

Recovery Protocols for Shared Memory Database Systems

Lory D. Molesky and Krithi Ramamritham

Department of Computer Science

University of Massachusetts

Amherst MA 01003-4610

e-mail: lory@cs.umass.edu, krithi@cs.umass.edu

Abstract

Significant performance advantages can be gained by implementing a database system on a cache-coherent shared memory multiprocessor. However, problems arise when failures occur. A single node (where a *node* refers to a processor/memory pair) crash may require a reboot of the entire shared memory system. Fortunately, shared memory multiprocessors that isolate individual node failures are currently being developed. Even with these, because of the side effects of the cache coherency protocol, a transaction executing strictly on a single node may become dependent on the validity of the memory of *many* nodes thereby inducing unnecessary transaction aborts. This happens when *database objects*, such as records, and *database support structures*, such as index structures and shared lock tables, are stored in shared memory.

In this paper, we propose crash recovery protocols for shared memory database systems which avoid the unnecessary transaction aborts induced by the failure dependencies that occur due to the use of shared physical memory. Our recovery protocols guarantee that if one or more nodes crash, all effects of active transactions running on the crashed nodes will be undone, and no effects of transactions running on nodes which did not crash will be undone. In order to show the practicality of our protocols, we discuss how existing features of cache-coherent multiprocessors can be utilized to implement these recovery protocols. Specifically, we demonstrate that (1) for many types of database objects and support structures, volatile (in-memory) logging is sufficient to avoid unnecessary transaction aborts, and (2) a very low overhead implementation of this strategy can be achieved with existing multiprocessor features.

Contents

1	Introduction	1
2	System and Transaction Model	2
3	Failure Dependencies Induced by Cache Coherency	3
4	Transaction Recovery in Shared Memory Database Systems	5
4.1	Recovery of Database Objects	6
4.2	Recovery of Database Management Structures	8
4.2.1	A Shared Memory Implementation of B-trees	9
4.2.2	A Shared Memory Implementation of Locking	11
5	Implementing LBM Policies in Shared Memory Systems	13
5.1	Enforcing the Volatile LBM Policy	14
5.2	Enforcing the Stable LBM Policy	15
6	Integrating the Recovery Protocol with other Transaction Processing Mechanisms	17
7	Related Work	19
8	Summary and Conclusions	20

1 Introduction

Shared memory systems offer significant performance advantages for applications which share data. But, when we consider how applications that have failure resilience requirements can benefit from cache-coherent shared memory (SM) systems, we find that current crash recovery mechanisms are insufficient to support these applications. With current mechanisms, a single node (where a *node* refers to a processor/memory pair) crash is likely to require a reboot of the entire shared memory system. Although the crash of a single node in an SM system should be infrequent, in very large systems if one node crash implies system failure, then the system will be down most of the time. Clearly, it will be beneficial to have *independent node failures* whereby the crash of one node not cause the failure of other nodes.

One class of applications which will benefit immensely from the provision of *independent node failures* is database applications. Primarily because of the absence of support for independent node failures, shared memory database implementations suffer from recovery problems (cf. Sigmod94 panel). This is unfortunate because database applications can easily exploit the performance advantages of shared memory architectures.

While the support for independent node failures in an SM system requires certain low-level mechanisms to be in place, such as mechanisms to detect and isolate hardware faults [11], these alone are not enough to guarantee that unnecessary failure dependencies do not form between nodes. Specifically, when we utilize the performance features of an SM database system, two or more *independent* transactions executing on different nodes may incur failure dependencies due to the side effects of the cache coherency protocol. When a node crash occurs (whereby the contents of the physical memory on the failed node are destroyed), the resolution of these system induced failure dependencies may require aborting otherwise independent transactions which execute on *other* nodes. These unnecessary transaction aborts must be avoided by appropriate design of recovery protocols. This is our goal.

Our recovery protocols guarantee the failure atomicity¹ of transactions, yet avoid unnecessary transaction aborts. Specifically, our recovery protocols guarantee that if one or more nodes crash, all effects of active transactions running on the crashed nodes will be undone, and no effects of transactions running on nodes which did not crash will be undone. The key features of these recovery protocols are (a) *logging-before-migration* (LBM) policies, which enforce specific

¹Failure atomicity of a transaction, also called the *all-or-nothing* property, means that either all or none of the transaction's operations are performed even when there are failures.

logging policies prior to the migration of uncommitted data from one node to another, and (b) pairing the migration of (a small amount of) information required for undo with the migration of uncommitted updates. The *volatile* LBM policy helps achieve our recovery goals with very low runtime overheads and can be implemented with existing features of cache-coherent multiprocessors.

This paper is structured as follows. Section 2 discusses our system and transaction model, and section 3 discusses how cache coherency can complicate recovery. In section 4, we present our recovery protocols and discuss how these protocols can be applied to database objects, such as records, and database support structures, such as index structures and lock tables. Techniques for implementing the LBM policies on a cache-coherent shared memory multiprocessor are discussed in section 5. Issues which arise when integrating our recovery protocols with other transaction processing components are discussed in section 6. Related work is discussed in section 7 and our conclusions are presented in section 8.

2 System and Transaction Model

In an SM database system, each node is connected to all disks in the system, and each node has access to shared memory. We consider SM database implementations where the shared memory is made coherent using a hardware-based cache coherency protocol. Hardware-based cache coherency provides low latency, high bandwidth access to shared memory, and ensures that each node reads the most recent version of the data in the shared address space. By utilizing these features of an SM system, it is possible to construct high performance multiple node database systems.

The coherency protocol, implemented in hardware, ensures that any read operation sees the most recently written value for any data item. Each node has its own cache, and before an operation is performed on a data item, the data item must first be brought into the cache. In general, operation execution time is minimal if the data item is already in the cache, more expensive if the data item is in another node's cache, and the most expensive if the data item must be fetched from disk. Typically, the hardware elements implementing the cache coherency scheme include a cache controller, a cache directory, and the cache itself. The cache contains the cached data, while the cache directory contains the addresses of all cached data. Our discussions in this paper assume a *write-invalidate* cache coherency protocol [1, 12] where, before a write to a cache line by one node occurs, all other cached copies of the line are first invalidated².

²Our results also apply to a *write-broadcast* cache coherency protocol, but a discussion of this has been omitted due to space limitations.

However, a cache line could be valid in multiple nodes after a series of read requests to that line have been issued. When a cache line l is resident on node x and another node y desires to update line l , l *migrates* from node x to node y . While the unit of I/O is a page, the unit of coherency is a *cache line*, and is typically smaller than a page.

