A) write(A) unlock(A)

S.P. 2PL

X 2PL

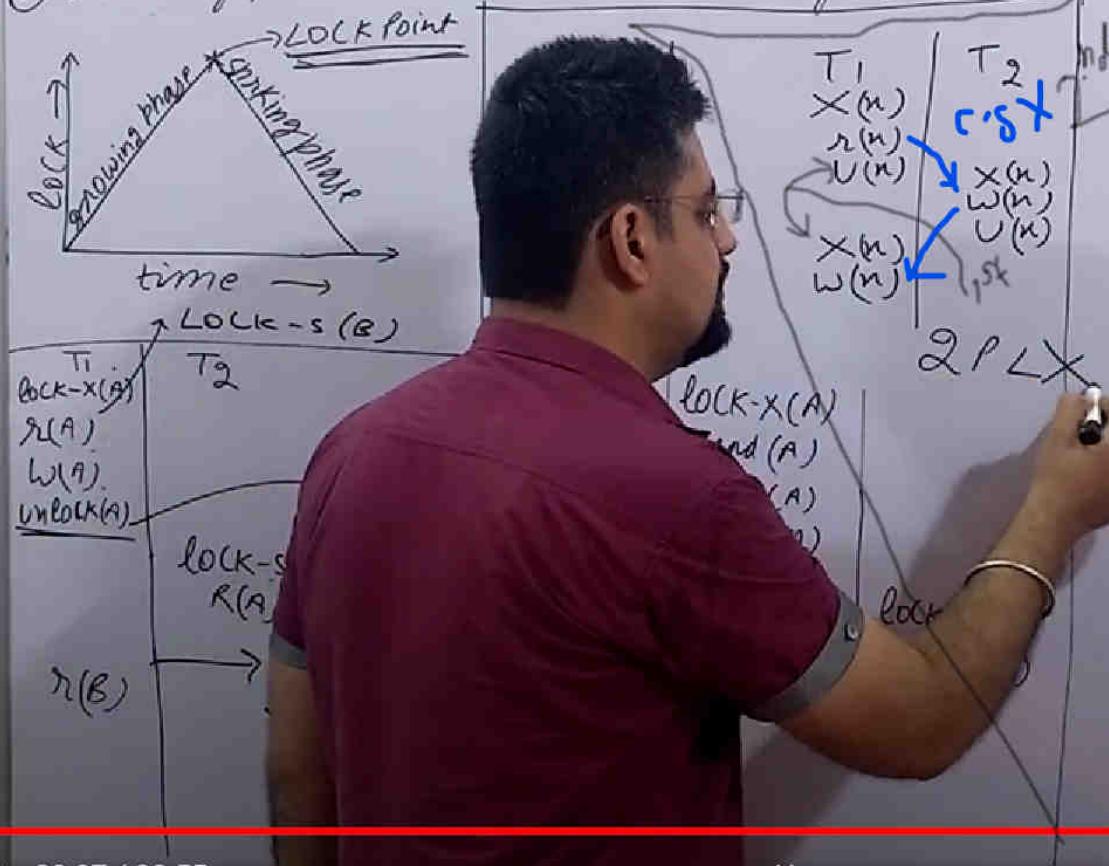
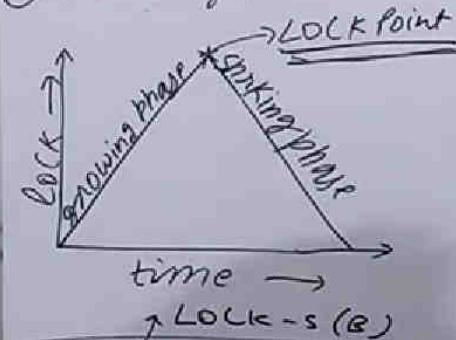
T ₁	T ₂	T ₃	T ₁	T ₂	T ₃
LOCK-X(A) read(A) LOCK-S(B) read(B) write(A) unlock(A)					
			LOCK-X(A) read(A) write(A) unlock(A)		w(B)
					w(C)
				R(A)	
				R(B)	
				R(C)	
					w(A)

2PL ✓

2 Phase locking

→ As the name suggest here we have 2 phases

- ① Growing phase → Transaction can only acquire lock
- ② Shrinking phase → Transaction can only release lock

T₁X(B)
X(C)T₂S(A)
R(A)
U(A)T₃

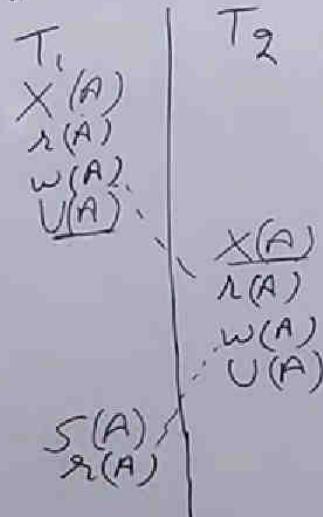
2 P.L

S(A)
W(B)
W(C)
U(B)
U(C)S(B)
R(B)
U(B)R(A)
Commit
U(A)U(K)
X(A)
W(A)
CommitS(B)
R(B)
U(B)

U(A)

Drawbacks of 2PL

① It guarantees conflict serializability of schedules hence Serializability is guaranteed.



Disadvantages:-

② 2PL doesn't allow all C.S. schedules.

③ 2PL is not deadlock free.

④ 2PL doesn't guarantee recoverable, cascading & strict schedules.

⑤ 2PL limits the amount of concurrency.

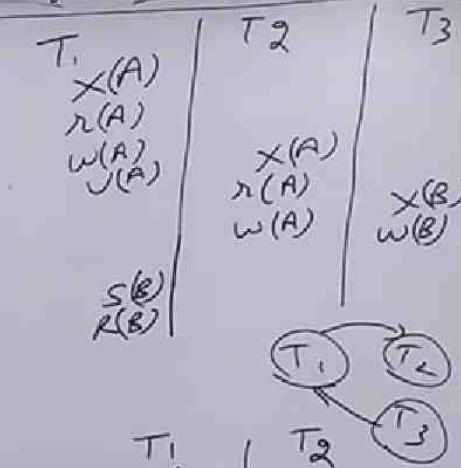
T → X, Y, Z

↓
SPL

T₁
X(A)
R(A)
W(A)
U(A)

T₂
X(A)
R(A)
W(A)
U(A)

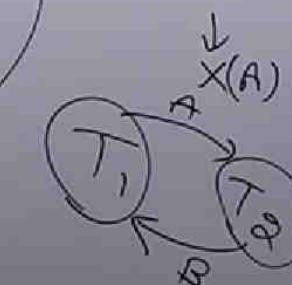
Commit
2PL ✓
Recoverable X



T₁
X(B)
R(B)

W(B)

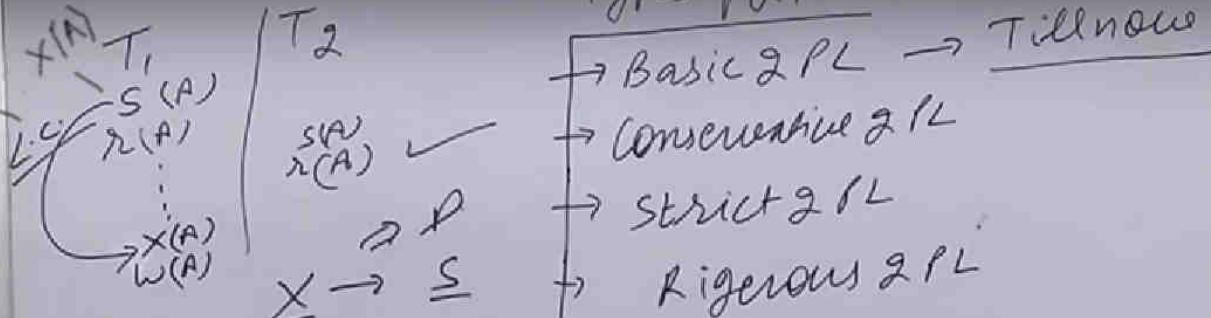
S(A)
R(A)
S(B)



X(A)



lock conv.

Basic 2 PL \rightarrow TillnowConservative 2 PLStrict 2 PLRigorous 2 PL

L.C.

Basic 2 PL will lock conversions

→ In basic 2 PL if transaction wants to read any item & then after some time it wants to write that item. It will demand exclusive lock on that item by lock manager directly i.e. read lock is not converted to exclusive lock.

→ But in basic 2 PL with lock conversion read lock can be converted to exclusive lock.

→ conversion of shared to exclusive lock is known as upgrading of lock & is done in growing phase only.

→ conversion of exclusive to shared lock is known as downgrading of lock & is done in shrinking phase only.

Conservative 2 PL :-

- Each transaction must precede its read set & write set (Assumption)
- Transaction will not start till it acquire lock on all items in read set & write set
- Hence transaction starts at lock point & hence enters the shunting phase.

Advantages :-

- ① It guarantees C-S of schedules
- ② It is deadlock free ✓

Disadvantages :-

- ① It limits the amount of concurrency →
- ② It doesn't guarantee recoverable, cascades & strict schedules
- ③ Not practical to use in DBMS
- ④ All C-S schedules are not allowed.

$T_1 \rightarrow u, y \rightarrow \text{read}$
 $3, w \rightarrow \text{write}$



T_1
 $r(u)$
 $w(y)$

T_2

$r(u)$
 $w(w)$
 con.

162. Basic 2PL, Conservative 2PL, Strict 2PL and Rigorous 2PL in DBMS | 2 Phase Locking in DBMS

→ Transaction will not release any of its originally locked row or lock.

T ₁	T ₂
X(A)	
R(A)	
W(A)	
Commit	✓
U(A)	
	X(A)
	R(A)
	W(A)
	Commit
	U(A)

T ₁	T ₂
S(A)	
X(B)	??
R(A)	
U(A)	
R(B)	
W(B)	
Commit	✓
U(B)	
	X(B)
	W(B)
	Commit
	U(B)

T ₁	T ₂
R(n)	
W(n)	
W(n)	
	X
	X

T ₁	T ₂
R(n)	
W(n)	
	R(n)
	X

Advantages:-

- ① It guarantees c-s. schedules
- ② It guarantees strict schedules (hence recoverable & cascadelock also)

Disadvantages:-

- | | |
|---------------------------------------|----------------------------------|
| ① It is not deadlock free | ③ All c-s. schedules not allowed |
| ② It limits the amount of concurrency | |



14:04 / 19:11



162. Basic 2PL, Conservative 2PL, Strict 2PL and Rigorous 2PL in DBMS | 2 Phase Locking in DBMS



rowing
phase

T ₁	T ₂
X(A)	
R(A)	
W(A)	
Commit	✓
U(A)	
X(A)	
R(A)	
W(A)	
Commit	
U(A)	

T ₁	T ₂
S(A)	
X(B)	✗
R(A)	
U(B)	
R(B)	
W(B)	
Commit	
U(B)	
U(A)	
X(B)	
W(B)	
Commit	
U(B)	

T ₁	T ₂
r(n)	
w(n)	
w(n)	
✓	

T ₁	T ₂
r(n)	
w(n)	
r(n)	
✗	

DBMS
→ Strict
→ Rigorous

Advantages:-

- ① It guarantees c-s. schedules
- ② It guarantees strict schedules (hence recoverable & cascadeable also)
- ③ Easy to implement

Disadvantages:-

- | | |
|---------------------------------------|----------------------------------|
| ① It is not deadlock free | ③ All c-s. schedules not allowed |
| ② It limits the amount of concurrency | |



17:04 / 19:11



Timestamp ordering protocol (2nd C.C.T)

- Timestamp is a unique value assigned by DBMS to each transaction.
- We can think timestamp as the start time of transaction.
- $TS(T_i) \Rightarrow$ Timestamp of T_i
- $TS(T_i) < TS(T_j) \Rightarrow T_j$ is younger & T_i is older
Junior Senior

Q1 How to implement timestamps?

Ans1 2 ways $\begin{cases} \text{System clock} \\ \text{Counter} \end{cases}$

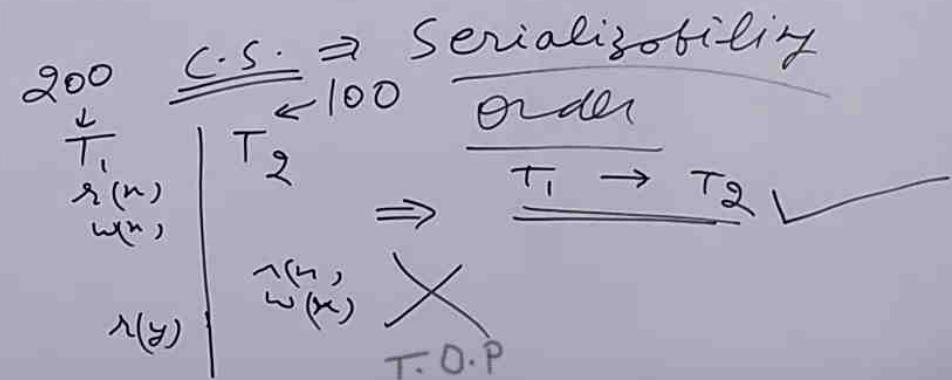
Q2 How to guarantee uniqueness?

Ans2 It's the responsibility of D.B.M.S

Q3 What is timestamp ordering?

Ans3 The main aim of T.O.P is to implement FIFO order i.e. system will allow only those schedules whose equivalent serial schedule contains transactions according to their time of entering the system.

But in 2PL FIFO order was not followed.



Timestamp ordering protocol (2nd C.C.T)

- Timestamp is a unique value assigned by DBMS to each transaction.
- We can think timestamp as the start time of transaction.
- $TS(T_i) \Rightarrow$ Timestamp of T_i
- $TS(T_i) < TS(T_j) \Rightarrow T_j$ is younger & T_i is older
Junior Senior

Q1 How to implement time

Ans1 2 ways $\begin{cases} \text{System clock} \\ \text{Counter} \end{cases}$

Q2 How to guarantee unique

Ans2 9 ts the responsibility of

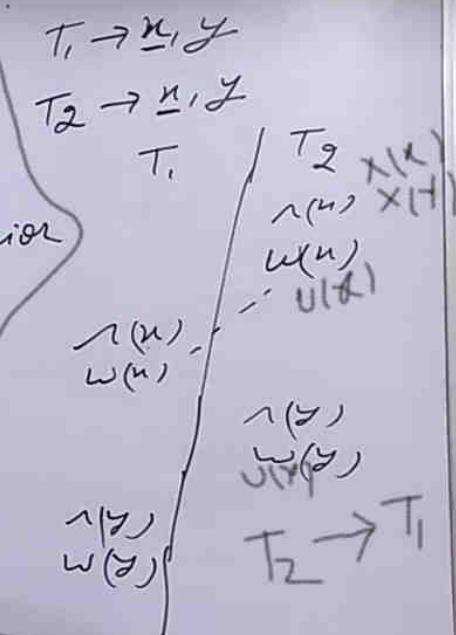
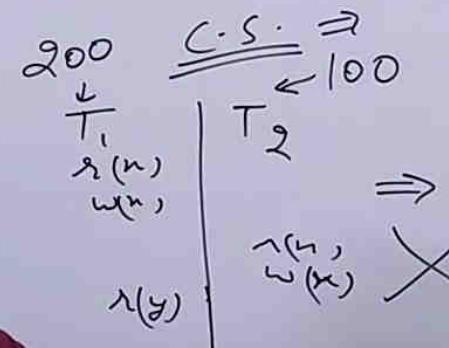
Q3 What is timestamp

Ans3 The main aim

those schedu

to their ti

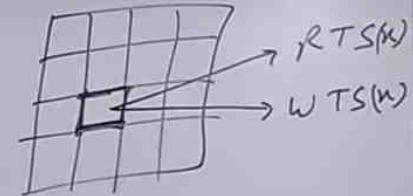
But in



present FIFO order i.e. sun will allow only
transactions according
not followed

Types of Time Stamp ordering protocol

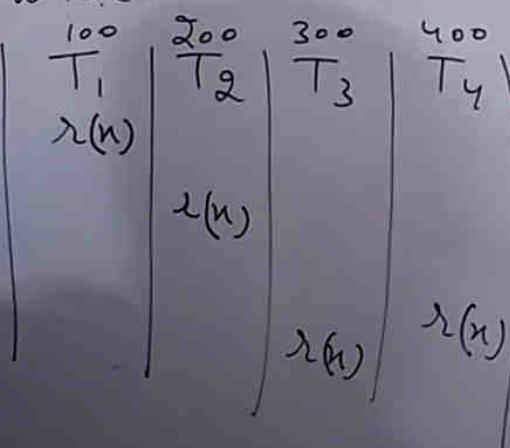
- Basic T.S.O.P —
- Strict T.S.O.P —
- Total T.S.O.P — }
Total T.S.O.P — }
- MultiVersion T.S.O.P — ??



