

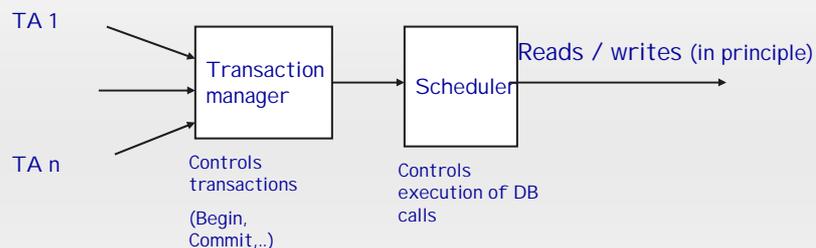
Transactions and Concurrency control

Transactions
Concepts
Serializability
Concurrency control
Locking
Multiversion cc
Optimistic cc

Concurrency control

...and serializability

- ▶ **Wanted:** effective real-time scheduling of operations with guaranteed serializability of the resulting execution sequence.



- ▶ **Concurrency control in DBS:** methods which schedule the operations of different TAs in a way which guarantees serializability of all transactions ("between system start and shutdown")

- ▶ Three principle concurrency control methods
 - ▶ Locking (most important)
 - ▶ Optimistic concurrency control
 - ▶ Time stamps
- ▶ No explicit locking in application programs,
 - error prone,
 - responsibility of scheduler (and lock manager)

In most DBS also explicit locking allowed in addition to implicit locking by scheduler

- ▶ Not considered here: transaction independent locking, e.g. writing a page p to disk requires a short term lock on p

hs / FUB dbs03-22-TAConCtrl-2-3

- ▶ Locking is **pessimistic**:
 - ▶ Assumption: during operation $op[x]$ of TA1 a (potentially) conflicting operation $op'[x]$ of TA2 will access the same object x
 - ▶ This has to be avoided by locking x before accessing x

S: $r1[y], r3[u], r2[y], w1[y], w2[x], w1[x], w2[z], c2, w3[x]$

↑
TA2 read x , wait until $c2$ has committed

- ▶ An **optimistic** strategy would be:
 - ▶ Perform all operations on a copy of the data. Check at the end - before commit - if there were any conflicts.
 - ▶ If no: commit, else abort (rollback) - more or less

hs / FUB dbs03-22-TAConCtrl-2-4

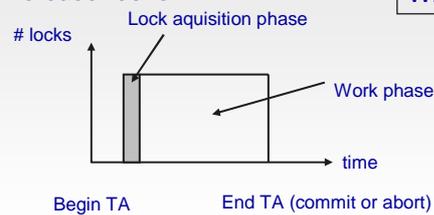
▶ Lock protocols

▶ Simplistic: Lock each object before writing, unlock afterwards ⇒ schedule will not be serializable (why?)

▶ Preclaiming

- Acquire all locks you need before performing an operation,
- release, if you do not get all of them and try again race condition!
- Execute transaction
- Release locks

Preclaiming serializable?
Why (not) ?



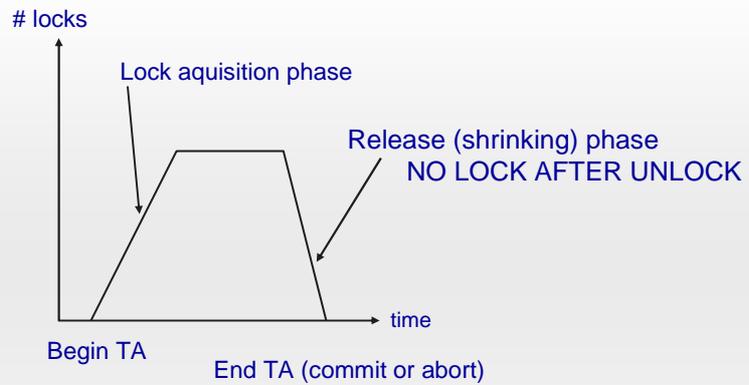
Bad: objects to be processed may not be known in advance.
Not used in DBS.

▶ Two phase locking (2PL)

Important method

1. Each object referenced by TAI has to be locked before usage
2. Existing locks of other TA's will to be respected
3. No lock is requested by a TAI , if a lock has been released by the same transaction TAI (no lock after unlock)
4. Locks are released at least at commit time
5. A requests of a lock by a TA which it already holds, has no effect.