Motivated by the proposed enhancements to SM architectures whereby individual node failures are isolated [11], we assume that nodes in the shared memory multiprocessor fail independently. When node x crashes while holding the only copy of a cache line, we assume that this cache line is unavailable for use by other nodes until a recovery procedure is executed. Node crashes which do not destroy the only copy of a cache line do not render these cache lines unavailable. The recovery procedure is used to reconstruct some or all of the dirty data of cache x . We do not distinguish between a node’s cache memory and other volatile memory, since, if a node crashes, the contents of all volatile physical memory are destroyed.

Consistent with most commercial database implementations, we assume that:

- Locking-based concurrency control is used. The basic lock modes are shared and exclusive. An exclusive lock on a record r guarantees that no other transaction will read or modify r , while a shared lock on r ensures that no other transaction will modify r . Note that several shared requests on r can be granted concurrently.
- No-force/steal buffer management policy is used [7].
- Each node maintains a log, and in-place updating is used in conjunction with the write-ahead log protocol (WAL) [2]. All update operations to this log take place in the node’s cache. This in-cache log is volatile, but can be made stable by writing it to one of the shared disks. Since updates to the log are performed only by the local system, with proper memory alignment (i.e. a cache line which contains local log information stores no other sharable information), we can ensure that local logs do not migrate between systems.

We focus on transaction workloads where independent transactions execute entirely on a single node. Although an SM database system is well suited for applications where a single (parallelized) transaction executes on multiple nodes, the presentation of the recovery strategies is simplified under the assumption that each transaction executes on a single node. However, our results easily generalize to address transactions which execute on multiple nodes.

3 Failure Dependencies Induced by Cache Coherency

Any data structure implemented in shared memory can be adversely affected by the failure dependencies induced by a cache-coherency protocol. In database systems, these data structures

include *database objects*, such as records, and *database support structures*, such as index structures and shared tables. Consider the storage requirements of records. Due to typical cache line sizes of current shared memory multiprocessors, it is likely (unless a lot of space is wasted) that multiple records will be stored in a cache line. For example, on the KSR-1 and KSR-2 multiprocessors [21], and on Stanford's FLASH [11] distributed shared memory machine, the cache line size is 128 bytes. Under a write-invalidate cache-coherency protocol, when two or more records stored in the same cache line are concurrently updated by transactions running on different nodes, the node where the last update occurred will hold the *only copy* of these records. Due to this access pattern, two types of failure dependencies may form: First, the crash of one node may not result in the complete annulment of transactions which execute strictly on the failed node. Second, the crash of one node may destroy updates performed by transactions which execute on other nodes.

The following example illustrates one case where a transaction executing on one node may become dependent on the validity of memory on some other node. Transaction t_x (executing entirely on node x) locks, then updates record $r1$. Then, transaction t_y (executing on node y) locks, then updates record $r2$. As a side effect of the write-invalidate cache coherency protocol, *the only copy* of cache line l now resides on node y . If node x crashes, the control state (registers, stack, etc.) of t_x will be destroyed; but since l migrated to node y , the uncommitted update to record $r1$ will remain intact. On the other hand, if node y crashes, l will be destroyed and even though the node t_x is executing on has not failed, the uncommitted update to $r1$ performed by t_x will be destroyed. In these two cases, failure dependencies have formed between transaction t_x and node y , and between transaction t_y and node x . Thus, if either node x or node y crashes, the failure atomicity of transactions may be compromised.

These failure dependencies can also form on database support structures, such as database indexes, or tables used by lock managers. Consider a shared memory implementation of a lock manager, where, in order to properly acquire and release locks, each node examines and updates information stored in shared memory. In this situation, multiple entries which describe the acquisition of database locks may be stored in a single cache line. For example, suppose two active transactions, running on different nodes, have acquired a lock on the same record (in shared mode). Further suppose that the lock control block, the shared data structure which contains the information describing the two holders of this lock, is stored in one cache line. The last node to acquire this record lock will hold the only copy of the lock control block. Without sufficient recovery provisions, a crash of this node will result in the loss of some of the information about locks granted to transactions running on other nodes. This can also lead to a violation of the failure atomicity of transactions.

One way to ensure the failure atomicity of transactions in these cases would be to abort all transactions which had formed failure dependencies on the memory of remote nodes. But this method is overkill, since unnecessary transaction aborts are incurred. Instead, our recovery protocols guarantee *Isolated Failure Atomicity* (IFA), by ensuring the failure atomicity of transactions, yet avoiding unnecessary transaction aborts. IFA ensures that if one or more nodes crash, *all* effects of active transactions running on the crashed nodes will be undone, and *no* effects of transactions running on nodes which did not crash will be undone.

4 Transaction Recovery in Shared Memory Database Systems

In this section, we propose SM crash recovery protocols which ensure IFA. While these recovery protocols can make use of stable storage to log recovery information, they incur substantially lower overheads when volatile logging is done. This is an important efficiency consideration, since, in most cases, utilizing stable storage requires use of the more expensive I/O operations.

Recovery protocols are generally comprised of two main mechanisms. The first takes the necessary steps to facilitate recovery at some future point in time. It ensures that sufficient information is available for the recovery procedure to re-establish a consistent database state. The second undertakes the actual recovery when needed, i.e., it restores the database to a consistent state after a node crash. In the interest of space, we focus on the first mechanism. It is made up of LBM policies, which perform *Logging Before Migration*, and a mechanism which pairs the migration of undo information with the migration of uncommitted updates. Prior to the migration of uncommitted data, LBM policies log sufficient information to allow the recovery procedure to ensure the failure atomicity of transactions without unnecessarily aborting transactions. *Undo* information is logged in order to ensure that if a node crashes, *all* effects of transactions running on the crashed node will be undone. *Redo* information is logged in order to ensure that if a node crashes, *no* effects of transactions running on nodes which did not crash will be lost.

We examine two different LBM policies, one based on volatile logging, and the other based on stable logging. Volatile logging is implemented by logging information into (volatile) memory, while stable logging is implemented by first logging into memory, then flushing the log to stable storage. While stable logging provides support for both undo and redo, volatile logging provides support only for redo. However, by augmenting volatile logging with a mechanism for tagging the uncommitted data with a small amount of information needed for undo, the undo requirements can be met without the use of stable logging. In section 4.1, we discuss the details of providing recovery support to guarantee IFA for database objects, while section 4.2 deals with ensuring IFA for database management structures.