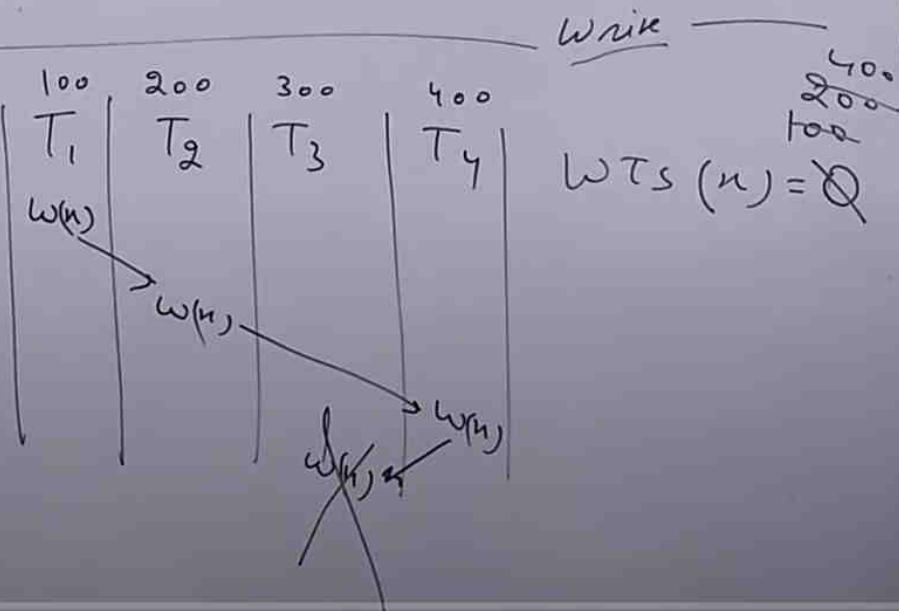
→ 2 timestamps associated with item 'n' i.e. RTS(n) & WTS(n)

RTS(n) → Time stamp of youngest transaction that has performed read on 'n'

WTS(n) →



$$\text{RTS}(n) = \underline{\underline{200}}$$



$$\text{WTS}(n) = \underline{\underline{300}}$$

164. basic timestamp ordering protocol

basic T.S.O.P

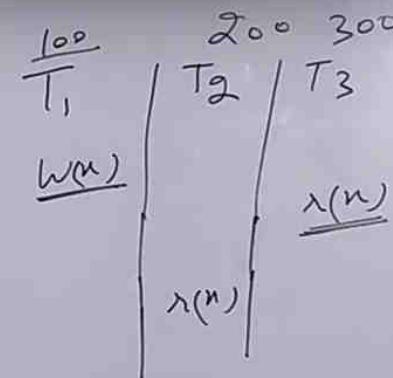
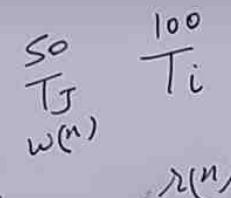
If T_i wants to read item 'n'

If $TS(T_i) < WTS(n)$

Abort T_i & restart it with new timestamp.

else
read operation is granted.

$$RTS(n) = \max(TS(T_i), RTS(n))$$



T_2 request for read

$TS(T_2) < WTS(n)$

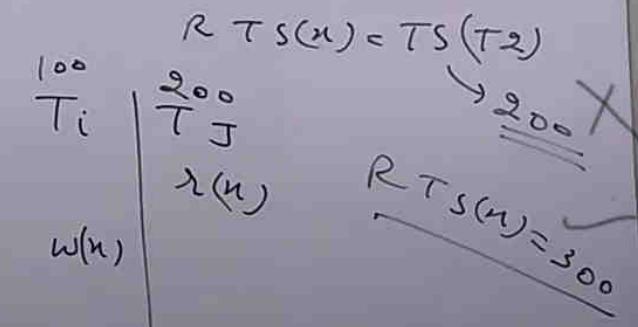
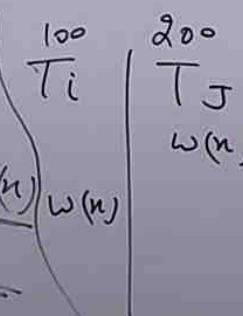
If T_i wants to write item 'n'

If $TS(T_i) < RTS(n)$ OR $TS(T_i) < WTS(n)$

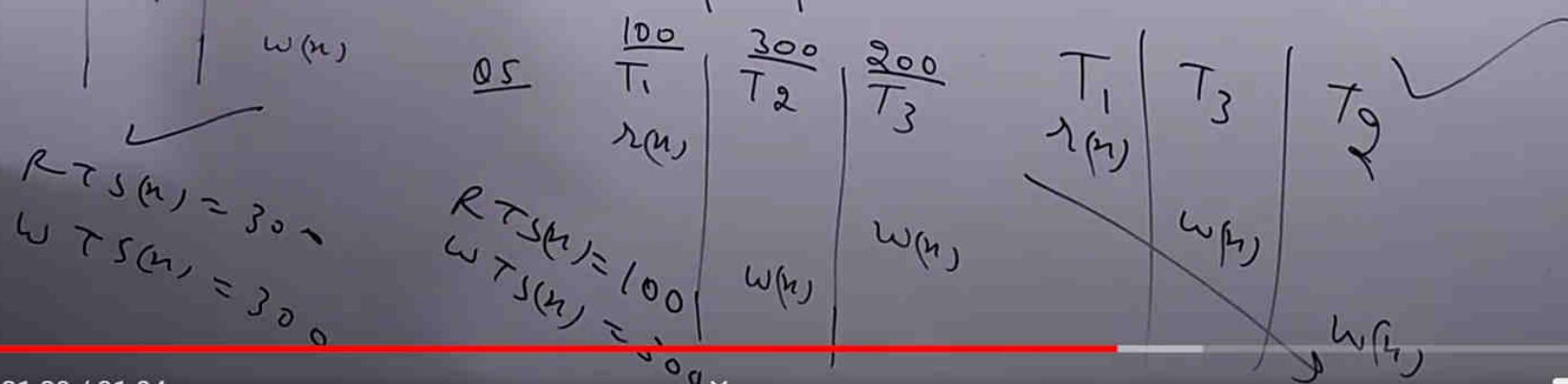
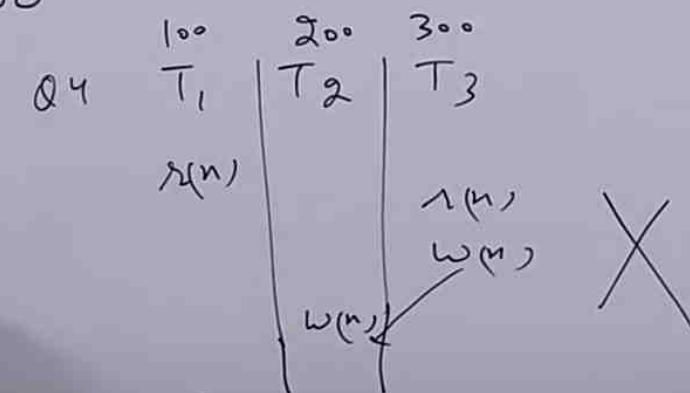
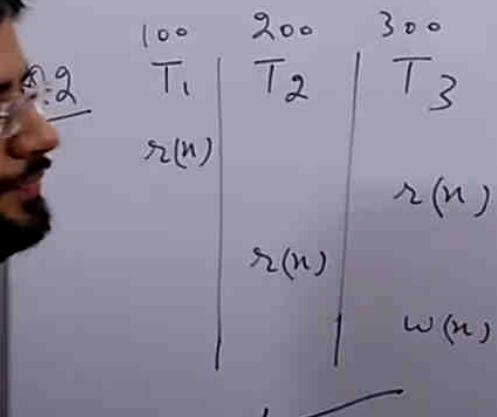
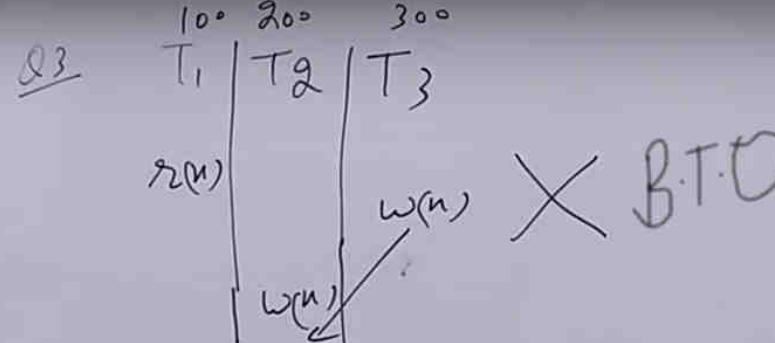
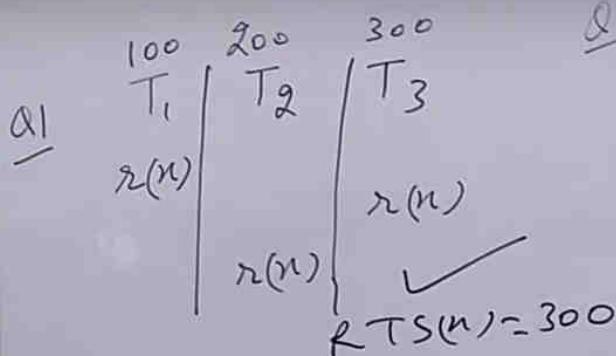
Abort T_i & restart it with new T.S.
else

Write operation granted

$$WTS(n) = TS(T_i)$$



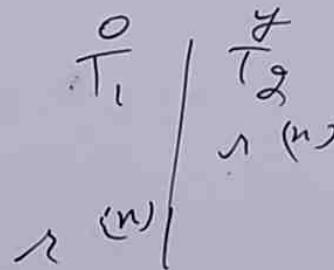
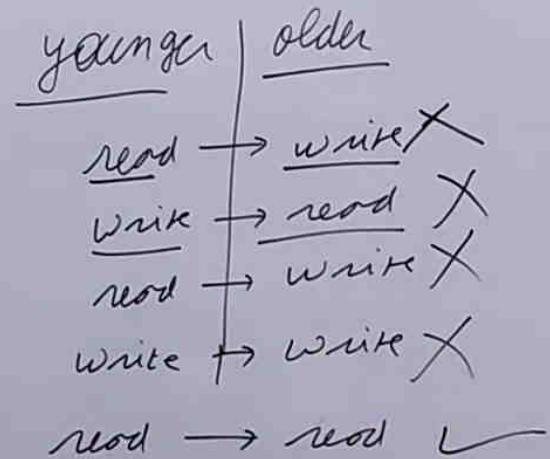
164. basic timestamp ordering protocol



21:29 / 31:04



164. basic timestamp ordering protocol



Advantages:-

- It guarantees conflict serializability of schedules.
- It is deadlock free technique.

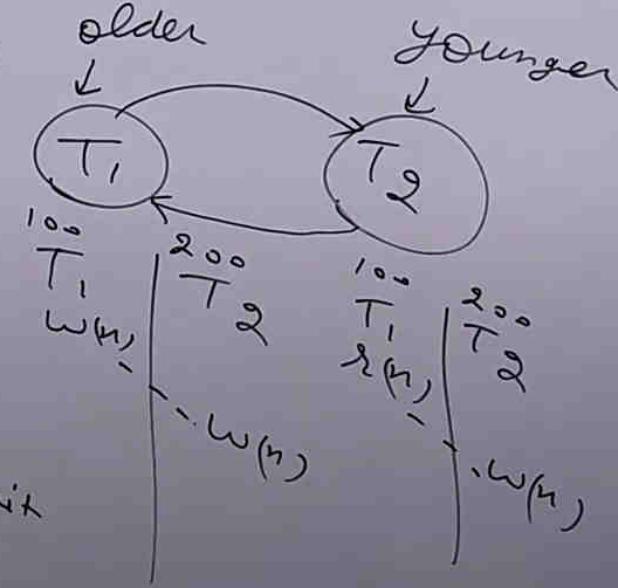
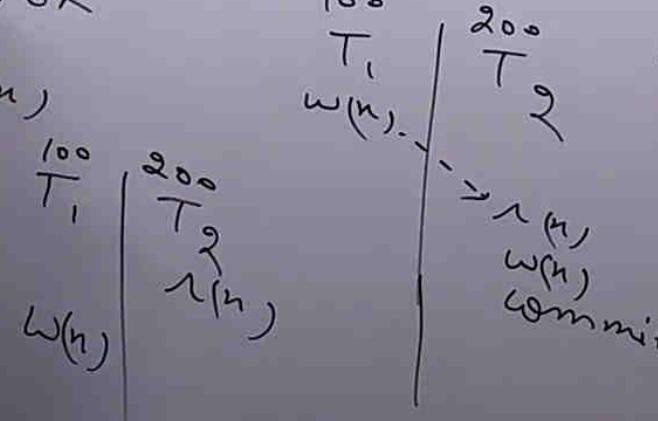
T_i T_j
 $w(n)$ $w(n)$

Disadvantages:-

- It does not guarantee recoverable, cascaded & strict schedules.
- All C-S schedules are not allowed (like 2PL)
- It is not starvation free (unlike 2PL)

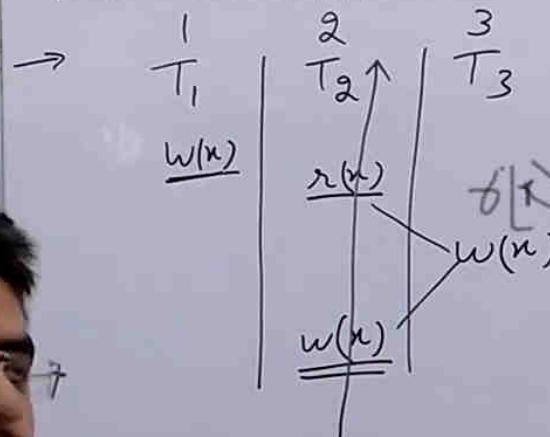
$T_i \rightarrow w(n)$
 $T_s(T_i) < R T_s(n) \text{ OK}$

$T_s(T_i) < w T_s(n)$



Thomas Write rule :-

→ Each transaction must read latest value of item (Not outdated value)



→ Here T2 has read latest value of n [written by T1]

→ T3 has read value of n written by T1 & when T3 read item n then system assumes T2 doesn't want to write n & hence for T3 value of n is latest.

→ Now after some time when T2 wants to write n then its against assumption of system hence it will reject w(n) of T2 & roll back it [\because the given schedule is not serializable]

Q What if younger has done blind write??

Any If T3 has done blind write then w(n) of T2 can be ignored & T2 can continue to execute

\because any Transaction Ti with $TS(T_i) < WTS(n)$ will be rolled back according to T.S.O.P

& any Transaction Tj with $TS(T_j) > WTS(n)$ will read value of n written by T3 ^{if it wants to read n}.

Thomas Write rule:-

→ If younger transaction has done $r(n)$ then $w(n)$ of older transaction will be repeated & it will be rolled back.

→ If younger transaction has done blind write then write operation (on same item) of older transaction can be ignored & it can continue execution.

<u>If $TS(T_2) > TS(T_1)$</u>	T_1	T_2	T_1	T_2	T_1	T_2	T_1	T_2
	<u>$r(n)$</u>	<u>$r(n)$</u>	<u>$r(n)$</u>	<u>$w(n)$</u>	<u>$w(n)$</u>	<u>$r(n)$</u>	<u>$w(n)$</u>	<u>$w(n)$</u>
$TSOP \rightarrow$	✓	X	X	X	X	X	X	X
$TWR \rightarrow$	✓	X	X	X	X	X	X	X

ignore
 $w(n)$ of T_1
& continue T_1

Thomas write rule:-

→ If younger transaction has done $r(n)$ then $w(n)$ of older transaction will be rejected & it will be rolled back.

→ If younger transaction has done blind write then write operation (on same item) of older transaction can be ignored & it can continue execution.

<u>If $TS(T_2) > TS(T_1)$</u>	<u>T_1</u>	<u>T_2</u>	<u>T_1</u>	<u>T_2</u>	<u>T_1</u>	<u>T_2</u>	<u>T_1</u>	<u>T_2</u>
		<u>$r(n)$</u>		<u>$w(n)$</u>		<u>$w(n)$</u>		<u>$w(n)$</u>

T_i wants to write item n

① If $TS(T_i) < RTS(n)$ then reject write request of T_i & roll back it

X

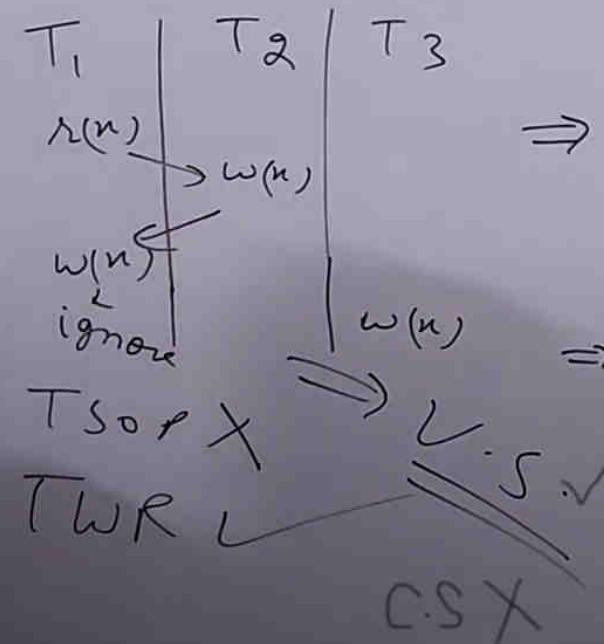
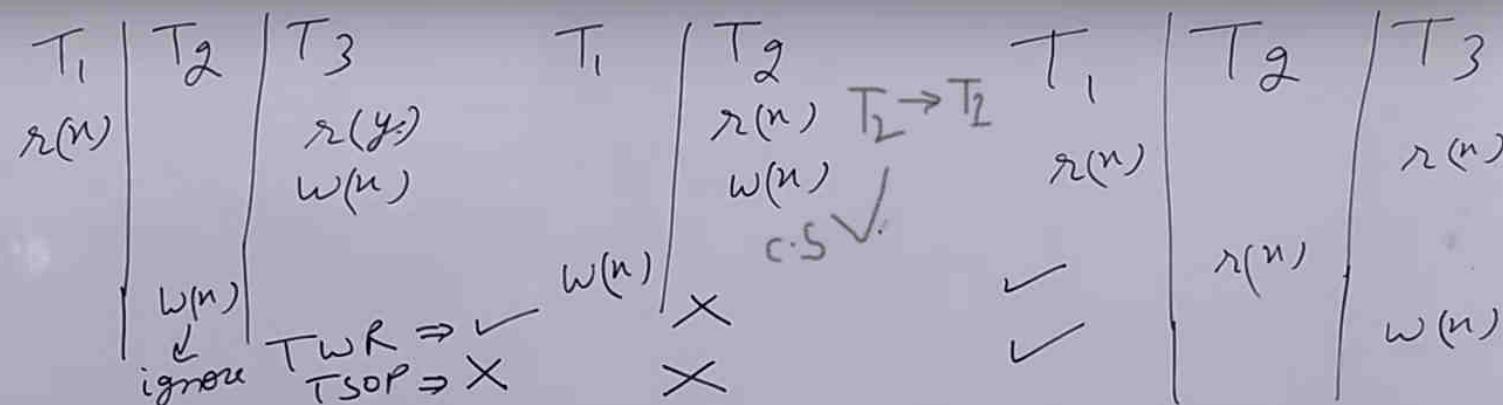
② else if $TS(T_i) < WTS(n)$ then ignore write request of T_i & continue it.