... but a read lock may be "upgraded" to a write lock



- ▶ Locked objects may be read / written already in lock aquisition phase

hs / FUB dbs03-22-TAConCtrl-2-7

- ▶ Why no lock after unlock?

- ▶ Example:

```
lock1(x), r1[x], x:=x*10, unlock1(x)
lock2(x), r1[x], x:=x+1, w2[x] unlock2(x)
lock1(x), w1(x), unlock1(x)
```

results in a lost update



⇒ Rule 3 is essential

hs / FUB dbs03-22-TAConCtrl-2-8

▶ 2-Phase locking theorem

If all transactions follow the 2-phase locking protocol, the resulting schedule is serializable

▶ Proof sketch:

- ▶ Suppose a resulting schedule is not serializable. when using 2PL \Rightarrow conflict graph contains a cycle \Rightarrow there are transactions TA1 and TA2 with conflict pairs (p, q) and (q', p') ,
 p, p' atomic operations of TA1, q, q' of TA2, p, q access the same object x , and q', p' an object y (assuming a cycle of length 1, induction for the general case)



▶ Proof (cont.)

- ▶ Let e.g. $(p, q) = (r1[x], w2[x]),$
 $(q', p') = (w2[y], w1[y])$

Analyze all of the possible sequences of Execution:

p, q, q', p

p, q', q, p'

q', p, q, p'

q', p, p', q

q', p', p, q

T2: Lock y, T1: Lock x, T2: Lock x, T1: Lock y

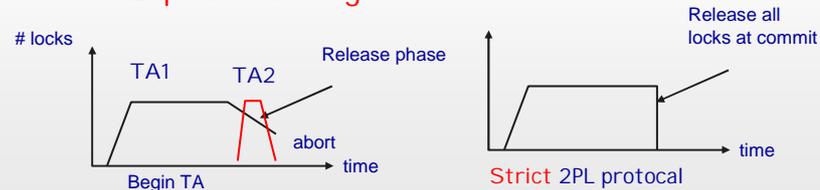
T1 must have released lock on x and aquired one on y (or T2 must have aquired after release)

Hurts 2-phase rule! Contradiction to assumption that all follow 2PL protocol

Same holds for the other possible sequences \Rightarrow Theorem

- ▶ Note: serializability does not imply 2PL, i.e. there are serializable schedules which cannot result from a scheduler employing 2PL

▶ **Strict 2-phase locking**



- ▶ Another transaction TA2 could have used an object x which was unlocked by TA1 in the release phase
 - not a problem, if TA1 commits
 - if TA1 aborts, TA2 has used a wrong state of x
TA2 has to be aborted by the system
- ▶ May happen recursively: cascading abort, **bad ...**
- ▶ **Strict 2PL: Release all locks at commit point.**

hs / FUB dbs03-22-TAConCtrl-2-11

▶ **Lock conflict**

- ▶ Two or more processes request an exclusive lock for the same object

▶ **Deadlock**

- ▶ Locking: threat of deadlock
 - No preemption
 - No lock release in case of lock conflicts
- ▶ Two-Phase locking may cause deadlocks

$L_i[x]$ = Transaction i requests lock on x

$U_i[x]$ = Transaction i releases lock on x

Lock sequence: $L_1[x], L_2[y], \dots, L_1[y], L_2[x]$ causes deadlock

hs / FUB dbs03-22-TAConCtrl-2-12

Concurrency control

Secondary goals

- ▶ Primary goal
 - ▶ no harmful effects (lost update, ...)
- ▶ Secondary goal
 - ▶ Degree of parallelism should be as high as possible, even when locking is used
 - ▶ Low deadlock probability, if any
- ▶ Ways to increase parallelism
 - ▶ Compatible locks (read versus write semantics)
 - ▶ Different lock granularity
 - ▶ Application semantics
 - ▶ No locks, optimistic cc

hs / FUB dbs03-22-TAConCtrl-2-13

Concurrency control

Lock modes

- ▶ Lock modes and lock compatibility
 - ▶ RX - model: read (R) and write(X) locks

	holder		
	R	X	
requester	R	+	-
	X	-	-

Lock compatibility matrix

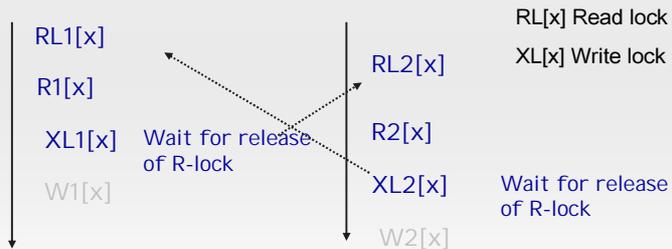
- ▶ Lock compatibility in the RX model:
 - ▶ Objects locked in R-mode may be locked in R-mode by other transactions(+)
 - ▶ Objects locked in eXclusively mode (X-mode) may not be locked by any other transaction in any mode.
- Requesting TA must wait.

