

# 8. Concurrency Control for Transactions Part Two

CSEP 545 Transaction Processing  
Philip A. Bernstein

Copyright © 2005 Philip A. Bernstein

2.45.05

1

## Outline

1. A Model for Concurrency Control
2. Serializability Theory
3. Synchronization Requirements for Recoverability
4. Two-Phase Locking
5. Implementing Two-Phase Locking
6. Locking Performance
7. Multigranularity Locking (revisited)
8. Hot Spot Techniques
9. Query-Update Techniques
10. Phantom s
11. Shared Disk System s
12. B-Trees
13. Tree locking

2.45.05

2

## 8.6 Locking Performance

- Deadlocks are rare
  - up to 1% - 2% of transactions deadlock
- The one exception to this is lock conversions
  - r-lock a record and later upgrade to w-lock
  - e.g.,  $T_i = \text{read}(x) \dots \text{write}(x)$
  - if two txns do this concurrently, they'll deadlock (both get an r-lock on x before either gets a w-lock)
  - To avoid lock conversion deadlocks, get a w-lock first and down-grade to an r-lock if you don't need to write.
  - Use SQL Update statement or explicit program hints

2.45.05

3

## Conversions in M S SQL Server

- Update-lock prevents lock conversion deadlock.
  - Conflicts with other update and write locks, but not with read locks.
  - Only on pages and rows (not tables)
- You get an update lock by using the UPD LOCK hint in the FROM clause

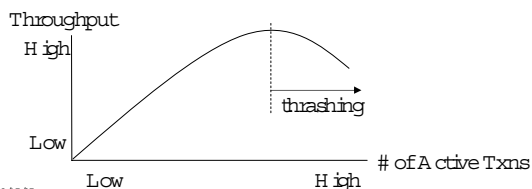
```
Select Foo A
From Foo (UPD LOCK)
Where Foo B = 7
```

2.45.05

4

## Blocking and Lock Thrashing

- The locking performance problem is too much delay due to blocking
  - little delay until locks are saturated
  - then major delay, due to the locking bottleneck
  - thrashing - the point where throughput decreases with increasing load



2.45.05

5

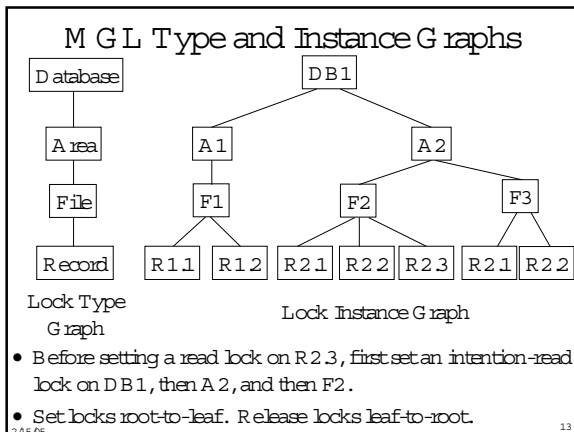
## More on Thrashing

- It's purely a blocking problem
  - It happens even when the abort rate is low
- As number of transactions increase
  - each additional transaction is more likely to block
  - but first, it gathers some locks, increasing the probability others will block (negative feedback)

2.45.05

6



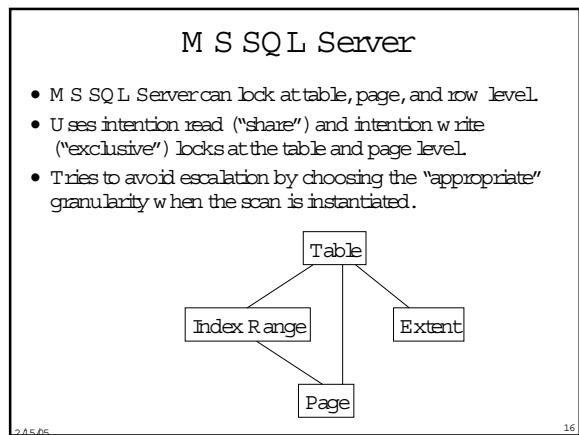
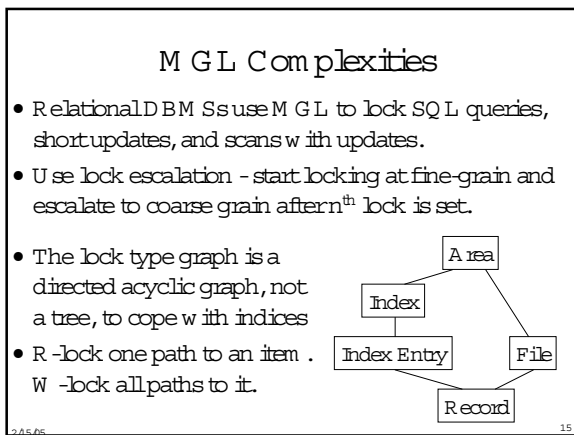


