

Concurrency Control Issues in Nested Transactions

Theo Härder and Kurt Rothermel

Received October 26, 1990; revised version received May 16, 1991; accepted July 20, 1992.

Abstract. The concept of nested transactions offers more decomposable execution units and finer-grained control over concurrency and recovery than “flat” transactions. Furthermore, it supports the decomposition of a “unit of work” into subtasks and their appropriate distribution in a computer system as a prerequisite of intra-transaction parallelism. However, to exploit its full potential, suitable granules of concurrency control as well as access modes for shared data are necessary. In this article, we investigate various issues of concurrency control for nested transactions. First, the mechanisms for cooperation and communication within nested transactions should not impede parallel execution of transactions among parent and children or among siblings. Therefore, a model for nested transactions is proposed allowing for effective exploitation of intra-transaction parallelism. Starting with a set of basic locking rules, we introduce the concept of “downward inheritance of locks” to make data manipulated by a parent available to its children. To support supervised and restricted access, this concept is refined to “controlled downward inheritance.” The initial concurrency control scheme was based on S-X locks for “flat,” non-overlapping data objects. In order to adjust this scheme for practical applications, a set of concurrency control rules is derived for generalized lock modes described by a compatibility matrix. Also, these rules are combined with a hierarchical locking scheme to improve selective access to data granules of varying sizes. After having tied together both types of hierarchies (transaction and object), it can be shown how “controlled downward inheritance” for hierarchical objects is achieved in nested transactions. Finally, problems of deadlock detection and resolution in nested transactions are considered.

Key Words. Nested transactions, concurrency control, locking, object hierarchies.

1. Introduction

When multiple users access a database simultaneously, their data operations have to be coordinated in order to prevent incorrect results and to preserve the consistency

Theo Härder, Ph.D., is Professor of Computer Science, University of Kaiserslautern, D-6750 Kaiserslautern, Federal Republic of Germany. Kurt Rothermel, Ph.D., is Professor of Computer Science, University of Stuttgart, D-7000 Stuttgart 80, Federal Republic of Germany.

of the shared data. This activity is called concurrency control and should provide each concurrent user with the illusion that he is referencing a dedicated database. The classical transaction concept (Eswaran et al., 1976) defines a transaction as the unit of concurrency control, that is, the database management system (DBMS) has to guarantee *isolated execution* for an entire transaction. This implies that its results as derived in a multi-programming environment should be the same as if obtained in a serial execution schedule. Other important transaction properties are atomicity, consistency, and durability as defined by Härder and Reuter (1983). In a DBMS, the component responsible for achieving these properties is transaction management which includes concurrency control as a major function.

In current DBMSs, transaction management is typically designed with a single-level control structure; its implementation is optimized to execute short transactions with only a few data references (Anon et al., 1985). Two-phase locking is, by far, the most common method for controlling concurrency among transactions and has been accepted as a standard solution (Gray, 1978; Bernstein and Goodman, 1981). When running on a centralized DBMS, transaction granularity as well as locking protocols usually obtain satisfactory performance; for high-performance transaction systems, special concurrency control methods are considered to be mandatory in order to increase the level of parallelism (Reuter, 1982; Gawlick, 1985), however, without requiring changes to the transaction concept.

When executing more complex transactions involving, for example, sequences of joins and sort operations in a relational DBMS, it turns out that single-level transactions do not achieve optimal flexibility and performance. Especially in distributed systems, it is highly desirable to have more general control structures to support reliable and distributed computing more effectively. Major concerns are more decomposable and finer grained control of concurrency and recovery. As a solution to these problems, the concept of nested transactions was proposed by Moss (1985) where single-level transactions are enriched by an inner control structure. Such a mechanism allows for the dynamic decomposition of a transaction into a hierarchy of subtransactions, thereby preserving all properties of a transaction as a unit and assuring *atomicity* and *isolated execution* for every individual subtransaction. As a consequence, subtransactions may be distributed in a system among various (processor) nodes performing subtasks of the entire transaction. These prime aspects of nested transactions—*decomposition of a "unit of work" into subtasks and their distribution*—lead to the following advantages in a computing system and, in particular, in a distributed DBMS.

1.1 Intra-transaction Parallelism

The larger a transaction, the more inherent parallelism may be anticipated during its execution. To take advantage of this inherent concurrency in the application, suitable granules of concurrency control as well as access modes (e.g., locking modes) are necessary. In environments enabling parallel execution, the nested-transaction concept embodies an appropriate control structure to support supervised, and therefore safe, intra-transaction concurrency, thereby increasing efficiency and decreasing response time.

1.2 Intra-transaction Recovery Control

An uncommitted subtransaction can be aborted and rolled back without any side-effects to other transactions outside its hierarchy. Hence, the concept of nested transactions contributes to a considerable refinement of the scope of in-transaction UNDO as compared to single-level transactions where UNDO-recovery necessarily yielded the state “begin of transaction” (BOT). It may be further refined by adding an appropriate savepoint concept to nested transactions (Härder and Rothermel, 1987; Rothermel and Mohan, 1989).

1.3 Explicit Control Structure

When parallel and asynchronous activities are to be coordinated for a single unit of work (from an external point of view), the introduction of a powerful explicit control structure allowing for the delegation of pieces of work and their atomic execution appears to be mandatory. Such a structure will greatly reduce the complexity of programming and enhance the reliability of transaction processing.

1.4 System Modularity

Subtransactions facilitate a simple and safe composition of a transaction program whose modules may be designed and implemented independently. This system modularity serves other design goals as well: encapsulation (information hiding), failure limitation, and security.

1.5 Distribution of Implementation

The concept of the nested transaction supports the implementation of distributed algorithms by a flexible control structure for concurrent execution. Distribution of

data and processing, in turn, has a major impact on overall efficiency, in terms both of cost-effective use of hardware (special processors, I/O devices) and responsiveness. Distribution also affects availability (replication of data). Hence, the robustness of the system may be improved in various ways.

In a centralized DBMS, nested transactions have some uses, however, they do not exploit their full potential due to the lack of resources. An obvious advantage is a clearer control structure for the execution of complex transactions supporting the design of more reliable programs. It also allows for the isolated rollback of an uncommitted subtransaction in the case of forced abort or transaction failure. When serializability of transactions controlled by strict 2-phase locking protocols (or equivalent methods) is required, neither lock granules nor lock duration are affected by such an approach. Subtransactions do not release their locks; they are inherited by their parent transaction.

For multi-layered centralized DBMSs, some kind of multi-level transaction management was provided where the subtransactions serve as control structures in the various layers. To gain a higher degree of concurrency and more flexible control of lock granules, a so-called multi-level concurrency control was introduced. Furthermore, isolated rollback of subtransactions can be guaranteed. In System R (Astrahan et al., 1976) this concept is used for two layers: locking is applied twice: on tuples until EOT (long tuple locks) and on pages for the duration of each tuple operation (short page locks for actions). Since tuple operations can be regarded as subtransactions and page locks are released before EOT of the parent transaction, this technique has been called "open-nested transaction." The problems involved were discussed by Gray (1981).

The generalization of open nested transactions for centralized systems—called *multi-level transactions*—was proposed to allow early release of locks at lower levels of control (Weikum and Schek, 1984; Moss et al., 1986; Weikum, 1986; Beerli et al., 1989); however, they rely on compensation operations for subtransactions to be applied in the case of rollback recovery. A detailed description of all aspects of multi-level transaction management including a discussion of performance issues is given by Weikum (1991). We don't want to consider this kind of multi-level structure and concurrency control for centralized DBMS operations.

Due to the salient properties supported by nested transactions, many researchers have focused attention on their design and implementation in distributed systems. Our approach is based on the proposals, results, and experiences of distributed systems design (Jessop, 1982; Allchin, 1983; Mueller et al., 1983; Spector and Schwarz, 1983; Walter, 1984; Liskov, 1985); it tries to adjust the concept of nested transactions and improves its use for distributed DBMSs. Our prime goal is its

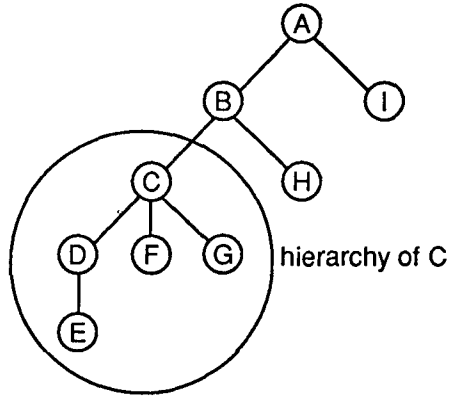
investigation and its conceivable extension for flexible intra-transaction parallelism. Due to space limitations we restrict our discussion to concurrency control and deadlock detection issues. Recovery problems are dealt with by Moss (1987), Härder and Rothermel (1987), and Rothermel and Mohan (1989). To facilitate our discussion, we introduce a model for nested transactions; it is designed so as to not prohibit parent/child- or sibling-parallelism.

In Section 3, the basic concurrency control model (Moss, 1985) is discussed. In some systems, it has been extended and refined by downward inheritance, enabling transactions to pass on locks to their child transactions. In Section 4, we propose a number of generalizations and extensions for concurrency control in nested transactions. Controlled downward inheritance enables a parent to give its child access to shared data and at the same time to restrict its mode of usage. Another refinement allows the use of more general lock modes as compared to the simple S-X lock model. Hence, applications may better adjust their synchronization needs. So far, all efforts are directed towards enhancement of concurrency control in transaction hierarchies operating on “flat” objects. Since every practical DBMS is forced to use an object hierarchy to provide fine as well as coarse lock granules at reasonable cost, we design a concurrency control protocol which combines object and transaction hierarchies as well as supports controlled downward inheritance. Section 5 discusses deadlock detection issues in nested transactions, and Section 6 compares concurrency control schemes of some system implementations. We conclude and summarize our results in the final section.