hs / FUB dbs03-22-TAConCtrl-2-14

Concurrency control

Reduce deadlock threat

- ▶ Deadlocks caused by typical read / write sequences
- ▶ TA1: read account_record x; incr(x.balance); write account_record
- ▶ TA2: read account_record x; incr(x.balance); write account_record



- ▶ Gives rise to Read-Update-Exclusive Model (RUX)

hs / FUB dbs03-22-TAConCtrl-2-15

Concurrency control

Lock modes

- ▶ **RUX Lock protocol**
 - ▶ Transactions which read and subsequently update an object y request a U-lock
 - ▶ U-locks are incompatible with U-locks \Rightarrow deadlock-thread avoided
 - ▶ U / R-lock compatibility asymmetric, why?

		holder		
		R	U	X
requester	R	+	-	-
	U	+	-	-
	X	-	-	-

How does DBS know, that update is intended?

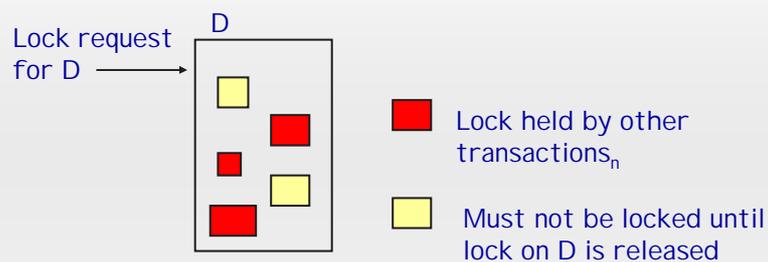
hs / FUB dbs03-22-TAConCtrl-2-16

▶ Hierarchical locking

- ▶ One single lock granularity (e.g. records) insufficient overhead of many record locks large compared to a table lock
- ▶ Most DBS have at least two lock granularities
- ▶ Problem: TA_i wants to lock table R
 - some rows of R are locked by different transactions
 - Different lock conflict than before: TA_i is waiting for the release of all record locks
 - No other TA should be able to lock a record, otherwise TA_i could starve

hs / FUB dbs03-22-TAConCtrl-2-17

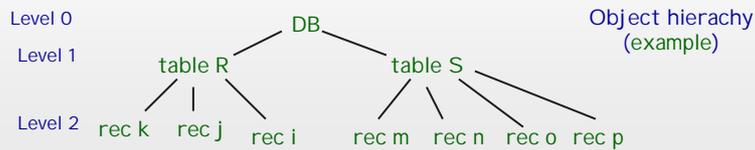
▶ Locks of different granularity



Efficient implementation of this type of situation??

hs / FUB dbs03-22-TAConCtrl-2-18

▶ Intention locks

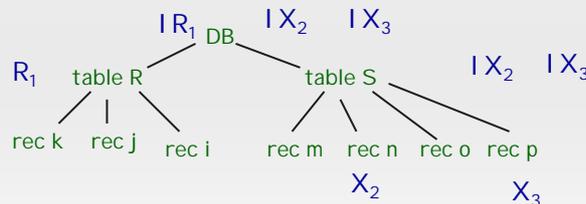


▶ Important feature of hierarchical locking: **intention locks**: for each lock mode, there is an intention lock, e.g. for RX-lock modes: IR and IX

▶ Semantics:
A TA holds a IM-lock on an object D on level i, if and only if it holds an M-lock on an object D' on level j > i subordinate to D

▶ Hierarchical locks :

- ▶ An object O on level i contains all objects x on level i+1
- ▶ Locks of O lock all subordinate objects x
- ▶ If a subordinate object x (level i+1) is locked, this is indicated by an intention lock on level i



▶ Lock escalation

If too many objects x on level i+1 are locked by a transaction, it may be converted into one lock on level i