4.1 Recovery of Database Objects

We assume that the possible operations on database objects (records) are *read* and *write*. Our recovery protocols assume many of the mechanisms used in existing commercial database systems, including certain logging techniques, in-place updating, the WAL protocol, and the use of *strict* 2PL to ensure serializability. By utilizing mechanisms and structures which are already part of the database management system, the incremental overheads associated with adopting our recovery protocols are minimized. For example, our LBM policies exploit the existence of undo and redo log records which are part of most database systems [7]. When a database record is updated, a redo log record (containing the value of the updated database record) is written to the volatile log. When a transaction first performs an update to a database record, an undo log record is written to the volatile log. This undo log record contains the before image (the last committed value) of the database record. To ensure the WAL protocol, prior to updating the *disk version* of a database record, the associated undo log record must be forced to stable store.

In the context of the WAL protocol, the assumption of *strict* 2PL allows transaction aborts to be implemented by simply replacing all the data touched by a transaction with their before images. Under strict 2PL, record locks are not released until after a transaction either commits or aborts. These protocols simplify recovery by ensuring that, at any time, only one transaction is associated with a particular uncommitted data item, and all before images exist in stable store. These assumptions also allow simple extensions to our LBM policies to implement transaction undos.

As we have shown earlier, when uncommitted data migrates between nodes, failure dependencies form between transactions running on one node and the volatile memory of another node. Thus, our discussions focus on recovery provisions for uncommitted transactions – the migration of committed data does not pose a recovery problem since the commit process ensures the durability of this data. Uncommitted data may migrate between nodes when two or more records are stored in the same cache line, and transactions executing on different nodes concurrently update these different records.

As mentioned earlier, this introduces two problems for recovery in a cache-coherent shared memory system. First, the crash of one node may not result in the complete annulment of transactions which execute strictly on the failed node. Second, the crash of one node may destroy updates performed by transactions which execute on other nodes. Consider any transaction t_{active} , which was active at the time one or more nodes crash: With respect to database objects, in order to ensure IFA, our recovery protocols guarantee that if one or more nodes crash, for all t_{active} ,

1. If t_{active} was running on a node that crashed, then all its updates must be undone.

Thus, in case some of these updates migrated to other nodes, sufficient information must be available to undo them.

2. If t_{active} was running on a node that did not crash, then none of its updates must be undone.

Thus, in case some of these updates migrated to a node that crashed, sufficient information must be available to redo them.

To provide this guarantee, we can employ either the volatile or stable LBM policies. For clarity of the presentation, we examine a scenario when *one* transaction becomes dependent on the memory of *one* remote node. Consider active transaction t_x with the *only copy* of a record (r) updated by t_x having migrated to another node (y). When we consider how the failure atomicity of a transaction may be compromised, we must consider two basic crash scenarios:

1. x , the node executing t_x crashes.
2. y crashes.

These crash scenarios necessitate the undo and redo requirements imposed above. For case 1, the update performed by t_x on r (which now resides on node y) must be undone. For case 2, the update performed by t_x to r must be redone, since node y crashed, destroying the contents of r .

Under the stable LBM policy, prior to the migration of record r , undo and redo log records for r are written to stable storage. For case 1, the stable undo log records are used by the restart recovery procedure to undo the update to r by t_x , ensuring the failure atomicity of t_x . For case 2, the stable redo log records are used by the restart recovery procedure to redo the update to r by t_x . Thus, under the stable LBM policy, both undo and redo information are durably maintained prior to cache line migration to ensure the failure atomicity of transactions – without incurring unnecessary transaction aborts – when nodes fail. The obvious disadvantage of this approach is that the runtime overheads of stable logging are very high (unless non-volatile RAM is assumed). These runtime overheads can be reduced by adopting the volatile LBM policy.

The volatile LBM policy is sufficient to ensure the *redo* requirements of our recovery protocol, since, if some node other than the one executing the transaction in question crashes, that transaction’s volatile log remains intact.

However, by itself, the volatile LBM policy does not provide support for undo. For example, in case 1, when the node executing t_x crashes, it could easily be the case that the transaction

management system *left no trace* of ever running t_x (as would be the case if neither stable logging or checkpointing had occurred after the start of t_x). Specifically, note that if the node executing t_x crashes, the volatile log which recorded the undo information corresponding to t_x 's update to any records which may have migrated is also destroyed.

Since the volatile LBM policy provides no guarantee that any part of a node's log will be stable while a transaction is active, we need an additional mechanism to identify the active records, i.e., records updated by active transactions, which require undo. For these records, which had been updated by t_x and which subsequently migrated to another node, the abort can be performed by installing each record's before image. A simple way to identify these active lines is to associate a node identifier with each data object. The node ID is stored in the *same cache line* as the active data object:

Tagging Rule: If multiple database objects are stored in a cache line, each active object is tagged with a node identifier. The node identifier indicates which node was executing the transaction that had updated the object. In the event of a crash of node x , all objects which are in the cache and tagged with node x are candidates for undo.

Note that because of the strict 2PL assumption, only one node will be associated with each active record. Once the data is no longer active, the node ID is assigned a null value. Thus, the node ID enables active cache lines (those cache lines which contain active data) to be identified for undo purposes. Of course, this approach requires additional space in order to maintain the per object node ID.

We have shown how (a) the stable LBM policy, or (b) the volatile LBM policy combined with record tags, can ensure IFA for database objects. These policies ensure that under any crash scenario, when a transaction becomes dependent on the memory of any other node, sufficient information will be available for the restart recovery procedure to mask this dependency. In order to guarantee IFA for (independent) transactions, we must also consider the failure dependencies which may form on database management structures, such as index structures and those related to database locking. This is the topic of the next subsection.

4.2 Recovery of Database Management Structures

By implementing database management structures in shared memory, the performance advantages of a shared memory multiprocessor can be fully exploited. Examples of these database management structures include hash tables, index structures such as B-trees, and tables used for lock management. In this section we examine the last two in detail.

For many transaction management issues, it is important to distinguish between structural and non-structural changes to database management structures. Examples of structural changes include B-tree page splits and the dynamic allocation of space used to store lock management information. Even in a multi-programmed uniprocessor database implementation, the subsequent use of this newly altered or created space by other transactions can cause transaction dependencies to form. To avoid this, when one transaction performs a structural change, it is customary to commit these changes immediately (regardless of the future commit or abort of the transaction that caused the change), prior to being released to other transactions [16, 13, 14]. For these reasons, we also assume this convention of the independent commit of structural changes done to database management structures.

