

Performance Characteristics of Epsilon Serializability
with Hierarchical Inconsistency Bounds[†]

Mohan Kamath and Krithi Ramamritham

Department of Computer Science
University of Massachusetts
Amherst, MA 01003, USA

Abstract

Epsilon serializability (ESR) is a weaker form of correctness designed to provide more concurrency than classic serializability (SR) by allowing, for example, query transactions to view inconsistent data in a controlled fashion *i.e.* limiting the inconsistency within the specified bounds. In the previous literature on ESR, inconsistency bounds have been specified with respect to transactions or with respect to objects. In this paper, we introduce the notion of hierarchical inconsistency bounds that allows inconsistency to be specified at different granularities. The motivation for this comes from the way data is usually organized, in hierarchical groups, based on some common features and interrelationships. Bounds on transactions are specified at the top of the hierarchy, while bounds on the objects are specified at the bottom and on groups in between. We also discuss mechanisms needed to control the inconsistency so that it lies within the specified bounds. While executing a transaction, the system checks for possible violation of inconsistency bounds bottom up, starting with the object level and ending with the transaction level.

Thus far, to our knowledge, no work has been done to determine the quantitative performance improvement resulting from ESR. Hence in this paper we report on an evaluation of the performance improvement due to ESR incorporating hierarchical inconsistency bounds. The tests were performed on a prototype transaction processing system that uses timestamp based concurrency control. For simplicity, our implementation uses a two level hierarchy for inconsistency specification - the transaction level and the object level. We present the results of our performance tests and discuss how the behavior of the system is influenced by the transaction and object level inconsistency bounds. We make two important observations from the tests. First, the thrashing point shifts to a higher multiprogramming level when transaction inconsistency bounds are increased. Further, for a particular multiprogramming level and a particular transaction inconsistency bound, the throughput does not increase with increasing object inconsistency bounds but peaks at some intermediate value.

[†] This material is based upon work supported by the National Science Foundation under grant IRI - 9109210

Contents

1	Introduction	1
2	Review of Epsilon Serializability	2
3	Specifying Inconsistency	3
3.1	Hierarchical Inconsistency Bounds	3
3.2	Two Level Specification used in our System	5
3.2.1	Transaction Level	5
3.2.2	Object Level	7
4	Time Stamp based ESR	8
5	Controlling Inconsistency	10
5.1	Import Inconsistency	10
5.2	Export Inconsistency	12
5.3	Other Issues	13
5.3.1	Controlling Inconsistency in a Multi-level Hierarchy	13
5.3.2	Controlling Inconsistency in Different Types of Queries	13
6	Details of the Prototype System	14
7	Objectives of the Performance Tests	15
8	Results and Discussion	16
9	Conclusion	20

1 Introduction

Epsilon Serializability (ESR) is a weaker form of correctness that allows transactions to view inconsistent data in a limited fashion. ESR does this by permitting some non-SR execution schedules, thereby enhancing concurrency. Potential benefits from ESR include increased system availability and greater autonomy [16,17,21]. A transaction can be classified as a *query* transaction or an *update* transaction, depending upon the types of operations involved. In typical scenarios, lengthy epsilon transactions (ETs) that just perform queries may execute in spite of ongoing concurrent update transactions. In such a situation, query ETs may view inconsistent data. Thus an update ET may *export* some inconsistency when it updates a data item while query ETs are in progress; a query ET may *import* some inconsistency when it reads possibly inconsistent data items while update ETs are in progress. The correctness notion in ESR is based on limiting the amount of imported and exported inconsistency within the specified bounds. If the bounds are set to zero, ESR reduces to SR. We review some key aspects of ESR further in section 2, where we also describe the various notations used in this paper.

We first introduce the notion of hierarchical inconsistency bounds which allows inconsistency specification at different levels. From a data organization point of view, data objects can be grouped hierarchically based on some common features and interrelationships. Depending upon the relative importance of the groups, there could be a limit on the inconsistency resulting from update to each of the groups. This forms the motivation for the concept of hierarchical inconsistency bounds. Bounds on transactions are specified at the top of the hierarchy, while bounds on the objects are specified at the bottom and groups in between. It is also necessary to look at ways of specifying these hierarchical inconsistency bounds. Section 3 discusses these aspects in more detail using a practical example. It also discusses the two-level inconsistency bounds used in our implementation — the transaction level and the object level.