2. A Model of Nested Transactions

The concurrency control techniques we present in this article are based on the nested transaction model introduced by Moss (1985). A transaction may contain any number of subtransactions, which again may be composed of any number of subtransactions—conceivably resulting in an arbitrarily deep hierarchy of nested transactions. The root transaction which is not enclosed in any transaction is called the *top-level transaction* (TL-transaction). Transactions having subtransactions are called *parents*, and their subtransactions are their *children*. We also speak of *ancestors* and *descendants*. The ancestor (descendant) relation is the reflexive transitive closure of the parent (child) relation. We use the terms *superior* or *inferior* for the non-reflexive version of the ancestor or descendant. The set of descendants of a transaction together with its parent/child relationships is called the transaction's *hierarchy*. In the following, we use the term “transaction” to denote both TL-transactions and subtransactions.

Figure 1. Example of a Transaction Tree.



The hierarchy of a TL-transaction can be represented by a so-called *transaction tree*. The nodes of the tree represent transactions, and the edges illustrate the parent/child relationships between the related transactions. In the transaction tree shown in Figure 1, the root is represented by TL-transaction A. The children of subtransaction C are D, F, and G; and the parent of C is B. The inferiors of C are D, E, F, and G, and the superiors are B and A. Of course, the descendant and ancestor sets of C additionally contain C itself. The hierarchy of C is depicted as the subtree spanned by C's descendants.

The properties defined for flat transactions are *atomicity*, *consistency*, *isolated execution*, and *durability* (ACID-properties) (Härder and Reuter, 1983). In the nested transaction model, the ACID-properties are fulfilled for TL-transactions, while only a subset of them are defined for subtransactions. A subtransaction appears atomic to the other transactions and may commit and abort independently. Aborting a subtransaction does not affect the outcome of the transactions not belonging to the subtransaction's hierarchy, and hence subtransactions act as firewalls, shielding the outside world from internal failures. If the concurrency control scheme introduced by Moss (1985) is applied, isolated execution is guaranteed for subtransactions. However, to increase intra-transaction parallelism the enhanced schemes proposed in this article allow transactions belonging to the same TL-transaction hierarchy to share data in a controlled manner. The durability of the effects of a committed subtransaction depends on the outcome of its superiors—even if it commits, aborting one of its superiors will undo its effects. A subtransaction's effects become permanent only when its TL-transaction commits. The consistency property for subtransactions seems to be too restrictive, because sometimes a parent transaction needs the results of several child transactions to perform some consistency preserving actions.

To exploit the inherent potential of nested transactions and their advantages as stated in Section 1, the degree of intra-transaction parallelism should be as high as possible. Two kinds of intra-transaction parallelism can be defined, parent/child-parallelism and sibling-parallelism. If the first kind of parallelism is supported, then a transaction may run in parallel to its children, while in the second kind, siblings are allowed to run concurrently. Using both these definitions, we are able to characterize four levels of intra-transaction parallelism:

1. *Neither parent/child- nor sibling-parallelism:* At any point in time there is at most one transaction active in a TL-transaction hierarchy, i.e. there is no intra-transaction parallelism at all. Since all transactions in a hierarchy are executed serially, no concurrency control among them is needed. For example, if each transaction is executed by a single process and processes communicate only by means of a (synchronous) remote procedure call mechanism, only this level of "parallelism" can be provided.
2. *Only sibling-parallelism:* If only siblings may be performed concurrently, then a transaction never runs in parallel with its superiors. This kind of restricted parallelism enables a transaction to share objects with its ancestors without further concurrency control. For example, in the ARGUS system (Liskov, 1985), the intra-transaction parallelism is restricted to this level.
3. *Only parent/child-parallelism:* Since a transaction and its children may run concurrently but siblings may not, in a TL-transaction hierarchy only the transactions along one path of the hierarchy may run in parallel. This kind of restriction simplifies intra-transaction concurrency control in the sense that only transactions residing in the same path have to be synchronized with each other. (This reason hardly justifies such a system design).
4. *Parent/child- as well as sibling-parallelism:* This level permits arbitrary intra-transaction parallelism, i.e., in principle, all transactions of a TL-transaction hierarchy may be executed concurrently. Of course, compared to the degrees of parallelism described above, this degree requires the most sophisticated concurrency control scheme. For example, LOCUS (Mueller et al., 1983) and CLOUDS (Allchin, 1983) support this level of parallelism.

As discussed so far, our transaction model does not contain any essential restrictions. Transactions may be performed either entirely on a single processor site or may be distributed over multiple processors located at one or more sites. Moreover, the model does not restrict the kind of data distribution implemented by

the underlying system, and hence our considerations apply for data sharing as well as for data distribution approaches (Rahm, 1992). Since we focus on concurrency control concepts, introduction of further refinements or implementation issues would only burden our discussion.

3. Basic Locking Rules for Nested Transactions

Locking as the standard method of concurrency control in DBMS has been used successfully for a variety of applications over the past decade and longer. Therefore, it is reasonable to choose conventional locking protocols as our starting point of investigation for nested transactions. Conventional locking protocols offer two modes of synchronization: read, which permits multiple transactions to Share an object at a time, and write, which gives the right to a single transaction for eXclusively accessing an object (Gray, 1978). As far as concurrency control is concerned, our data model initially consists of disjoint (“flat”) objects, O_i , which are the *lockable units*.

In the next part of this section, we summarize the locking scheme for nested transactions proposed by Moss (1985). This scheme only allows for upward inheritance of locks, i.e., a transaction can inherit locks from its children, but not vice versa. In the last part, we extend this scheme such that it supports upward as well as downward inheritance. Both schemes presented in this section have been implemented in several systems.

3.1 Upward Inheritance of Locks

Before describing the locking rules proposed by Moss (1985), we have to introduce some terminology. Possible lock modes of an object are NL-, S-, and X-mode. The *null mode* (NL) represents the absence of a lock request for or a lock on the object. A transaction can acquire a lock on object, O , in some mode, M ; then it *holds* the lock in mode M until its termination. Besides holding a lock, a transaction can *retain* a lock. When a subtransaction commits, its parent transaction inherits its locks and then retains them. If a transaction holds a lock, it has the right to access the locked object (in the corresponding mode), which is not true for retained locks. A retained lock is only a place holder. A retained X-lock, denoted by $r:X$ (as opposed to $h:X$ for an X-lock held), indicates that transactions outside the hierarchy of the retainer cannot acquire the lock, but that descendants of the retainer potentially can. That is, if a transaction, T , retains an X-lock, then all non-descendants of T cannot hold the lock in either X- or in S-mode. If T is a retainer of an S-lock,

it is guaranteed that a non-descendant of T cannot hold the lock in X-mode, but potentially can in S-mode. As soon as a transaction becomes a retainer of a lock, it remains a retainer for that lock until it terminates.

Having introduced this terminology, we can formulate the locking rules now:

Rule 1: Transaction T may acquire a lock in X-mode if

- no other transaction holds the lock in X- or S-mode, and
- all transactions that retain the lock in X- or S-mode are ancestors of T.

Rule 2: Transaction T may acquire a lock in S-mode if

- no other transaction holds the lock in X-mode, and
- all transactions that retain the lock in X-mode are ancestors of T.

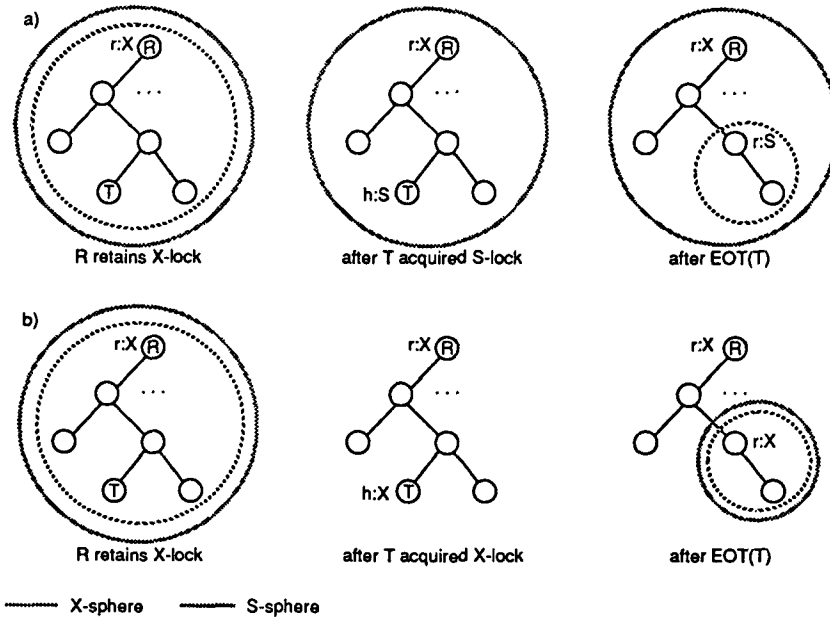
Rule 3: When subtransaction T commits, the parent of T inherits T's locks (held and retained). After that, the parent retains the locks in the same mode (X or S) in which T held or retained the locks previously.¹

Rule 4: When a transaction aborts, it releases all locks it holds or retains. If any of its superiors holds or retains any of these locks they continue to do so.