► Hierarchical locking (cont)

- Advantage: **one lookup** is sufficient to check **if a lock on higher level (say on a table) can be granted**
- **Protocol:** if a TA wants to lock an object on level *i* in mode *<M>* (X or R), lock all objects on higher level (on the path to root) in I<M> - mode
- Easy to check, if the locks on all subordinate objects are released: implement I<M>-lock as a counter

		holder				
		IR	IX	R	X	
requester	IR	+	+	+	-	Compatibility matrix How to combine with U-lock mode?
	IX	+	+	-	-	
	R	+	-	+	-	
	X	-	-	-	-	

hs / FUB dbs03-22-TAConCtrl-2-21

► Deadlocks

... can happen with 2PL protocol (see above)

- Release of a lock could break rule 4
 XL1[x] , XL2[y], XL1[y] -> TA1: WAIT for XU2[y] , XL2[x] -> TA2: WAIT for XU1[x]

- Note: deadlock different from lock conflicts:
 XL1[x] , XL2[y], XL1[y] -> TA1: WAIT for XU2[y] XL2[z], w2[y], w2[z], XU2[y],...



Not schedules, but sequences of operations including lock / unlock, e.g. in a schedule

hs / FUB dbs03-22-TAConCtrl-2-22

▶ Resolving deadlocks**▶ Cycle check in Wait-for-graph**

- Waiting of TA1 for release of lock on x by TA2 is indicated by an arc from TA1 to TA2 labeled "x"
- Cycles indicate deadlock
- In a distributed environment, deadlocks may involve different systems. How to detect cycles?
- One of the waiting transaction ("victim") is rolled back
- Which one??

▶ Timeout

- If TA has been waiting longer than the time limit, it is aborted.
- Efficient but may roll back innocent victims (deadlock does not exist)

Oracle: WF-graph in central systems, timeout in distributed

hs / FUB dbs03-22-TAConCtrl-2-23

▶ Avoiding deadlocks

- ▶ Deadlocks only occur, if no lock holder is preempted, i.e. loses his lock
- ▶ Preemption may be forced by the lock manager
- ▶ If TA t is preempted, it is forced to rollback
- ▶ Preemption \Rightarrow no deadlocks, but living transactions may be killed

▶ Wait/Die - Wound/Wait : Basic idea

- ▶ Solve lock conflicts by rollback of one of the conflicting transactions....
- ▶ but not always
- ▶ Rollback dependent on the relative age of the transactions
- ▶ Time stamp for each transaction

hs / FUB dbs03-22-TAConCtrl-2-24

Concurrency control

Deadlock avoidance

▶ Wound/Wait – Wait / Die methods

- ▶ Each transaction has an initial timestamp $TS(i)$
- ▶ If TA_2 requests a lock on x and there is a lock conflict with TA_1 , then

▶ WOUND / WAIT

if $ts(TA_1) < ts(TA_2)$ then $TA_2.WAIT$ else $TA_1.ABORT$



If $TS(TA_1) < TS(TA_2)$ then $TA_2.wait$ else $TA_1.abort$

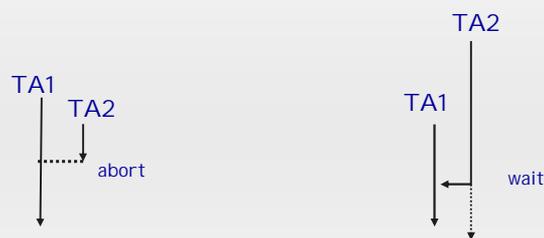
hs / FUB dbs03-22-TAConCtrl-2-25

Concurrency control

Deadlock avoidance

▶ WAIT / DIE

if $ts(TA_1) < ts(TA_2)$ then $TA_2.ABORT$ else $TA_2.WAIT$



If $TS(TA_1) < TS(TA_2)$ then $TA_2.abort$ else $TA_2.wait$

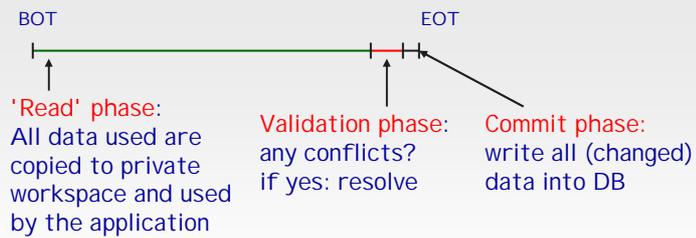
No deadlocks! Why?