X

③ else write operation granted

$$WTS(n) = TS(T_i)$$

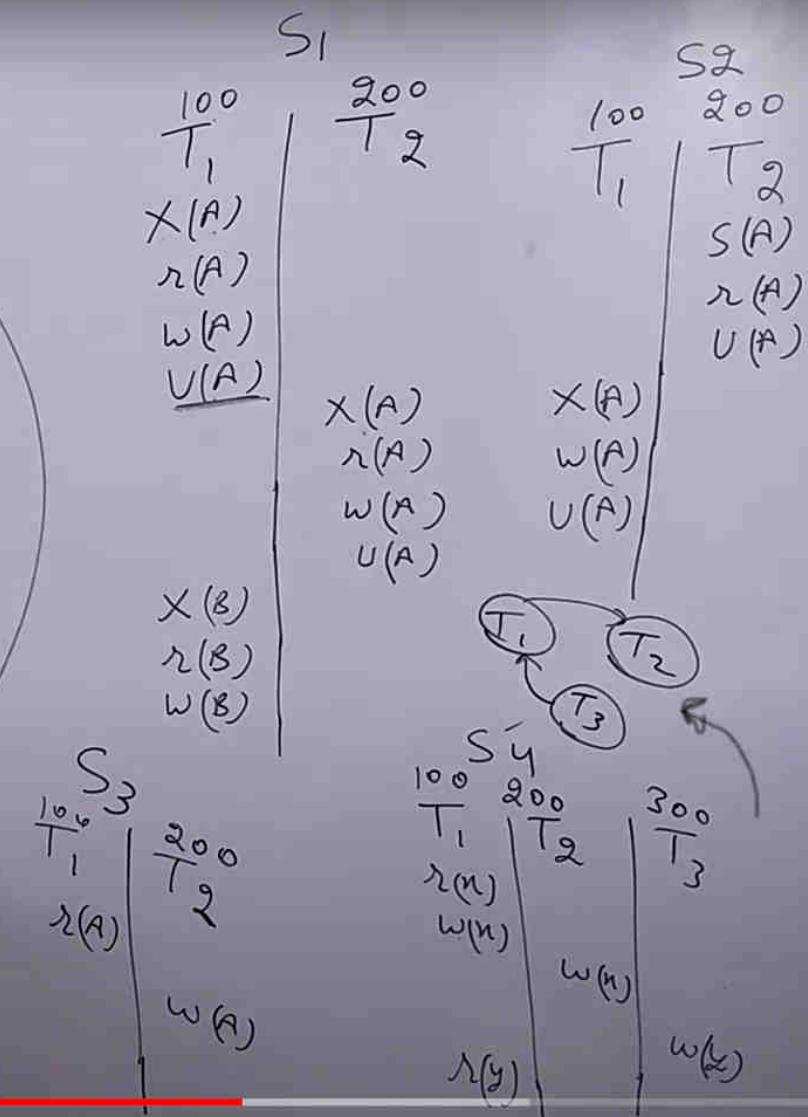
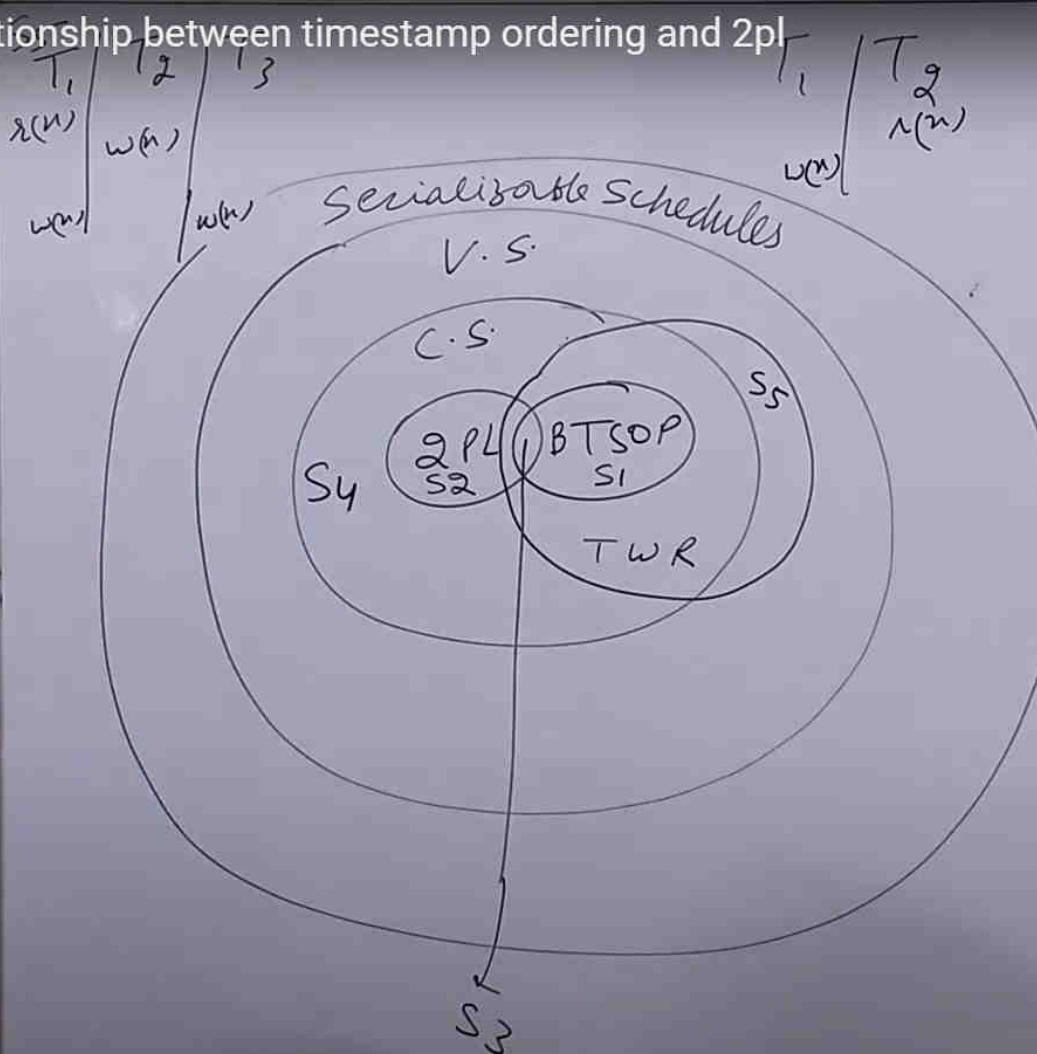
ignore
 $w(n)$ of T_i ,
& continue it)



⇒ TWR doesn't guarantee conflict serializability rather it is based on concept of view serializability. So it will result in fewer number of roll backs.

⇒ But all C.S. schedules are not allowed by T.W.R.

166. relationship between timestamp ordering and 2pl



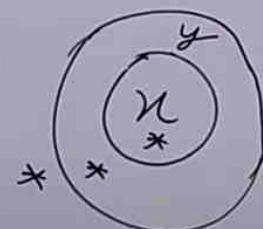
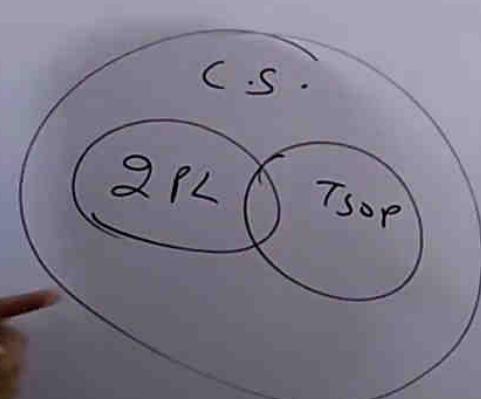
166. relationship between timestamp ordering and 2pl

Q1 2 PL is sufficient condⁿ for conflict serializability ✓ $\frac{VS}{VS}$

Q2 TSOP is sufficient condⁿ for ✓ X V.S. ✓

Q3 TWR

Q4 TWR is necessary condⁿ for conflict serializability X V.S. X



n is sufficient condⁿ ✗

Total Timestamp ordering protocol

- With each database item n we associate only one timestamp ie $TS(n)$
- $TS(n) \rightarrow$ Timestamp of youngest transaction that has performed read or write on item n .

T_i wants to read/write item n :

If $TS(T_i) < TS(n)$ then

abort T_i & restart it with new timestamp

else

read or write operation is granted

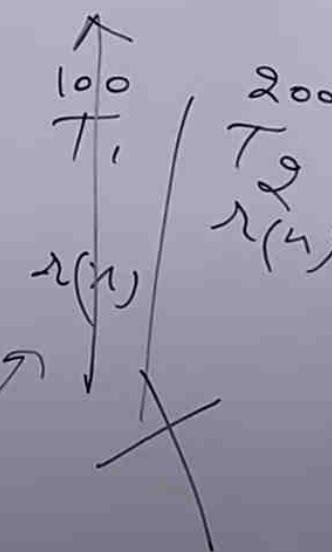
$$TS(n) = TS(T_i)$$

Advantage

① Implementation is easy & memory is saved.

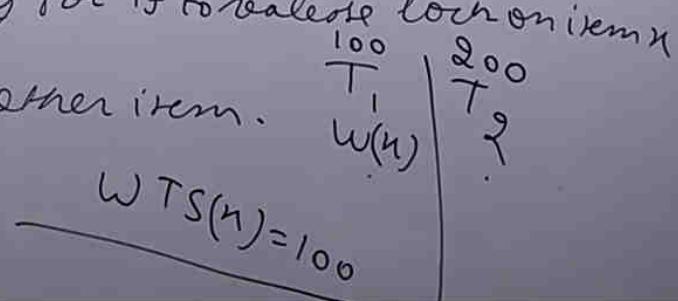
Disadvantage:-

① It will have unnecessary roll backs.



Strict Timestamp ordering protocol

- BTOP doesn't guarantee recoverable, cascadeless & strict schedules
- So we have a variant of BTOP called Strict Timestamp ordering protocol.
- As the name suggest it guarantees strict schedules hence recoverable & cascadeless is also guaranteed.
- If T_i wants to read or write item n such that $TS(T_i) > WTS(n)$ then this read or write operation is delayed till transaction T_j which wrote the value of n ($i.e. WTS(n) = T_j$) commits or aborts.
- To implement this scheme we must use concept of locking i.e. T_j must not release lock on item n till it commits or aborts.
- This technique is deadlock free ∵ if T_i is waiting for T_j to release lock on item n then T_j can't wait for T_i to release lock on any other item.

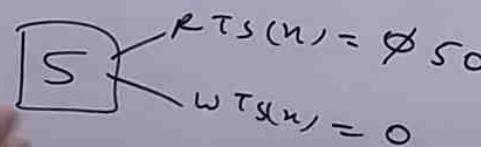


167. total, strict, multi version timestamp ordering protocol

Multi version timestamp ordering protocol

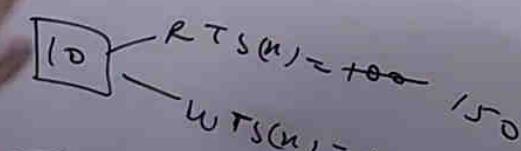
- As the name suggest we maintain multiple versions of same item 'n'
- Transaction which wants to read any item will read latest value of 'n' according to its timestamp.
- Lot of memory wasted but it will have least no. of rollbacks.
Rollbacks:-

TTSOP > B TSOP > m VTSOP

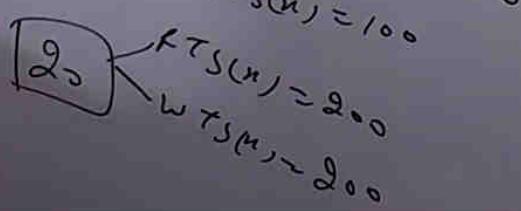


T_1^{100}
 $n^{(n)}$
 $U^{(n)}$

T_2^{200}
 $n^{(n)}$
 $U^{(n)}$



$T_3^{TS(T_3) = 50}$
 $T_4^{TS(T_4) = 150}$



16:09 / 16:46



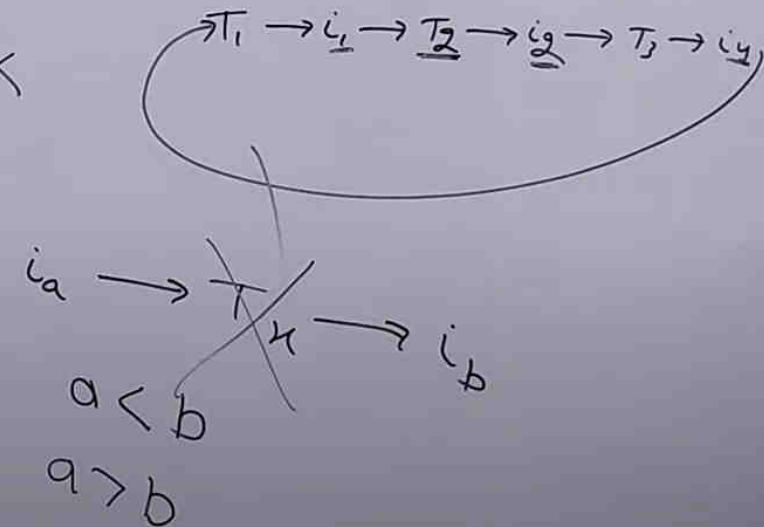
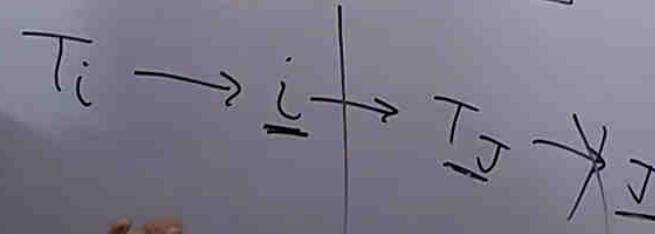
Deadlock Handling

- Deadlock Prevention
- Deadlock detection & recovery

- Use of conservative 2 PL → X
- * → Ordering all database items $i_1 \leq i_2 \leq i_3 \leq i_4$ ~~X~~
- Time out X

Premptive Scheme

- No waiting algorithm X
- Contious waiting algo
- Use of time stamp.



use of time stamp:-Wait & die

→ T_i wants an item held by T_j then:-

If $TS(T_i) < TS(T_j)$ then T_i is allowed to wait
else roll back T_i & restart it
with old time stamp.

→ Non pre-emptive

→ Both techniques are starvation free & Junior is always roll back.

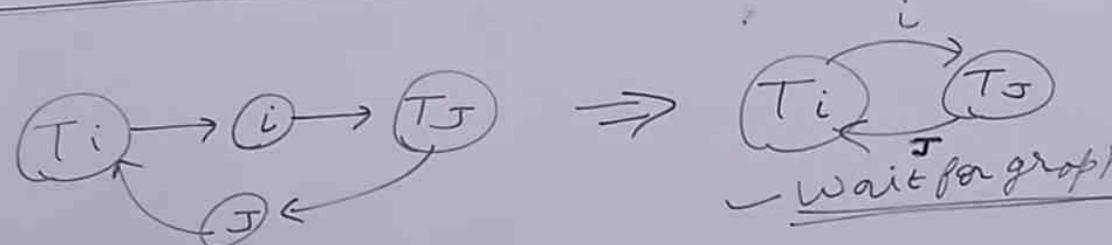
Wound & Wait

→ T_i wants an item held by T_j then:-

If $TS(T_i) < TS(T_j)$ then T_i bumps the item from T_j & T_j is rolled back & started with old time stamp

else T_i is allowed to wait.

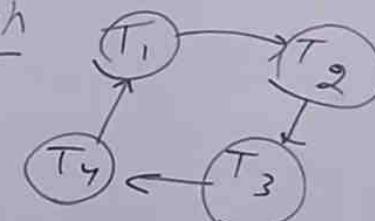
→ Pre-emptive

Deadlock Detection & RecoveryWait for graph is used for deadlock detection

→ cycle in wait for graph means deadlock

Recovery:-

- simplest way of deadlock recovery is to rollback any transaction.
- Transaction which is chosen for rollback is known as victim transaction
- Many criteria for choosing victims like youngest, few locks,



Consider the following 2PL. Let transaction T access set of objects $\{O_1, \dots, O_k\}$. This is done in following manner.