Obviously, the rules stated above only allow for upward inheritance of locks, i.e., a transaction can only inherit its children's locks, but not vice versa. The principle of upward inheritance is exemplified in Figure 2, where we use the notions X- and S-sphere for describing the implications of this principle. The X-sphere (S-sphere) of an object is defined to be the set of transactions that can potentially lock this object in X-mode (S-mode). In Figure 2a, the X-sphere of an object disappears entirely when transaction T acquires an S-lock on this object, i.e., no transaction may acquire an X-lock on this object, while each transaction in R's hierarchy may lock the object in S-mode. After commit of T, a new X-sphere is established, which consists of the descendants of T's parent transaction. In Figure 2b, the X- as well as the S-sphere disappear when T acquires an X-lock. When T commits, a new X- and S-sphere are established. In general, a transaction acquiring a lock on an object may cause the object's X- or S-sphere to shrink, while the termination of a transaction may cause them to grow.

The rules stated above only allow for upward inheritance at commit time, i.e., a transaction may not inherit a child's locks before the latter commits. This restriction

1. Note that the inheritance mechanism may cause a transaction to (conceptually) retain several locks on the same object. Of course, the number of locks retained by a transaction should be limited to one by only retaining the most restrictive lock.

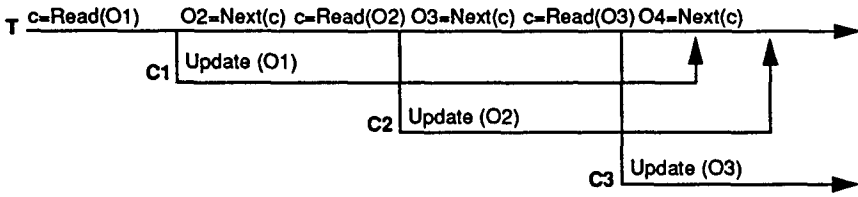
Figure 2. Upward inheritance of locks.

guarantees that transactions can see the effects of committed children only, and hence are not affected by failures of children. Furthermore, this restriction ensures that the subtransactions of a transaction tree are serializable. Allowing upward inheritance before commit time would cause transactions to become dependent on the outcome of child transactions, i.e., subtransactions would not act as firewalls anymore (Härder and Rothermel, 1987) such that application code within a subtransaction had to cope with concurrency and recovery issues.

3.2 Downward Inheritance of Locks

We feel that the restrictions caused by allowing only upward lock inheritance especially prevent desirable decompositions of transactions into a set of cooperating subtransactions. For example, assume an application that navigates through an object base and updates each of the accessed objects. A desirable decomposition of the above task is depicted in Figure 3. Transaction T reads an object and determines the next object to be accessed by applying the Next operation on the current object's content. Then T creates a new child transaction, which asynchronously performs an Update operation on the current object, while T reads in the next object, on which it acts as described above. This decomposition has some appealing characteristics: (1) Update operations are performed in parallel. (2) If an update operation fails,

Figure 3. Decomposition of an Application



it does not affect the other operations. A failed operation can be restarted at a later point in time. (3) The update operations are performed in isolation from each other. This is of particular importance if the update of an object may imply updates on other objects. For example, the update of two different objects may imply two updates of the same access path.

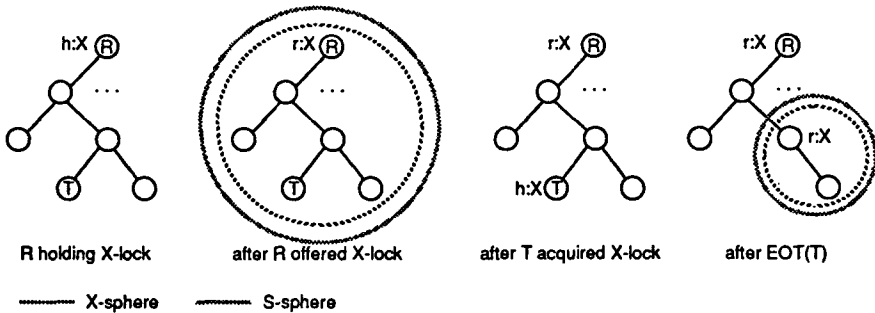
Unfortunately, this decomposition is impossible if the basic locking rules proposed by Moss (1985) are applied. To be able to perform the *Next* operation on an object, T must hold an R-lock on this object. Since T must hold this R-lock until it commits, no child of T can ever acquire an X-lock on this object. In other words, once an object has been read by T, it cannot be updated by T's children anymore.

The decomposition required in the example of Figure 3 is possible as soon as downward inheritance of locks is supported by the underlying locking scheme. In such a scheme, subtransactions may inherit locks from superiors, where inheritance of a lock can only take place after the superior holding this lock has explicitly offered this lock for downward inheritance. A transaction can offer a lock it holds to the transactions in its hierarchy, which can then acquire the lock according to the locking rules stated above. Consequently, the concept of downward inheritance allows a transaction to make all or a subset of its locks available to its hierarchy. The locking rules proposed by Moss (1985) can easily be extended to support downward inheritance by adding a new rule:

Rule 5: Transaction T, holding a lock, can offer the lock (to the transactions in its hierarchy). After offering the lock, T retains the lock in the same mode it held the lock before.

A transaction offering a lock disclaims (temporarily) the right to access the locked object and gives the transactions in its hierarchy the opportunity to lock this object in any mode. Of course, in its hierarchy there might be either at most one transaction holding the lock in X-mode or a number of transactions holding the lock in S-mode. Since the transaction offering a lock still retains the lock in the

Figure 4. Downward inheritance of locks.



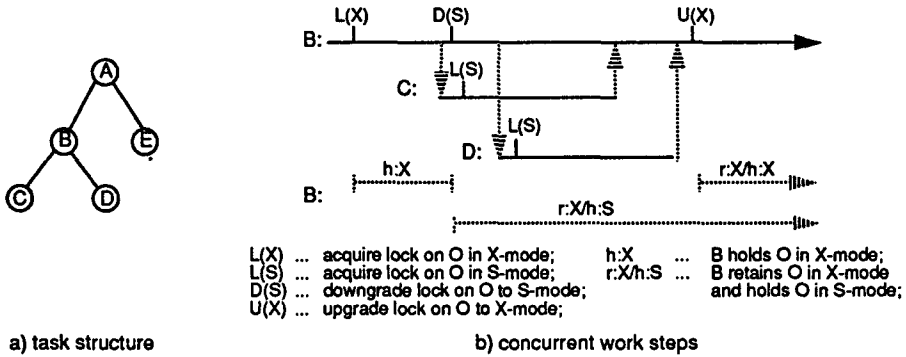
mode it held the lock before, no transaction from outside its hierarchy can lock the object in a mode that conflicts with the mode of the retained lock. To become a holder again, the transaction must acquire the lock anew, which only succeeds if rules R1 and R2 stated above are fulfilled.

An example for applying the lock-offering mechanism is illustrated in Figure 4. When transaction R offers the X-lock it holds, an S- and X-sphere comprising R's hierarchy is established for the corresponding object, i.e., all descendants of R have the opportunity to lock the object either in S- or in X-mode. As depicted in the example, the object's X- and S-sphere disappear when a descendant of R locks the object in X-mode.

If downward inheritance of locks is possible, the isolation property of transactions may be violated. While transactions belonging to different TL-transaction hierarchies still cannot interfere, transactions of the same hierarchy may share data. As a consequence, a transaction may see uncommitted data of superiors. This, however, cannot lead to inconsistencies because the effects of the transaction are undone when a superior aborts. On the other hand, a transaction may never see uncommitted data of inferiors, i.e., subtransactions act as firewalls even if downward inheritance of locks is allowed.

A lock-offering mechanism similar to the one described above has been implemented in the LOCUS system (Mueller et al., 1983). Some kind of automatic downward inheritance is provided in the ARGUS system (Liskov, 1985). In this particular approach, concurrency control is considerably simplified, since conflicts among transactions in a hierarchical path are prevented by allowing only sibling-parallelism. Automatic downward inheritance is then implicitly obtained by the rule that a transaction may acquire a lock if each transaction holding this lock is a superior of it.

Figure 5. Decomposition of a design task.



4. Enhanced Concurrency Control for Nested Transactions

By using the idea of downward inheritance we gain more flexibility of lock inheritance in a given transaction hierarchy. We have poor control over its specific usage, however. For this reason, this kind of offering concept still has some shortcomings in situations where a transaction offering a lock desires to control the mode in which its inferiors can hold the lock. For example, consider the access sequence shown in Figure 3 once more. With the additional rule R5 it is possible to make the desired decomposition. When transaction T offers the X-lock for O1, O2, O3, etc., then its children C1, C2, C3, etc., can acquire and hold the lock in any mode later on. However, it would be helpful if T could prevent some child Ci from being able to hold the lock in X-mode in order to make sure that Ci cannot change the respective Oi.

4.1 Controlled Downward Inheritance

The need for controlling the lock mode in which inferiors can access an offered object becomes more obvious if we consider an example from a cooperative design environment (Kim et al., 1984; Bancilhon et al., 1985). Figure 5 shows a design task that is structured as a three-level transaction hierarchy. Assume that transaction B generates an object, O, describing the interface of a work piece. Transactions C and D, which are children of B, design subparts of the work piece and therefore require read access to the interface description. To allow its children to read O, B must offer the lock that it holds on O. If there is no way to control the mode

in which children can hold the lock, one of the children may acquire the lock in X-mode which has two undesirable consequences: First, the child can change O, and second, the child blocks its siblings by preventing them from reading O.

To overcome these problems we suggest an extension of the locking rules introduced in the previous section. In the scheme discussed previously, a transaction can offer locks that it holds to its inferiors. In the extended scheme, we will replace the lock-offering mechanism by primitives supporting the upgrading and downgrading of locks:

Downgrade. Transaction T, holding a lock in mode M, can downgrade the lock to a less restrictive mode M' . After downgrading the lock, transaction T holds the lock in mode M' and retains the lock in mode M. For example, a transaction holding a lock in X-mode can downgrade the lock to mode S or NL.²

Upgrade. Transaction T, holding a lock in mode M, can upgrade the lock to a more restrictive mode M' if the following condition is satisfied: No other transaction holds the lock in a mode conflicting with M' and all transactions that retain the lock in a mode conflicting with M' are ancestors of T.³ For example, transaction T, holding a lock in S-mode, can upgrade the lock to mode X if no other transaction holds the lock in X- or S-mode, and all transactions retaining the lock in X- or S-mode are ancestors of T.

