

CSE 562 Database Systems

More About Transaction Processing

Some slides are based or modified from originals by
Database Systems: The Complete Book,
 Pearson Prentice Hall, 2nd Edition
 ©2008 Garcia-Molina, Ullman, and Widom

cse@buffalo

More on Transaction Processing

Topics:

- Cascading rollback, recoverable schedule
- Deadlocks
 - Prevention
 - Detection
- View serializability
- Distributed transactions
- Long transactions (nested, compensation)

Concurrency Control & Recovery

Example:

Tj	Ti
⋮	⋮
Wj(A)	⋮
⋮	ri(A)
⋮	Commit Ti
⋮	⋮
Abort Tj	⋮

☛ Non-Persistent Commit (Bad!)

Concurrency Control & Recovery

Example:

Tj	Ti
⋮	⋮
Wj(A)	⋮
⋮	ri(A)
⋮	Commit Ti
⋮	⋮
Abort Tj	⋮

☛ Non-Persistent Commit (Bad!)

avoided by recoverable
schedules

Concurrency Control & Recovery

Example:

T_j	T_i
⋮	⋮
$W_j(A)$	⋮
⋮	$ri(A)$
⋮	$w_i(B)$
⋮	⋮
Abort T_j	⋮
	[Commit T_i]

↪ Cascading rollback (Bad!)

Concurrency Control & Recovery

Example:

T_j	T_i
⋮	⋮
$W_j(A)$	⋮
⋮	$ri(A)$
⋮	$w_i(B)$
⋮	⋮
Abort T_j	⋮
	[Commit T_i]

↪ Cascading rollback (Bad!)

avoided by
avoids-cascading-rollback (ACR)
schedules

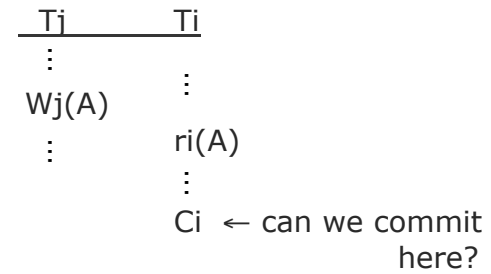
- Schedule is conflict serializable
- $T_j \longrightarrow T_i$
- But not recoverable

- Need to make "final" decision for each transaction:
 - **commit decision** – system guarantees transaction will or has completed, no matter what
 - **abort decision** – system guarantees transaction will or has been rolled back (has no effect)

To model this, two new actions:

- Ci - transaction Ti commits
- Ai - transaction Ti aborts

Back To Example:



Definition

Ti reads from Tj in S ($Tj \Rightarrow_S Ti$) if

- (1) $wj(A) <_S ri(A)$
- (2) $aj \not<_S ri(A)$ ($\not<$: does not precede)
- (3) If $wj(A) <_S wk(A) <_S ri(A)$ then $ak <_S ri(A)$

Definition

Schedule S is recoverable if
whenever $Tj \Rightarrow_S Ti$ and $j \neq i$ and $Ci \in S$
then $Cj <_S Ci$

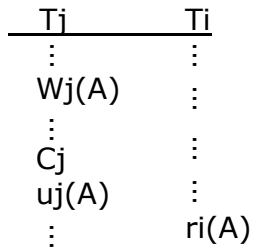
Note: in transactions, reads and writes precede commit or abort

- If $C_i \in T_i$, then $ri(A) < C_i$
 $w_i(A) < C_i$
- If $A_i \in T_i$, then $ri(A) < A_i$
 $w_i(A) < A_i$

- Also, one of C_i, A_i per transaction

How to achieve recoverable schedules?

- With 2PL, hold write locks to commit (strict 2PL)



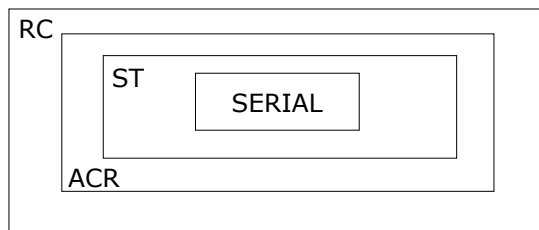
- With validation, no change!

