

## Concurrency control in Homogeneous Distributed Databases (2)

- Timestamp ordering
- Basic implementation
- Optimistic CC in distributed DB
- Distributed deadlock detection

based on slides by Weikum / Vossen: Transactional Information Systems; H. Garcia Molina

## Non-locking concurrency control

### Time stamp ordering

Basic idea:

- assign timestamp when transaction starts
- if  $ts(t_1) < ts(t_2) \dots < ts(t_n)$ , then scheduler has to produce history equivalent to  $t_1, t_2, t_3, t_4, \dots t_n$

Timestamp ordering rule:

If  $p_i[x]$  and  $q_j[x]$  are conflicting operations, then  $p_i[x]$  is executed before  $q_j[x]$  ( $p_i[x] < q_j[x]$ ) iff  $ts(t_i) < ts(t_j)$

Issue: how to find out that x has been modified by a younger / older TA ?? Timestamp for each data item !

hs / FUB dbsII-03-17DDBCC1-2

### Timestamp ordering

- TO concurrency control guarantees conflict-serializable schedules:

If not: cycle in conflict graph

cycle of length 2:  $ts(t_1) < ts(t_2) \wedge ts(t_2) < ts(t_1) \quad \#$

induction over length of cycle  $\Rightarrow \#$

$\Rightarrow$  No cycle in conflict graph  $\checkmark$

hs / FUB dbsII-03-17DDBCC1-3

### Example: Distributed case

(Node S1)	(Node S2)
(t1) $a \leftarrow X(S1)$	(t2) $d \leftarrow Y(S2)$
(t1) $X \leftarrow a+100$	(t2) $Y \leftarrow 2d$
(t2) $c \leftarrow X(S1)$	(t1) $b \leftarrow Y(S2)$
(t2) $X \leftarrow 2*c$	(t1) $Y \leftarrow b+100$

read: t1 reads X@S1 into a ... etc

$ts(t_1) < ts(t_2)$

Abort t1 at S1

Cascading abort of t2

Abort t2 at S2

hs / FUB dbsII-03-17DDBCC1-4

### Strict timestamp ordering

- Strict TO
- Lock the items changed until ta has been committed (or aborted)

(Node S1)	(Node S2)
(t1) $a \leftarrow X(S1)$	(t2) $d \leftarrow Y(S2)$
(t1) $X \leftarrow a+100$	(t2) $Y \leftarrow 2d$
LOCK X	(t1) $b \leftarrow Y(S2)$
abort t1 at S1	(t1) $Y \leftarrow b+100$

UNLOCK T1

(t2)  $c \leftarrow X(S1)$

(t2)  $X \leftarrow 2*c$

$ts(t_1) < ts(t_2)$   
abort t1

hs / FUB dbsII-03-17DDBCC1-5

### TO Scheduler

- Basic principle:
- Abort transaction if its operation is "too late"
- Remember timestamp of last write of X:  $\max W[X]$
- and last read  $\max R[X]$

Transaction i:  $t_i$  with timestamp  $ts(t_i)$

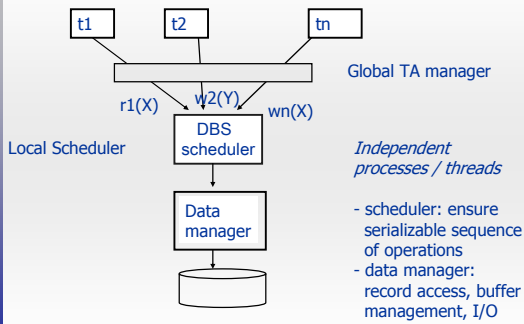
Operations:  $ri(X) / wi(X)$  -  $t_i$  wants to read / write X

Scheduler state:  $\max R[X] / \max W[X]$

timestamp of youngest TA which read X / has written X

hs / FUB dbsII-03-17DDBCC1-6

## Scheduler / Data manager architecture



hs / FUB dbsII-03-17DDBCC1-7

## TO Scheduler: read

**Read:** TA  $t_i$  with timestamp  $ts(t_i)$  reads  $X$  :  $ri(X)$

$\max W[X] > ts(t_i)$ :  
there is a younger TA which has written  $X$   
 $\Rightarrow$  contradicts timestamp ordering:  
 $t_i$  reads too late  
 $\Rightarrow$  abort  $TA\ t_i$ , restart  $t_i$

$\max W[X] < ts(t_i) \Rightarrow$  go ahead

Example:  $-----|-----|----->$   
 $ri(X) \quad wj(X) \quad ts(t_i) < ts(t_j)$   
 but data manager may execute  $wj(X)$  before  $ri(X)$

hs / FUB dbsII-03-17DDBCC1-8

## TO Scheduler: read

**Write:** TA  $t_i$  with timestamp  $ts(t_i)$  writes  $X$  :  $wi(X)$

$\max W[X] > ts(t_i) \vee \max R[X] > ts(t_i)$ :  
*/\* but X has been written or read by younger transaction.*