In the *extended scheme*, holding and retaining a lock have exactly the same semantics as in Moss's scheme (1985). After downgrading a lock from mode M to mode M' , a transaction holds the lock in mode M' and retains it in mode M. Since the transaction retains the lock in M-mode, it prevents transactions outside its hierarchy from holding the lock in a mode conflicting with M. On the other hand, since it holds the lock in mode M' it keeps its inferiors from holding the lock in a mode conflicting with M' . That is, in contrast to the offering of locks described in the basic scheme, downgrade allows a transaction to control how its inferiors can hold a lock. For example, if transaction T downgrades a lock from X- to S-mode, it prevents transactions outside its hierarchy from holding the lock in any mode, and precludes its inferiors from holding the lock in X-mode, but allows its inferiors to hold the lock in S-mode. Of course, downgrading a lock to NL-mode is equivalent to offering a lock in the basic scheme. As stated above, the holder of a lock can

2. Downgrading to NL does not correspond to a general release of the lock. Release of such locks is limited to the sphere of the downgrader.

3. This condition is equivalent to the condition that must be satisfied for a transaction acquiring a lock.

upgrade the lock to a mode which is more restrictive than its current hold mode. This feature allows a transaction to upgrade a lock which it downgraded previously, e.g., a transaction that downgraded a lock to S-mode could again upgrade the lock to mode X as soon as its children have committed. Of course, transactions can also upgrade a lock without having downgraded it before. In the following, we will describe the extended locking rules. Italics will be used to point out the extensions added to Moss's scheme:

Extended rule 1: Transaction T may acquire a lock in X-mode or *upgrade a lock it holds to mode X* if

- no other transaction holds the lock in X- or S-mode, and
- all transactions that retain the lock in X- or S-mode are ancestors of T.

Extended rule 2: Transaction T may acquire a lock in S-mode if

- no other transaction holds the lock in X-mode, and
- all transactions that retain the lock in X-mode are ancestors of T.

Extended rule 3: When subtransaction T commits, the parent of T inherits T's locks (held and retained). After that, the parent retains the locks in the same mode (X or S) as T held or retained them before.

Extended rule 4: When a top-level transaction commits, it releases all locks it holds or retains.

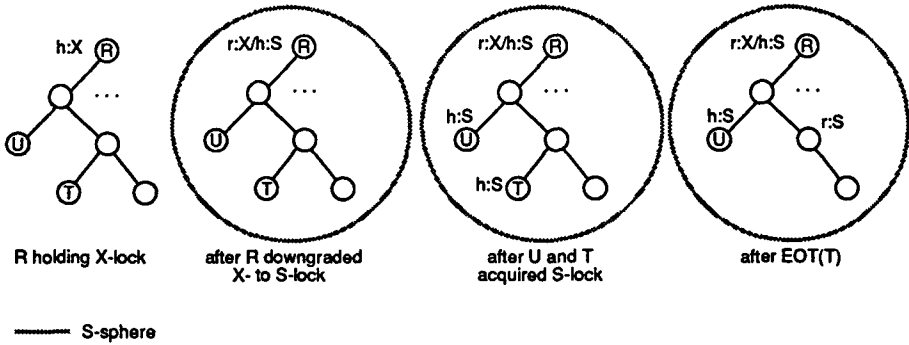
Extended rule 5: When a transaction aborts, it releases all locks it holds or retains. If any of its superiors hold or retain any of these locks, they continue to do so.

Extended rule 6: *Transaction T, holding a lock in X-mode, can downgrade the lock to mode S or NL. After performing the downgrade operation, T retains the lock in X-mode.*

Extended rule 7: *A transaction holding a lock in S-mode can downgrade the lock to mode NL. After performing the downgrade operation, T retains the lock in S-mode.*

The mode to which a T transaction downgrades a lock determines the modes in which the transactions of T's hierarchy cannot hold the lock. If the downgraded mode is S, the transactions of T's hierarchy cannot hold the lock in X-mode (since S conflicts with X). If the downgraded mode is NL, then the transactions in T's hierarchy can potentially hold the lock in any mode.

Figure 6. Controlled downward inheritance.



Some examples may help to clarify the key issue of controlled downward inheritance. The effect of offering an X-lock can be depicted in the scenario in Figure 4. A similar scenario in Figure 6 illustrates the downgrading of an X-lock to mode S. (Downgrading of S-locks is handled in an analogous manner). The essential issue observed in this example is that only S-locks may be granted within R's hierarchy, i.e., no X-sphere is established when the lock is downgraded to mode S.

Given these extended locking rules, the problem described in the design environment example above can be solved very easily (see Figure 5). After having generated object O, transaction B downgrades the X-lock it holds on O to mode S. Since it then holds the lock in S-mode, C and D are prevented from holding the lock in X-mode which guarantees that they cannot change O or block each other. Note that since B retains the lock in X-mode after downgrading the lock, transactions A and E cannot hold the lock in any mode, i.e., A and E can neither read nor write O. After commit of C and D, B can upgrade the lock once again.

4.2 Correctness Concerns

As stated in Section 3.1, upward lock inheritance at commit ensures that subtransactions act as firewalls in case of a failure and that all subtransactions of a TL-transaction remain isolated. Since we have proposed the concept of controlled downward inheritance, we now discuss the impact of this concept on the correctness of concurrent executions.

TL-transactions are serializable because

- each transaction in a TL-transaction tree locks each data object before accessing it
- all locks held by transactions in a transaction tree are released, but not before the TL-transaction commits.

This locking protocol corresponds to strict 2-phase locking for TL-transactions. It determines their serialization order by the time of their commit, as it holds for single-level transactions.

Now let us discuss the visibility of data changes and their induced dependencies within a transaction tree. In Moss's nested transaction model (1985) the following holds:

- A transaction may see changes only of those transactions that are committed and on which it depends.⁴ We say transaction T depends on a transaction T', if undoing the effects of T' causes the abortion of T.
- Once transaction T has seen a state of an object, this state never will be seen or changed by another transaction before T commits.

In contrast, our model allows for controlled downward inheritance which makes uncommitted data available to inferiors. For this reason, we observe the following properties:

- Transaction T may see changes only of
 - those transactions that are committed and on which T depends,
 - those transactions that are superiors of T
- An object state seen by transaction T may be changed by inferiors of T.

A transaction may see changes of superiors only if these transactions have downgraded the corresponding locks explicitly. That is, whether or not a transaction may see the effects of superiors can be controlled by the application logic. In terms of failures a transaction seeing changes of superiors causes no problems because, if a superior aborts, this transaction is aborted also. Note that once a transaction has seen an object, the object cannot be changed by a superior again before this transaction commits. If a transaction downgrades a lock, it must be aware of the

4. Remember that the effects of a committed (sub) transaction only become permanent when the top-level transaction commits. The reason why a transaction may only see data of committed transactions is to ensure that it is aborted when the effects of one of these transactions is wiped out due to a failure. In the transaction hierarchy of Figure 5a, transaction E may see changes from transaction C, but not before B has committed. When B commits, E becomes dependent on C as a failure, wiping out effects of C, which causes A and, of course, E to be aborted.

consequences of the reduced isolation. Downgrading from S- to NL-mode may cause unrepeatable reads from the downgrading transaction's point of view. With X-locks two cases must be considered: Downgrading from X- to NL-mode and from X- to S-mode. In the first case, from the downgrader's point of view lost updates and unrepeatable reads are possible in principle. However, a much more flexible cooperation is enabled where the correctness of execution has to be enforced by application level protocols. In CSCW (computer-supported cooperative work)-like applications, it is even conceivable that this kind of high-level control is based on so-called social protocols between end users. In the latter case, which prevents the inferiors of the downgrading transaction from keeping the downgraded lock in X-mode, neither unrepeatable reads nor lost updates can occur.

An important question is whether the firewall property of nested transactions is in some way affected by the downgrading mechanism. A transaction downgrading a lock does not become dependent on the outcome of its inferiors. When a child fails, its updates (possibly on objects with downgraded locks) are rolled back. Therefore, the downgrading transaction is not affected and may create another child to do the work.

In summary, the fact that a transaction may see changes from superiors causes no problems. The firewall property is not affected by the downgrading mechanism. Lost updates may only happen when locks are downgraded from X- to NL-mode. In this case, which provides the highest degree of flexibility in terms of cooperation, application-level concurrency control mechanisms are needed to ensure the required form of correctness. Since the application itself can decide how and when to use the downgrading mechanism, it can adapt the level of system-supported isolation to its cooperation needs and its facilities for application-specific concurrency control.

4.3 Generalization of Lock Modes

Thus far, we have described and refined a concurrency control scheme for S-X locks on "flat," non-overlapping objects (e.g., tuples or relations); in particular, we have developed a mechanism for controlled downward inheritance of locks in nested transactions. Closer consideration reveals that the lock modes (comprising only S and X so far) may be enriched by special modes to better adapt concurrency control to access patterns in practical applications. For example, tailored lock modes for frequent kinds of object access could be helpful to more effectively exploit the inherent parallelism of concurrent transactions. Furthermore, the use of semantic knowledge could greatly optimize some contention patterns of data access. However, this requires enhanced lock modes; in particular, it presupposes the ability

to introduce user-defined lock modes (Allchin, 1983; Schwarz and Spector, 1984).