### M G L Compatibility Matrix

	r	w	ir	iw	riw
r	y	n	y	n	n
w	n	n	n	n	n
ir	y	n	y	y	y
iw	n	n	y	y	n
riw	n	n	y	n	n

riw = read with intent to write, for a scan that updates some of the records it reads

- E.g., ir conflicts with w because ir says there's a fine-grained r-lock that conflicts with a w-lock on its parent
- To r-lock an item, need an r-, ir- or riw-lock on its parent
- To w-lock an item, need a w-, iw- or riw-lock on its parent



### 8.8 Hot Spot Techniques

- If each txn holds a lock for t seconds, then the max throughput is 1/t txns/second for that lock.
- Hot spot - A data item that's more popular than others, so a large fraction of active txns need it
  - Summary information (total inventory)
  - End-of-file marker in data entry application
  - Counter used for assigning serial numbers
- Hot spots often create a convoy of transactions. The hot spot lock serializes transactions.

### Hot Spot Techniques (cont'd)

- Special techniques are needed to reduce it
  - Keep the hot data in main memory
  - Delay operations on hot data till commit time
  - Use optimistic methods
  - Batch up operations to hot spot data
  - Partition hot spot data

## Delaying Operations Until Commit

- Data manager logs each transaction's updates
- Only applies the updates (and sets locks) after receiving Commit from the transaction
- IMS FastPath uses this for
  - Data Entry DB
  - Main Storage DB
- Works for write, insert, and delete, but not read

2.45.05

19

## Locking Higher-Level Operations

- Read is often part of a read-write pair, such as Increment( $x, n$ ), which adds constant  $n$  to  $x$ , but doesn't return a value.
- Increment (and Decrement) commute
- So, introduce Increment and Decrement locks

	r	w	inc	dec
r	y	n	n	n
w	n	n	n	n
inc	n	n	y	y
dec	n	n	y	y

- But if Inc and Dec have a threshold (e.g. a quantity of zero), then they conflict (when the threshold is near)

2.45.05

20

## Solving the Threshold Problem

A nother IMS FastPath Technique

- Use a blind Decrement (no threshold) and Verify( $x, n$ ), which returns true if  $x \neq n$
- Re-execute Verify at commit time
  - If it returns a different value than it did during normal execution, then abort
  - It's like checking that the threshold lock you didn't set during Decrement is still valid.

```
bEnough = Verify(iQuantity, n);
If (bEnough) Decrement(iQuantity, n)
else print ("not enough");
```

2.45.05

21

## Optimistic Concurrency Control

- The Verify trick is optimistic concurrency control
- Main idea - execute operations on shared data without setting locks. At commit time, test if there were conflicts on the locks (that you didn't set).
- Often used in client/server systems
  - Client does all updates in cache without shared locks
  - At commit time, try to get locks and perform updates

2.45.05

22

## Batching

- Transactions add updates to a mini-batch and only periodically apply the mini-batch to shared data.
  - Each process has a private data entry file, in addition to a global shared data entry file
  - Each transaction appends to its process' file
  - Periodically append the process file to the shared file
- Tricky failure handling
  - Gathering up private files
  - Avoiding holes in serial number order

2.45.05

23

## Partitioning

- Split up inventory into partitions
- Each transaction only accesses one partition
- Example
  - Each ticket agency has a subset of the tickets
  - If one agency sells out early, it needs a way to get more tickets from other agencies (partitions)