- S is recoverable if each transaction *commits* only after all transactions from which it read have committed.
- S avoids cascading rollback if each transaction may *read* only those values written by committed transactions.

- S is strict if each transaction may *read and write* only items previously written by committed transactions.



Where are serializable schedules?



Examples

- Recoverable:
 - $w_1(A) w_1(B) w_2(A) r_2(B) c_1 c_2$
- Avoids Cascading Rollback:
 - $w_1(A) w_1(B) w_2(A) c_1 r_2(B) c_2$
- Strict:
 - $w_1(A) w_1(B) c_1 w_2(A) r_2(B) c_2$

Assumes $w_2(A)$ is done without reading

Deadlocks

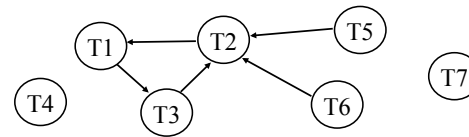
- Detection
 - Wait-For Graph
- Prevention
 - Resource Ordering
 - Timeout
 - Wait-Die
 - Wound-Wait

UB CSE 562 Spring 2009

21

Deadlock Detection

- Build Wait-For Graph
- Use lock table structures
- Build incrementally or periodically
- When cycle found, rollback victim



UB CSE 562 Spring 2009

22

Resource Ordering

- Order all elements A_1, A_2, \dots, A_n
- A transaction T can lock A_i after A_j only if $i > j$

Problem: Ordered lock requests not realistic in most cases

UB CSE 562 Spring 2009

23

Timeout

- If transaction waits more than L sec., roll it back!
- Simple scheme
- Hard to select L

UB CSE 562 Spring 2009

24

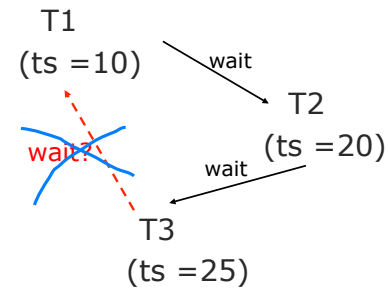
Wait-Die

- Transactions given a timestamp when they arrive ... $ts(T_i)$
- T_i can only wait for T_j if $ts(T_i) < ts(T_j)$...else die

UB CSE 562 Spring 2009

25

Example



UB CSE 562 Spring 2009

26

Starvation with Wait-Die

- When transaction dies, re-try later with what timestamp?
 - original timestamp
 - new timestamp (time of re-submit)

UB CSE 562 Spring 2009

27

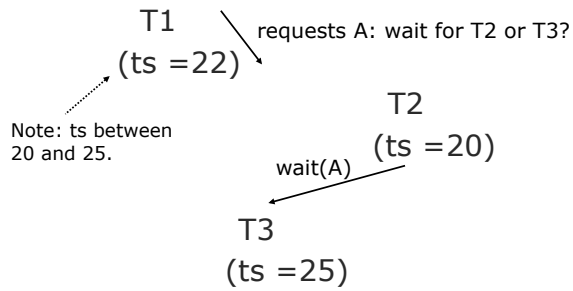
Starvation with Wait-Die

- Resubmit with original timestamp
- Guarantees no starvation
 - Transaction with oldest ts never dies
 - A transaction that dies will eventually have oldest ts and will complete...

UB CSE 562 Spring 2009

28

Second Example

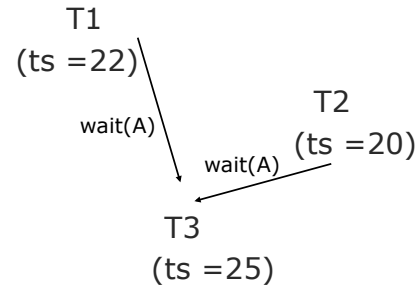


UB CSE 562 Spring 2009

29

Second Example (cont'd)

One option:
T1 waits just for T3, transaction holding lock.
But when T2 gets lock, T1 will have to die!

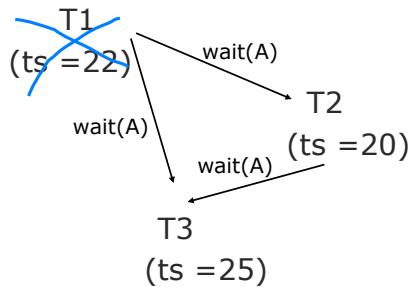


UB CSE 562 Spring 2009

30

Second Example (cont'd)

Another option:
T1 only gets A lock after T2, T3 complete,
so T1 waits for both T2, T3 \Rightarrow T₁ dies right away!