Such a refinement of lock modes may be easily integrated into our model presented so far. Assume that the data model remains unchanged. Then our locking rules stated in Section 4.1 for S-X schemes can be generalized for basic and/or user-defined lock modes as follows:

Generalized Rule 1: Transaction T may acquire a lock in mode M or upgrade a lock it holds to mode M if

- no other transaction holds the lock in a mode that conflicts with M, and
- all transactions that retain the lock in a mode conflicting with M are ancestors of T.

Generalized Rule 2: When subtransaction T commits, the parent of T inherits T's locks (held and retained). After that, the parent retains the locks in the same mode as T held or retained them before.

Generalized Rule 3: When a top-level transaction commits, it releases all locks it holds or retains.

Generalized Rule 4: When a transaction aborts, it releases all locks it holds or retains. If any of its superiors hold or retain any of these locks, they continue to do so.

Generalized Rule 5: Transaction T, holding a lock in mode M, can downgrade the lock to a less restrictive mode, M'. After downgrading the lock, T retains it in mode M.

The locking rules stated above allow upward as well as controlled downward inheritance for arbitrary lock modes. If generalized rule 5 were omitted, we would get a generalization of Moss's scheme (1985) which only provides for upward inheritance.

4.4 Use of Hierarchical Locks in Nested Transactions

Let us now reconsider our underlying data model which has some serious drawbacks for realistic concurrency control situations. In particular, the flat object structure that requires disjoint lockable units of a given granule makes it impractical for large databases when small granules are needed for some transactions and larger ones for others. To improve selective access to data granules of varying sizes, hierarchical locking schemes have been proposed. In our context, hierarchically structured objects introduce a certain complexity, due to orthogonal transaction and data hierarchies.

As mentioned earlier, locking of disjoint partitions of a given size is insufficient for performance reasons in most applications. The choice of lockable units affects locking overhead of a transaction (space for lock control blocks, time to request and release locks) as well as concurrency among transactions. Hence, it implies a dichotomy of increased concurrency using fine lockable units and higher cost for lock management. While small granules are appropriate for “simple” transactions accessing a few tuples, they are intolerable (and hard to implement) for “complex” transactions accessing a large fraction of the database. Assume, for example, a sequential scan of a relation with 10^6 tuples; having only tuples as lockable units would require 10^6 consecutive lock requests and storing of just as many lock control blocks (of course, in main memory for performance reasons). Hence, coarser granularity locks are sometimes more natural and efficient (e.g., when sorting or reorganizing a relation).

These arguments should convince every DBMS designer that an object hierarchy for locking purposes has to be provided. In fact, every “practical” DBMS supports such a hierarchy of typically 2, 3, or 4 levels (e.g., System R has a generic 4-level hierarchy: database, segment, relation, tuple; Astrahan et al., 1976).

An appropriate hierarchical locking scheme was proposed for flat transactions (Gray et al., 1976). Two key ideas allowed for the design of a scheme that could be adapted to a transaction’s needs for either locking a few items using a fine lockable unit or locking larger sets of items with larger lock granules:

- A node, R, in a hierarchy can be locked explicitly. As a result, its entire subtree is implicitly locked, too.
- A transaction, locking part of the hierarchy, places “Intention mode” locks along the path to R to avoid a situation where an ancestor node of R is locked in an incompatible mode as compared to R. I-locks merely serve as place holders, signalling the fact that locking of a subtree is done at a lower level of the hierarchy, thereby preventing incompatible locks from being granted for the corresponding nodes.

Besides the known modes, S and X, an *Intention Share* mode (IS) and an *Intention eXclusive* mode (IX) were introduced to express a transaction’s intent to read and to update or read an object at a lower level of the hierarchy, respectively. A further refinement is the *Share and Intention eXclusive* mode (SIX) which grants an S-lock for the entire subtree to a transaction. In addition, it indicates the transaction’s intention to request X-locks explicitly for “finer” object granules later on. Table 1 (from Gray et al., 1976) shows the compatibilities among request/lock modes which derive from these semantics.

Table 1. Compatibility modes for hierarchical locking.

Compatibility Mode of request	Mode of lock					
	NL	IS	IX	S	SIX	X
NL	yes	yes	yes	yes	yes	yes
IS	yes	yes	yes	yes	yes	no
IX	yes	yes	yes	no	no	no
S	yes	yes	no	no	no	no
SIX	yes	no	no	no	no	no

For a comprehensive discussion of the precise effects of lock modes and their compatibilities see Gray (1978).

4.4.1 Basic Locking Rules for Object Hierarchies. We have now introduced the essential ingredients of both generalized locking rules for nested transactions and appropriate lock modes for an object hierarchy. How can we combine both together? We start with the basic concurrency control model where only upward inheritance is allowed. For the transaction hierarchy, our generalized rules 1–4 apply. Furthermore, when acquiring a lock on an object, O , we have to consider additional rules resulting from the object hierarchies.

As opposed to flat objects, an approach for controlling concurrent access to an object hierarchy has to obey the following rules (Gray et al., 1976):

1. Instead of locking an object directly, every transaction has to observe a strict hierarchical protocol requesting appropriate locks from root to leaf in the object hierarchy (denoted in the following as *root-to-leaf rule*.) A lock is granted at each level according to the compatibilities expressed in Table 1. As soon as a lock is obtained, a transaction may request another appropriate lock at the same or at the next lower level.
2. Level-to-level transitions should obey the following constraints called *level-to-level rules*:
 - IS held at a node only allows IS and S to be requested on descendant nodes.
 - IX granted for a node carries the privilege to request IS, IX, S, SIX and X at the next level.

Table 2. Locks in an object hierarchy.

Object hierarchy	Before EOT(T ₁)	After EOT(T ₁)		
		P	T ₂	T ₃
Database DB	h:IX	r:IX	h:IX	h:IS
Segment S	h:IX	r:IX	h:IX	h:IS
Relation R	h:X	r:X	h:IX	h:IS
Tuples;			h:X on t ₁	h:S on t ₃
			h:X on t ₂	h:S on t ₄

- S and X allow read and write access (respectively) to all descendants of the node without further locking.
- SIX carries the privileges of S and IX; hence, while S mode allows read only access to all descendants, write access at lower levels may be requested by IX or X at the next level.

As far as acquiring locks is concerned, the rules obtained for the transaction hierarchy and the object hierarchy must be satisfied independently. Following the root-to-leaf rule, transactions must request their locks from root to leaf in the object hierarchy. Whether or not a lock for an object may be granted in a particular mode is decided according to the level-to-level rules, the generalized locking rules (1-4), and the lock mode compatibilities depicted in Table 1. Since the rules introduced for the object hierarchy are independent of the underlying transaction model and the rules for both hierarchies are applied independently, our protocol and that proposed by Gray et al. (1976) for flat transactions only differ in the rules implied by the transaction model.

An example may clarify the issues involved in lock retainment for object hierarchies. In the scenario of Table 2, T₁ passes on its hierarchical locks for X-access on relation R to its parent P at EOT(T₁). After having retained the locks, P and its inferiors T_i are qualified to acquire read or write access on R or to tuples of R. Table 2 shows the locks of T₂ and T₃ for obtaining write and read access on tuples of R.

4.4.2 Upgrading and Downgrading Hierarchical Locks. Although we have succeeded in tying together both hierarchy types (transaction and object), we have so far obtained only a more economical and efficient solution of the concurrency problem than compared to the basic approach in Section 3.1. Since we cannot

Table 3. Inconsistent downgrading of a lock.

Object hierarchy	P	P after downgrade	T using downgraded lock
Database DB	h:IX	h:IX	h:IX
Segment S	h:IX	h:IX	h:IX
Relation R	h:SIX	r:SIX/-	h:SIX
Tuples _i	h:X on t ₁ h:X on t ₂	h:X on t ₁ h:X on t ₂	h:X on t ₃ h:X on t ₄

make a transaction's objects available to its inferiors, all arguments discussed earlier apply. Therefore, it is desirable to enable controlled downward inheritance in the presence of object hierarchies, too.

Assume, for example, a transaction P holds an SIX-lock on a relation, R, and wants to permit write access to tuples of R by its inferiors. Using the same kind of inheritance mechanism as in Section 4.1, P has to downgrade its lock on the object to an appropriate mode. To do so, P retains the SIX-lock on R (r:SIX) and holds R in IX-mode (h:IX). Note that r:SIX only prevents incompatible locks on R from being granted to non-descendants, but not to inferiors.

Let us examine whether such a straightforward approach may be applied. In the scenario depicted by Table 3, P holds R in SIX-mode and some tuples of R in X-mode. We assume that P downgrades the SIX-lock on R. Requesting a lock by an inferior T implies that T obeys the root-to-leaf and the level-to-level rules. Hence, as soon as T has acquired appropriate locks for the ancestors of R, it can request a compatible lock for R. The presented scenario is meant to serve as a counter-example for arbitrary inheritance of hierarchical locks and aims at clarifying a new issue: Inheritance of objects in data hierarchies. It shows that P holds some locks at the tuple level while it has downgraded the corresponding lock at the relation level to NL-mode, i.e., without particular protection.

In Table 3, T acquires an SIX-lock on R giving read access to all tuples of R. On the other hand, P still has some tuples locked in X-mode, namely t₁ and t₂. These exclusively locked tuples would be read by T, since read access to tuples of R need not be checked by T anymore. Even worse, write-write interference on the same tuples could occur, if T had locked R in X-mode. Of course, the sketched examples may cause severe consistency problems. These anomalies would not occur if the lock on relation R, together with all locks on its tuples, had been downgraded.

Control given by the hold-mode alone would not guarantee the desired consistency as exemplified by only downgrading R to S-mode.