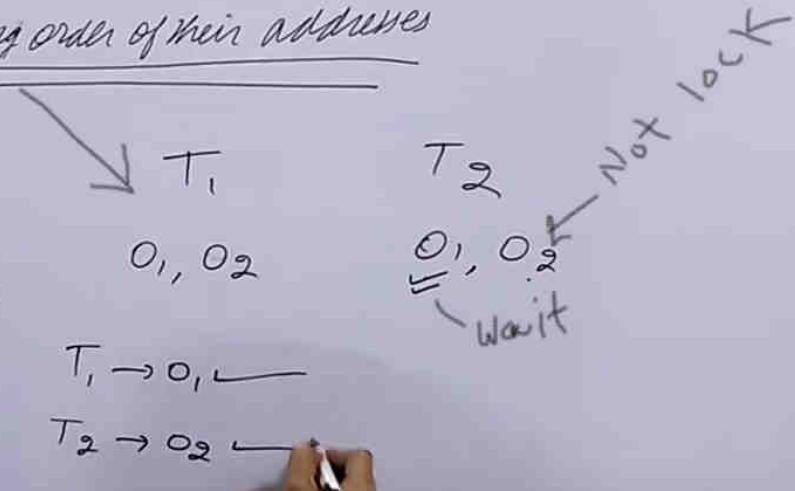
Step 1 T acquires locks on O_1, \dots, O_k in increasing order of their addresses

Step 2 The required operations are performed.

Step 3 All locks are released.

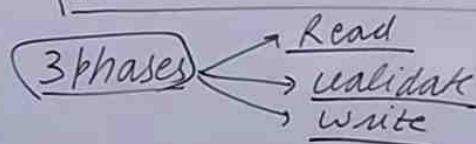
Guarantees serializability,

Guarantees deadlock



Optimistic Concurrency Control technique

- In pessimistic technique we have a common buffer in which transactions read or write data.
- But no concept of common buffer in optimistic technique.
- Hence each transaction has its own private buffer in which it read & write data.

Validation based protocol:-

T_1
→ $r(n)$

Read phase → Transaction read items from database & store them in buffer & then perform read & write operations blindly in buffer itself.

validate phase → When transaction completes then it enters validation phase where test for serializability is performed. If transaction passes validation test then it enters write phase else transaction is restarted.

write phase → All updated done in buffer are applied to database.

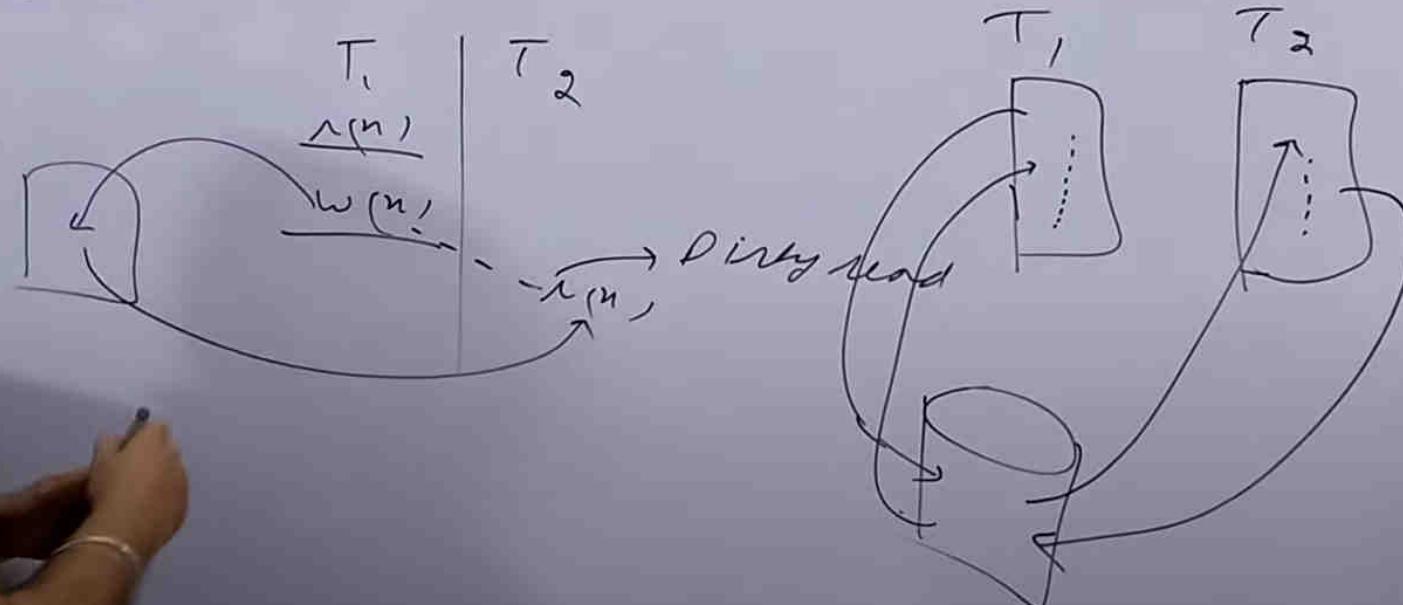
Exit full screen (f)

169. optimistic concurrency control technique

Advantages:-

- ① Optimistic techniques are very fast because transactions are executed blindly.
But in pessimistic techniques some check is performed before transaction read or write data.

It guarantees recoverable & cascading schedules \therefore No dirty read.



Log based Recovery

Recovery → Recovery is a technique to restore the database to most recent consistent state if any kind of failure occurs.

failure → disk failure, catastrophic failure, transaction failure

Logical error → no item found

System error → deadlock.

→ Hence we have to undo all the changes done by uncommitted transactions & redo all changes done by committed transactions.

→ Now DBMS remembers all this will be studied in log based recovery.

Log → Log is a special file maintained by DBMS to keep track of all transactions operations that affect the values of database items.

→ Log file is stored on stable storage so that data never loss

170. log based recovery

Log file :-

$\langle T_1, \text{start} \rangle \rightarrow \text{Start record}$
 $\langle T_1, x_{\text{old}} = 500, x_{\text{new}} = 1000 \rangle \rightarrow \text{Write record}$
 $\langle T_1, y_{\text{old}} = 1000, y_{\text{new}} = 1500 \rangle \rightarrow \text{Commit}$
 $\langle T_1, \text{commit} \rangle \rightarrow \text{Record lost line}$
 $\langle T_2, \text{start} \rangle \rightarrow \text{Start record}$
 $\langle T_2, z_{\text{old}} = 5000, z_{\text{new}} = 1500 \rangle \rightarrow \text{Write record}$
 $\langle T_2, \text{abort} \rangle \rightarrow \text{Record lost line}$

Deferred updates

- If failure occurs then Redo (Set values of all items updated by transaction to new value) those transactions for which commit is found.
- Simply ignore those transactions for which commit is not found i.e. No need of undo [∴ database is not modified].
- Hence name of recovery algo is No Undo/Redo algo.

- Whenever transaction performs the write, the log record for that write must be created before the database is modified otherwise recovery may not be possible.
- Whenever $\langle T_i, \text{commit} \rangle$ record is written in the log then this guarantees that all modifications done by transaction is successfully recorded in the log file in form of write records in stable storage.
- 3 possibilities
 - Database is not modified at all → Deferred update
 - Few writes are done in database (but can't guarantee all writes). → Restricted immediate update
 - All writes are done in database. → Immediate update

immediate update

before fail
case



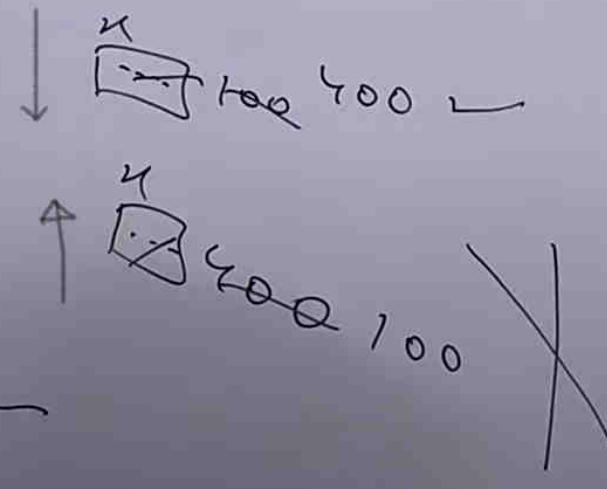
19:16 / 39:19



No undo/redo Algo:-

- This recovery algo is used in deferred updates scheme.
- After failure scan the log in forward direction.
- If $\langle T_i, \text{start} \rangle$ record is found but $\langle T_i, \text{commit} \rangle$ record is missing then ignore such transactions.
- If $\langle T_i, \text{start} \rangle$ & $\langle T_i, \text{commit} \rangle$ both records are present then redo T_i .
- Why forward ??

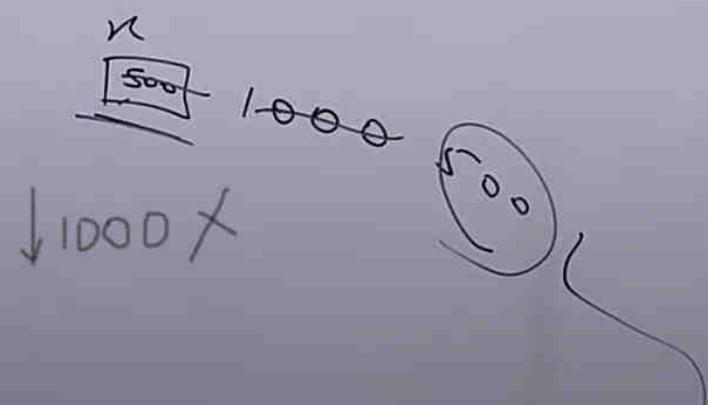
$\langle T_1, \text{start} \rangle$
 $\langle T_1, n_{new}=100 \rangle$
 $\langle T_1, \text{commit} \rangle$
 $\langle T_2, \text{start} \rangle$
 $\langle T_2, n_{new}=400 \rangle$
 $\langle T_2, \text{commit} \rangle$



Restricted immediate update

- Before writing $\langle T_i, \text{commit} \rangle$ record in log file DBMS guarantee that all changes done by transaction T_i are already written in database
- Hence if any failure occurs then scan the log file & if $\langle T_i, \text{start} \rangle \& \langle T_i, \text{commit} \rangle$ both log records are present then ignore T_i .
- If $\langle T_i, \text{start} \rangle$ is present but $\langle T_i, \text{commit} \rangle$ is absent then undo T_i .
Hence name of algo is
- Direction ??
No redo/undo

$\langle T_1, \text{start} \rangle$
 $\langle T_1, n_{\text{old}}=500 \rangle$
 $\langle T_2, \text{start} \rangle$
 $\langle T_2, n_{\text{old}}=1000 \rangle$



Immediate updates

- Before writing $\langle T_i, \text{commit} \rangle$ record in log file No guarantee that all changes done by transaction T_i are already written in database.
- Hence if any failure occurs then scan the log file & if $\langle T_i, \text{start} \rangle$ & $\langle T_i, \text{commit} \rangle$ both log records are present then redo T_i .
- If $\langle T_i, \text{start} \rangle$ is present but $\langle T_i, \text{commit} \rangle$ is absent then undo T_i .
- Hence name of algo is undo/redo algo.
- So in log file for each write record both old & new value must be stored.
- If failure occurs then make 2 lists of transactions ie redo list containing transactions ID for which commit operation is found in log & undolist containing transactions ID for which commit operation is missing
1st do undo in backward direction & then redo in forward direction.

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

$\langle T_1, y_{\text{old}} = 500, y_{\text{new}} = 1000 \rangle$

$\langle T_2, \text{start} \rangle$

$\langle T_2, y_{\text{old}} = 1000, y_{\text{new}} = 2000 \rangle$

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

No undo/redo (Deferred updates)

→ Redo T_1

No redo/undo (Restricted immediate updates)

→ Undo T_2

redo/undo (Immediate updates)

→ Redo T_1 , Undo T_2

1. T₁: start
 2. T₁: bold = 12000, new = 10000
 3. T₁: old = 0, new = 2000
 4. T₁: commit
 5. T₂: start
 6. T₂: bold = 10500
 7. T₂: GO
- Database had just been written before log was written

T₁ - 3
6

then (let)
2, 3
undo

Q1 T₁ ⇒ 2, 3, 4, 5 & commit is present

T₂ ⇒ 7, 8, 9 & commit is absent

1st Undo T₂ ⇒ 9 → 8 → 7 then

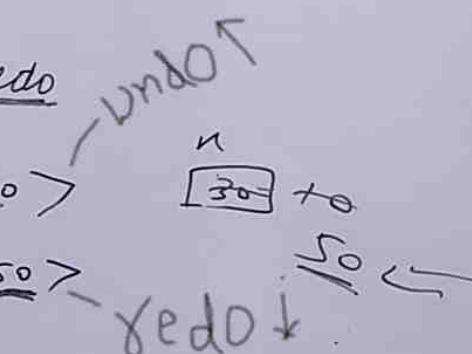
Redo T₁ ⇒ 2 → 3 → 4 → 5

Q2 Why 1st undo & then redo

↑ <T₁, n_{old} = 10, n_{new} = 30>

↓ <T₂, n_{old} = 10, n_{new} = 50>

↓ <T₂, commit>

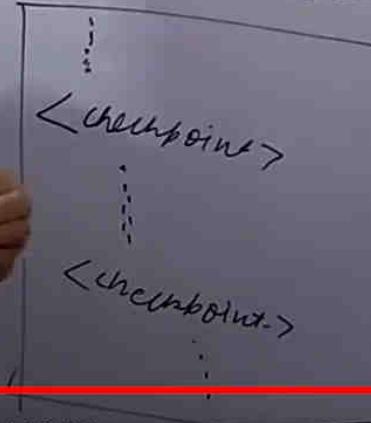


Q3 Why Undo in backward direction & Redo in forward direction

✓ <T₁, n_{old} = 10, n_{new} = 20>
✗ <T₁, n_{old} = 20, n_{new} = 30>

Checkpoints:-

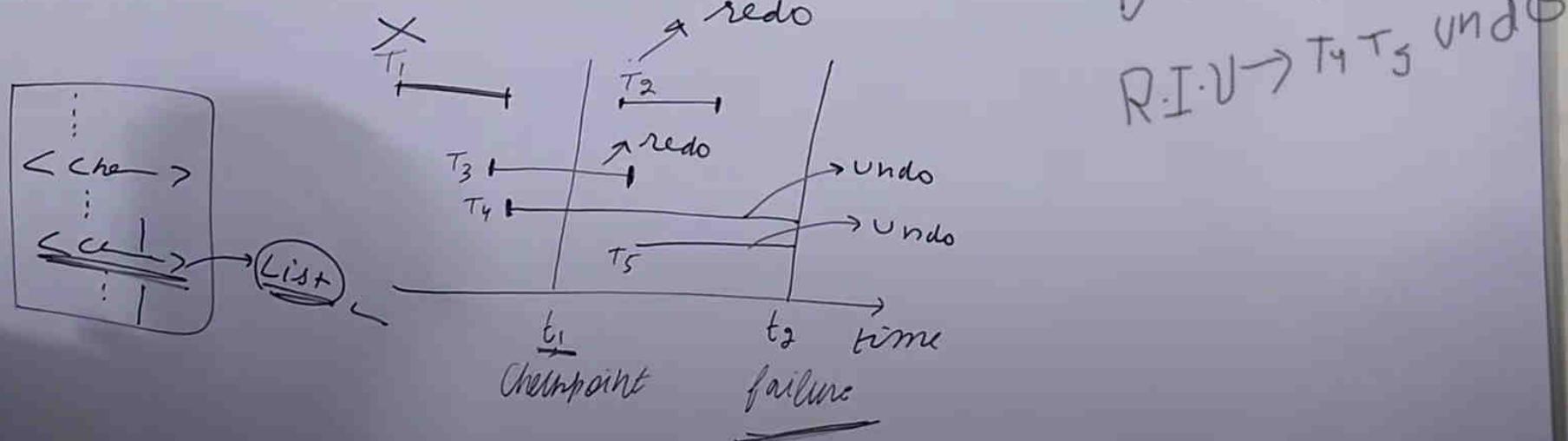
- We discuss checkpoints for immediate updates (undo/redo algo).
- If size of log file is very large then we have 2 problems:-
 - ① Search process to make redo & undo list is time consuming
 - ② Most of the transactions updates may be written to the database so if we redo them again then it is wastage of time.
- Hence system do checkpoints from time to time which require 3 steps.
 - ① Store all log records currently in RAM to stable storage.
 - ② All modified buffer blocks are written to Database.
 - ③ <checkpoint> log record is written to stable storage.



Note:- Transactions are not allowed to perform any update actions like writing to a buffer block or writing a log record while checkpoint is in progress.

Recovery in case of checkpoints.