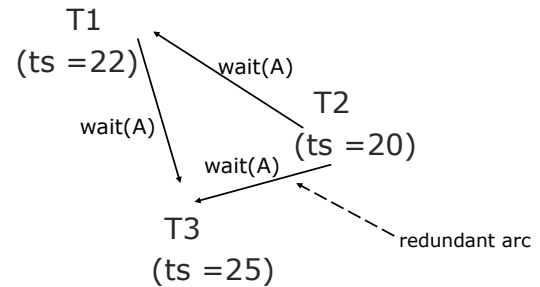


UB CSE 562 Spring 2009

31

Second Example (cont'd)

Yet another option:
T1 preempts T2, so T1 only waits for T3;
T2 then waits for T3 and T1... \Rightarrow T2 may starve?



UB CSE 562 Spring 2009

32

Wound-Wait

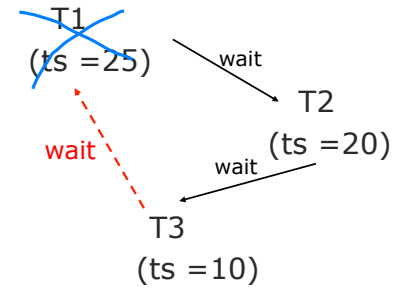
- Transactions given a timestamp when they arrive ... $ts(T_i)$
- T_i wounds T_j if $ts(T_i) < ts(T_j)$
else T_i waits

“Wound”: T_j rolls back and gives lock to T_i

UB CSE 562 Spring 2009

33

Example



UB CSE 562 Spring 2009

34

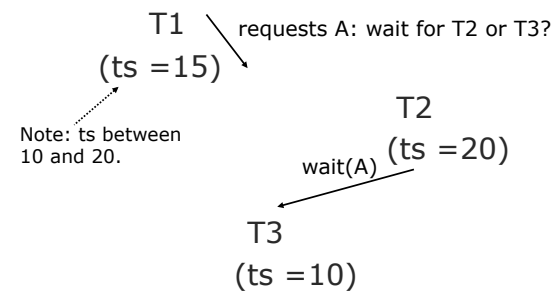
Starvation with Wound-Wait

- When transaction dies, re-try later with what timestamp?
 - original timestamp
 - new timestamp (time of re-submit)

UB CSE 562 Spring 2009

35

Second Example

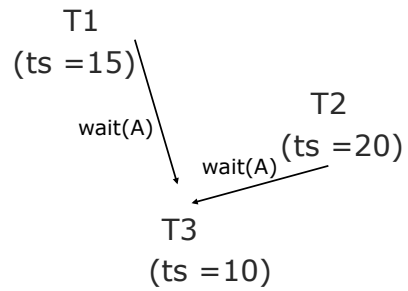


UB CSE 562 Spring 2009

36

Second Example (cont'd)

One option:
T1 waits just for T3, transaction holding lock.
But when T2 gets lock, T1 waits for T2 and wounds T2.

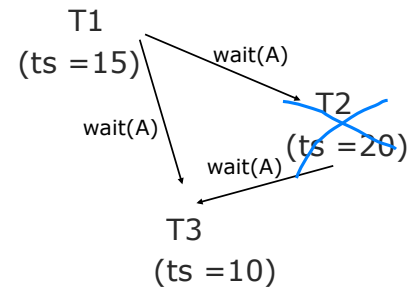


UB CSE 562 Spring 2009

37

Second Example (cont'd)

Another option:
T1 only gets A lock after T2, T3 complete,
so T1 waits for both T2, T3 \Rightarrow T2 wounded right away!

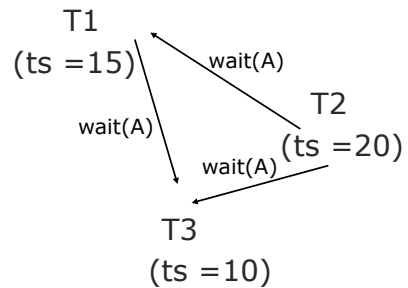


UB CSE 562 Spring 2009

38

Second Example (cont'd)

Yet another option:
T1 preempts T2, so T1 only waits for T3;
T2 then waits for T3 and T1... \Rightarrow T2 is spared!