Though update ETs can view inconsistent data the same way query ETs do, in this paper we focus our attention on the situation where query ETs run concurrently with *consistent* update ETs in a centralized system.

Just like SR, ESR can be implemented using one of the many concurrency control mechanisms available. In this paper, we discuss the various control aspects of ESR based on timestamp ordering as it forms the basis of concurrency control in our experiments. Section 4 gives further details of time stamp based ESR where we identify the situations where inconsistent operations will be allowed to execute, provided they do not violate the inconsistency bounds.

We present detailed mechanisms for controlling the inconsistency seen by the queries within the specified bounds at different levels of a hierarchy in section 5. Wu et. al [21] suggest some general ways by which the amount of inconsistency can be controlled in ETs. However our notion of hierarchical inconsistency control attempts to provide a finer-grained approach to controlling inconsistencies in ETs.

ESR is relatively a new concept and its performance has been predicted in the literature [16,21] but never evaluated. So another contribution of our work is the evaluation of the quantitative performance improvement resulting from ESR by changing the various settings of the inconsistency bounds.

There has been considerable work in the area of performance evaluation of concurrency control protocols. Notably, Agrawal et. al [1] discuss the performance of concurrency control protocols and present results from exhaustive performance tests conducted on various cases. Cordon and Garcia-Molina [6] have studied the performance of a concurrency control mechanism, based on semantic knowledge about transactions. Badrinath and Ramamritham [2] evaluated the performance of multilevel concurrency control protocols to enhance concurrency by exploiting the synchronization properties and structure of operations. While most of these use simulation models in evaluating the performance, we have actually built a prototype system, to do a more realistic appraisal of ESR. Our prototype transaction processing system is based on the client-server model and uses timestamp ordering for concurrency control as mentioned earlier. Further details of our prototype system are presented in Section 6.

Section 7 discusses the objectives of the performance tests. We look at some of the performance metrics and the various settings for the tests. Results of the performance tests are presented in Section 8. Here, we also discuss the behavior of the system and the effects of the various inconsistency bounds on system performance. Section 9 concludes with a summary of the paper.

2 Review of Epsilon Serializability

In this section we just touch upon some of the key aspects of ESR that are necessary for our discussion. While [21] gives the details of ESR, a more detailed and formal description can be found in [19].

To an application designer and transaction programmer, an ET is a classic transaction with the addition of inconsistency limits. A query ET has an *import-limit*, which specifies the maximum amount of inconsistency that can be imported by it. Similarly an update ET has an *export-limit* that specifies the maximum amount of inconsistency that can be exported by it. These notions have been used widely in [16,17,19,21]. Henceforth these will be referred to as TIL (*transaction-import-limit*) and TEL (*transaction-export-limit*) respectively. For a query ET Q_i and an update ET U_j they are denoted as TIL_{Q_i} and TEL_{U_j} respectively. In addition to the overall inconsistency experienced by a transaction, bounds could be placed, for example on inconsistency due to each object. These are called OIL (*object-import-limit*) and OEL (*object-export-limit*). Given a query Q , for an object x they are denoted as OIL_x^Q and OEL_x^Q respectively. OIL_x^Q specifies the maximum amount of inconsistency the reads of query Q can view with respect to object x and similarly OEL_x^T specifies the maximum amount of inconsistency the writes of an update transaction T can export with respect to object x . OIL_x^Q is represented as $import_limit_x^Q$ in [19]. In the rest of this paper OIL_x^Q will be assumed to be the same for all Q and is denoted by OIL_x . Similarly OEL_x^T will be assumed to be the same for all T and will be denoted by OEL_x . The notations just introduced will be used in the rest of the paper to denote the various bounds. A more detailed discussion on the bounds appears in sections 3 and 5. An application designer specifies the limits for each ET and the system ensures that these limits are not exceeded during execution.

A database state is a set of data values. A database state space is a set of all database states. ESR is applicable to a database state space S_{DB} if it is a metric space *i.e.* it has the following

properties :

- A distance function $distance(u, v)$ is defined over every pair of states u, v in the state space S_{DB} . This gives the absolute value of the difference between two states of data items.
- Has symmetry: $\forall u, v \in S_{DB}, distance(u, v) = distance(v, u)$
- Satisfies the triangle inequality. $\forall u, v \in S_{DB}, distance(u, v) + distance(v, w) \geq distance(u, w)$.