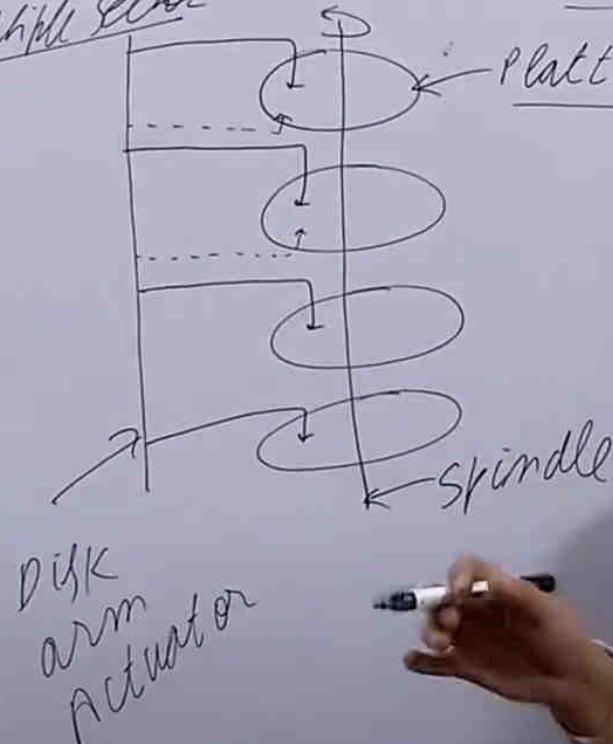
- If T_i commits before <checkpoint> then we can ignore it.
- Scan the log file in backward direction & find 1st checkpoint (actually it is lost or most recent checkpoint)
- DBMS will make list of transactions that were active during checkpoint.
- From this list make redo & undo lists



$D \rightarrow T_1 T_3 \text{耶} \log$
 $R.I.U \rightarrow T_4 T_5 \text{ undo}$

Block → 1 sector
or
multiple sectors

0 & 1
→
magnetic material



Seek time → track
Latency time → sector time



SIT

Platter

Tracks

Spindle

Sector → 512 B

M-T-U

▼

Q1 Disk contains 100 tracks & 500 sectors find Average seek time & Average latency time if disk is rotating @ 3600 rpm & time needed by read/write head to jump from one track to another is 6 msec.

Ans Avg seek time = $\frac{1}{3} \times \cancel{\max^n \text{ seek time}} = \frac{1}{3} \times 99 \times 6 \text{ msec}$

$$\frac{1}{2} \times \cancel{(n-1)} \times \cancel{\text{Time Jump}} \leftarrow \frac{1}{2} \times \cancel{\max^n \text{ seek time}} = \frac{1}{2} \times 99 \times 6 \text{ msec}$$

$$\frac{1 \rightarrow 1 + 1 \rightarrow 2 + 1 \rightarrow 3 + \dots + 1 \rightarrow 100}{100} =$$

$$\frac{\frac{1}{2} \times \cancel{\text{Resolution time}}}{3600} \rightarrow 60 \text{ Sec} \quad \frac{0 + 1 + 2 + \dots + 99}{100} = \frac{9.9 \times 100}{2 \times 100} = 49.5$$

$$\frac{1}{2} \rightarrow \frac{60 \times \frac{1}{2} \text{ Sec}}{3600}$$

Q2 If track size is 50 KB ^{bytes} disk rotation speed is 3600 rpm. Find transfer rate

Ans. Amount of data transferred from HD to RAM in 1sec is known as transfer rate.

$$3600 \rightarrow 60 \text{ sec}$$

$$\textcircled{1} \rightarrow \frac{60}{3600} \text{ sec}$$

$$\frac{60}{3600} \text{ sec} \rightarrow 50 \text{ KB}$$

$$1 \text{ sec} \rightarrow \frac{50 \times 60}{60} \text{ KB} \quad \nearrow 3000 \text{ KB per sec}$$

24 MBPS

3 MBPS

3 MB per sec

Q3 If track size is 50 KB ^{bytes} disk rotation speed is 3600 rpm . Find transfer rate & Block transfer time if block size is 1 KB

Ans. Once read/write head is placed at beginning of correct block some time is needed to transfer data of block from RAM that time is block transfer time.

$$\text{Transfer rate} = 3 \text{ MB/s}$$

$$3000 \text{ KB} \rightarrow 1 \text{ sec}$$

$$1 \text{ KB} \rightarrow \frac{1}{3000} \text{ sec}$$

Block Size

Transfer rate

Time needed to find &
Transfer block from RAM

$$S + L + \frac{B}{T.R.}$$

Q4 If disk contain 200 tracks with track size of 50KB & rotating at speed of 3600 rpm. What is the time needed to find & transfer 1 block of size 1KB from HDD to RAM. Jump time = 6 msec

Ans

$$\text{find} \Rightarrow A.S.T + A.L.T$$

$$3600 \rightarrow 60 \\ 1 \rightarrow \frac{60}{3600}$$

$$\text{find} = \frac{1}{2} \times \frac{199 \times 60}{1000} + \frac{1}{2} \times \frac{60}{3600} \quad \frac{1}{2} \rightarrow$$

$$\text{Transfer time} = \frac{\text{Block size}}{\text{T.R.}} = \frac{1 \text{ KB}}{3000 \text{ KB}}$$

$$\text{Ser} = \boxed{\frac{1}{3000} \text{ sec}}$$

$$\frac{60}{1000} \times \frac{1}{2} \times 199 + \frac{1}{2} \times \frac{60}{3600} + \frac{1}{3000}$$

173. spanned vs unspanned strategy

File organization

Basic concepts:-

- Database → Collection of files → Collection of records (tuples) → Collection of fields
- A table may be stored as single file or file may contain multiple tables or table may be stored as multiple files
- Q. How records are stored on H.D.?



name varchar(20):



Fixed length records

→ Size of record is fixed

← Block

R₁ R₂ R₃

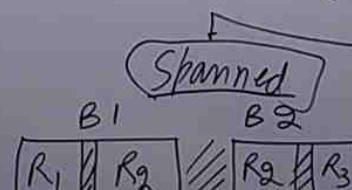
Internal fragmentation

Inter record gap

Unspanned strategy is used

Variable length records

→ Size of record is not fixed ∵ size of field is not fixed like varchar(20) or File contain multiple tables or we don't want to store repeating values



Unspanned

Spanned

Interblock gap



10:02 / 11:32



174. pile or heap file organization

- Pile or heap
- Sequential file
- Induced sequential file
- B⁺ Tree
- B Tree

Pile or Heap :-

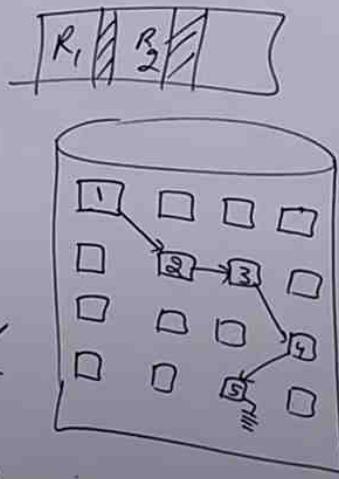
- File is a Link list of blocks
- Block may contain multiple records
- If records in a file are unsorted
Then we call it as Pile
- Here records are stacked up
or piled up one after another.

→ Insertion of record is done by stacking

→ Searching of record is brute force ie linear search
so we have to read in block one by one & search
each & every record in it.

→ If we want to delete any record then 1st we have
to search it & then we can delete it.

→ Deletion is not physical.



174. pile or heap file organization

Q1 If file is stored in 'B' blocks then on an average how many blocks are searched for particular record?

Ans $\frac{B}{2}$

Q2 If we want to search records on the basis of field which can be repeated then on an average how many blocks are searched? B blocks

Ans B

Q3 In file searching is block by block (block search) or record by record (record search)

Ans Record search

Hence file can't be used for larger databases \therefore search is major operation
which is inefficient in file.

Sequential file

- Records are kept in sequential order i.e in sorted fashion on some attribute(s) & usually this attribute(s) is primary key.
- Strict sequential file

<u>bid</u>	<u>title</u>	<u>YR-PUB</u>
1	A	N
3	B	Y
5	C	Z

b1

8	-	-
10	-	-
15	-	-

b2

17	-	-
----	---	---

b3

(2, -, -)

Inserion → If we want to insert a new record then 1st we have to find its position & then it is stored at its proper place.

→ In worst case we have to read & write entire file (if insertion is at begining)

Deletion → Deletion is physical i.e records are shifted

Searching → If we assume all blocks are present in RAM then apply binary search $\Rightarrow \lceil \log_2 B \rceil$ blocks are searched in worst case.

→ If blocks are not in RAM then still searching is efficient than file because here we have block search & not record search.

175. sequential file organization

Normal sequential file

- In Normal sequential file, we insert tuples in the end by changing 2 new pointer values.
- Hence records are logically sorted but physically unsorted

	<u>bid</u>	<u>title</u>	<u>YR-PUB</u>	<u>Next</u>
B ₁	1	A	X	Start
	3	B	Y	
	5	C	Z	

B ₂	8	-	-	
	10	-	-	
	15	-	-	

B ₃	17	-	-	
	2	-	-	



To physically sort the records in file time to time reorganization is done

→ Deletion is not physical & we change only 2 next pointer values for deletion.

→ 2 Special pointers are used i.e start & available

→ Start always points to 1st record of file & it forms the link list of sorted records

→ Available points to 1st record which is deleted

↳ of deleted Record

Indexed Sequential file :-

- Main disadvantage of sequential file is that it is always sequential access ∵ of Link List
- Hence binary search cannot be applied.
- If we assume entire file is in RAM or if we assume file header contain address of all blocks of file then we can apply binary search.

Q How to make sequential file as random access? file descriptor

Ans search key

bid	BLOCK pointer
1	2000
4	1000
7	4000

Index file

bid	title	YR-PUB	Nent
2	-	-	\$
3	-	-	\$
4	-	-	\$

bid	title	YR-PUB	Nent
5	-	-	\$
6	-	-	\$
7	-	-	\$

bid	title	YR-PUB	Nent
8	-	-	\$
9	-	-	\$

ordering field

Index file + sequential file
build on ordering field.

- We can make this sequential file as random access by maintaining a small file called index file which contains only 2 columns i.e. search key(field) & block(record) pointer
- This file is always sorted according to search field value.

Indexed Sequential file :-Search key :-

- It is an attribute of table which is present in index file. It may or may not be Primary key
- Search key is a field on which index is built.
- It may or may not be ordering field.

What is the maximum number of disk access required to fetch a record from indexed sequential file assuming index file is in RAM?

1

In above question if index file is in HD & all pointers to index block are available then how many disk access?

$$\lceil \log_2 B_i \rceil + 1$$

if not index file
 $\lceil \log B \rceil$

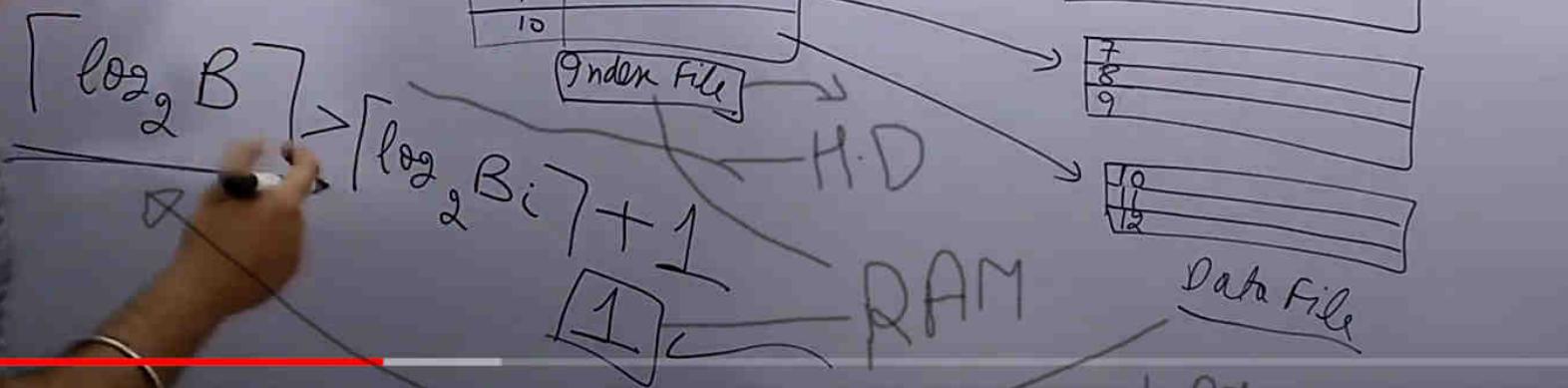
$$B \gg B_i$$

Indexes

- Indexes in database are similar to indexes at last of books.
- They provide faster way of accessing records of file.
- For a file we can have more than one index also & in fact for every field of file we can build index (but not good idea)
- If lots of query on book table involve conditions on title & bid then make 2 index files 1st having search key as title & 2nd has search key as bid.

SELECT *
FROM book
WHERE bid=5

$$\lceil \log_2 B \rceil$$



Types of Indexes

single level Ordered Indexes

Index file always

Sorted acc. to Search key

Multilevel Indexes

Dynamic Multilevel Indexes (B-tree, B+tree)

Primary Index

- Index is build on ordering field which is candidate key of data file.

- File is sorted acc. to C.K. & on that C.K. index is build

clustered Index

- Index is build on ordering field which is not C.K. of data file (i.e. repeating)

- File is sorted acc. to non key & that non key is used to build index.

Secondary Index

- Index is build on non ordering field which may or may not be C.K.

B.I.D	B.P.

unsorted &

like

→ sorted

Types of Indexes

Single level ordered Indexes

A file can have almost one primary index

file can have almost one clustered index

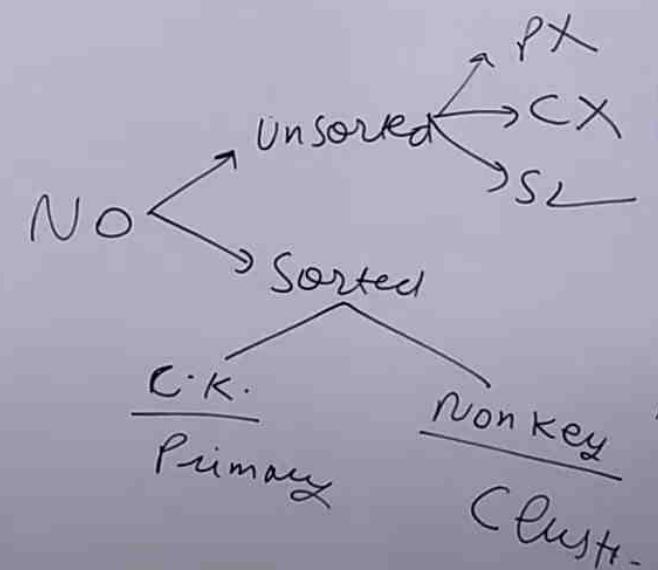
file can have any NO. secondary index

file can have both primary & clustered index

A file

Multilevel Indexes

Dynamic Multilevel Indexes (Btree, B+tree)



Types of Indexes

single level ordered indexes

multilevel indexes

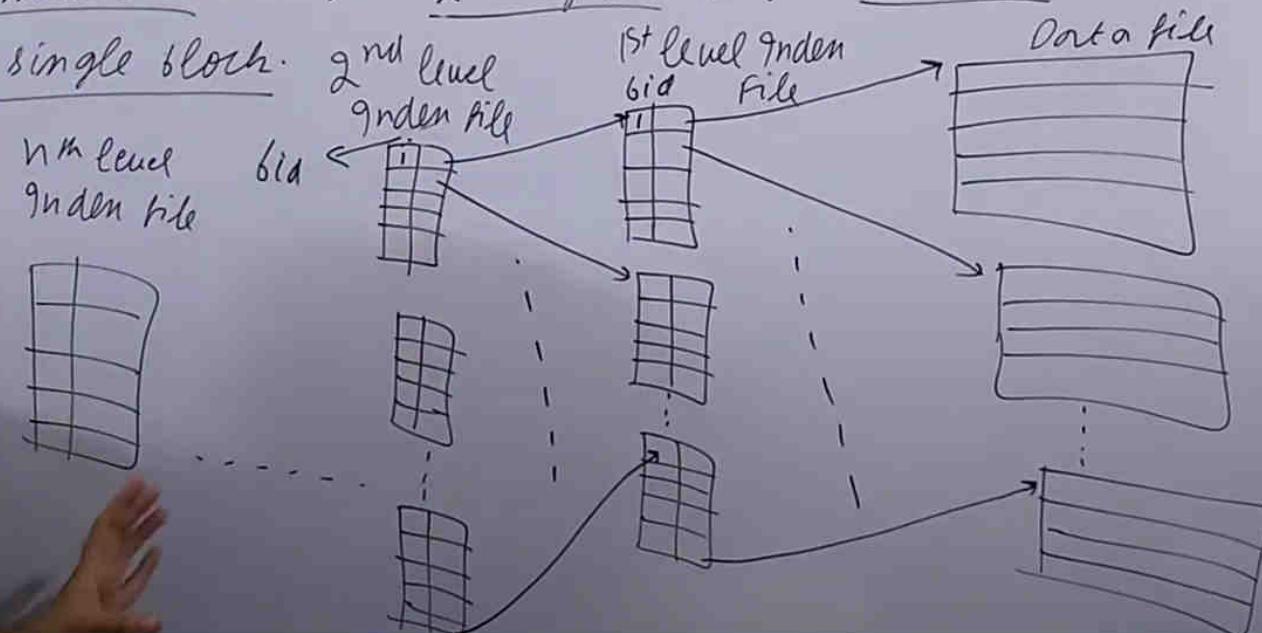
Dynamic multilevel indexes (Btree, B+tree)

multilevel indexes :-

Since index file at 1st level is sorted according to Search key

we can make index of index file itself till we get index file which can be stored in single block.

sorted





Types of Indexes

Pense

Index entry for every search

Key value (hence for every record)

$$\frac{\text{No. of records}}{\text{in index file}} = \frac{\text{No. of records in Data file}}{\text{Data file}}$$