2.45.05

24

## 8.9 Query-Update Techniques

- Queries run for a long time and lock a lot of data – a performance nightmare when trying also to run short update transactions
- There are several good solutions
  - Use a data warehouse
  - Accept weaker consistency guarantees
  - Use multiversion data
- Solutions trade data quality or timeliness for performance

2.45.05

25

## Data Warehouse

- A data warehouse contains a snapshot of the DB which is periodically refreshed from the TPDB
- All queries run on the data warehouse
- All update transactions run on the TPDB
- Queries don't get absolutely up-to-date data
- How to refresh the data warehouse?
  - Stop processing transactions and copy the TPDB to the data warehouse. Possibly run queries while refreshing
  - Treat the warehouse as a DB replica and use a replication technique

2.45.05

26

## Degrees of Isolation

- Serializability = Degree 3 Isolation
- Degree 2 Isolation (a.k.a. cursor stability)
  - Data manager holds read-lock (x) only while reading x, but holds write locks till commit (as in 2PL)
  - E.g., when scanning records in a file, each get-next-record releases lock on current record and gets lock on next one
  - read (x) is not "repeatable" within a transaction, e.g.,  
 $r_1[x] r_1[x] w_1[x] w_2[x] w_2[x] w_2[x] r_1[x] r_1[x] w_1[x]$
  - Degree 2 is commonly used by ISAM file systems
  - Degree 2 is often a DB system's default behavior! And customers seem to accept it!!!

2.45.05

27

## Degrees of Isolation (cont'd)

- Could run queries Degree 2 and updates Degree 3
  - Updates are still serializable w.r.t. each other
- Degree 1 - no read locks; hold write locks to commit
- Unfortunately, SQL concurrency control standards have been stated in terms of "repeatable reads" and "cursor stability" instead of serializability, leading to much confusion.

2.45.05

28

## ANSI SQL Isolation Levels

- Uncommitted Read - Degree 1
- Committed Read - Degree 2
- Repeatable Read - Uses read locks and write locks, but allows "phantom s"
- Serializable - Degree 3

2.45.05

29

## MS SQL Server

- Lock hints in SQL FROM clause
  - All the ANSI isolation levels, plus ...
  - UPDLOCK - use update locks instead of read locks
  - READPAST - ignore locked rows (if running read committed)
  - PAGLOCK - use page lock when the system would otherwise use a table lock
  - TABLOCK - shared table lock till end of commit and/or transaction
  - TABLOCKX - exclusive table lock till end of commit and/or transaction

2.45.05

30

## Multiversion Data

- Assume record granularity locking
- Each write operation creates a new version instead of overwriting existing value.
- So each logical record has a sequence of versions.
- Tag each record with transaction id of the transaction that wrote that version

Tid	Previous	E#	Name	Other fields
123	null	1	Bill	
175	123	1	Bill	
134	null	2	Sue	
199	134	2	Sue	
227	null	27	Steve	

2.45.05

31

## Multiversion Data (cont'd)

- Execute update transactions using ordinary 2PL
- Execute queries in snapshot mode
  - System keeps a commit list of tids of all committed txns
  - When a query starts executing, it reads the commit list
  - When a query reads x, it reads the latest version of x written by a transaction on its commit list
  - Thus, it reads the database state that existed when it started running

2.45.05

32

## Commit List Management

- Maintain and periodically recompute a tid  $T-O_{list}$ , such that
  - Every active txn's tid is greater than  $T-O_{list}$
  - Every new tid is greater than  $T-O_{list}$
  - For every committed transaction with tid  $\in T-O_{list}$ , its versions are committed
  - For every aborted transaction with tid  $\in T-O_{list}$ , its versions are wiped out
- Queries don't need to know tids  $\in T-O_{list}$ 
  - So only maintain the commit list for tids  $> T-O_{list}$

2.45.05

33

## Multiversion Garbage Collection

- Can delete an old version of x if no query will ever read it
  - There's a later version of x whose tid  $\leq T-O_{list}$  (or is on every active query's commit list)
- Originally used in Prime Computer's CODASYL DB system and Oracle's Rdb/VM S

2.45.05

34

## Oracle Multiversion Concurrency Control

- Data page contains latest version of each record, which points to older version in rollback segment.
- Read-committed query reads data as of its start time.
- Read-only isolation reads data as of transaction start time.
- "Serializable" query reads data as of the txn's start time.
  - An update checks that the updated record was not modified after txn start time.
  - If that check fails, Oracle returns an error.
  - If there isn't enough history for Oracle to perform the check, Oracle returns an error. (You can control the history area's size.)
  - What if  $T_1$  and  $T_2$  modify each other's readset concurrently?