The key observation in the example above is that downgrading a lock without considering the whole object hierarchy may lead to inconsistencies. The same can be shown for upgrading locks in object hierarchies. For example, if a transaction has locked a database in IS-mode and upgrades an S-lock that it holds on a segment of this database to X-mode, similar inconsistencies may occur. Obviously, to prevent violations of the level-to-level rules, upgrading or downgrading of a lock may enforce upgrade or downgrade operations on other locks held in the object hierarchy.

When a transaction T upgrades a lock held on an object, O, within an object hierarchy, it might be necessary to also upgrade locks of T held on superior objects of O in order to satisfy the level-to-level rules. For example, if T holds an IS-, IS-, and S-lock on a database, a segment of this database and a relation of this segment, respectively, the level-to-level rules enforce the upgrading of both IS-locks to IX-mode before the relation lock can be upgraded to X-mode. Because upgrading the locks on an object and superior objects is not performed in an atomic manner, upgrading should be done in a root-to-leaf direction. Of course, an upgrade operation can only take place if the generalized locking rules 1-4 are fulfilled. Otherwise, upgrading is blocked which may cause deadlocks to occur (see Section 5).

Because upgrading a lock on an object, O, converts the mode of the locks to a more restrictive one, the level-to-level rules are not violated as far as locks on inferior objects of O are concerned. However, due to the upgrade operation, locks held by the upgrading transaction on inferior objects of O may become useless. For example, when a lock on relation R is upgraded from SIX- to X-mode (lock escalation; Bernstein et al., 1987), all locks held by the upgrading transaction on individual tuples of R are not needed anymore. A clean approach to handling those useless locks is to release them as part of the upgrade operation. An actual implementation may optimize this cleanup process using pragmatic arguments (e.g., see System R; Astrahan et al., 1976).

Downgrading a lock held by transaction T on object O is confined to the subhierarchy with O as the root object. Superiors of O in the object hierarchy are not involved because downgrading cannot violate the level-to-level rules as far as superiors of O are concerned. However, with respect to the objects in its subhierarchy, downgrading the lock on O may cause a violation of the level-to-level rules: if T holds a lock on subobject O' of O, then after downgrading, the mode of the lock held on O' may violate the level-to-level rules. For example, assume T holds an IX-lock on a relation, R, and an X-lock on a tuple, t of R. If T downgrades

Table 4. Consistent lock modes of subobjects.

Object O of transaction T downgraded to mode	Consistent modes for locks of T on subobjects of O
NL	NL
IS	NL, IS, S
IX	NL, IS, IX, SIX, S, X
SIX	NL, IX, X
S	NL

Table 5. Consistent downgrading of a lock.

Object hierarchy	P	P	T
Database DB	h:IX	h:IX	h:IX
Segment S	h:IX	h:IX	h:IX
Relation R	h:SIX	r:SIX/h:IX	h:SIX
Tuples; 	h:X on t_1 h:X on t_2	r:X/h:S on t_1 r:X/- on t_2	h:S on t_1 h:X on t_2

its lock on R from IX to IS, then T's X-lock on t is not consistent anymore with the lock mode of its parent object. As a consequence, downgrading the lock on O may require downgrading locks held by T on objects in the subhierarchy of O, such that the level-to-level rules are satisfied. In Table 3, this would require downgrading T's lock on tuple t to S- or NL-mode.

Table 4, derived from the level-to-level rules, lists for each possible mode to which the lock on O can be downgraded, the modes in which T can hold locks on the objects in the subhierarchy of O without violating the level-to-level rules. For example, if the lock is downgraded to IS-mode, T can hold subobjects of O either in NL-, IS- or S-mode. If subobjects are held in more restrictive modes, the locks on these objects must be downgraded to one of the listed modes. Note, since downgrading an entire subhierarchy cannot be done atomically, downgrading should be performed in a leaf-to-root direction.

By observing these rules, consistency-preserving downward inheritance of locks may be easily achieved by P in the example in Table 3 by downgrading the tuples t_1 and t_2 before downgrading relation R. Control of lock usage is then possible by

downgrading to the appropriate modes. In Table 5, the locks in the subhierarchy of relation R have been downgraded to different modes which allows for selective control of access to R's subobjects.

Downgrade of an intention mode (IS, IX, SIX) implies subsequent downgrades of locks on subobjects in order to satisfy the level-to-level rules. This, however, can be avoided by restricting downgrade operations to S- and X-locks. If T holds a lock on object O in S- or X-mode, then the entire subhierarchy of O is locked implicitly for T, that is, T does not hold any locks on subobjects of O. Hence, downgrade of O does not involve downgrading locks on lower levels of O's subhierarchy.

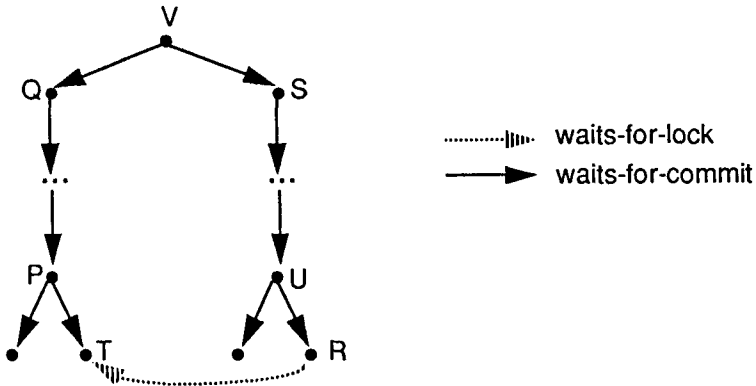
Let us summarize our findings for controlled downward inheritance of locks in a data hierarchy. In general, downgrading of entire subtrees is necessary for hierarchical objects to guarantee consistency of downward inheritance in nested transactions. That is, if an M-lock held by transaction T on object O is downgraded, it might be necessary to downgrade locks held by T on inferiors of O in order to satisfy the level-to-level rules. If downgrading is allowed only on X- and S-locks, then the downgrading of a lock never involves locks held on lower levels of the object hierarchy, which simplifies the downgrade mechanism substantially.

5. Deadlock Detection in Nested Transactions

Lock protocols are pessimistic, that is, they block lock requests of data currently granted to another transaction in a conflicting mode, and therefore are not immune to deadlocks. Deadlocks may occur among transactions belonging to various TL-transactions and even among subtransactions within a single transaction hierarchy. For deadlock detection, we mainly follow the basic approach sketched in Moss (1985), which allows identification of existing deadlocks. In addition, we propose the maintenance of further information (waits-for-retained-locks relation) to detect opening-up deadlocks as early as possible.

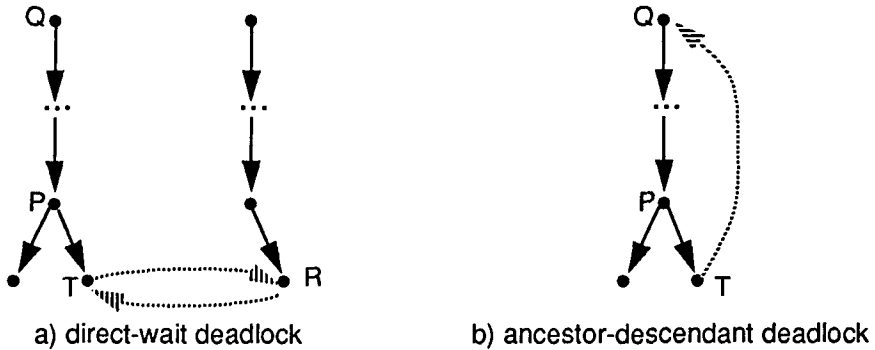
Deadlocks in nested transactions can be resolved by the concepts known for single-level transactions extended by some mechanisms tailored to the properties of the nested structure (Moss, 1985; Rukoz, 1991). When a transaction acquires a lock for data which is incompatible with a lock held by another transaction, the requesting transaction is deactivated: a direct wait for the lock holder occurs. All direct waits are maintained in a waits-for-lock relation in order to detect deadlocks. Using this waits-for-lock relation, deadlock detection can be performed immediately when a transaction is blocked or after some elapsed time. A deadlock exists if, and only if, a cycle is found in the waits-for-lock relation. For single-level transactions, such a cycle is composed of direct waits (or waits for lock) only.

Figure 7. Lock and commit waits.



As we have seen in Section 3.1, nested transactions have an inner structure which determines along which paths locks are inherited and whether retained locks can be acquired. Assume a subtransaction, R, waits for a lock held by a subtransaction, T. After commit of T, its locks are inherited and retained by its parent transaction, P (Figure 7). Now lock requests from transactions in P’s X- or S- sphere can be served. Outside P’s sphere, however, transaction R cannot acquire retained locks; for this reason, it has to wait for the retained locks of P. Such waits for retained locks are indirect waits. They propagate along the ancestor hierarchy of P. In the following, we will introduce two different waiting relationships:

1. *Waits-for-retained-locks:* A lock requestor, R, directly waits for a lock holder, T, if the mode of the requested lock is in conflict with the lock mode held by T. Let Q be the highest ancestor of T that is not an ancestor of R. Then R indirectly waits for all ancestors of T up to Q until they commit. Those wait relationships are called waits-for-retained-lock. This wait rule implies that if requestor R is not in the same TL-transaction as T, then R must wait for the retained locks of T’s TL-transaction, which are released at its commit.
2. *Waits-for-commit:* Since a waiting lock requestor, R, cannot proceed with its work, all ancestors of R may have to wait. In Figure 7, U cannot commit until R does, and U’s parent cannot commit until U does, that is, all ancestors of R cannot commit before R does. We denote this kind of wait relationship by waits-for-commit, which can be represented by the parent-child relationships, as outlined in Figure 7. Due to this dependency, an ancestor of R may wait

Figure 8. Deadlock situations.

for all transactions for which R directly or indirectly waits. Of course, waiting may be broken up as soon as T or one of its ancestors aborts. As illustrated in Figure 7, R directly waits for T and indirectly waits for retained locks for P, ..., Q. Furthermore, since S, ..., U wait for commit of R, they also wait for T, P, ... Q.