UB CSE 562 Spring 2009

39

User/Program Commands

Lots of variations, but in general

- Begin_work
- Commit_work
- Abort_work

UB CSE 562 Spring 2009

40

Nested Transactions

User program:

```
⋮  
Begin_work;  
⋮  
⋮  
If results_ok, then commit work  
  else abort_work
```

UB CSE 562 Spring 2009

41

Nested Transactions

User program:

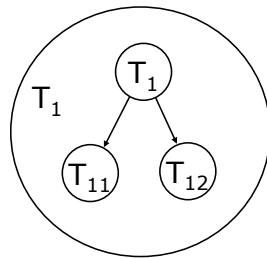
```
⋮  
Begin_work;  
  Begin_work;  
  ⋮  
  If results_ok, then commit work  
    else {abort_work; try something else...}  
  ⋮  
If results_ok, then commit work  
  else abort_work
```

UB CSE 562 Spring 2009

42

Parallel Nested Transactions

```
T1: begin_work  
    ⋮  
    parallel:  
      T11: begin_work  
          ⋮  
          commit_work  
      T12: begin_work  
          ⋮  
          commit_work  
    ⋮  
    commit_work
```



UB CSE 562 Spring 2009

43

Locking

Locking

What are we really locking?



UB CSE 562 Spring 2009

44

Example

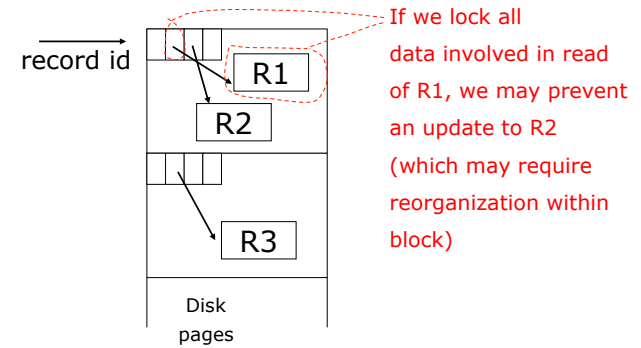
Ti :
 :
 Read record r1
 :
 Read record r1
 :
 Modify record r3
 :
 :
 :

do record locking

UB CSE 562 Spring 2009

45

But Underneath



UB CSE 562 Spring 2009

46

Solution: View DB At Two Levels

Top level: record actions
 record locks
 undo/redo actions — logical

e.g., Insert record(X,Y,Z)
 Redo: insert(X,Y,Z)
 Undo: delete

UB CSE 562 Spring 2009

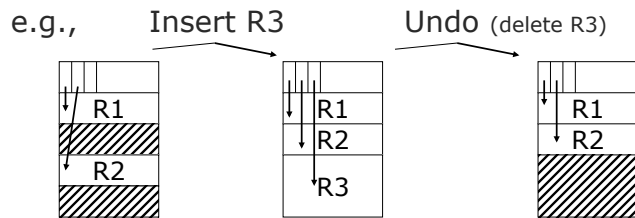
47

Low level: deal with physical details
 latch page during action
 (release at end of action)

UB CSE 562 Spring 2009

48

Note: undo does not return physical DB to original state; only same logical state



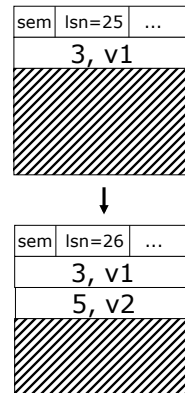
Logging Logical Actions

- Logical action typically span one block (physiological actions)
- Undo/redo log entry specifies undo/redo logical action
- Challenge: making actions idempotent
 - Example (bad): redo insert \Rightarrow key inserted multiple times!

