## CHAPTER 16

## **Concurrency Control**

## **Solutions to Practice Exercises**

**16.1** Suppose two-phase locking does not ensure serializability. Then there exists a set of transactions  $T_0$ ,  $T_1 \dots T_{n-1}$  which obey 2PL and which produce a non-serializable schedule. A non-serializable schedule implies a cycle in the precedence graph, and we shall show that 2PL cannot produce such cycles. Without loss of generality, assume the following cycle exists in the precedence graph:  $T_0 \rightarrow T_1 \rightarrow T_2 \rightarrow \dots \rightarrow T_{n-1} \rightarrow T_0$ . Let  $\alpha_i$  be the time at which  $T_i$  obtains its last lock (i.e.  $T_i$ 's lock point). Then for all transactions such that  $T_i \rightarrow T_j$ ,  $\alpha_i < \alpha_j$ . Then for the cycle we have

$$\alpha_0 < \alpha_1 < \alpha_2 < \dots < \alpha_{n-1} < \alpha_0$$

Since  $\alpha_0 < \alpha_0$  is a contradiction, no such cycle can exist. Hence 2PL cannot produce non-serializable schedules. Because of the property that for all transactions such that  $T_i \rightarrow T_j$ ,  $\alpha_i < \alpha_j$ , the lock point ordering of the transactions is also a topological sort ordering of the precedence graph. Thus transactions can be serialized according to their lock points.

**16.2 a.** Lock and unlock instructions:

 $T_{31}: lock-S(A)$ read(A) lock-X(B) read(B) if A = 0then B := B + 1write(B) unlock(A) unlock(B)

```
T_{32}: lock-S(B)
read(B)
lock-X(A)
read(A)
if B = 0
then A := A + 1
write(A)
unlock(B)
unlock(A)
```

- **b.** Execution of these transactions can result in deadlock. For example, consider the following partial schedule:
  - $\begin{array}{c|c} T_{31} & T_{32} \\ \hline \text{lock-S}(A) & \\ & \text{lock-S}(B) \\ \text{read}(A) \\ \text{lock-X}(B) & \\ \hline \end{array}$

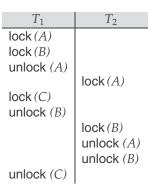
The transactions are now deadlocked.

- **16.3** Rigorous two-phase locking has the advantages of strict 2PL. In addition it has the property that for two conflicting transactions, their commit order is their serializability order. In some systems users might expect this behavior.
- **16.4** The proof is in Buckley and Silberschatz, "Concurrency Control in Graph Protocols by Using Edge Locks," Proc. ACM SIGACT-SIGMOD Symposium on the Principles of Database Systems, 1984.
- 16.5 Consider the tree-structured database graph given below.



Schedule possible under tree protocol but not under 2PL:

Exercises 79



Schedule possible under 2PL but not under tree protocol:

$T_1$	$T_2$
lock(A)	
	lock(B)
lock(C)	
	unlock (B)
unlock (A)	
unlock (C)	

- **16.6** The proof is in Kedem and Silberschatz, "Locking Protocols: From Exclusive to Shared Locks," JACM Vol. 30, 4, 1983.
- **16.7** The proof is in Kedem and Silberschatz, "Controlling Concurrency Using Locking Protocols," Proc. Annual IEEE Symposium on Foundations of Computer Science, 1979.
- **16.8** The proof is in Kedem and Silberschatz, "Controlling Concurrency Using Locking Protocols," Proc. Annual IEEE Symposium on Foundations of Computer Science, 1979.
- **16.9** The access protection mechanism can be used to implement page level locking. Consider reads first. A process is allowed to read a page only after it read-locks the page. This is implemented by using mprotect to initially turn off read permissions to all pages, for the process. When the process tries to access an address in a page, a protection violation occurs. The handler associated with protection violation then requests a read lock on the page, and after the lock is acquired, it uses mprotect to allow read access to the page by the process, and finally allows the process to continue. Write access is handled similarly.
- **16.10** The proof is in Korth, "Locking Primitives in a Database System," JACM Vol. 30, 1983.
- **16.11** It would make no difference. The write protocol is such that the most recent transaction to write an item is also the one with the largest timestamp to have done so.

- **16.12** If a transaction needs to access a large a set of items, multiple granularity locking requires fewer locks, whereas if only one item needs to be accessed, the single lock granularity system allows this with just one lock. Because all the desired data items are locked and unlocked together in the multiple granularity scheme, the locking overhead is low, but concurrency is also reduced.
- **16.13** In the concurrency control scheme of Section 16.3 choosing  $\mathbf{Start}(T_i)$  as the timestamp of  $T_i$  gives a subset of the schedules allowed by choosing  $\mathbf{Validation}(T_i)$  as the timestamp. Using  $\mathbf{Start}(T_i)$  means that whoever started first must finish first. Clearly transactions could enter the validation phase in the same order in which they began executing, but this is overly restrictive. Since choosing  $\mathbf{Validation}(T_i)$  causes fewer nonconflicting transactions to restart, it gives the better response times.
- **16.14** Two-phase locking: Use for simple applications where a single granularity is acceptable. If there are large read-only transactions, multiversion protocols would do better. Also, if deadlocks must be avoided at all costs, the tree protocol would be preferable.
  - Two-phase locking with multiple granularity locking: Use for an application mix where some applications access individual records and others access whole relations or substantial parts thereof. The drawbacks of 2PL mentioned above also apply to this one.
  - The tree protocol: Use if all applications tend to access data items in an order consistent with a particular partial order. This protocol is free of deadlocks, but transactions will often have to lock unwanted nodes in order to access the desired nodes.
  - Timestamp ordering: Use if the application demands a concurrent execution that is equivalent to a particular serial ordering (say, the order of arrival), rather than *any* serial ordering. But conflicts are handled by roll-back of transactions rather than waiting, and schedules are not recoverable. To make them recoverable, additional overheads and increased response time have to be tolerated. Not suitable if there are long read-only transactions, since they will starve. Deadlocks are absent.
  - Validation: If the probability that two concurrently executing transactions conflict is low, this protocol can be used advantageously to get better concurrency and good response times with low overheads. Not suitable under high contention, when a lot of wasted work will be done.
  - Multiversion timestamp ordering: Use if timestamp ordering is appropriate but it is desirable for read requests to never wait. Shares the other disadvantages of the timestamp ordering protocol.
  - Multiversion two-phase locking: This protocol allows read-only transactions to always commit without ever waiting. Update transactions follow 2PL, thus allowing recoverable schedules with conflicts solved by waiting rather than roll-back. But the problem of deadlocks comes back, though read-only transactions cannot get involved in them. Keeping multiple ver-