In Section 4.4, hierarchical locking was employed to nested transactions. The key observation exhibited is that object and transaction hierarchies are orthogonal. As a consequence, further aspects are not added to deadlock detection when hierarchically composed objects are used. As illustrated by Figure 7, waits-for relations occur among transactions; thus, the rules of the hierarchical locking protocol do not interfere with the waits-for-lock and waits-for-commit relationships as long as the root-to-leaf and the level-to-level rules are observed.

5.1 Detection of Existing Deadlocks

In order to handle deadlock detection in nested transactions successfully, we have to combine the various waits-for relations. Considering the waits-for-lock relation, only direct-wait deadlocks can be found, as indicated in Figure 8a, whereas other kinds of deadlocks cannot be detected. This is true no matter whether a deadlock occurs within a TL-transaction or among subtransactions of various TL-transactions. The cycle in Figure 8a consists of direct waits only, that is, all transactions in the cycle cannot proceed any further. In Figure 8b, however, another kind of cycle is encountered; this situation does not mean that progress has stopped everywhere in the cycle. T waits for a lock of Q; Q, ..., P may proceed for some time, but cannot commit without aborting T. Since T must be rolled back anyway, the best decision is to detect and resolve this ancestor-descendant deadlock immediately. A

request of T, causing a lock wait on its ancestor Q, can be detected only by using the combined waits-for-lock and waits-for-commit relation information.

Situations such as those illustrated by Figure 8b would frequently be caused when descendants refer to exclusively used data in an uncoordinated way. Controlled downgrading of locks, however, provides a mechanism to avoid such cycles, that is, application knowledge is applied to reduce the possibility of a deadlock involving lock and commit waits. Coordinated work requires that a parent, P, should downgrade the lock on an object, O, currently granted to P, before it creates a child, T, to do some work on O. Then T can acquire the lock for O in a non-conflicting mode without causing a blocking situation. Downgrading enables deadlock-free cooperation, but cannot enforce it; if T requests the lock in a mode more restrictive than the offered one, a deadlock may arise.

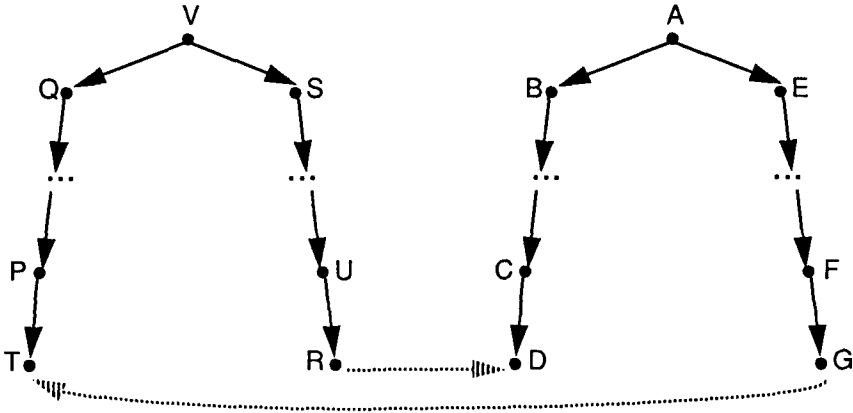
Upgrading a lock may lead to wait situations and therefore to deadlocks as they occur in single-level transactions. Assume in Figure 7 that T and R already hold an S-lock on object O. Now, if R upgrades the lock to mode X, R has to wait until a direct ancestor of R retains the lock or, if T and R are not in the same TL-transaction, until T's TL-transaction has committed. Hence, our wait rule applies to lock upgrades, too.

5.2 Detection of Opening-up Deadlocks

The combined use of waits-for-lock and waits-for-commit relations turned out to be sufficient for nested transactions to detect existing cycles embodying direct-wait or ancestor-descendant deadlocks. Since a waits-for-lock relationship is only represented between the requestor and the holder of a lock (or, after commit of the lock holder, the current retainer), waits-for-retained-lock relationships between the requestor and all ancestors of the holder (retainer) are not explicitly established in the waits-for information. For nested transactions, however, these waits-for-retained-lock relationships should be taken into account to provide for early deadlock detection. This may save a lot of useless work as shown in the scenario of Figure 9. Figure 9 represents a deadlock-free situation, since transactions T, D, and possibly others can proceed with their work. R waits for D, and G for T to obtain the requested locks. All other waits indicated are waits-for-commit. R indirectly waits for A, which is the oldest ancestor of D that is not an ancestor of R. On the other hand, G indirectly waits for V. If we evaluate this information ($R \rightarrow A$, $G \rightarrow V$), we can immediately detect a cycle opening up.

Optimistically, one may not care about such an opening-up deadlock, since an abort of any transaction involved would eventually avoid the actual deadlock. For

Figure 9. Opening-up deadlock among nested transactions.



example, the abort of any transaction in Figure 9 resolves the opening-up deadlock before all progress ceases within the TL transactions V and A. However, transaction aborts are regarded as exceptions and should not be considered a remedy to break opening-up deadlock cycles.

In contrast, a pessimistic approach usually saves work. If we use the transitive waits-for-commit and waits-for-retained-lock relationships of all ancestors, e.g., of V on R and A on G, as well as R on A and G on V, we can construct a direct (future) cycle between V and A and can roll back either V or A. However, deadlock detection and resolution at the level of the highest non-common ancestors of the transactions which have caused the conflict may not be appropriate. Deadlock resolution is typically based on transaction rollback and should affect minimal data granules or work lost.

For this reason, special measures should be used to determine an opening-up deadlock as early as possible and at the suitable level in the nested transaction hierarchies. In addition to waits-for-lock and waits-for-commit relations, the waits-for-retained-lock relationships have to be included in the waits-for information. For example in Figure 9, the relationships $R \rightarrow C, \dots, R \rightarrow B, R \rightarrow A$, as well as $G \rightarrow P, \dots, G \rightarrow Q, G \rightarrow V$ have to be represented in order to successfully search for opening-up cycles. Once an opening-up deadlock is detected, all transactions involved have to be considered to determine a low-cost victim for rollback. Since rollback of a parent transaction implies rollback of all its inferiors (committed and uncommitted), rollback of a child transaction is always cheaper than that of the corresponding parent transaction. For this reason, rollback of a lock holder (retainer) or a lock requestor is always cheaper than a rollback of their ancestors

in a potential cycle; hence, the set of transactions from which to choose a rollback victim is the set of lock holders (retainers) and lock requestors. In Figure 9, this set of candidates is D,T and R,G, respectively.

Note that, in contrast to Moss (1985), in our transaction model these transactions must not be leaves in the current transaction tree. However, the same methods and cost measures could be applied for breaking up the cycle, in principle. Since candidate transactions can occur arbitrarily in the transaction hierarchy, resource estimation involving the evaluation of subtrees may become much more complicated.

To summarize, waits-for-retained-lock relationships are evaluated only to detect opening-up deadlocks as early as possible. Since candidate transactions for breaking up the cycle are lock holders and lock requestors, the mechanisms for deadlock resolution can be derived from those provided for single-level transactions.

Early detection of opening-up deadlocks saves transaction work. However, as discussed above, the additional representation and management of waits-for-retained-lock relationships require some overhead. If deadlocks are infrequent, a particular system implementation has to take this trade-off into account.

6. Comparison of Some System Implementations

In the following, we will compare some systems that implement nested transactions with regard to the degree of parallelism supported, the applied concurrency control schemes, and the way deadlocks are treated. In particular, we will consider ARGUS (Liskov et al., 1987; Liskov, 1988), Camelot (Spector et al., 1988; Epinger et al., 1991), Clouds (Ahamad et al., 1987; Dasgupta et al., 1989), Eden (Almes et al., 1985; Pu and Noe, 1988), and LOCUS (Mueller et al., 1983; Weinstein et al., 1985). Table 6 summarizes the results.

While Camelot, Clouds, Eden and LOCUS allow for parent/child as well as sibling parallelism, ARGUS does not permit a parent transaction to run in parallel with its children, resulting in simpler locking rules (Liskov et al., 1987). All the systems considered are based on two-phase locking; however, only Clouds and LOCUS support downward inheritance. While the downward inheritance scheme in LOCUS requires a lock holder to explicitly state when downward inheritance may potentially take place, transactions in CLOUDS are allowed to share the locks of their ancestors in a totally uncontrolled manner (Allchin, 1983). When a transaction closes a file in LOCUS, the lock held on this file becomes a retained lock. Although it supports downward inheritance by means of an explicit offering mechanism, the LOCUS scheme is uncontrolled in the sense that a lock holder

Table 6. Comparison of system implementations.

	Parent/child parallelism	Sibling parallelism	Downward inheritance	Controlled downward inheritance	Object hierarchy support	Deadlock avoidance & detection
Argus	no	yes	no	no	no	timeout-based resolution
Camelot	yes	yes	no	no	no	timeout-based resolution
Clouds	yes	yes	yes	no	no	timeout-based resolution
Eden	yes	yes	no	no	no	wound-wait avoidance scheme
LOCUS	yes	yes	yes	no	no	neither resolution nor avoidance

offering a lock cannot control in which mode its descendants may acquire this lock. None of the five systems supports controlled downward inheritance, neither do they support object hierarchies. ARGUS, Camelot, and Clouds implement deadlock resolution based on a timeout mechanism, whereas Eden applies a wound-wait deadlock avoidance scheme (Rosenkrantz et al., 1978). LOCUS neither performs deadlock detection, nor implements an avoidance scheme. However, it provides an interface to operating system data, permitting a system process to detect deadlock by constructing a wait-for-graph. In this manner, different deadlock resolution strategies may be implemented (Weinstein et al., 1985).