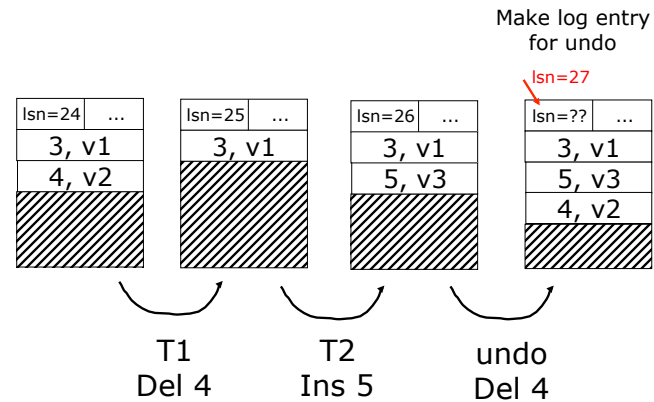
Solution: Add Log Sequence Number

Log record:

- LSN=26
- OP=insert(5,v2) into P
- ...



Still Have a Problem!



Compensation Log Records

- Log record to indicate undo (not redo) action performed
- Note: Compensation may not return page to exactly the initial state

UB CSE 562 Spring 2009

53

At Recovery: Example

Log:

...	lsn=21	...	lsn=27	...	lsn=35	...
	T1		T1		T1	
	a1		a2		a2 ⁻¹	
	p1		p2		p2	

UB CSE 562 Spring 2009

54

What To Do With p2 (During T1 Rollback)?

- If $\text{lsn}(p2) < 27$ then ... ?
- If $27 \leq \text{lsn}(p2) < 35$ then ... ?
- If $\text{lsn}(p2) \geq 35$ then ... ?

Note: $\text{lsn}(p2)$ is lsn of p copy on disk

UB CSE 562 Spring 2009

55

Recovery Strategy

- [1]** Reconstruct state at time of crash
- Find latest valid checkpoint, Ck , and let ac be its set of active transactions
 - Scan log from Ck to end:
 - For each log entry $[lsn, page]$ do:
 - if $\text{lsn}(page) < lsn$ then redo action
 - If log entry is start or commit, update ac

UB CSE 562 Spring 2009

56

Recovery Strategy

[2] Abort uncommitted transactions

- Set ac contains transactions to abort
- Scan log from end to Ck :
 - For each log entry (not undo) of an ac transaction, undo action (making log entry)
- For ac transactions not fully aborted, read their log entries older than Ck and undo their actions

UB CSE 562 Spring 2009

57

Example: What To Do After Crash

Log:

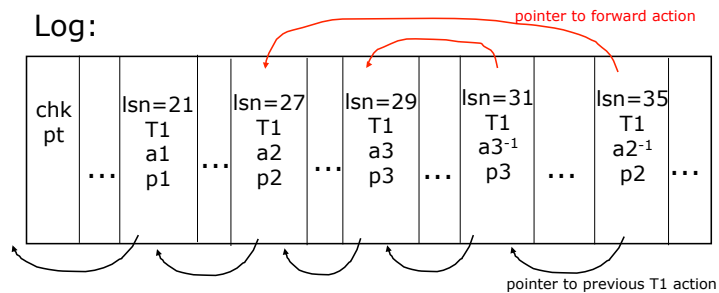
chk		lsn=21		lsn=27		lsn=29		lsn=31		lsn=35	
pt	...	T1	...	T1	...	T1	...	T1	...	T1	...
		a1		a2		a3		a3 ⁻¹		a2 ⁻¹	
		p1		p2		p3		p3		p2	

UB CSE 562 Spring 2009

58

During Undo: Skip Undo's

Log:



UB CSE 562 Spring 2009

59

Related idea: Sagas

- Long running activity: T_1, T_2, \dots, T_n
- Each step/transaction T_i has a compensating transaction T_{i-1}
- Semantic atomicity: execute one of
 - T_1, T_2, \dots, T_n
 - $T_1, T_2, \dots, T_{n-1}, T_{n-1}^{-1}, T_{n-2}^{-1}, \dots, T_1^{-1}$
 - $T_1, T_2, \dots, T_{n-2}, T_{n-2}^{-1}, T_{n-3}^{-1}, \dots, T_1^{-1}$
 - ⋮
 - T_1, T_1^{-1}
 - nothing

UB CSE 562 Spring 2009

60

Summary

- Cascading rollback
Recoverable schedule
- Deadlock
 - Prevention
 - Detection
- Nested transactions
- Multi-level view