2.45.05

35

## Oracle Concurrency Control (cont'd)

$r_1 [x] r_1 [y] r_2 [x] r_2 [y] w_1 [x] c_1 w_2 [y] c_2$

- The result is not serializable!
- In any SR execution, one transaction would have read the other's output

2.45.05

36

## 8.10 Phantom s

- Problem s when using 2PL w ith inserts and deletes

A ccounts			A ssets	
A cct#	Location	Balance	Location	Total
1	Seattle	400	Seattle	400
2	Tacom a	200	Tacom a	500
3	Tacom a	300		

$T_1$ : Read A ccounts 1, 2, and 3  
 $T_2$ : Insert A ccounts [4, Tacom a, 100] ← The phantom record  
 $T_2$ : Read A ssets (Tacom a), returns 500  
 $T_2$ : Write A ssets (Tacom a, 600)  
 $T_1$ : Read A ssets (Tacom a), returns 600  
 $T_1$ : Com m it

2.45.05

37

## The Phantom Phantom Problem

- It looks like  $T_1$  should lock record 4, which isn't there!
- Which of  $T_1$ 's operations determined that there were only 3 records?
  - Read end-of-file?
  - Read record counter?
  - SQL Select operation?
- This operation conflicts w ith  $T_2$ 's Insert A ccounts [4, Tacom a, 100]
- Therefore, Insert A ccounts [4, Tacom a, 100] shouldn't run until after  $T_1$  com m its

2.45.05

38

## A voiding Phantom s - Predicate Locks

- Suppose a query reads all records satisfying predicate P. For example,
  - Select \* From A ccounts Where Location = "Tacom a"
  - Normally would hash each record id to an integer lock id
  - And lock control structures. Too coarse grained.
- Ideally, set a read lock on P
  - which conflicts w ith a write lock Q if some record can satisfy (P and Q)
- For arbitrary predicates, this is too slow to check
  - Not w ithin a few hundred instructions, anyway

2.45.05

39

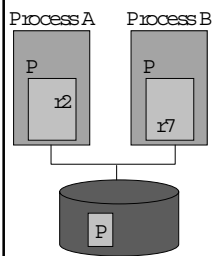
## Precision Locks

- Suppose update operations are on single records
- Maintain a list of predicate Read-locks
- Insert, Delete, & Update write-lock the record and check for conflict w ith all predicate locks
- Query sets a read lock on the predicate and check for conflict w ith all record locks
- Cheaper than predicate satisfiability, but still too expensive for practical implementation.

2.45.05

40

## 8.11 Shared Disk System s



- Can cache a page in two processes that write-lock different records
- Only one process at a time can have write privilege
- Use a global lock manager
- When setting a write lock on P, may need to refresh the cached copy from disk (if another process recently updated it)

- Use a version number on the page and in the lock

2.45.05

41

## Shared Disk System

- When a process sets the lock, it tells the lock manager version number of its cached page.
- A process increments the version number the first time it updates a cached page.
- When a process is done w ith an updated page, it flushes the page to disk and then increments version number in the lock.

2.45.05

42

## Logging

- Since updates are happening on different systems, where is the log?
- A single log server is simplest, but makes logging more expensive.
- Be careful not to flush to the log until necessary.
  - This requires locally-assigned LSNs
  - Must flush the log before flushing an updated page

2.45.05

43

## 8.12 B-Trees

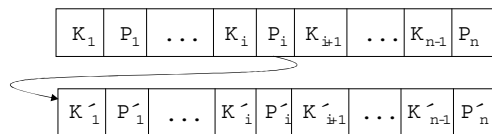
- An index maps field values to record ids.
  - Record id = [page-id, offset-within-page]
  - Most common DB index structures: hashing and B-trees
  - DB index structures are page-oriented
- Hashing uses a function  $H: V \rightarrow B$ , from field values to block numbers.
  - $V$  = social security numbers,  $B = \{1 \dots 1000\}$
  - $H(v) = v \bmod 1000$
  - If a page overflows, then use an extra overflow page
  - At 90% load on pages, 1.2 block accesses per request!
  - BUT, doesn't help for key range access ( $10 < v < 75$ )