Many database state spaces have such a regular geometry and examples of this includes dollar amount of bank account and airplane seats in airline reservation systems. Hence the distance can be measured from any database state ‘ a ’ to any database state ‘ b ’. During the execution of each ET, the system needs to maintain the amount of inconsistency the ET has imported so far. The amount of inconsistency is given by the distance function and the incremental accumulation of inconsistency depends upon the triangle inequality property of metric spaces. Without triangle inequality, we would have to recompute the distance function for the entire history each time a change occurs. The system should ensure that all the inconsistencies viewed while ETs execute are within the limits. As mentioned earlier, when these limits are set to zero ESR reduces to SR.

3 Specifying Inconsistency

Inconsistency in a data item x with respect to a query q is defined as the difference between the current value of x and the value of x if no updates on x were allowed to execute concurrently. We will denote this by $Inconsistency_x$. When an inconsistency bound, $Limit_x$, is specified for a data item x with respect to query q then $Inconsistency_x$ should be less than or equal to $Limit_x$. In the previous literature on ESR, inconsistencies have been specified with respect to transactions [21] or with respect to objects [19]. Here we first present the concept of hierarchical inconsistency bounds and then focus on specifying inconsistency in a two level system - *i.e.* at the root and the leaf levels (as we have done in our tests using the prototype system).

3.1 Hierarchical Inconsistency Bounds

In a banking system or an airline reservation system, data can be grouped hierarchically based on some commonality. Hence while executing a transaction by reading objects, apart from specifying a bound on the overall inconsistency experienced by the transaction, bounds can also be specified on the inconsistency arising from groups and data items. Essentially the constraints are hierarchically ordered, the inconsistency limits being associated with the *nodes*¹ of the tree. It should be noted that, the data items (objects) which are accessed are actually present only at the leaf level and the intermediate nodes just represent groups. Using the transitive relationship between the nodes and its children, it can be seen that, if a node m has a inconsistency limit \mathcal{L} , then the sum of the inconsistencies viewed by all the leaves (data items) of the subtree rooted at node m is less than or equal to \mathcal{L} . The inconsistency limit on a node at level i , places a bound on

¹a node could represent a data item (object) or a group of data items.

the inconsistencies viewed by nodes at level $i + 1$. Also at a particular level j , the inconsistency limit on a node places bounds on the inconsistency tolerable by its siblings. The information flow is top-down during the inconsistency specification stage and bottom-up during the control stage. During the control stage if the inconsistency bounds are violated at any level of the tree, then the transaction has to be aborted.

To motivate the need for hierarchical inconsistency bounds we look at a practical situation. *Figure 1* pertains to a banking system. For accounting purposes, the bank categorizes the accounts broadly into some classes such as *company*, *preferred customer*, *personal accounts* etc.. Let us consider the case when the bank needs to estimate the *overall* amount held by the bank. Apart from specifying the inconsistency bound on the *overall* estimate ($Limit_{overall}$), it can specify inconsistency bounds for each of the categories ($Limit_{company}, Limit_{preferred}, Limit_{personal}$ etc..). Further, each category could be subclassified and can be assigned an inconsistency bound ($Limit_{com1}, Limit_{com2}$ etc..). This will proceed till the inconsistency is specified at the leaf level ($Limit_{div1}, Limit_{div2}$, etc..). It should be noted that all these bounds are with respect to an individual transaction (query).