Given this assumption, structural changes will not result in dependencies between active transactions and the memory of some other node. Thus, for structural changes, no additional recovery provisions are necessary in a SM implementation. However, for non-structural changes, uncommitted data migration can occur, potentially violating IFA. We now focus on applying our recovery protocols to database management structures for which non-structural changes are likely.

4.2.1 A Shared Memory Implementation of B-trees

Consider a shared memory implementation of an index with the B+-tree, where records are stored only in leaf nodes. Examples of non-structural B+-tree updates include the insertion and deletion of records. Consistent with existing commercial database implementations, we assume that strict write locks are obtained on the record to be inserted or deleted. We also assume that the associated logical log records are written for insert and delete operations, enabling the logical undo and redo of these operations [16, 15]. In conjunction with the WAL protocol, these assumptions ensure that each B+-tree update is associated with at most one transaction, simplifying the implementation of transaction aborts.

When the B+-tree index is implemented in shared memory, non-structural updates such as insert and delete can trigger the migration of uncommitted data between nodes. For example, one transaction may perform an insert operation, in the process adding a record to an existing leaf page p , and installing a pointer to this leaf page in an internal node (page) p' . If other records may also be stored in the same cache line l as where the newly inserted record is stored, it is possible for the newly inserted (uncommitted) record to migrate to some other node. Note that this migration may be triggered by a transaction running on some other node which either updates an existing record stored in l , or inserts or deletes a record in l . In this case, just as for

record updates, failure dependencies may form between an independent transaction executing entirely on one node and the memory of some other node. Likewise, if multiple pointers (to leaf pages) are stored in the same cache line, pointers to newly inserted (uncommitted) records may also migrate to some other node, again forming failure dependencies.

For such updates to B+-trees, these failure dependencies can be avoided by employing our recovery techniques. As with record updates, stable LBM alone will ensure that sufficient undo and redo information is available for the restart recovery procedure to ensure IFA for non-structural updates to B+-trees. Likewise, IFA for non-structural updates to B+-trees can be ensured with volatile LBM combined with tagging. Since the details of the stable LBM and the redo mechanism under volatile LBM are similar to those for record updates, next we discuss the support for undo under volatile LBM. Recall that this undo requirement arises when the node a transaction is executing on crashes, but some of the uncommitted data of this transaction has migrated to another node.

Just as for record updates, to enable the undo of active insert and delete operations, a node identifier can be tagged to each active record. If a node crash does not result in the complete annulment of a transaction, the recovery procedure can identify all cache lines which need undo based on the node identifier. However, since the steps required to perform an undo of an insert or delete operation are different than those required to undo a record update, an *undo operation code* is also stored in the same cache line as the record. During crash recovery, once a cache line which requires undo is identified, the appropriate undo operation is performed based on the undo operation code. Since the undo operation code need only distinguish between record updates, inserts, and deletes, two bits are sufficient to store this operation code.

Issues related to space management require that a subtle implementation of record delete operations be employed. To ensure that the space freed by a delete is not used until the transaction which performed the delete commits, it is customary to perform record deletes logically, by marking the record as deleted [14]. Once the transaction which performed the delete commits, the space freed by the deleted record can be used by other transactions. This strategy also enables an efficient implementation of the undo requirement for volatile LBM – since any migrating cache line which contains an uncommitted delete will also contain the original record, the undo of a delete is effected by merely “unmarking” this record.

No such special provisions need be made for space management for the undo of an insert, since allocated space can always be freed.

4.2.2 A Shared Memory Implementation of Locking

Next, we consider how a shared memory implementation of database locking may benefit from our recovery protocols. For a lock table implemented in shared memory, almost all operations on a lock table are non-structural (space allocation operations are the exception). Because of the likelihood of many non-structural operations to a lock table during transaction execution, to guarantee IFA for transactions, it is important to apply our recovery protocols to this database management structure.

One strategy for implementing a lock manager in a multi-node system is to designate some node as being responsible for managing each database object, and allow remote nodes to access locks by using message passing. This is the approach of many shared-disk (SD) systems [22, 17, 19]. The presence of shared memory in an SM database system allows a more efficient approach to be taken for the implementation of database locking. In this strategy, which we call *SM locking*, LCB's (lock control blocks) are stored in shared memory, and transactions acquire and release locks via operations on these LCB's. The performance gains of SM locking stem from the elimination of all inter-process communication [18].

Consider acquiring a record lock using SM locking. A lock request consists of a lock *name* and a lock *mode*. Using a hash function, the name is translated to an LCB address specific to one lock. An LCB stores the current mode of the lock, plus two transaction lists, one containing the current holder(s) of the lock, the other containing any transaction(s) waiting for the lock. All updates to the LCB are performed inside a critical section. If the requested mode is compatible with the mode stored in the LCB, and there are no conflicting waiters, an entry containing the requesting transaction and requested mode is added to the holder list, and the lock is granted. This entry is called the *lock acquisition record*. Otherwise an entry is added to the wait list, and a not-granted flag is returned to the requestor. The strategy for releasing a lock is similar. After finding the appropriate LCB, the tuple identified by the transaction is deleted from the holder list, and any lock requests in the wait list which become compatible due to the release are granted.

When lock information pertaining to two or more transactions is stored in a single cache line, recovery issues similar to those for record updates arise. For example, after two transactions running on different nodes have acquired a compatible lock, the LCB will be valid at the node which last acquired the lock. In this case, a node crash may lose some but not all of a transaction's lock information. Note that this scenario is only applicable to uncommitted transactions, since committed transactions have no effect on the lock space (once a transaction

has committed, all its locks are released)³. In contrast, each lock acquired by an uncommitted transaction will have a corresponding entry in the lock space.

Next, we consider recovery issues for acquired locks⁴. Consider any transaction t_{active} , which was active at the time that one or more nodes crash. To ensure IFA for SM locking, we guarantee that if one or more nodes crash, for all t_{active} ,

1. All locks acquired by transaction t_{active} running on a node which crashes will be released.
2. No locks acquired by transaction t_{active} running on a node which did not crash will be released.