2.45.05

44

## B-Tree Structure

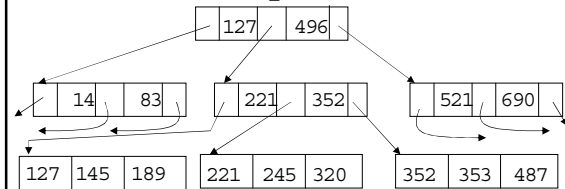
- Index node is a sequence of [pointer, key] pairs
- $K_1 < K_2 < \dots < K_{n-2} < K_{n-1}$
- $P_1$  points to a node containing keys  $< K_1$
- $P_i$  points to a node containing keys in range  $[K_{i-1}, K_i)$
- $P_n$  points to a node containing keys  $> K_{n-1}$
- So,  $K'_1 < K'_2 < \dots < K'_{n-2} < K'_{n-1}$



2.45.05

45

## Example n=3



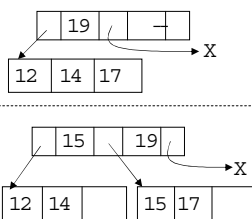
- Notice that leaves are sorted by key, left-to-right
- Search for value  $v$  by following path from the root
- If key = 8 bytes, ptr = 2 bytes, page = 4K, then  $n = 409$
- So 3-level index has up to 68M leaves ( $409^3$ )
- At 20 records per leaf, that's 136M records

2.45.05

46

## Insertion

- To insert key  $v$ , search for the leaf where  $v$  should appear
- If there's space on the leaf, insert the record
- If no, split the leaf in half, and split the key range in its parent to point to the two leaves



- To insert key 15
- split the leaf
  - split the parent's range  $[0, 19)$  to  $[0, 15)$  and  $[15, 19)$
  - if the parent was full, you'd split that too (not shown here)
  - this automatically keeps the tree balanced

2.45.05

47

## B-Tree Observations

- Delete algorithm merges adjacent nodes < 50% full, but rarely used in practice
- Root and most level-1 nodes are cached, to reduce disk accesses
- Secondary (non-clustered) index - Leaves contain [key, record id] pairs.
- Primary (clustered) index - Leaves contain records
- Use key prefix for long (string) key values
  - drop prefix and add to suffix as you move down the tree

2.45.05

48



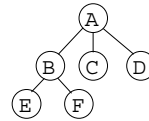
### Key Range Locks

- Lock on B-tree key range is a cheap predicate lock
- 
- Select DeptW here (Budget > 250 and Budget < 350)
  - lock the key range [221, 352) record
  - only useful when query is on an indexed field
- Commonly used with multi-granularity locking
    - Insert/delete locks record and intention-write locks range
    - MGL tree defines a fixed set of predicates, and thereby avoids predicate satisfiability

2.45/65 49

### 8.13 Tree Locking

- Can beat 2PL by exploiting root-to-leaf access in a tree
- If searching for a leaf, after setting a lock on a node, release the lock on its parent



wl(A) wl(B) wu(A) wl(E) wu(B)

- The lock order on the root serializes access to other nodes

2.45/65 50

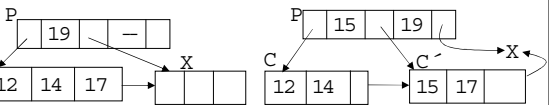
### B-tree Locking

- Root lock on a B-tree is a bottleneck
  - Use tree locking to relieve it
  - Problem: node splits
- 
- If you unlock P before splitting C, then you have to back up and lock P again, which breaks the tree locking protocol.
- So, don't unlock a node till you're sure its child won't split (i.e. has space for an insert)
  - Implies different locking rules for different ops (search vs. insert/update)

2.45/65 51

### B-link Optimization

- B-link tree - Each node has a side pointer to the next
- After searching a node, you can release its lock before locking its child



- Searching has the same behavior as if it locked the child before releasing the parent... and ran later (after the insert)

2.45/65 52