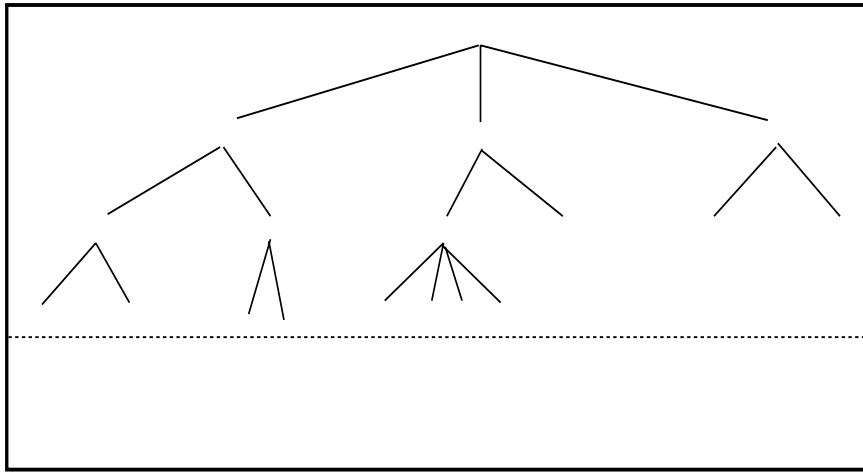


Figure 1: Hierarchy of objects in banking systems

Thus the inconsistency limits specified in the above hierarchy imply the following :

$$\begin{aligned}
 Inconsistency_{company} + Inconsistency_{preferred} + \dots + Inconsistency_{personal} &\leq Limit_{overall} \\
 \sum_{i=1}^n Inconsistency_{com_i} &\leq Limit_{company} \\
 \sum_{j=1}^n Inconsistency_{div_j} &\leq Limit_{company_k} \text{ etc...}
 \end{aligned}$$

There are a many different ways of specifying inconsistency limits at multiple levels. Inconsistency bounds could also be specified using relative weights for the nodes in the tree and thus the inconsistency viewed at a level will be the weighted sum of the inconsistency of the nodes at that level. Examples of hierarchical data organization also occur in other environments where transactions are commonly used, for e.g., airline reservation.

Having introduced the concept of hierarchical inconsistency bounds, we outline a method to specify it. One such generic specification is shown in *Figure 2*. Here a transaction accesses objects,

some of which may be independent and some could be part of a group.

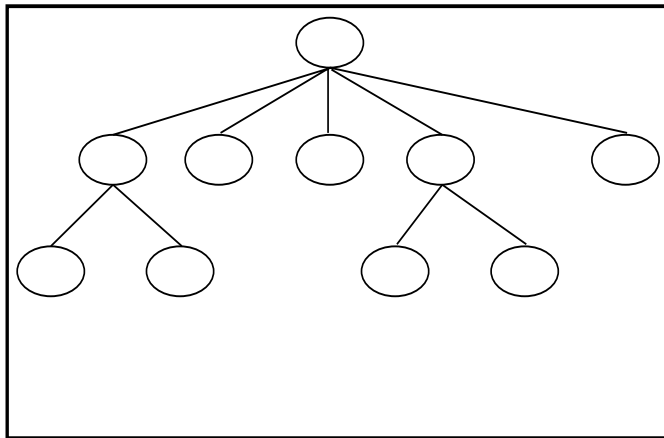


Figure 2: Specifying inconsistency at multiple levels.

GIL_i specifies the total inconsistency that can be contributed from operations that access objects within a particular group i . Groups in turn can contain subgroups. Thus the total inconsistency that is viewed by the transaction is the sum of the inconsistencies from the operations that access the independent objects and objects that are parts of groups. Here is a query specifying hierarchical inconsistency bounds (the syntax of Begin, Read etc.. are explained in later sections).

```
BEGIN Query TIL 10000
  LIMIT company 4000
  LIMIT preferred 3000
  LIMIT personal 3000
  LIMIT com1 200
  ...
  t1=Read com2745
  t2=Read com4639
  ...
END
```

Thus each transaction could have an inconsistency specification part at the beginning before the actual operations are submitted. It should be mentioned that hierarchical specification and control does not come free of charge and a small price is to be paid for it as we discuss in section 5. We now discuss a two level system used in our implementation.

3.2 Two Level Specification used in our System

3.2.1 Transaction Level

Let us first consider inconsistency at the transaction level. A large bank may wish to know how many dollars, in millions, are there in the corporate accounts of its clients. If this query is executed directly on the accounts during banking hours, there would be substantial interference from the ongoing updates. Most of the interference is irrelevant, however, since typical updates refer to small amounts compared to the query's results, which is in millions of dollars. This can

be exploited using ESR, to execute the query during banking hours with the desired accuracy. Specifically, under ESR, if we specify a TIL for the query ET, for example, of \$ 100,000, for this query, the result also would be guaranteed to be within \$ 100,000 of a consistent value (produced by the serial execution of the same transactions). For example, if the ET returns the value \$ 450,500,000 then at least one of the serial transaction executions would have yielded a serializable query result in the \$ 450,500,000 \pm \$ 100,000 interval.