Because of (1), any lock acquired by t_{active} running on a node which *had crashed* and stored in LCB's which *survived* the crash must be released by the restart recovery procedure. Because of (2), any lock acquired by t_{active} running on a node which *did not crash* and stored in LCB's for which *no copy survived* the crash must be restored by the restart recovery procedure. As with database objects, guaranteeing condition 1 requires undo information to be maintained, while guaranteeing condition 2 requires redo information to be maintained.

Let us first consider the guarantee of condition 2, since either stable or volatile LBM can be used to provide this guarantee. To guarantee condition 2, logged redo information is used. Suppose a node crash destroys the only copy of a lock acquisition record pertaining to a transaction running on a node which did not crash. During restart, the recovery procedure will restore these lock acquisitions as the LCB's are being reconstructed. Since these locks were acquired by nodes which did not crash, the log records contained in the volatile logs of these nodes are used by the restart recovery procedure. Thus, either the volatile and stable LBM policy is sufficient to ensure condition 2.

Consider how the LBM policies can ensure condition 1 in the case where node x crashes. Prior to acquiring (or releasing) a lock on node x , a logical log record [6] is written to the log on node x . Storing the transaction ID in the LCB is sufficient to guarantee condition 1 with the stable LBM policy. Since the stable portion of the log on node x will survive the crash of node x , the restart recovery procedure will be able to determine which locks were held by active transactions running on node x . Any LCB's, unaffected by the crash of node x , which contain lock acquisition entries for transactions running on node x will be repaired by the restart

³An exception to this rule occurs in some SD systems [17] where since lock acquisition is expensive in SD systems, locks are sometimes *retained* on local nodes after the transaction commit has occurred.

⁴The recovery issues for transactions which are waiting for locks are similar, but the discussion has been omitted due to space limitations.

recovery procedure. This repair involves releasing any locks acquired by transactions on node x .

As with access to database records, we can ensure that the undo requirement (condition 1) is met with volatile LBM if, as before, the record is tagged. In this case, the node ID is tagged to each lock acquisition record of the LCB. Note that, if the transaction ID encodes the node ID, then storing an *additional* node ID tag would not be necessary. For example, suppose transaction t_x (running on node x) acquires a lock, the LCB containing the lock acquisition information migrates to another node, then node x crashes. This lock acquisition record can be identified by its node ID of x available in the migrated lock acquisition record. This information is used by the restart recovery procedure to ensure that all locks which were acquired by transactions running on nodes which had crashed are released.

Some issues related to ensuring the failure atomicity of database management structures are covered in [23, 22, 10], where crash recovery issues for an SD lock manager implementation on a VAXcluster are discussed. When a node crash is detected, all locking activity in the database system is stopped. Then, any updates performed by failed transactions are undone. This is accomplished by the installation of the before images of the associated records. After this restart recovery procedure is complete, user activity may proceed.

In contrast, in a cache-coherent SM database system, if certain mechanisms, discussed presently, are available, locking activity does not need to be stopped when a node crash is detected. In an implementation of SM locking, problems of ensuring a consistent lock space may arise if a node holding the only copy of a LCB crashes, but other nodes, not detecting the existence of this LCB, create a new LCB and incorrectly release an acquired lock. This will not be a problem if the underlying hardware support of the SM multiprocessor ensures that (just the) references made to cache lines residing on crashed nodes are stalled.

This section showed how our recovery protocols can be applied both to database objects and database management structures to achieve independent failure atomicity (IFA) for transactions. In the next section, we discuss how a low overhead implementation of our recovery protocols can be achieved in a cache-coherent SM multiprocessor.

5 Implementing LBM Policies in Shared Memory Systems

The LBM policies require that prior to the migration of a cache line, either stable logging or volatile logging is performed. Here, we discuss the implementation of these Logging Before Migration policies on a shared memory multiprocessor. We will show that sufficient primitives are already available on existing multiprocessor hardware to efficiently enforce the volatile LBM

policy, but not the stable LBM policy.

To enforce the volatile and stable LBM policies, it is sufficient to construct the appropriate log record at any point after line l is updated and before l migrates. For a number of reasons, it is best to perform volatile logging immediately after an update is performed. This is both a logical and efficient point to perform volatile logging, since at the time of an update most of the relevant information is already cached locally, and performing a few additional local memory references to write the log record minimizes the additional overheads. However, to reduce the overheads of stable logging, it is wise to minimize the frequency of log forces. Thus, while it is best to enforce volatile LBM immediately, it is best to delay enforcing stable LBM as long as possible. We discuss the enforcement of the volatile LBM policy in section 5.1, and the stable LBM policy in section 5.2.

5.1 Enforcing the Volatile LBM Policy

A (cache) *line lock* [21], commercially available on the KSR-1 multiprocessor, is an example of a mechanism useful for efficiently implementing critical sections in a cache-coherent multiprocessor. Thus, they can be utilized to achieve a very low overhead implementation of the volatile LBM policy. Once a line lock is acquired on cache line l by a process running on node x , the underlying hardware ensures the following properties:

- Line l is held *exclusively* in cache x .
- No other process, whether it be on the same or a different node, can read or write to l until l is explicitly released by the holding process.

The `getline(l)` instruction is used to obtain and hold a cache line l in a mutually exclusive (ME) state. The semantics of the `getline` primitive are such that, if the cache line is not already in ME state in any cache, the local cache acquires it in ME state. The `releaseline(l)` primitive releases the cache line from ME state⁵. The advantages of the line lock are that (a) locking and unlocking the line each requires only a single instruction, and (b) in the process of locking the line, the line also becomes resident in the local cache.

The volatile LBM policy can be efficiently enforced with the use of the line lock as follows. Whenever an update to a database object or database support structure is performed, a log record is written describing this update. To ensure that a cache line l does not migrate between

⁵On the KSR-1, these primitives are called `gsp` and `rsp`, get subpage and release subpage. We have renamed these primitives to be consistent with the literature on cache coherency.

the time it is updated and the time the log record is written, a line lock can be used. Thus, prior to performing an update to data stored in cache line l , `getline(l)` is issued to lock the line in cache. The update is performed, and the log record is written prior to releasing the line with `releaseline(l)`.

Our experience in the implementation and empirical performance evaluation of database mechanisms on the KSR-1 confirms the expected performance gains provided by the line lock primitive. In [18], we implemented a prototype lock manager, using the line lock to ensure mutually exclusive updates to the shared memory implementation of the lock space. Our empirical performance studies have shown that under low contention, the mean execution time to obtain a line lock is less than 10 μ s, and under high contention (32 processors simultaneously attempting to acquire the *same* line), the mean execution time to obtain a line lock is less than 40 μ s.