Sparse

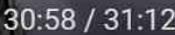
Index entry for only some of search key values (not for every record)

$$\frac{\text{No. of records}}{\text{in index file}} < \frac{\text{No. of records in data file}}{\text{Data file}}$$

$$\underline{B} \gg \underline{B_i}$$

$$\underline{\underline{B}} \gg \underline{\underline{B_i}}$$

$$m < n$$

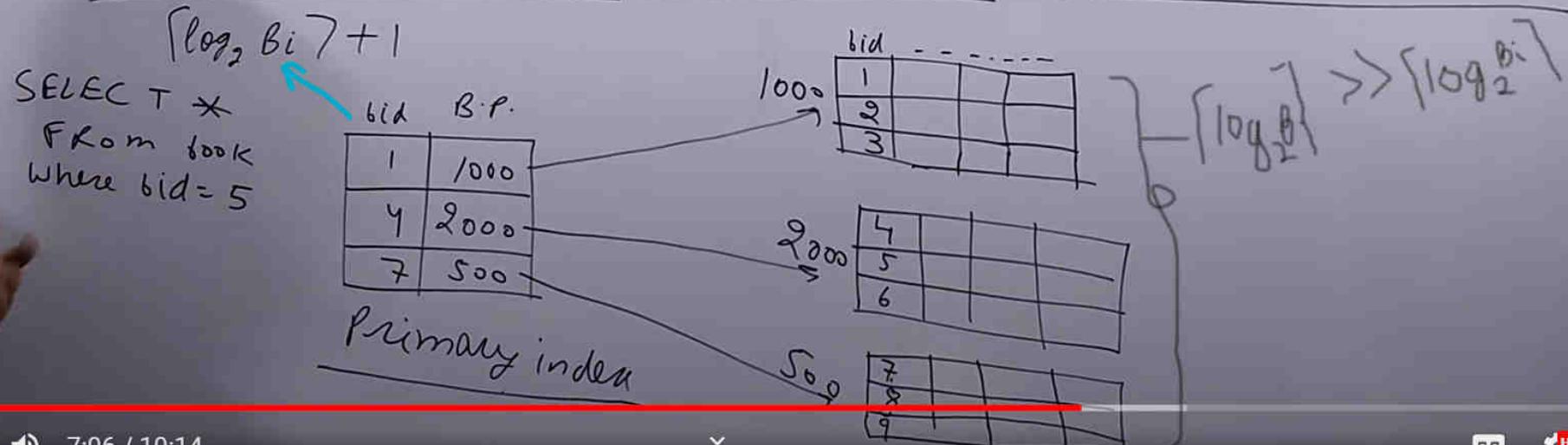


30:58 / 31:12



Primary Index

- Primary index is an ordered file with 2 fields i.e. Search key & Block pointer
- Search key is an ordering field of data file & this ordering field must be C.K.
- There is one index entry (index record) for each block of data file
- 1st record in each block of data file is called anchor record or block anchor
- Primary index is non dense or sparse index : index file has one entry for each block rather than each record.
- $B \gg B_i$: size of each record is smaller in index file & No. of records are also less.
- No. of disk access required to fetch a record given a search key value will be



179. numerical on primary indexing

A file consists of 50,000 records, each record is of 100 bytes. Block size is 1024 bytes. If primary indexing is used & block pointer is of 6 bytes & primary key is of 2 bytes

- Q1 What is blocking factor for data file & index file
 - Q2 No. of blocks needed to store data file & index file
 - Q3 No. of block accesses to answer data with & w/o indexing assuming all index & data block addresses are available but their blocks are in H.P.
- How many block accesses if multilevel indexing is used.

Blocking Factor → No. of records which can be stored in each block

Q1

$$\left\lceil \frac{\text{Block size}}{\text{Record size}} \right\rceil = \left\lceil \frac{1024}{8} \right\rceil = \left\lceil \frac{2^10}{2^3} \right\rceil = 2^7$$

$$\text{Data file} = \left\lceil \frac{1024}{100} \right\rceil = 10$$

= Index file

128



3:21 / 15:49



179. numerical on primary indexing
Total no. of records = 30,000, each record is of 100 bytes. Block size is 1024 bytes.

If primary indexing is used & block pointer is of 6 bytes & primary key is of 2 bytes

Q1 What is blocking factor for data file & index file

Q2 No. of blocks needed to store data file & index file

Q3 No. of block accesses to answer data with & w/o indexing assuming all index & data block addresses are available but their blocks are in H.P.

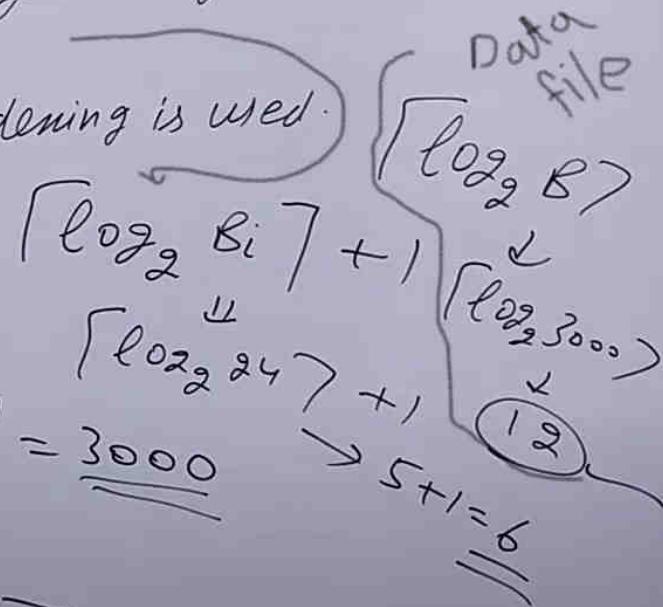
Q4 How many block accesses if multilevel indexing is used.

Blocking factor for index file $BFR_i = 128$

., ., ., data file $BFR_D = 10$

$$\text{No. of blocks for data file} = \left\lceil \frac{30,000}{10} \right\rceil = \underline{\underline{3000}}$$

$$\text{., ., ., index file} = \left\lceil \frac{3000}{128} \right\rceil = 24$$



8:25 / 15:49



179. numerical on primary indexing
A file contains 30,000 ^{UNSPANNED} records, each record is of 100 bytes. Block size is 1024 bytes.

If primary indexing is used & block pointer is of 6 bytes & primary key is of 2 bytes

Q1 What is blocking factor for data file & index file

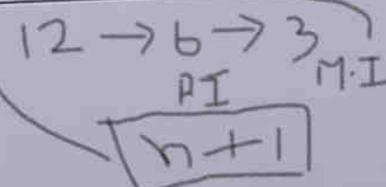
Q2 No. of blocks needed to store data file & index file

Q3 No. of block accesses to answer data with & w/o indexing assuming all index & data block addresses are available but their blocks are in H.P.

Q4 How many block accesses if multilevel indexing is used.

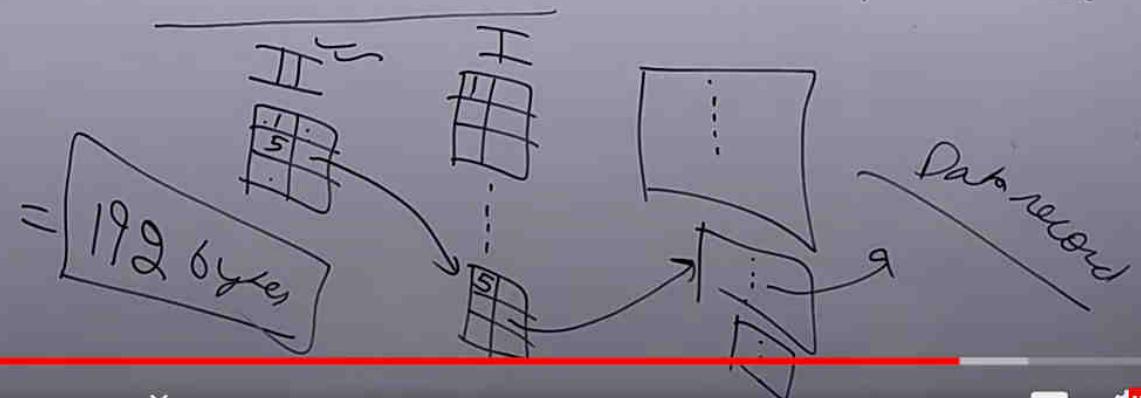
$$\text{No. of blocks for data file} = 3000$$

$$\dots \text{index} \dots = 24$$



→ Since index file is stored in more than one block so we can do indexing of index file itself which is known as multilevel indexing & we will stop when we get index file of one block.

$$3 \leftarrow \underbrace{(1 + 1 + 1)}_{\substack{\text{I} \\ \text{II}}} \quad 24 \times 8 = 192 \text{ bytes}$$



clustered Index :-

- A clustered index is an ordered file with 2 fields i.e. search key & block pointer
- Search key is an ordering field of file & it does not have unique value for each record i.e. it must be non key.
- We call this ordering field as clustering field :: it is used to make clusters of records.
- There is one index entry (index record) for each distinct value of clustering field.
- Block pointer points to the block containing 1st record with respective search key value.
- Clustering index is also wondense or sparse :: it has index entry for every distinct value of clustering field rather than for every record of file.

SK	B.P.
1	1000
2	1000
3	2000
4	1500
5	1500

Clustered Index
File

1	
1	
1	
2	

1000 clustering field

2	
2	
3	
3	

2000

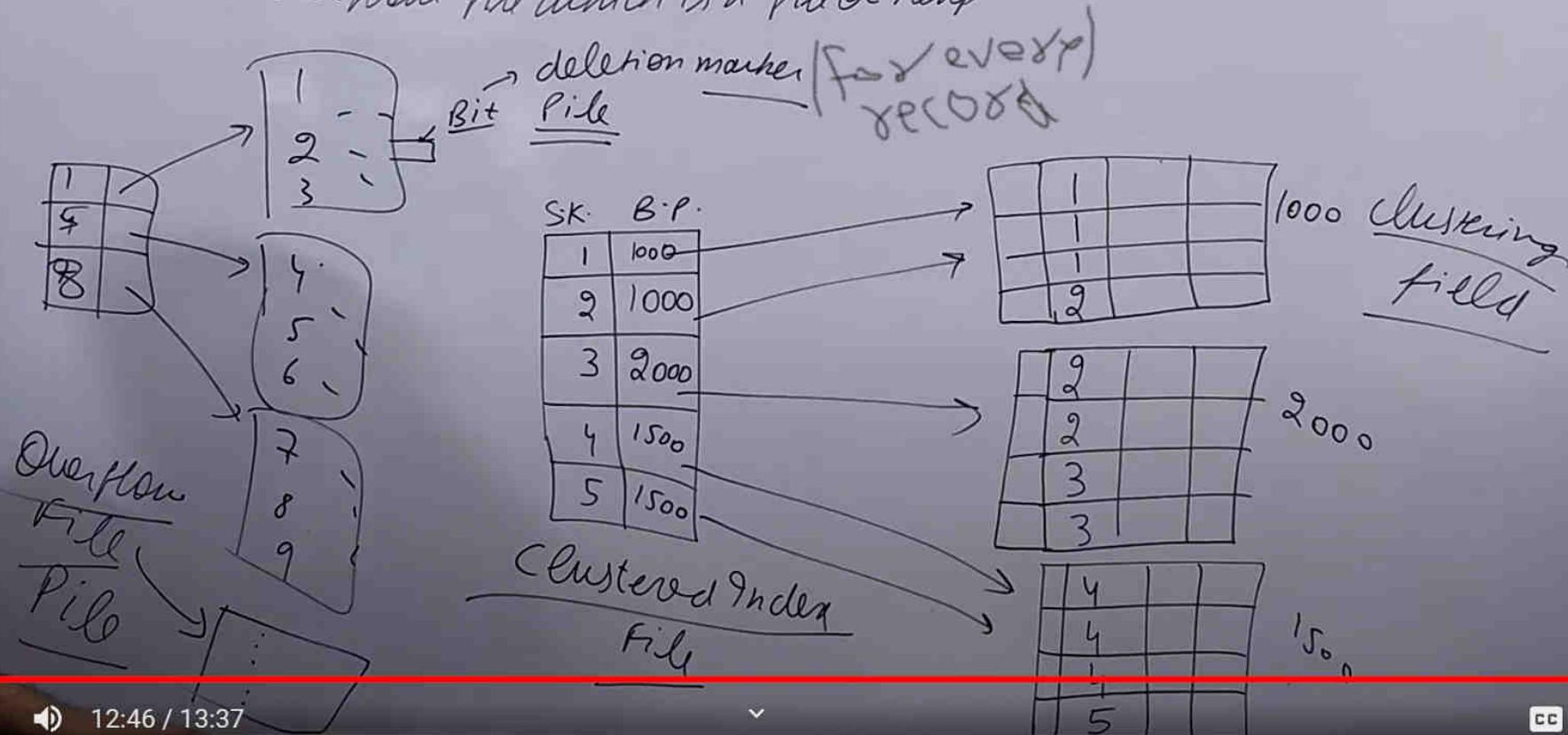
4	
4	
4	
5	

1500

180. clustered indexing

clustered index :-

- Insertion & deletion of records in both primary & clustered indexing is time consuming.
- In clustered indexing we maintain a separate block or cluster of contiguous blocks for each unique clustering field value.
- In primary indexing deletion is handled by deletion markers & for insertion we can use overflow file which is a pile or heap.



181. numerical on clustered indexing

If clustered indexing is used on a field such that every value is repeated 6 times. Size of block pointer is 6 bytes & clustering field value is of 9 bytes

~~30000
6~~

Q1 What is blocking factor for data file & index file

No. of blocks needed to store index file & data file

No. of block accesses to access data assuming index file is in RAM & all block addresses are available

How many block accesses if multilevel indexing is used.

$$BFR_D = 10 \checkmark$$

$$BFR_I = 68 \checkmark$$

$$\left\lceil \frac{\text{No. of records}}{BFR_D} \right\rceil$$

$$\left[\frac{30,000}{10} \rightarrow 3000 \right]$$

No. of blocks needed to store Data file

" " " " " " index file

$$\left[\frac{56000}{68} \right] \rightarrow 824 \text{ unique}$$



5:04 / 16:46



181. numerical on clustered indexing
no of records = 747, no of unspanned records each of size 100 bytes - block size is 1024 bytes

If clustered indexing is used on a field such that every value is repeated 6 times. Size of block pointer is 6 bytes & clustering field value is of 9 bytes

Q1. What is blocking factor for data file & index file

Q2. No. of blocks needed to store index file & data file

No. of block accesses to access data assuming index file is in RAM & all block addresses are available.

Now many block accesses if multilevel indexing is used.

$$\lceil \log_2 747 + K \rceil^2$$

$$7 + 2 = 9$$

Blocking factor = 10

No of blocks in which records having respective value of search key are present



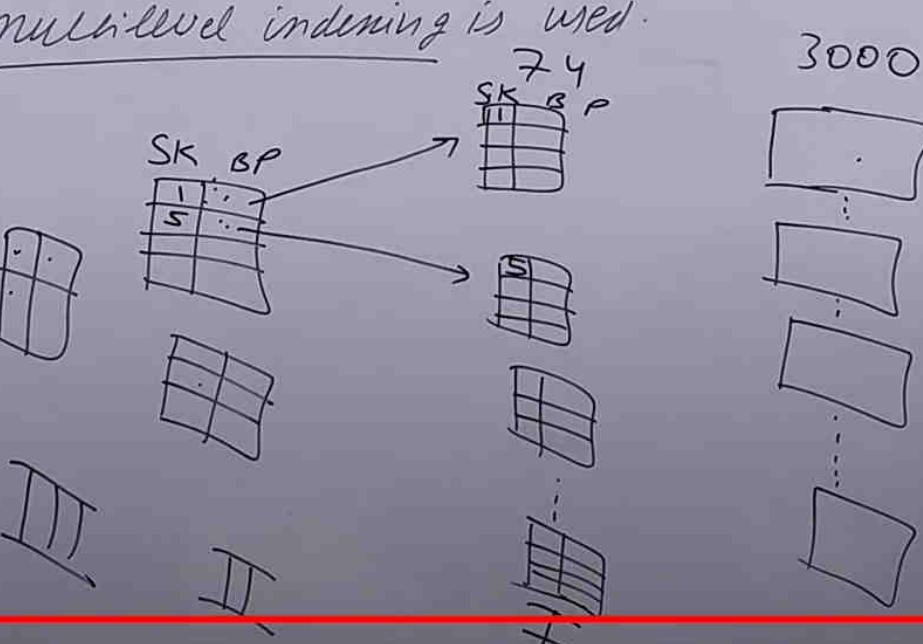
181. numerical on clustered indexing

If clustered indexing is used on a field such that every value is repeated 6 times. Size of block pointer is 6 bytes & clustering field value is of 9 bytes

- Q1 What is blocking factor for data file & index file
- Q2 No. of blocks needed to store index file & data file
- Q3 No. of block accesses to access data assuming index file is in RAM & all block addresses are available.
- Q4 How many block accesses if multilevel indexing is used.

$$\text{BFR}_i \geq 6.8$$

$$\begin{array}{c} \text{III} \quad \text{II} \quad \text{I} \quad \text{D} \\ \downarrow \quad \downarrow \quad \downarrow \quad \downarrow \\ 1+ \quad 1+ \quad 1+ \quad k \\ \downarrow \quad \downarrow \quad \downarrow \\ 1+1+k \end{array}$$
$$\left[\frac{24}{68} \right] = 2$$
$$= 5$$



Secondary Index :-

- Secondary index is an ordered file having 2 fields i.e. Search key & block or record pointer
- Search key is a non ordering field & it may or may not be candidate key.
- Hence there can be many secondary indexes for a file
- Record pointer is a block pointer appended with some offset
- Secondary index on nonordering key field (unique) is always dense i.e. index file has one index entry for each record of data file.
- Index entry contains value of this key field & record or block pointer
- Why dense? ↗

$$BP = 7 \text{ bytes}$$

Block \rightarrow 256 records

$$R.P. = 7 + 1 = 8 \text{ bytes}$$

Secondary Index

SK	BP
1	
2	
3	
5	
6	
8	
9	
10	

5		
9		
3		

10		
6		
8		

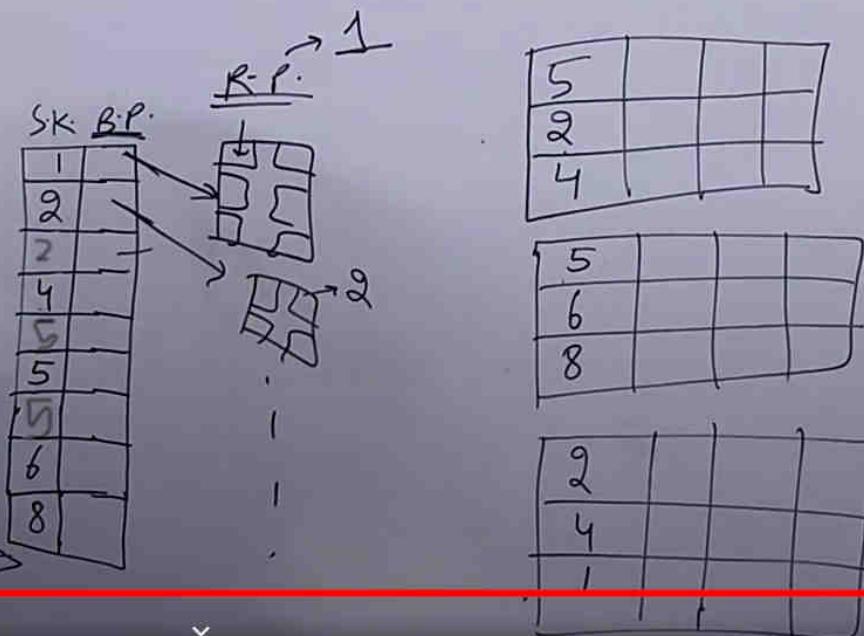
2		
4		
1		

Secondary Index :-

- If secondary index is built on non ordering field which is repeated ie non key then we have many options of creating secondary index.
- But we can't skip any record.