$\Rightarrow$  # timestamp ordering  
 $\Rightarrow$  abort  $TA\ t_i$

otherwise:  $\Rightarrow$  schedule  $wi(X)$  for execution

► Same issue as with 'read'

► Solution: 'lock variables'

number of Readers of  $X$  :  $nR[X]$   
 number of Writers of  $X$  :  $nW[X]$

$nR[X] == 0 \wedge nW[X] == 0$  : schedule  $wi(X)$   
 otherwise enqueue  $wi(X)$

similar for reads...

hs / FUB dbsII-03-17DDBCC1-9

## TO scheduling

Note:

- This kind of blocking reads / writes not necessary with 2PL since 2PL lock manager blocks anyway
- Different from waiting for EOT to enforce strict TO protocol

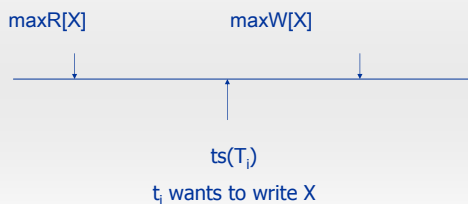
Can be implemented with  $nW[X]$ : if  $t_i$  has incremented  $nW[X]$   $0 \rightarrow 1$   
 wait for commit ( $t_i$ ) before decrementing  $nW[X]$

- Enqueue operations  $ri(X)$  /  $wi(X)$   
 if  $nW[X] > 0$  /  $nW[X] > 0 \vee nR[X] > 0$
- Release appropriate element from queue if  $nW[X]$  and/or  $nR[X]$  decreased

hs / FUB dbsII-03-17DDBCC1-10

## Thomas Write Rule

- Idea: younger write overwrites older write without changing effect of timestamp ordering



hs / FUB dbsII-03-17DDBCC1-11

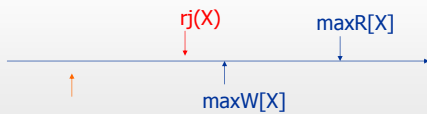
## Modified Time ordered scheduling of writes

Queue  $q$ ; ....

```
void write ( TA ti, Object X ) { /* write wi(X)
    if ts(ti) < maxR[X] ti.abort();
    else if ts(ti) < maxW[X];
        /* IGNORE WRITE (tell ti "success")
    else { /* process write as before...
        maxW[X] ← ts(ti);
        if ( q.isEmpty() ∧ nW[x]==0 ∧ nR[X]==0 ) {
            nW[X] ← 1; dataMgr.write(X);
            WAIT ( commit (ti) );
        }
        else q.add ( W, ti );
    }
}
```

hs / FUB dbsII-03-17DDBCC1-12

## Another rule?



$ts(t_i)$

TA  $t_i$  with  $ts(t_i) < maxW[X] < maxR[X]$  wants to write

Can  $t_i$  go ahead and ignore  $w_i(X)$  ??

If there is a read  $r_j(X)$  before the write  $\rightarrow maxW[X]$   
then TA  $t_j$  has read a wrong value:  
wrong order if  $t_j < t_i$

hs / FUB dbsII-03-17DDBCC1-13

## Management of timestamps

data	maxR	maxW	nR / nW
X1			
X2			
⋮	⋮	⋮	
Xn			

$maxR / maxW$  must be kept for each data item in DB  
**Space! I/O traffic !!**

hs / FUB dbsII-03-17DDBCC1-14

## Timestamp cache

data	maxR	maxW	nR / nW
X1			
X2			
⋮	⋮	⋮	
Xn			

tsMin

- if TA reads / writes data item: make cache entry
- flush cache periodically: purge rows with timestamp  $\leq tsMin$  e.g.  $currentTime - \delta$
- ts for data on disk: tsMin

hs / FUB dbsII-03-17DDBCC1-15

## Timestamp cache

- ▶ Enforce timestamp order rule for  $w_i(X)$  :
  - Use  $maxW(X)$  if  $X$  is in cache
  - Assume  $maxW(X) = tsMin$  otherwise
  - Same for reads
- ▶ Hash table as data structure (like lock table)

hs / FUB dbsII-03-17DDBCC1-16

## Time stamp order and distribution

- ▶ Distributed TO scheduling
  - Basic prerequisite: *total order of TA timestamps*  
e.g. local counter ++ server#
  - TO schedulers independent
  - Total order guarantees serializable read / write at all sites
  - At TA commit: release  $nW(X) / nR(X)$  locks
- Tricky detail:  
suppose only a few TA at site  $S_1$ , many at site  $S_2$   
 $\Rightarrow counter[S_1] \ll counter[S_2]$   
 $\Rightarrow$  TAs originating at  $S_2$  will be frequently aborted
- Any idea to solve this problem?

hs / FUB dbsII-03-17DDBCC1-17

## Time stamp cc and 2PL

TO  $\Rightarrow$  conflict serializable  
2PL  $\Rightarrow$  conflict serializable

2PL = TO ??