Aborted transaction restarts with old timestamp in order to avoid starvation

hs / FUB dbs03-22-TAConCtrl-2-26

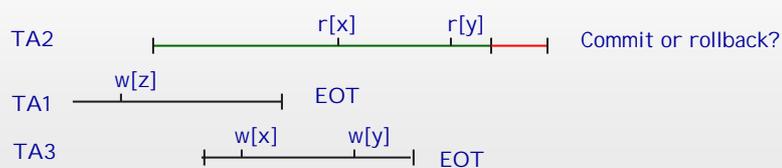
▶ Optimistic concurrency control

- ▶ Locks are expensive
- ▶ If conflicts are rare, retrospective check for conflicts and – if any – abort are cheaper, hopefully
- ▶ Basis idea: let all transactions work on copies, at the end, check for conflicts, commit (or abort)



hs / FUB dbs03-22-TAConCtrl-2-27

▶ Backward oriented concurrency control (BOCC)



- ▶ ReadSet $R(T)$ = set of data, transaction T read in read phase
- ▶ WriteSet $W(T)$ = set of data, T has changed during read phase
Assumption: $W(T) \subseteq R(T)$
- Example above: $x, y \in R(T2)$, $x, y \in W(T3)$, $z \in W(T1)$

hs / FUB dbs03-22-TAConCtrl-2-28

Concurrency control

Optimistic CC

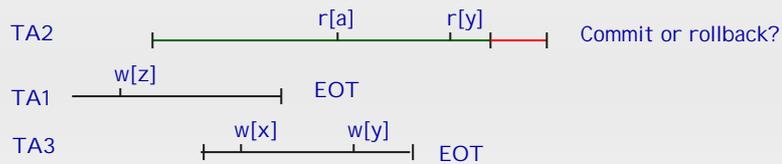
▶ BOCC_validate(T) :

if for all transactions (T') which committed after BOT(T) :

$$R(T) \cap W(T') = \emptyset$$

then T.commit // successful validation

else T.abort



Note :

Conflict caused by r[y] is NOT harmful,
but rollback of TA2 because $R(TA2) \cap W(TA3) \neq \emptyset$

=> Unnecessary Aborts

hs / FUB dbs03-22-TAConCtrl-2-29

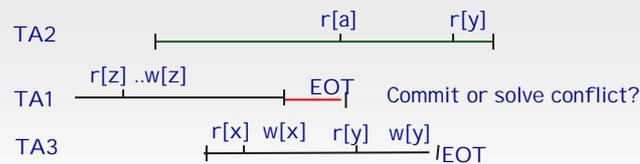
Concurrency control

Optimistic CC

▶ Forward oriented optimistic Concurrency control (FOCC)

▶ Forward looking validation phase:

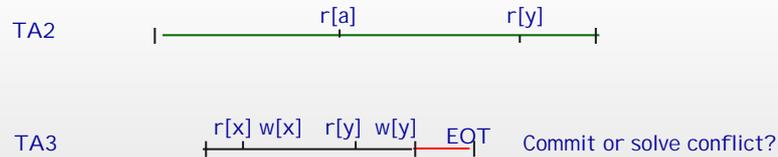
if there is a running transaction T' which read data written by the validating transaction T then solve the conflict (e.g. kill T'), else commit



hs / FUB dbs03-22-TAConCtrl-2-30

Concurrency control

Optimistic CC



FOCC_validate(T) : if for all running transactions (T')

$$R^n(T') \cap W(T) = \emptyset$$

then T.commit // successful validation
 else solve_conflict (T, T')

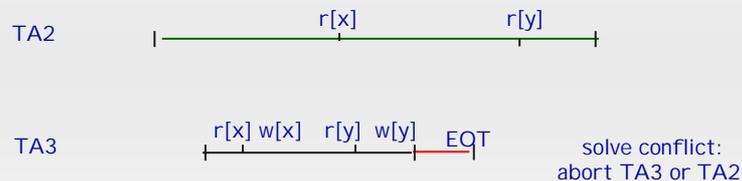
$R^n(T')$: Read set of T' at validation time of T (current read set)

hs / FUB dbs03-22-TAConCtrl-2-31

Concurrency control

Optimistic CC

- ▶ Validation of read only transactions T:
FOCC guarantees successful validation !
- ▶ FOCC has greater flexibility
Validating TA may decide on victims!