Secondary Index → key ✓
 non key ✗

Secondary Index



183. numerical on secondary indexing

File consists of 30,000 unspinned records, each of size 100 bytes, block size = 1024 bytes
secondary index is build on a field which is unique of size 6 bytes & block pointer = 9 bytes

Q1 What is blocking factor for data file & index file

Q2 No. of blocks needed to store data file & index file

Q3 No. of block accesses to access data with or without indexing assuming all
index & data block addresses are available but these blocks are in H.D.

How many block accesses if multilevel indexing is used.

$$BFR_i = 68$$

$$\frac{BFR_D}{BFR_i} = 10$$

$$\left\lceil \frac{30,000}{10} \right\rceil = 3000$$

$$\left\lceil \frac{30,000}{68} \right\rceil = 442$$

183. numerical on secondary indexing

File consists of 30,000 Unspanned records, each of size 100 bytes, block size = 1024 bytes
secondary index is build on a field which is unique of size 6 bytes & block pointer = 9 bytes

- Q1 What is blocking factor for data file & index file
- Q2 No. of blocks needed to store data file & index file
- Q3 No. of block accesses to access data with or without indexing assuming all index & data block addresses are available but these blocks are in H.D.
- Q4 How many block accesses if multi-level indexing is used.

$$\text{No. of blocks for data file} = 3000$$

$$\text{No. of blocks for index file} = 442$$

$$\begin{aligned} \text{With indexing} &\Rightarrow [022\ 442] + 1 \\ \text{W/o indexing} &\Rightarrow 3000 \Rightarrow 9 + 1 = 10 \\ &\Rightarrow 1500 \quad \text{Worst} \\ &\quad \text{Avg} \end{aligned}$$

183. numerical on secondary indexing

File consists of 50,000 unspanned records, each of size 100 bytes, block size = 1024 bytes
 Secondary index is build on a field which is unique of size 6 bytes & block pointer = 9 bytes

Q1 What is blocking factor for data file & index file

Q2 No. of blocks needed to store data file & index file

Q3 No. of block accesses to access data with or without indexing assuming all index & data block addresses are available but these blocks are in H.D.

How many block accesses if multi-level indexing is used.

$$FR_i = \frac{68}{68}$$

$$FR_D = \frac{3000}{10}$$

$$\left[\frac{442}{68} \right] = 7$$



$7 \times 10 = 70$ bytes

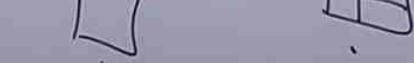
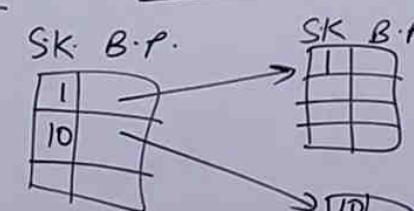
$$\left[\frac{7}{68} \right] = 1$$

III

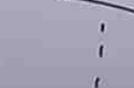
II

I

D = 4



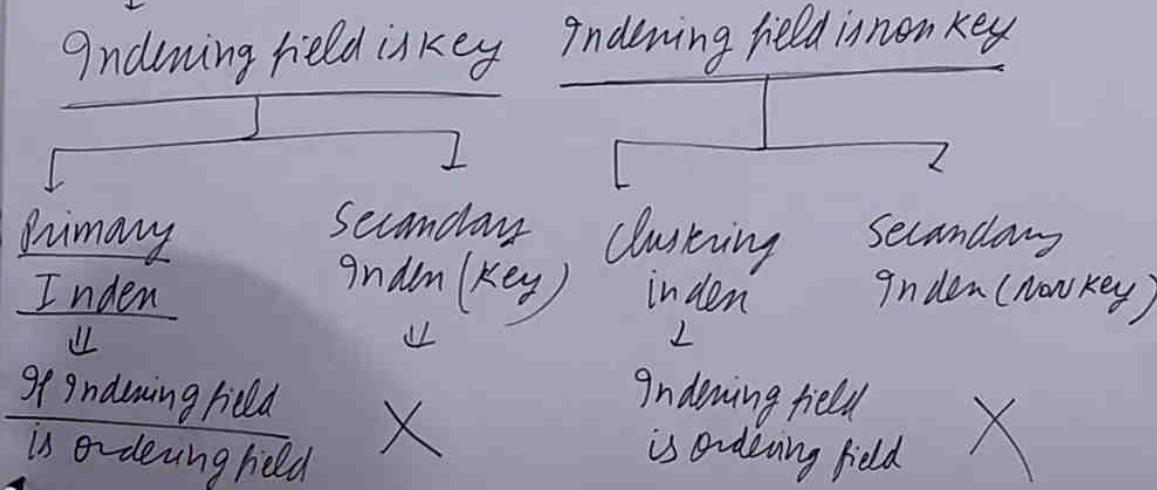
3000



Data file

SummaryType of IndexNo. of (1st level)
Index EntriesDense or SparseBlock anchoring on
Data filePrimaryNo. of blocks in
Data fileSparseYesclusteringNo. of distinct
index field valuesSparseYes/NoSecondary (key)No. of records in
Data fileDenseNoSecondary (Non key)=Dense / SparseNo

Types of indexes



→ If 1st level index file (Primary, Secondary, clustering) is very large then applying binary search on it is time consuming.

→ Hence multilevel indexing is used in such cases to decrease search time from $\log_2 n$ to $\log_{BFR_i} n$

No. of blocks at 1st level

$$\log_{BFR_i} n$$

Fanout

$$\frac{n}{BFR_i}$$

BFR_i

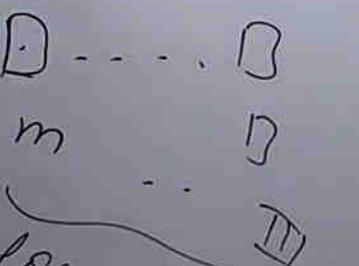
$$\log_{BFR_i} n$$

$$BFR_i > 2$$

BFR_i

$$\frac{n}{BFR_i}$$

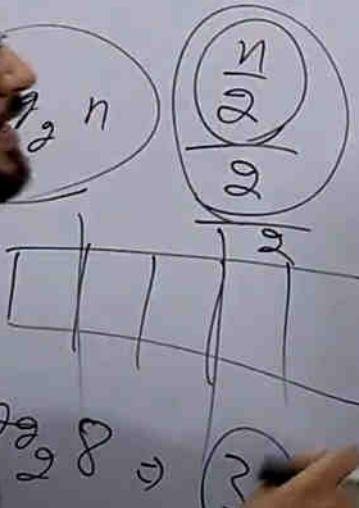
$$\square$$



$$\log_2 n \gg m$$

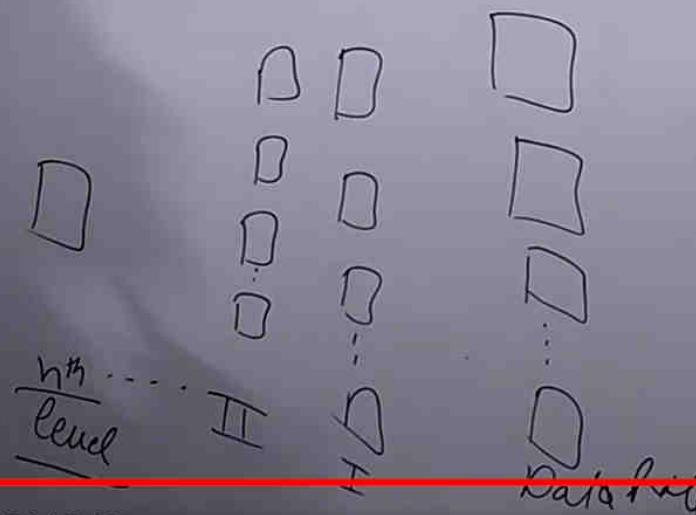
m

I



Static Multilevel Indexing

- Insertion & deletion of records from data file will lead to change in index files at various levels & making these changes is time consuming

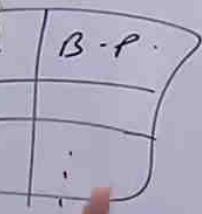
Dynamic Multilevel Indexing

- Insertion & deletion of records from data file may or may not lead to change in index file at various levels & we have efficient algorithms to handle these problems.
- B tree index files
- B⁺ tree index file

188. B-tree introduction



- B tree → Balanced → All leaf nodes are at same level
- Index file is stored in B tree D.S. so that searching, insertion & deletion of index entry becomes efficient.
- Assumption → Field on which index is built must be key i.e. non-repeating (unique)



Flag_i Bi?

used to store
index file



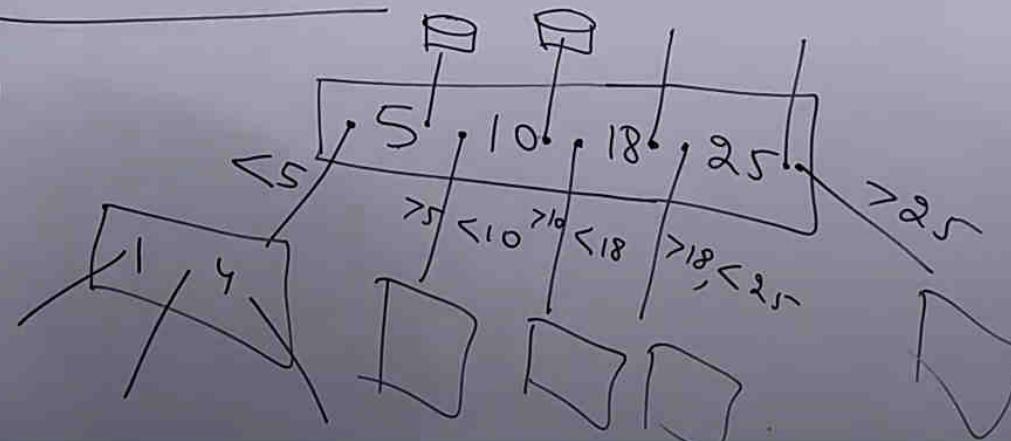
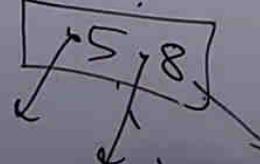
4:29 / 16:04



Properties of B Tree

- ① Root of B Tree can have children b/w $\frac{2}{P}$ & P where P is called order of tree
order means maximum No. of children a node can have
- ② Internal nodes can have children b/w $\lceil \frac{P}{2} \rceil$ & P & Keys b/w $\lceil \frac{P}{2} \rceil - 1$ & $P - 1$
- ③ Structure of internal nodes & leaf nodes are same
 $\langle P_1, \langle K_1, P_{11} \rangle, P_2, \langle K_2, P_{12} \rangle, P_3 \langle K_3, P_{13} \rangle \dots \langle K_{P-1}, P_{1P-1} \rangle, P_P \rangle$
- ④ Leaf nodes can have keys b/w $\lceil \frac{P}{2} \rceil - 1$ to $P - 1$ & 0 children.
- ⑤ Keys must be present in sorted order.

$$P = 5 \Rightarrow \lceil \frac{5}{2} \rceil = 3$$



Q. Find min. & max. No. of keys & children for root, internal, & leaf nodes of B Tree of order a) 5, b) 10

Order = 10

Root
internal
leaf

Children		Keys	
Min	Max	Min	Max
2	10	1	9
5	10	4	9
0	0	4	9

$$\left\lceil \frac{p}{2} \right\rceil$$



189. numerical on B tree

Consider a B Tree in which

Search key = 15 bytes

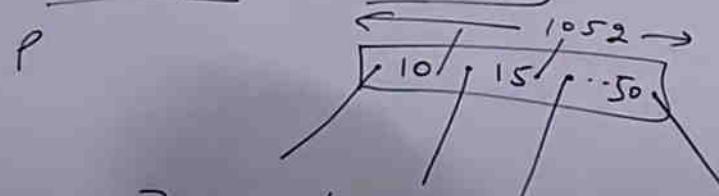
Block size = 1052 bytes

Record pointer = 14 bytes

Block pointer = 7 bytes

Calculate order \Rightarrow max^m No. of children a node can have

Node size = Block size



$$7P + (15+14)(P-1) \leq 1052$$

$$36P \leq 1052 + 29$$

$$36P \leq 1081$$

$$P \leq \frac{1081}{36} \Rightarrow P = 30$$

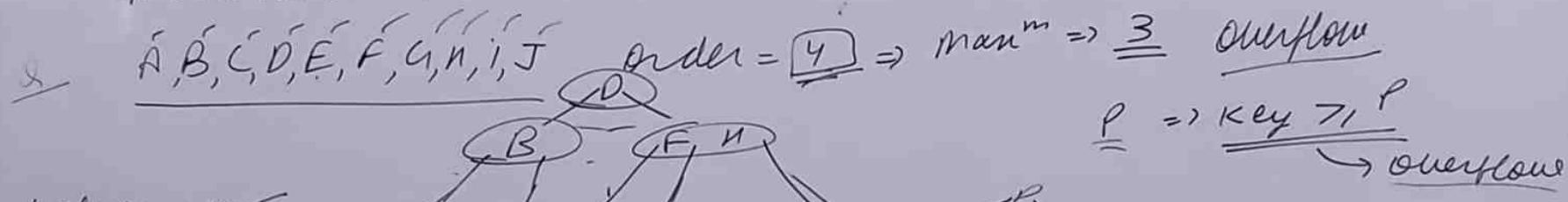
	Min Node	Min Key	Min Pointers
Root	1	$\left[\frac{P}{2}\right] + 1$	2
Level 0	2	14×2	15×2
Level 1	30	14×30	15×30
Level 2	450	14×450	15×450
Level 3			

	Max Nodes	Max Keys	Max Pointers
Level 0	1	$P-1$	$30 \leftarrow P$
Level 1	30	29×30	30×30
Level 2			
Level 3			

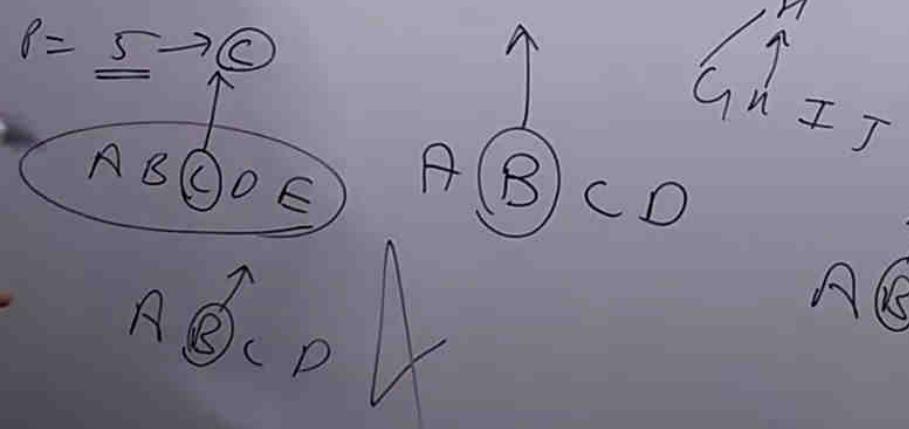
Subtitles/closed captions (c)

B tree insertion

→ Insertion is always done at leaf. Whenever we have overflow then split the node & push median to parent.