5.2 Enforcing the Stable LBM Policy

One approach to enforcing the stable LBM rule would be to flush the log as part of the update protocol. In this solution, line locks would be retained on the updated cache lines until the update is performed, the log record is written, and the log force is completed. Although this solution would guarantee the stable logged rule, it is also very inefficient – a log flush is performed on each update, regardless of whether the cache line ever migrates. In order to guarantee the stable LBM rule and minimize the frequency of log flushes, we must address the following:

- What is the latest point in a cache line use history where the log must be flushed?
- What are the appropriate enforcement mechanisms?

We first consider how read and write operations on cache lines affect the state transitions of the cache coherency protocol. Recall that cache line migration occurs when the state of a cache line changes from exclusive in one cache to exclusive in another cache. As the following histories indicate, cache line migration may occur *directly*, or *indirectly*:

$$\begin{aligned}
 H_{direct} &= w_x[l]; w_y[l]; \\
 H_{indirect} &= w_x[l]; (r_x[l];)^* r_x[l]; (r_z[l];)^* w_y[l];
 \end{aligned}$$

In H_{direct} , a cache line l migrates directly from node x to node y . In this case, a write by node x to line l ($w_x[l]$) occurs, causing line l to be valid only in node x 's cache (recall our assumption of a write-invalidate cache coherency protocol). The next operation issued on line l is a write by node y , causing l to be invalidated in all other caches (x), and held exclusively in node y 's cache.

$H_{indirect}$ shows how intermediate read operations can cause the state of l to be held in shared mode (in potentially many caches) between the $w_x[l]$ and $w_y[l]$ operations. This shared state of l may occur after $w_x[l]$ occurs, zero or more reads on l are issued by node x , one read is issued by some node other than x (\bar{x}), and zero or more read operations are issued by any other node (as indicated in the above grammar).

In order to minimize the frequency of log flushes, we would like to determine the latest point at which it is necessary to force the undo and redo logs. In H_{direct} , this point would be immediately prior to $w_y[l]$, since this is the operation which causes the transition from exclusive in cache x to exclusive in cache y . The analysis is not so straightforward for $H_{indirect}$: Following $w_x[l]$, any number of reads by node x to line l ($(r_x[l];)^*$) will not change the state of line l . However, the next read by some node other than x ($r_{\bar{x}}[l]$) will downgrade l from exclusive to shared mode on node x , allowing node \bar{x} to also hold l in shared mode. Although not pointed out in our earlier examples, this pattern of sharing poses problems for recovery. If node x crashes after $w_x[l]; r_{\bar{x}}[l]$, the crash recovery procedure would require undo information to complete the abort of active transactions that were running on node x . Clearly, to permit this, at least the undo portion of the log must be forced prior to $r_{\bar{x}}[l]$.

Thus, for a line l which has been updated by some node x , the latest point at which the stable LBM policies must be enforced corresponds to the downgrade or invalidation of l (for undo) and the invalidation of l (for redo). By triggering log flushes based on these cache line state changes, the number of log flushes can be minimized. This log flush would only be done if the cache line contains database related information for which the corresponding log records had not been flushed to stable store.

Unfortunately, triggers associated with the change of a cache line state are not a feature of any commercial multiprocessor that we know of. This extension to the cache coherency protocol can be implemented by dedicating one bit per cache line to indicate whether the line contains active data. Updates to the cache line would set this bit, and log flushes would clear the bits of all associated cache lines.

Summarizing this section, it is clear that even though the stable LBM policies are simpler to explain, they are more expensive to realize and in fact are not implementable with the features in today's multiprocessors. Volatile LBM policies, on the other hand, lend themselves to fairly efficient implementations with very little additional demands being placed on multiprocessors. The only requirement they have is a small amount of additional space for tagging active data.

6 Integrating the Recovery Protocol with other Transaction Processing Mechanisms

In this section we consider the important issues involved when our recovery protocols are integrated with other salient components of a database system. Although many of the implementation issues of interest to us have been addressed in the context of SD systems [17, 15, 23, 22, 19, 20], significant differences between SD and SM systems require that, for the most efficient implementation, different mechanisms be developed for SM. These differences stem from the different approaches used to achieving coherency. In SM, cache coherency is achieved transparently by the underlying hardware cache coherency protocol. In contrast, in an SD system, coherency is achieved entirely in software and is closely coupled with the lock and buffer managers [17, 19]. Two significant implications of these factors are (1) an SM database need not include the SD mechanisms used for ensuring coherency, and (2) an SM database can utilize the low latency access to shared memory to yield very efficient adaptations of other SD mechanisms.

In the rest of this section, we discuss how these implications affect the transaction processing components needed for SM systems in the context of our recovery mechanisms. We focus on the protocols and implementation mechanisms used in a recent SD system [17]. This system supports record level locking and uses in-place updating in conjunction with the WAL protocol. It uses the repeating history paradigm followed by undos to recover from failures. We discuss how WAL and the techniques for ensuring the repeating history paradigm can be efficiently implemented in SM, in light of the mechanisms for achieving LBM.

Since this system has also addressed some of the issues related to the migration of uncommitted data, we discuss this aspect of [17] first. In order to ensure inter-node coherency, [17] defines four inter-node page transfer schemes, two of which allow the migration of uncommitted data, called *fast*, and *super-fast* schemes. In the fast scheme, all updates to page p must be stable logged prior to p 's migration. In the super-fast scheme, all updates to page p must be volatile logged prior to p 's migration. In the super-fast scheme, because page updates are not necessarily stable logged prior to uncommitted data migration, enforcing the WAL protocol may require referencing merged log. To implement the WAL rule for SD, each updating system remembers an LSN (log sequence number) greater than or equal to its last update to page p [17]. Page p can be written to the StableDB only after all systems which have updated p have forced their logs up to this LSN.

Closely related to the page transfer schemes is the *ordered update logging* rule [17, 15], which is important for supporting the repeating history paradigm (even for a multiprogrammed uniprocessor database system). This rule guarantees that the order of logging of updates to a

page is the same as the order with which those updates are performed on the page. In [15], the *ordered update logging* rule is satisfied as part of the update protocol. This guarantee is provided by acquiring and holding a semaphore on the page to be updated for the duration of the update *and* the log write.