- ▶ **Issues** for both approaches:
fast validation – only one TA can validate at a time.
Fast and atomic commit processing,
- ▶ Useful in situation with few expected conflicts

hs / FUB dbs03-22-TAConCtrl-2-32

- ▶ Combining locks with optimism

- ▶ Example: high traffic reservation system
- ▶ typical TA: `check "seats_avail >0 ?" //seats_avail is hot spot`
`if yes, do this and that;`
`write seats_avail-1`
- ▶ `seats_avail` is a hot spot object
- ▶ Not the state per se is important but the predicate `"seats_avail >0 ?"`
- ▶ Optimism: if pred is true at BOT then it will be true with high probability at EOT
- ▶ But if not: abort

hs / FUB dbs03-22-TAConCtrl-2-33

- ▶ Additional operations Verify and Modify:

Verify P : check predicate P ("`seats_avail >0`")
 //like read phase

put "`seats_avail-1`" into to_do list

rest of TA

EOT:

Modify : for all operations on to_do list

```
{ lock; verify once more;
  if 'false' rollback else write updates;}
unlock all;
```

Roughly the approach taken in IMS / Fast Path

- ▶ Short locks, more parallelism
- ▶ If only decrement / increment operations: concurrent writing possible without producing inconsistencies

hs / FUB dbs03-22-TAConCtrl-2-34

Concurrency control

Multiversion concurrency

Arrows from
TA2-ops to
conflicting TA1-ops

▶ Multiversion CC:

r1[x] w1[x] r2[x] w2[y] r1[y] w1[z] c1 w2[a] c2
not serializable.

If r1[y] had arrived at the scheduler before w2[y] the schedule would have been serializable.

- ▶ Main idea of **multiversion concurrency control** : Reads should see a consistent (and committed) state, which might be older than the current object state.
- ▶ Necessary: **Different version** of an object
- ▶ Read and write locks compatible (!)
- ▶ Example: TA2 must not write its version of y before TA1 has release lock on y

hs / FUB dbs03-22-TAConCtrl-2-35

Concurrency control

Multiversion concurrency

▶ Lock based MVCC

- ▶ Read locks always granted, write lock if object not write locked
=>
two versions: consistent one and writable private copy
- ▶ When TA wants to write modified copy of x into DB it has to wait until all reader of x have released read lock
- ▶ write is delayed to ensure consistent read using a certify lock

	R	W	C
R	+	+	-
W	+	-	-
C	-	-	-

C = Certify

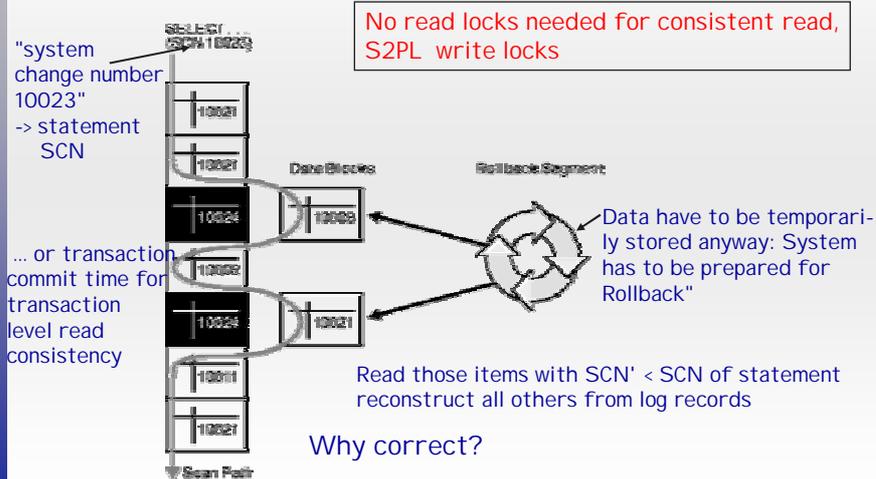
Read locks needed ??

hs / FUB dbs03-22-TAConCtrl-2-36

Concurrency control

Multiversion concurrency

- ▶ Read Only Multiple version CC (used in Oracle)