sions adds space and time overheads though, therefore plain 2PL may be preferable in low conflict situations.

**16.15** A transaction waits on *a*. disk I/O and *b*. lock acquisition. Transactions generally wait on disk reads and not on disk writes as disk writes are handled by the buffering mechanism in asynchronous fashion and transactions update only the in-memory copy of the disk blocks.

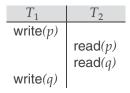
The technique proposed essentially separates the waiting times into two phases. The first phase – where transaction is executed without acquiring any locks and without performing any writes to the database – accounts for almost all the waiting time on disk I/O as it reads all the data blocks it needs from disk if they are not already in memory. The second phase — the transaction re-execution with strict two-phase locking — accounts for all the waiting time on acquiring locks. The second phase may, though rarely, involve a small waiting time on disk I/O if a disk block that the transaction needs is flushed to memory (by buffer manager) before the second phase starts.

The technique may increase concurrency as transactions spend almost no time on disk I/O with locks held and hence locks are held for shorter time. In the first phase the transaction reads all the data items required—and not already in memory—from disk. The locks are acquired in the second phase and the transaction does almost no disk I/O in this phase. Thus the transaction avoids spending time in disk I/O with locks held.

The technique may even increase disk throughput as the disk I/O is not stalled for want of a lock. Consider the following scenario with strict twophase locking protocol: A transaction is waiting for a lock, the disk is idle and there are some item to be read from disk. In such a situation disk bandwidth is wasted. But in the proposed technique, the transaction will read all the required item from the disk without acquiring any lock and the disk bandwidth may be properly utilized.

Note that the proposed technique is most useful if the computation involved in the transactions is less and most of the time is spent in disk I/O and waiting on locks, as is usually the case in disk-resident databases. If the transaction is computation intensive, there may be wasted work. An optimization is to save the updates of transactions in a temporary buffer, and instead of reexecuting the transaction, to compare the data values of items when they are locked with the values used earlier. If the two values are the same for all items, then the buffered updates of the transaction are executed, instead of reexecuting the entire transaction.

**16.16** Consider two transactions  $T_1$  and  $T_2$  shown below.



## 82 Chapter 16 Concurrency Control

Let  $TS(T_1) < TS(T_2)$  and let the timestamp test at each operation except write(*q*) be successful. When transaction  $T_1$  does the timestamp test for write(*q*) it finds that  $TS(T_1) < R$ -timestamp(*q*), since  $TS(T_1) < TS(T_2)$  and R-timestamp(*q*) =  $TS(T_2)$ . Hence the writeoperation fails and transaction  $T_1$  rolls back. The cascading results in transaction  $T_2$  also being rolled back as it uses the value for item *p* that is written by transaction  $T_1$ .

If this scenario is exactly repeated every time the transactions are restarted, this could result in starvation of both transactions.

- **16.17** In the text, we considered two approaches to dealing with the phantom phenomenon by means of locking. The coarser granularity approach obviously works for timestamps as well. The  $B^+$ -tree index based approach can be adapted to timestamping by treating index buckets as data items with timestamps associated with them, and requiring that all read accesses use an index. We now show that this simple method works. Suppose a transaction  $T_i$  wants to access all tuples with a particular range of search-key values, using a  $B^+$ -tree index on that search-key.  $T_i$  will need to read all the buckets in that index which have key values in that range. It can be seen that any delete or insert of a tuple with a key-value in the same range will need to write one of the index buckets read by  $T_i$ . Thus the logical conflict is converted to a conflict on an index bucket, and the phantom phenomenon is avoided.
- **16.18** Note: The tree-protocol of Section 16.1.5 which is referred to in this question, is different from the multigranularity protocol of Section 16.4 and the  $B^+$ -tree concurrency protocol of Section 16.9.

One strategy for early lock releasing is given here. Going down the tree from the root, if the currently visited node's child is not full, release locks held on all nodes except the current node, request an X-lock on the child node, after getting it release the lock on the current node, and then descend to the child. On the other hand, if the child is full, retain all locks held, request an X-lock on the child, and descend to it after getting the lock. On reaching the leaf node, start the insertion procedure. This strategy results in holding locks only on the full index tree nodes from the leaf upwards, until and including the first nonfull node.

An optimization to the above strategy is possible. Even if the current node's child is full, we can still release the locks on all nodes but the current one. But after getting the X-lock on the child node, we split it right away. Releasing the lock on the current node and retaining just the lock on the appropriate split child, we descend into it making it the current node. With this optimization, at any time at most two locks are held, of a parent and a child node.