On the surface, there are many similarities between the SD protocols described in [17] and our recovery protocols. However, there are important differences. Many of the protocols of [17], such as the page transfer protocols, are designed for the purpose of enforcing inter-system page coherency. In contrast, for cache-coherent SM, our LBM policies are designed specifically to isolate the crash of one node from affecting transactions which execute on other nodes. Thus, whereas the protocols of [17] are aimed at achieving page coherency, we are motivated by the need to eliminate the ill-effects of system-ensured cache coherency! Furthermore, as was discussed in section 5, efficiently enforcing the LBM policies on a cache-coherent SM multiprocessor requires a careful analysis of the coherency protocol, and a novel application of multiprocessor features.

Finally, the availability of shared memory in SM systems allows for more efficient implementations of many transaction processing mechanisms. To illustrate this, next we show how shared memory can be utilized in the adaptation of two mechanisms used in the SD system of [17]: the ordered update rule and enforcing the WAL protocol under the volatile LBM policy.

In section 4, we discussed how line locks could be used to enforce the volatile LBM policy as part of the update protocol. The line lock mechanism can also be employed in SM to efficiently enforce the ordered update logging rule. Consider updating record r stored in page p , when it is also necessary to update the page's Page-LSN field⁶. Once a record lock is obtained on r , a line lock is acquired on (a) the cache line (by convention, the first cache line of page p) containing the Page-LSN of p , and on (b) the cache line containing the record (assuming these cache lines are different). Once these line locks are acquired, r and Page-LSN(p) are updated. Finally, the log record for the update is written and the two line locks are released. By using line locks instead of semaphores to enforce this protocol, runtime overheads, as measured in terms of the number of instructions executed, are substantially reduced.

To enforce WAL under volatile LBM, we can adopt the same bookkeeping technique as done in SD, but exploit the available shared memory to minimize the runtime overheads. Each updating system remembers an LSN equal to its last update to page p . Page p can be written to the StableDB only after all nodes which have updated p have forced their logs up to this LSN.

⁶Each database page has a Page-LSN field which contains the LSN of the log record that describes the latest update to that page. The Page-LSN is used during restart and media-recovery to determine which logged updates have been applied to the page.

The determination of whether any other node is required to force their log can be computed very fast by maintaining this table of page/LSNs in shared memory.

7 Related Work

Many studies have demonstrated the performance advantages of using SM implementation platforms for database systems. Based on a TP1 benchmark performed on a Sequent Symmetry shared memory multiprocessor, [25] conclude that an SM database system can deliver very high performance. In [26], an analytical and simulation study compares SN (shared nothing), SD, and SIM (shared intermediate memory⁷). This comparison concludes that the data sharing architectures, especially SIM, are more resilient to transaction load surges. In [3, 4], simulation studies compare SN, SD, and SM (called SE (shared everything) in this reference), and concludes that SM outperforms SN and SD by a fairly wide margin. Our work exploits the performance advantages of SM systems, yet guarantees good failure properties for transactions.

In the previous sections, we discussed how our work is related to work in shared disk systems [17, 23, 22, 19]. Architectural differences between SD and cache-coherent SM, such as the unit of inter-system data sharing, how coherency is achieved, and whether shared memory is available, have a significant influence on the design and implementation of SM crash recovery protocols. Furthermore, our goal was to achieve IFA with minimal extra overheads while capitalizing on features available on or proposed for SM multiprocessor systems.

In [20], augmenting an SD system with a non-volatile global extended memory (GEM) is considered. System performance can be improved by adding GEM to a system that otherwise can only communicate by message passing. Failure atomicity for data structures can be ensured by propagating their updates to the GEM. However, non-volatile memory is much more expensive than volatile memory, and is a significant departure from a database implementation on off-the-shelf shared memory multiprocessors wherein the cache – where (parts of) data structures may reside – is volatile.

Transactional Memory [8] is another approach to supporting transactions on a cache-coherent shared-memory architecture. Transactional memory allows programmers to define customized read-modify-write operations that apply to multiple, independently-chosen words of memory. However, this approach is intended to replace short critical sections, i.e., it works well for short transactions with relatively small data sets. Our recovery protocols do not have these restrictions.

⁷The SIM model is slightly different than the shared memory model. In SIM, a shared intermediate memory serves as a global shared buffer for all nodes.

Volatile logging has been used in the context of *process checkpointing* schemes [9, 24] to ensure a consistent state of a distributed computation without necessitating the rollback of any processes other than ones that failed. In process checkpointing, information is periodically checkpointed to disk in order to ensure forward progress of a computation in the event that a processor and its associated memory fail. In this case, a checkpoint consists of the necessary process state for restarting execution, such as the program counter, process identifier, and register contents. Here, messages sent between processes trigger log records to be written to volatile memory. Recovery of a failed process is achieved by restarting the failed process from its checkpoint and replaying the message from the sender’s logs.

8 Summary and Conclusions

In this paper, we have presented crash recovery protocols for shared memory database systems which avoid unnecessary transaction aborts. For independent transactions, our recovery protocols guarantee IFA – that is, if one or more nodes crash in a system that isolates individual failures, all effects of active transactions running on crashed nodes will be undone, and no effects of active transactions running on nodes which did not crash will be undone. By applying our recovery protocols to database objects and database support structures, IFA is ensured for transactions under any crash scenario.

This guarantee can be provided by using either the volatile or the stable LBM policy:

- The volatile LBM has lower runtime overheads, and we showed how it can be efficiently supported with existing multiprocessor primitives. If the volatile LBM policy is chosen, it must be combined with a mechanism which tags the active database objects (or entries in the database support structures) with a small amount of additional information needed to perform undo, such as the node identifier and an undo operation code.
- By itself, the stable LBM policy ensures that sufficient information will be available during recovery to ensure IFA. The frequency of log flushes can be reduced with an extension to the cache coherency protocol which performs a log flush triggered by the migration of a cache line l containing uncommitted data. The undo log is forced by the invalidation of a cache line, while the redo log is forced by the invalidation or downgrade of the cache line.

Figure 1 summarizes these differences. It also shows that with features available in today’s SM multiprocessors it is easier to achieve volatile LBM than stable LBM.

We have also demonstrated how these protocols can be integrated with well established database design and recovery principles, such as the use of in-place updating in conjunction with

	<i>Runtime Overheads</i>	<i>Requires Tagging</i>	<i>Requires Cache Coherency Extensions</i>
<i>Volatile LBM</i>	Low	Yes	No
<i>Stable LBM</i>	High	No	Yes

Figure 1: Comparison of LBM Policies

the WAL (write-ahead logging protocol), the flexible no-force/steal buffer management policies, fine-granularity locking, and the repeating history paradigm. By exploiting mechanisms and structures which are already part of many databases, the incremental overheads associated with adopting our recovery protocols are minimized.