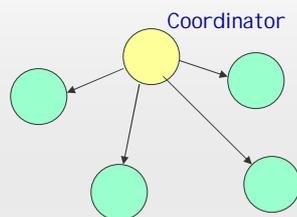


hs / FUB dbs03-22-TAConCtrl-2-37

Distributed Transactions

Intermezzo

- ▶ Configuration

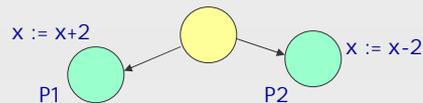


- ▶ Different resource managers involved in the transaction e.g: database systems, mail server, file system, message queues,...
- ▶ Asynchronous and independent
- ▶ One transaction coordinator can be a resource manager or not
- ▶ One or more participants
- ▶ Examples: Transfer of money/shares / ... from Bank A to B
ECommerce systems
All kinds of processing in decentralized organizations
- ▶ Frequently used in multi-tier architectures: middle tier accesses different databases
Just open two or more JDBC-connections

hs / FUB dbs03-22-TAConCtrl-2-38

► Problems

- No problem ... if all systems work reliable
- Deadlock? Difficult to detect: use optimistic locking
- Obvious inconsistencies, if one participant commits, another crashes:



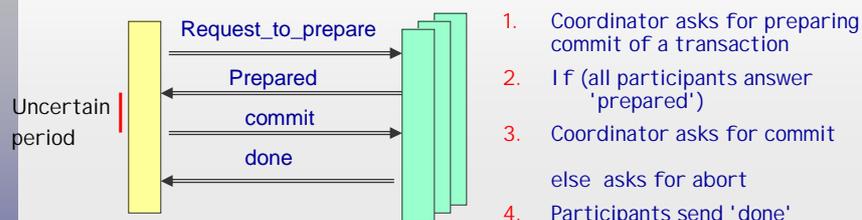
Coordinator: COMMIT
 P1 commits
 P2 crashes, undo??
 Introduces global inconsistency

► Assumptions

- Each resource manager has a transactional recovery system (log operations, commit, rollback)
- There is exactly one commit coordinator, which issues commit for a transaction exactly once
- A transaction has stopped processing at each site before commit is issued

hs / FUB dbs03-22-TAConCtrl-2-39

► The Two Phase Commit protocol (2PC)



1. Coordinator asks for preparing commit of a transaction
2. If (all participants answer 'prepared')
3. Coordinator asks for commit else asks for abort
4. Participants send 'done'

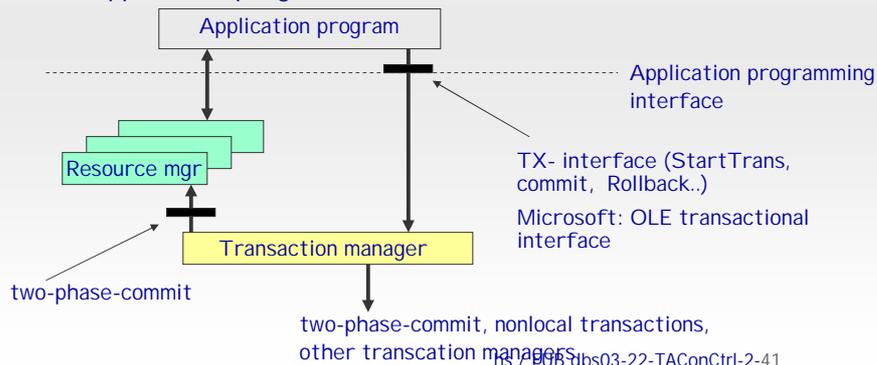
- After prepare phase: participants are ready to commit or to abort; they still hold locks
- If one of the participants does not reply or is not able to commit for some reason, the global transaction has to be aborted.
- Problem: if coordinator is unavailable after the prepare phase, resources may be locked for a long time

hs / FUB dbs03-22-TAConCtrl-2-40

Distributed Transactions

▶ Transaction managers: the X/Open transaction model

- ▶ Independent systems which coordinate transactions involving multiple resource managers as a service for application programs



Summary: Transactions and concurrency

- ▶ Transactions: Very important concept
- ▶ Model for consistent, isolated execution of TAs
- ▶ Scheduler has to decide on interleaving of operations
- ▶ Serializability: correctness criterion
- ▶ Implementation of serializability:
 - ▶ 2-phase-locking
 - ▶ hierarchical locks
 - ▶ variation in order to avoid deadlocks :
wound / wait, wait / die
 - ▶ optimistic concurrency control
 - ▶ multiversion cc