Throughout our discussion, for simplicity, we assume that an object is read or written once within a transaction and that the transaction calculates the sum of the read values. To allow multiple operations on the same object in a transaction, a simple extension to the control mechanism should suffice. This would calculate the inconsistency based on the values maintained for the minimum and maximum seen thus far by reads on an object from a query transaction (writes, in the case of an update transaction). This will take care of the worst case where, two reads would see the positive and negative extremes of the bound and hence could violate the inconsistency bounds. We discuss this further in section 5, where we also discuss the effect of inconsistency on transactions that compute results other than the sum.

Import Inconsistency:

Bounds on imported inconsistency is specified for query ETs. If a query ET Q_i has k read operations, we define its total inconsistency to be $\sum_{i=1}^k Inconsistency_from_Read_i$. Since TIL denotes the total inconsistency that can be tolerated for a query,

$$\sum_{i=1}^k Inconsistency_from_Read_i \leq TIL_{Q_i}$$

We define the *proper* value of an object that is being read by a query ET as the value that would have been seen by the read if no concurrent update ETs were allowed.

Inconsistency_from_Read is defined as the difference between the *present* value of the object and the *proper* value of the object. This definition is applicable only when there are concurrent update ETs. The concept of *proper* and *present* value of an object is discussed in more detail in the section on controlling inconsistency. In this section, wherever appropriate, we have included examples of the query/update ETs used in the evaluation. The following is a query ET. The Read operation takes the *id* of the object whose value is to be read. The BEGIN operation takes the *type* of the transaction and the corresponding inconsistency bound as the arguments.

```

BEGIN Query TIL = 100000
  t1 = Read 1863
  t2 = Read 1427
  t3 = Read 1912
  t4 = Read 1543
  t5 = Read 1657
  t6 = Read 1138
  t7 = Read 1729
  t8 = Read 1336
  output('Sum is : ',t1+t2+t3+t4+t5+t6+t7+t8)
COMMIT

```


Export Inconsistency:

Bounds on exported inconsistency is specified for update ETs. TEL denotes the total inconsistency that can be tolerated for an update ET. If an update ET U_j has k write operations, we define its total inconsistency to be $\sum_{j=1}^k Inconsistency_from_Write_j$. Since TEL denotes the total inconsistency that can be tolerated for a update,

$$\sum_{j=1}^k Inconsistency_from_Write_j \leq TEL_{U_j}$$

We define *new* value as the value that would be written by a write if it is allowed to execute. *Inconsistency_from_Write* is defined as the difference between the *new* value of the object and the *proper* value of the object. Our notion of this is slightly different from the one mentioned in [21] and we discuss this further in section 5. This definition is applicable when there are concurrent query ETs. An update ET submitted to the system for processing is shown below. Note that in this, the value of the writes are dependent upon the reads and hence reads from update ETs must be consistent. The argument for the Write operations are the *id* of the object and the *value* to be written.

```
BEGIN Update TEL = 10000
  t1 = Read 1923
  t2 = Read 1644
  Write 1078 , t2+3000
  t3 = Read 1066
  t4 = Read 1213
  Write 1727 , t3-t4+4230
  Write 1501 , t1+t4+7935
COMMIT
```

3.2.2 Object Level

We just saw how inconsistency limit can be specified at the transaction level. Now we go to a finer granularity by specifying the inconsistency limits on objects. We use a client-server model for our implementation and usually transaction limits are specified at the client side and object limits at the server side, but this could be overridden by explicitly specifying the object limits in the specification stage at the beginning of the transaction.