7. Conclusions

We have presented an investigation of concurrency control in nested transactions. The primary focus of our article has been on achieving a high degree of intra-transaction parallelism within nested transactions by using locking protocols.

Our initial concurrency control mechanism for nested transactions was based on S-X locking protocols on flat objects which seriously limit parent/child parallelism. Therefore, the concept of downward inheritance was introduced and refined to

controlled downward inheritance in order to enable a transaction to restrict the access mode of its inferiors for an object. Controlled downward inheritance turned out to be a useful concept for achieving safe parent/child cooperation on data structures to be read or written in a shared manner.

Practical applications sometimes have a need for specialized lock modes as well as multi-level object hierarchies offering efficient ways to lock granules of varying sizes. Therefore, we have generalized the locking rules for nested transactions to be applied for richer access modes on flat objects. Most importantly, this kind of generalization was a prerequisite for the integration of transaction and object hierarchies, since the appropriate use of object hierarchies implied suitable access modes beyond S- and X-locks. As a result, we could combine both types of hierarchies in a general concurrency control model and then could enhance the model again, using the concept of controlled downward inheritance, for the even richer set of access modes. Finally, we studied the principles of deadlock detection in nested transactions. In contrast to single-level transactions where the waits-for-lock relation is sufficient to search for waiting cycles among transactions, detection of all deadlocks in nested transactions further requires the maintenance of the waits-for-commit relation and its combined use with the waits-for-lock relation. If deadlocks are frequently anticipated, opening-up deadlocks, which may span transaction trees, should be detected as early as possible to save transaction work. For this purpose, we have additionally introduced the waits-for-retained-lock relation.

Acknowledgments

This research was carried out by the authors as visiting scientists at the IBM Almaden Research Center, San Jose, California. C. Mohan shared his great knowledge and experience on concurrency control with us. We thank him for his contributions, which led to essential simplifications and clarifications of the concepts proposed in the paper. We also thank J. Palmer and P. Schwarz and the VLDB referees for their helpful comments on this paper.

References

- Ahamad, M., Dasgupta, P., Le Blanc, R.J., and Wilkes, C.T. Fault tolerant computing in object based distributed systems. *IEEE 6th Symposium on Reliability in Distributed Software and Database Systems*, 1987.

- Allchin, J.E. An architecture for decentralized systems. Technical Report GIT-ICS83-23, School of Information and Computer Science, Georgia Institute of Technology, 1983.
- Almes, G.T., Black, A.P., Lazowska, E.D., and Noe, J.D. The Eden system: A technical review. *IEEE Trans. Software Engineering*, 11(1):43-58, 1985.
- Anon et al. A measure of transaction processing power. *Datamation*, April, 1985.
- Astrahan, M.M., Blasgen, M.W., Chamberlin, D.D., Eswaran, K.P., Gray, J.N., Griffith, P.P., King, W.F., Lorie, R.A., McJones, P.R., Mehl, J.W., Putzolu, G.R., Traiger, I.L., Wade, B., and Watson, V. System R: Relational approach to database management. *ACM TODS*, 1(2):97-137, 1976.
- Bancilhon, F., Kim, W., and Korth, H.F. A model of CAD transactions. *Proceedings of the 11th International Conference on VLDB*, Stockholm, 1985.
- Beeri, C., Bernstein, P.A., and Goodman, N. A model for concurrency in nested transaction systems. *Journal of the ACM*, 36(1):230-269, 1989.
- Bernstein, P.A. and Goodman, N. Concurrency control in distributed database systems. *ACM Computing Surveys*, 13(2):185-221, 1981.
- Bernstein, P.A., Hadzilacos, N., and Goodman, N. *Concurrency Control and Recovery in Database Systems*, Addison-Wesley: Reading, Pennsylvania, 1987.
- Dasgupta, P., Le Blanc, R., and Appelbe, W. The Clouds distributed operating system: Functional description, implementation details and related work. *IEEE 8th International Conference on Distributed Computing Systems*, San Jose, California, 1989.
- Eswaran, K.P., Gray, J.N., Lorie, R.A., and Traiger, I.L. The notions of consistency and predicate locks in a database system. *Communications of ACM*, 19(11):624-633, 1976.
- Eppinger, J.L., Mummert, L.B., and Spector, A.Z., eds. *Camelot and Avalon—A Distributed Transaction Facility*. Morgan Kaufmann: San Mateo, California, 1991.
- Gawlick, D. Processing "hot spots" in high performance systems. *Proceedings of the IEEE*, Spring CompCon, San Francisco, 1985.
- Gray, J.N. Notes on database operating systems. Operating Systems—An Advanced Course. In: Bayer, R., Graham, R.M., and Seegmueller, G., eds., *Lecture Notes in Computer Science 60*, Springer-Verlag: Berlin, 1978, pp. 393-481.
- Gray, J.N. The transaction concept: Virtues and limitations. *Proceedings of the Seventh International Conference on VLDB*, Cannes, 1981.
- Gray, J.N., Lorie, R.A., Putzolu, G.R., and Traiger, I. Granularity of locks and degrees of consistency in a shared data base. *Proceedings of the IFIP Working Conference on Modeling in Data Base Management Systems*, Freudenstadt, Germany, 1976.

- Härder, T. and Reuter, A. Principles of transaction-oriented database recovery. *ACM Computing Surveys*, 15(4):287-318, 1983.
- Härder, T. and Rothermel, K. Concepts for transaction recovery in nested transactions. *Proceedings of the ACM SIGMOD*, San Francisco, 1987.
- Jessop, W.H. The EDEN transaction-based file system. *Proceedings of the Second Symposium on Reliability in Distributed Software and Databases*, Pittsburgh, 1982.
- Kim, W., Lorie, R., McNabb, D., and Plouffe, W. Nested transactions for engineering design databases. *Proceedings of the Tenth International Conference on VLDB*, Singapore, 1984.
- Liskov, B. The ARGUS language and system. Distributed systems—Methods and tools for specification: An advanced course. In: Paul, M. and Siegart, H.J., eds., *Lecture Notes in Computer Science 190*, Springer-Verlag: Berlin, 1985, pp. 343-430.
- Liskov, B. Distributed programming in ARGUS. *Communications of the ACM*, 31:300-312, 1988.
- Liskov, B., Curtis, D., Johnson, P., and Scheifler, R. Implementation of ARGUS. *Proceedings of the 11th Symposium on Operating System Principles, ACM Operating Systems Review*, 21(5):111-122, 1987.
- Mohan, C., Lindsay, B., and Obermark, R. Transaction management in R* distributed data base management systems. IBM Research Report #RJ5037, San Jose, California, 1986.
- Moss, J.E.B. *Nested Transactions: An Approach to Reliable Distributed Computing*. M.I.T. Press, Cambridge, 1985.
- Moss, J.E.B., Griffeth, N.D., and Graham, M.H. Abstraction in recovery management. *Proceedings of the International Conference on Management of Data (SIGMOD)*, Washington, D.C., 1986.
- Moss, J.E.B. Log-based recovery for nested transactions. *Proceedings of the 13th VLDB Conference*, Brighton, 1987.
- Mueller, E.T., Moore, J.D., and Popek, G.A. Nested transaction mechanism for LOCUS. *Proceedings of the 9th Symposium on Operating Systems Principles, ACM/ SIGOPS*, Bretton Woods, 1983.
- Pu, C. and Noe, J.D. Nested transactions for general objects: The Eden implementation. Department of Computer Science, University of Washington, TR-85-12-03, 1988.
- Rahm, E. A Framework for Workload allocation in distributed transaction processing systems. *Journal of Systems and Software*, 18:171-190, 1992.
- Reuter, A. Concurrency on high-traffic data elements. *Proceedings on the Conference on Principles of Database Systems*, Los Angeles, 1982.

- Rosenkrantz, D.J., Stearns, R.E., and Lewis II, P.M. System level concurrency control for distributed database systems. *ACM Transactions on Database Systems*, 3(2): pages? 1978.
- Rothermel, K. and Mohan, C. ARIES/NT: A recovery method based on write-ahead logging for nested transactions. *Proceedings of the 15th International Conference on Very Large Data Bases*, Amsterdam, 1989.
- Rukoz, M. Hierarchical deadlock detection for nested transactions. *Distributed Computing*, 4:123-129, 1991.
- Schwarz, P.M. and Spector, A.Z. Synchronizing shared abstract types. *ACM TOCS*, 2(3):223-250, 1984.
- Spector, A.Z. and Schwarz, P.M. Transactions: A construct for reliable distributed computing. *Operating Systems Review*, 17(2):18-35, 1983.
- Spector, A.Z., Pausch, R.F., and Bruell, G. Camelot: A flexible, distributed transaction processing system. *Proceedings of the Spring IEEE COMPCON*, San Francisco, 1988.
- Walter, B. Nested transactions with multiple commit points: An approach to the structure of advanced database applications. *Proceedings of the Tenth International Conference on VLDB*, Singapore, 1984.
- Weikum, G. A theoretical foundation of multi-level concurrency control. *Proceedings of the ACM SIGACT-SIGMOD: Symposium on Principles of Database Systems*, Cambridge, Massachusetts, 1986.
- Weikum, G. Principles and realization strategies of multilevel transaction management, *ACM TODS*, 16(1):132-180, 1991.
- Weikum, G. and Schek, H.-J. Architectural issues of transaction management in layered systems. *Proceedings of the Tenth International Conference on VLDB*, Singapore, 1984.
- Weinstein, M., Page, T., Livezey, B., and Popek, G. Transactions and synchronization in a distributed operating system. *Proceedings of the Tenth Symposium on Operating Systems Principles, ACM/SIGOPS*, Orcas Island, Washington, 1985.