The recovery protocols developed in this paper assume that only read/write operations are performed on database objects. We are currently working on the extensions required to accommodate arbitrary operations on (abstract data type) objects.

In addition to providing support for SM database systems, our recovery protocols can also be applied to provide operating system support for handling independent node failures. For example, in multiprocessor systems where a *single* virtual address space is employed [5, 21], virtual memory allocation is most efficiently implemented by maintaining operating system tables in stored shared memory. When the data structures describing the allocation of virtual memory migrate between systems, just as in our database scenarios, failure dependencies will form between nodes. Recovery techniques similar to ours can be applied to these operating system data structures in order to ensure that the crash of one node does not necessarily affect the integrity of the process management information on other nodes.

References

- [1] J. Archibald and J. Baer. Cache Coherence Protocols: Evaluation Using a Multiprocessor Simulation Model. *ACM Transactions on Computer Systems*, 4(4):273–298, November 1986.
- [2] P. A. Bernstein, V. Hadzilacos, and N. Goodman. *Concurrency Control and Recovery in Database Systems*. Addison-Wesley, Reading, MA, 1987.
- [3] A. Bhide. An Analysis of Three Transaction Processing Architectures. *Proceedings of the 14th International Conference on Very Large Data Bases*, 14:339–350, September 1988.
- [4] A. Bhide and M. Stonebraker. A Performance Comparison of Two Architectures for Fast Transaction Processing. *IEEE Proc. 4th Intl. Conference on Data Engineering*, pages 536–545, February 1988.
- [5] J. Chase, H. Levy, E. Lazowska, and M. Barker-Harvey. Lightweight Shared Objects in a 64-Bit Operating System. *Proceedings of the ACM Conference on Object-Oriented Programming Systems, Languages, and Applications*, pages 397–413, 1992.
- [6] J. Gray and A. Reuter. *Transaction Processing: Concepts and Techniques*. Morgan Kaufmann, 1993.
- [7] T. Haerder and A. Reuter. Principles of Transaction-Oriented Database Recovery. *ACM Computing Surveys*, 15(4):287–317, December 1983.
- [8] M. Herlihy and E. Moss. Transactional Memory: Architectural Support for Lock-free Data Structures. *Proceedings of the 20th International Symposium on Computer Architecture*, May 1993.
- [9] D. Johnson and W. Zwaenepoel. Sender-Based Message Logging. *Proceedings of the 17th International Symposium on Fault-Tolerant Computing*, pages 14–19, 1987.
- [10] N. Kronenberg, H. Levy, and W. Streker. Vaxclusters: A Closely-Coupled Distributed System. *ACM Transactions on Computer Systems*, 4(2):130–146, May 1986.
- [11] J. Kuskin, D. Ofelt, M. Heinrich, J. Heinlein, R. Simoni, K. Gharachorloo, J. Chapin, D. Nakahira, J. Baxter, M. Horowitz, A. Gupta, M. Rosenblum, and J. Hennessy. The Stanford FLASH Multiprocessor. *Proceedings of the 21th International Symposium on Computer Architecture*, pages 302–313, April 1994.

- [12] D. Lilja. Cache Coherence in Large-Scale Shared-Memory Multiprocessors: Issues and Comparisons. *ACM Computing Surveys*, 25(3):303–338, September 1993.
- [13] D. Lomet and B. Salzberg. Access Method Concurrency with Recovery. *Proceedings of the 1992 ACM SIGMOD International Conference on Management of Data*, 21:351–360, June 1992.
- [14] C. Mohan and D. Haderle. Algorithms for Flexible Space Management in Transaction Systems Supporting Fine-Granularity Locking. *Proc. International Conference on Extending Data Base Technology*, March 1994.
- [15] C. Mohan, D. Haderle, B. Lindsay, H. Pirahesh, and P. Schwarz. ARIES: A Transaction Recovery Method Supporting Fine-Granularity Locking and Partial Rollbacks Using Write-Ahead Logging. *ACM Transactions on Database Systems*, 17:94–162, March 1992.
- [16] C. Mohan and F. Levine. ARIES/IM: An Efficient and High Concurrency Index Management Method Using Write-Ahead Logging. *Proceedings of the 1992 ACM SIGMOD International Conference on Management of Data*, 21:371–380, June 1992.
- [17] C. Mohan and I. Narang. Recovery and Coherency-Control Protocols for Fast Intersystem Page Transfer and Fine-Granularity Locking in a Shared-Disk Transaction Environment. *Proceedings of the 17th International Conference on Very Large Data Bases*, 17:193–207, 1991.
- [18] L. D. Molesky and K. Ramamritham. Efficient Locking for Shared-Memory Database Systems. Technical Report 94–10, University of Massachusetts Dept. of Computer Science, February 1994.
- [19] E. Rahm. Concurrency and Coherency Control in Database Sharing Systems. *Technical Report, University of Kaiserslautern, Germany*, December 1991.
- [20] E. Rahm. Use of Global Extended Memory for Distributed Transaction Processing. *Proceedings of the 4th Int. Workshop on High Performance Transaction Systems, Asilomar, CA.*, September 1991.
- [21] Kendall Square Research. *KSR1 Principles of Operation*. KSR Research, Waltham, Mass., 1992.
- [22] W. Snaman and D. Thiel. The VAX/VMS Distributed Lock Manager. *Digital Technical Journal*, (5):29–44, September 1987.

- [23] T. Rengarajan P. Spiro and W. Wright. High Availability Mechanisms of VAX DBMS Software. *Digital Technical Journal*, (8):88–98, February 1989.
- [24] R. Strom, D. Bacon, and S. Yemini. Volatile Logging in n-Fault-Tolerant Distributed Systems. *Proceedings of the 18th International Symposium on Fault-Tolerant Computing*, pages 44–49, 1988.
- [25] S. Thakkar and M. Sweiger. Performance of an OLTP Application on Symmetry Multiprocessor System. *Proceedings of the 17th International Symposium on Computer Architecture*, pages 228–238, May 1990.
- [26] P. Yu and A. Dan. Performance Evaluation of Transaction Processing Coupling Architectures for Handling System Dynamics. *IEEE Transactions on Parallel and Distributed Systems*, 5(2):139–153, June 1994.