Import Inconsistency:

OIL is defined as the maximum amount of inconsistency that can be viewed (imported) by read operations of a query ET on an object.

$$Inconsistency_from_Read(x) \leq OIL_x$$

Export Inconsistency:

OEL is defined as the maximum amount of inconsistency that can be exported by the write operations from an update ET on an object.

$$Inconsistency_from_Write(x) \leq OEL_x$$

Thus a system that implements ESR should ensure that the various inconsistencies are within the bounds. In our system which implements a two-level hierarchy, for an operation to succeed, at the object level

$$object_import_inconsistency \leq object_import_limit \text{ (Read from Query)}$$

$$object_export_inconsistency \leq object_export_limit \text{ (Write from Update)}$$

and at the transaction level

$$transaction_import_inconsistency \leq transaction_import_limit \text{ (Read from Query)}$$

$$transaction_export_inconsistency \leq transaction_export_limit \text{ (Write from Update)}$$

We had a fairly detailed look at specifying hierarchical inconsistency which will be usually done by the application programmer/user. The control mechanisms which ensures that the inconsistency limits are not violated are built into the system and are not accessible to the application programmer.

4 Time Stamp based ESR

Our goal was to build a simple prototype system to measure the performance of ESR — especially its impact on concurrency. Hence we chose timestamp ordering for concurrency control to avoid the problem of deadlock detection and recovery that is present in the case of 2PL. To simplify recovery, we enforce *strict* ordering [5] by using a wait based protocol for concurrent operations that are not able to execute. For late operations that are not able to execute we do aborts with immediate restarts. Hence we don't have the need to maintain detailed histories/logs since recovery is simple and rollbacks are not necessary. Recall that timestamps are assigned when transactions begin.

We do not go into the details of timestamp based SR, but all standard terminologies and techniques are applicable here. The algorithm used is an enhanced version of the regular time stamp ordering algorithm used to implement SR. The enhancements come in places where the algorithm would normally abort the transaction because the operation was late or some other conflicting operation is running concurrently that prevents the current operation from being executed. The detailed algorithm for implementing ESR using time stamp ordering is shown in *Figure 3*. The enhancements that have been incorporated to implement ESR have been highlighted. In the ESR case (highlighted), the operations are allowed to take place, provided the inconsistency is within the inconsistency (epsilon) bounds. In Section 5 we discuss in detail how to check whether the inconsistency viewed by an ET has crossed the inconsistency bounds. Referring to *Figure 3*, it can be seen that there are three cases where operations that will normally be rejected in SR, will be given a chance to perform under ESR. The first two cases correspond to reads from query ETs and the last one corresponds to a write from update ETs, as mentioned below:

1. A query read views committed data but the query ET's timestamp is older than the object's last-write timestamp.

2. A query read views uncommitted data from a concurrent update ET.
3. A write from an update ET arrives with a timestamp that is older than the object last-read timestamp and the last read was from a query ET

Case 1 can also be considered to occur due to a concurrent update because if the read operation has an older time-stamp and it read a committed value with a newer timestamp then the two transactions must have been concurrent at some point of time. A similar argument applies to case 3 where the last read is required to be from a query ET because we allow only consistent update transactions and hence, reads from one update transaction will conflict with writes from another update transaction.