**NO!**

$t_1: w_1[Y]$   
 $t_2: r_2[X] \ r_2[Y] \ w_2[Z]$   
 $t_3: w_3[X]$   
 $ts(t_1) < ts(t_2) < ts(t_3)$

$S: r_2[X] \ w_3[X] \ w_1[Y] \ r_2[Y] \ w_2[Z]$

$S$  could be produced with T.O. but not with 2PL

hs / FUB dbsII-03-17DDBCC1-18

## Pessimistic vs Optimistic

- ▶ Timestamp Order is pessimistic
  - All checks are made *before* operation is scheduled
  - Optimistic: work isolated on copy of data  
write back if no potential conflict has occurred

hs / FUB dbsII-03-17DDBCC1-19

## Optimistic protocol

- ▶ Optimistic CC in homogeneous DDBS
- ▶ Same protocol as centralized DBS
  - Read phase
  - Validation phase:  
ReadSet ( $t_j \mid t_j \text{ running}$ )  $\cap$  WriteSet (validation TA)  $== \emptyset$  ?  
Backward oriented optimistic CC (BOCC)
- ▶ Issue: Validation of TAs on different Servers in the same order

hs / FUB dbsII-03-17DDBCC1-20

## CC in homogeneous DDBS

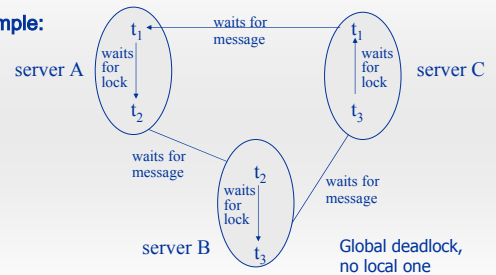
- ▶ **2PL**
  - Used frequently in commercial systems
  - Simple enhancement in distributed systems
  - Deadlocks possible – critical in distributed DBS
- ▶ **Timestamp ordering**
  - a reasonable alternative
  - aborts more likely
  - no deadlocks
  - useful in a distributed system
- ▶ **Optimistic**
  - become popular in widely distributed systems

hs / FUB dbsII-03-17DDBCC1-21

## Distributed Deadlocks

### Deadlock detection in DDB

Example:



Expl by Weikum / Vossen

hs / FUB dbsII-03-17DDBCC1-22

## Distributed Deadlocks

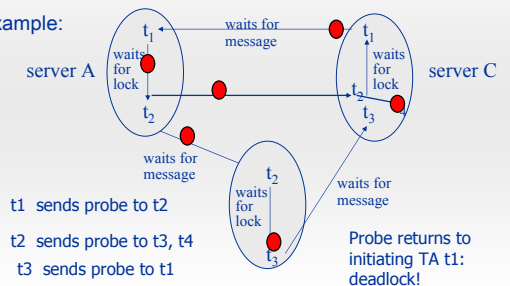
- Timeout**
  - used in most commercial systems
  - Abort transaction if TA waits longer than specified timeout. Victim?
- Centralized DL detection**
  - one site S is responsible for DL detection
  - other sites send periodically local Wait-For-graphs to S
  - S forms Global WF-graph and checks for cycles

hs / FUB dbsII-03-17DDBCC1-23

## Distributed Deadlock detection

- Distributed detection**
  - Edge chasing (probing)

Example:



hs / FUB dbsII-03-17DDBCC1-24

## Distributed deadlock detection

### ► Path pushing

#### Algorithm:

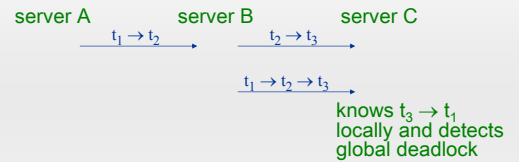
- each node that has a wait-for path from transaction  $t_i$  to  $t_j$  such that  $t_i$  has an *incoming wait-for-message edge* and  $t_j$  has an *outgoing wait-for-message edge*, sends the path to the server along the outgoing edge (if the id (or timestamp) of  $t_i$  is smaller than the id of  $t_j$ )
- upon *receiving a path*, each node *concatenates it with its local paths* and forwards it along outgoing edges
- if there is a cycle among  $n$  servers, *at least one server will detect the cycle* after at most  $n$  such rounds

2PL: no false deadlocks

hs / FUB dbsII-03-17DDBCC1-25

## Distributed deadlock detection: path pushing

### ► Example from above



hs / FUB dbsII-03-17DDBCC1-26

## Summary

### ► Homogenous concurrency control

- Slight extension of centralized protocols
- Always possible to introduce some kind of centralized control
- Contradicts principles of avoiding single point of failure and scalability
- TS ordering with little overhead
- Deadlock detection: expensive or time-out

hs / FUB dbsII-03-17DDBCC1-27