Left biased =
Right biased

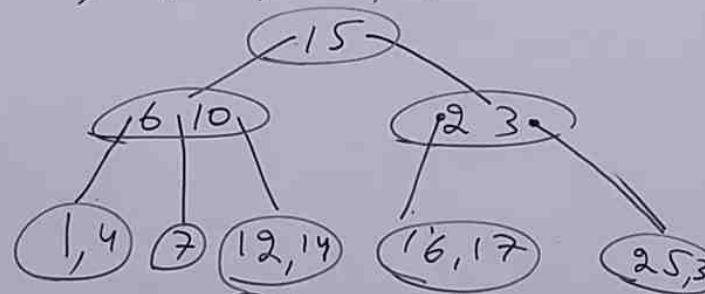


B tree Insertion

→ Insertion is always done at leaf. When ever we have overflow then split the node & push median to parent.

8 $7, 12, 15, 16, 1, 10, 14, 6, 4, 17, 23, 25, 30$ order = 4
Right biased

Underflow X



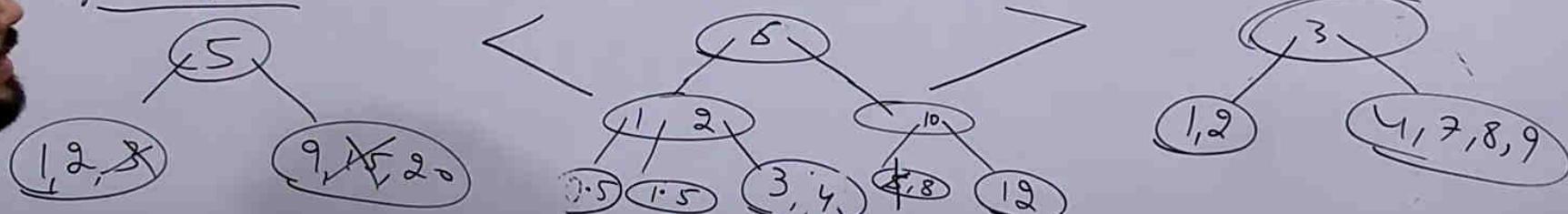
$$\left\lceil \frac{p}{2} \right\rceil = 2$$

P $\xrightarrow{\text{Odd}} 5 \Rightarrow$ $\left\lceil \frac{5}{2} \right\rceil = 3$

$\xrightarrow{\text{Even}} 6 \Rightarrow$ $\left\lceil \frac{6}{2} \right\rceil = 3$ — Key

Deletion in B tree

- Deletion is always done at leaf. If key to be deleted is in internal node then swap it with immediate successor or predecessor.
- If underflow occurs at leaf then take help from left or right siblings.
- If left & right siblings already have minimum No of Keys then merge given leaf with its left or right sibling along with one key from parent.
- If underflow in parent then do same process recursively.



Delete = 3, 15
Order = 4
 $\left(\frac{4}{2}\right) = 2 = 1$
Underflow

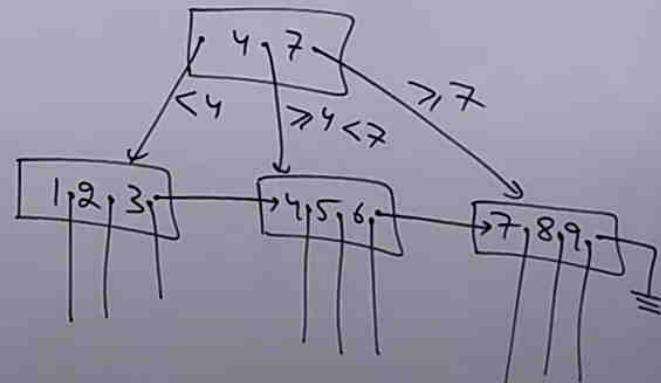
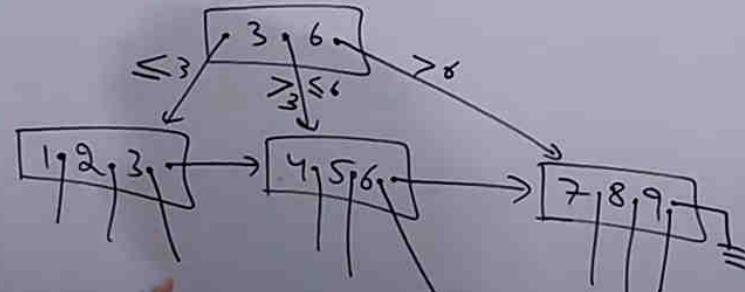
Delete = 5
Order = 4

$\left(\frac{2}{2}\right) = 1$
R $\left(\frac{2}{2}\right) = 2$
L $\left(\frac{2}{2}\right) = 2$

Delete = 6
Order = 5
 $\left(\frac{5}{2}\right) = 3 = 2$

B⁺ tree (Similar to B tree)

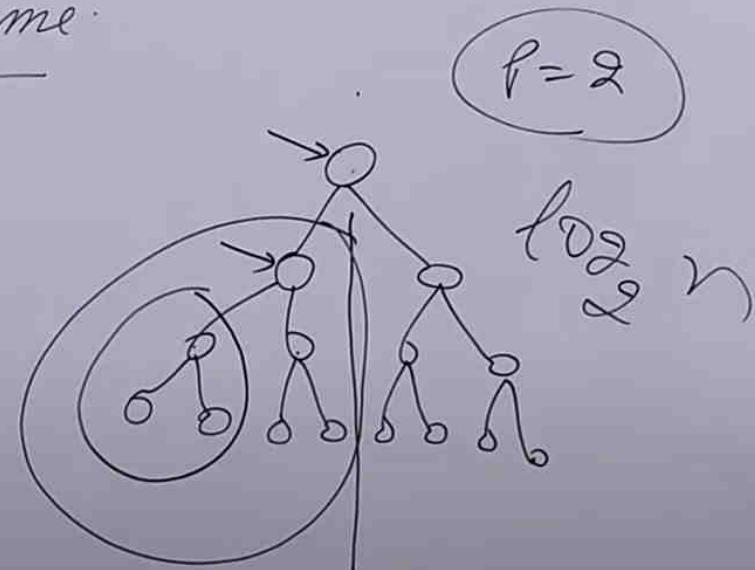
- Here record pointer is missing in internal node. Hence internal nodes can have more no. of children i.e. order ↑ & depth ↓
- Leaves have record pointer, so structure of leaf & internal nodes is different
- In leaf node last pointer is not record pointer but it is block pointer which points to next leaf node & it helps in range queries



B⁺ tree

- Order of leaf nodes & internal nodes are different in B⁺ tree (unlike B tree)
- Searching a key needs $O(\log_p n)$ where p is order. Order of B⁺ tree $>$ Order of B tree
Hence searching is efficient in B⁺ tree. So B⁺ tree is mostly used for multilevel indexing..
- Properties of B tree & B⁺ tree are same.

$$\lceil \frac{p}{2} \rceil \quad \log_{50} n \quad \log_{100} \sqrt{n}$$

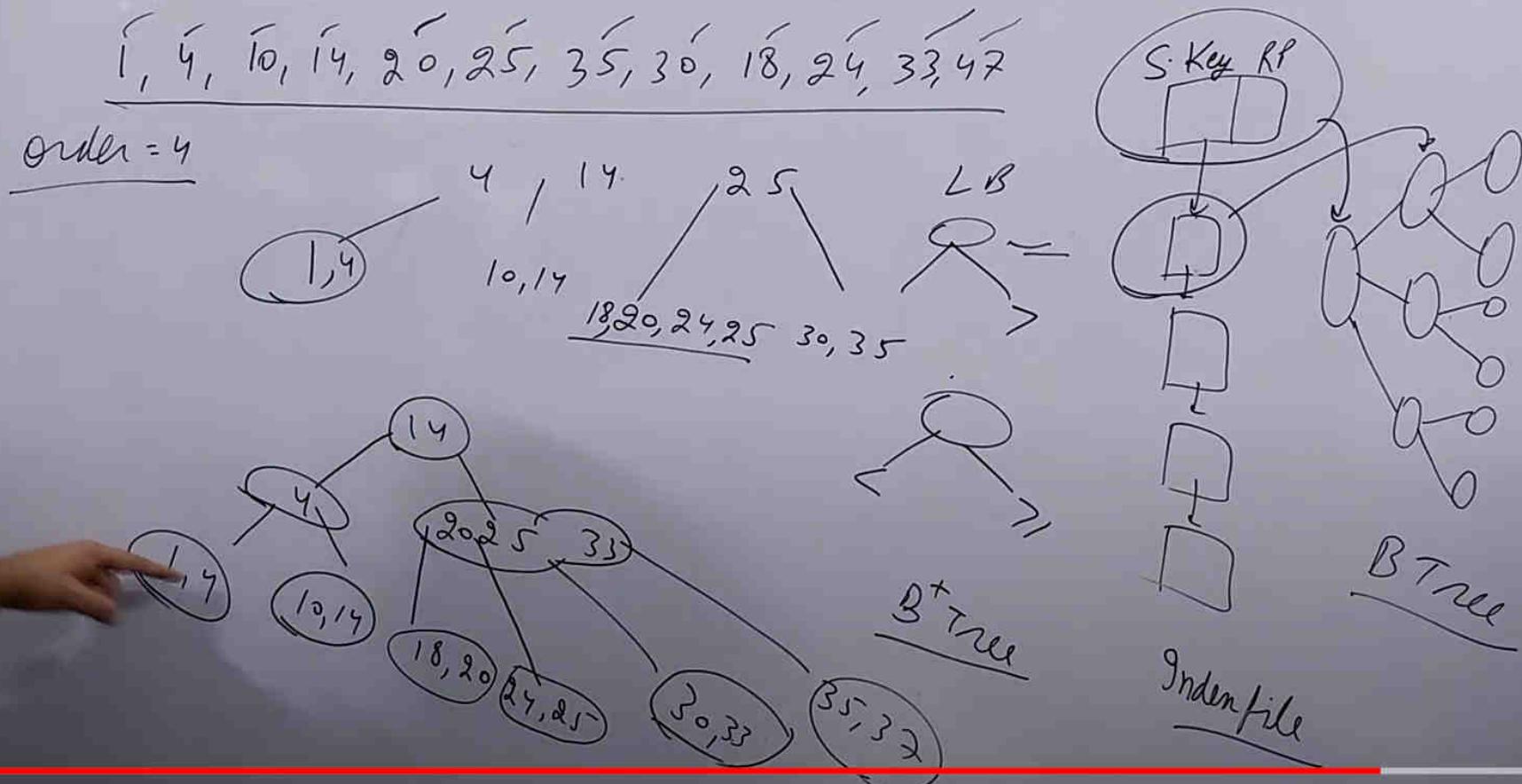


→ Like B tree in B⁺ tree we assume index is build on key field ie. index,
will be unique.



B⁺ Tree insertion

- Insertion done at leaf (like B Tree). But here difference is that in case of overflow in leaf node median go to parent & we keep it in leaf node also
- In case of overflow in non leaf node median go to parent but don't copy it.



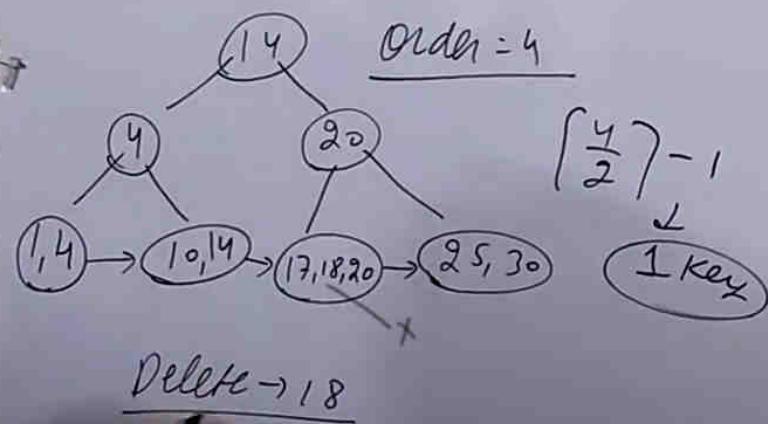
B⁺ Tree Deletion

→ Similar to B tree but key might appear in almost 2 nodes (leaf, internal or root)

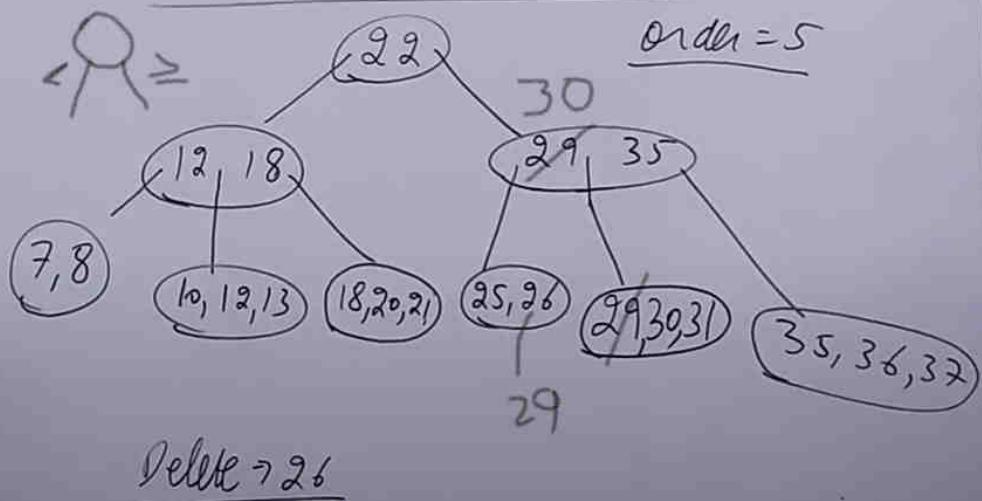
→ Delete entry from root & handle underflow like B tree

→ If entry present in internal node then replace it with immediate successor
or immediate predecessor $\Rightarrow \begin{cases} \text{Successor} \\ \text{Predecessor} \end{cases}$

Case 1 Key only in leaf & No underflow



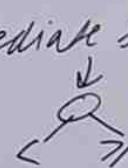
Key only in leaf & underflow



→ Similar to B tree but key might appear in almost 2 nodes (leaf, internal or root)

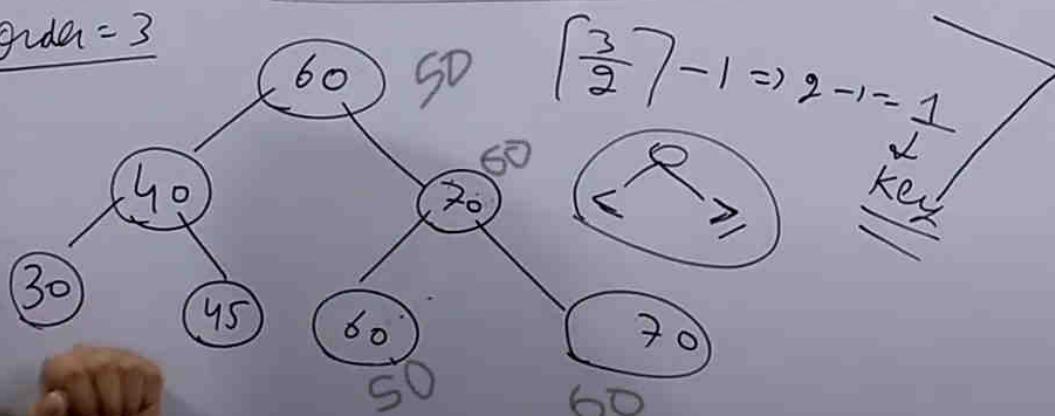
→ Delete entry from root & handle underflow like B tree

→ If entry present in internal node then replace it with immediate successor
or immediate predecessor $\Rightarrow \begin{cases} \leftarrow & \rightarrow \\ \leq & \geq \end{cases}$



Case 3 Delete key present in both leaf & internal node

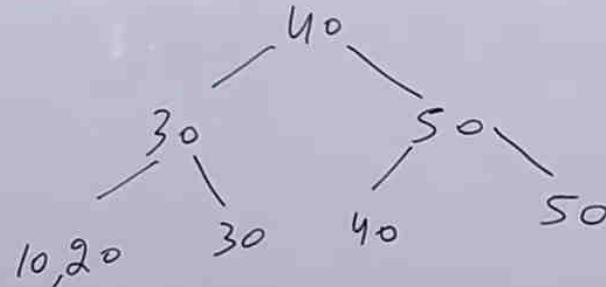
Order = 3



Delete $\rightarrow 50$

1st delete leaf
2nd ... internal

B⁺ Tree → at most 2 Keys + 3 links



Insert 15, 25

Delete 50

Table $T_1 \rightarrow 2000$ records stored in 80 blocks \Rightarrow 1 block 25 Records

Table $T_2 \rightarrow 400$ records stored in 20 blocks \Rightarrow 1 block 20 records

Main memory can store 1 block of T_1 & 1 block of T_2

Nested loop Join algo

Record
64 Records

for ($i = 1; i \leq 2000; i++$)

Block Nested loop Join algo

Block by
Block

for ($j = 1; j \leq 400; j++$)

Join operation

for ($i = 1; i \leq 80; i++$)

for ($j = 1; j \leq 20; j++$)

Join operation

T_1 T_2
1 20

for ($i = 1; i \leq 400; i++$)

for ($j = 1; j \leq 2000; j++$)

Join operation

{ for ($i = 1; i \leq 20; i++$)
 for ($j = 1; j \leq 80; j++$)

Join operation