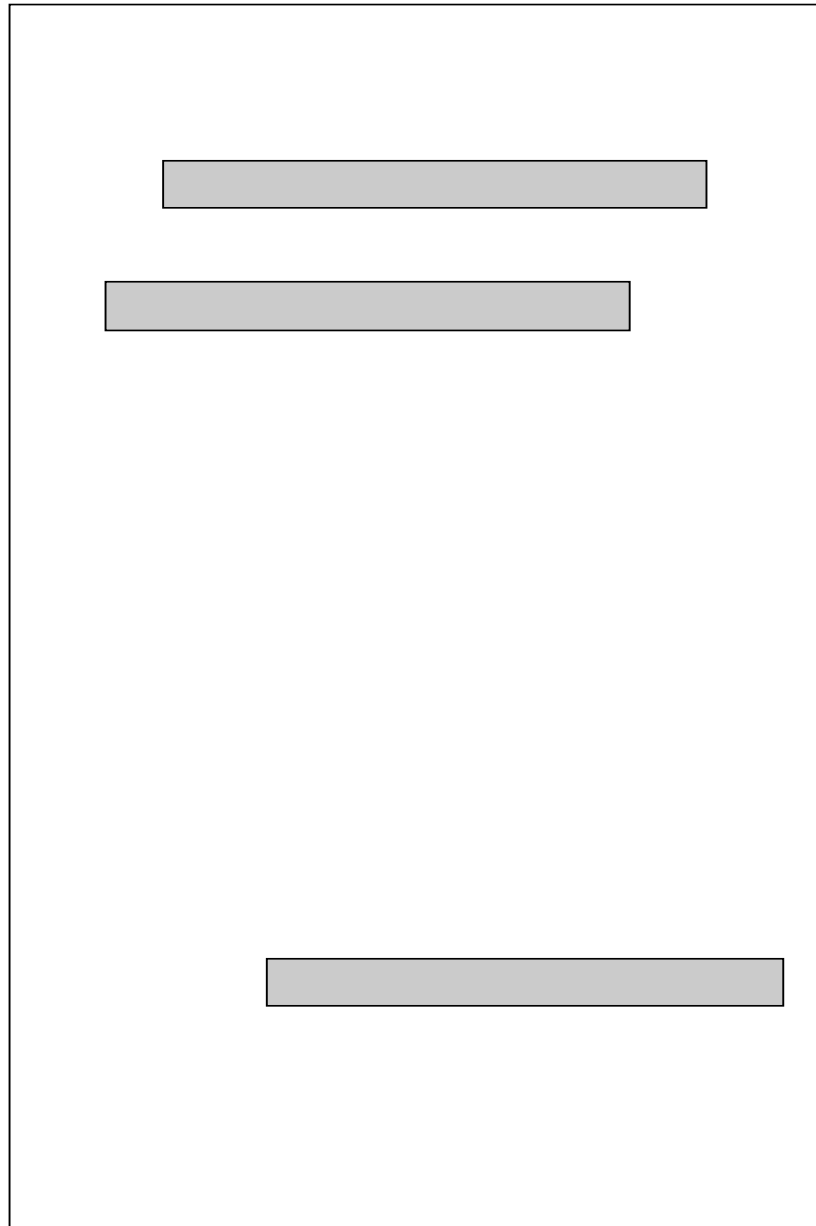


Figure 3: Algorithm for Time Stamp based ESR.

Figure 4 gives some examples of the three cases of concurrent updates discussed above. The

first occurs when $R_{Q1}(x)$ (read from query $Q1$) comes after the value written by $W_{U4}(x)$ has committed, the second when $R_{Q2}(x)$ views uncommitted data written by $W_{U3}(x)$ and the third when $W_{U2}(x)$ needs to write when $R_{Q3}(x)$ has already read it.

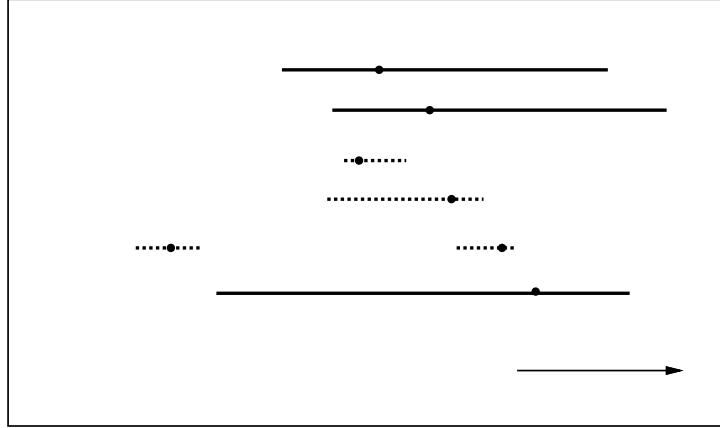


Figure 4: Time Stamp ordering for ESR - ETs with conflicts on object x

In all these cases the reads from query ETs view inconsistent data states because of out of order timestamps. The magnitude of this inconsistency must be controlled within the specifications. In the previous section we discussed ways to specify the allowable inconsistency. The following section deals with the mechanisms to meet these specifications.

5 Controlling Inconsistency

Throughout this discussion we have seen that in the ESR case, for an operation to succeed, the inconsistency should be controlled within the specified bounds. In this section we describe the timestamp based mechanisms we have used to control the inconsistencies. Some methods suggested in the literature [21], where the total inconsistency of a query ET is determined by multiplying the bound of an operation and the number of operations in the ET may result in the overestimation of the accumulated errors. We present detailed mechanisms for controlling the inconsistency at fine grain levels.

In keeping the prototype system simple and to perform controlled experiments, we have a two level specification and we look at only one type of query, those that compute the *sum*. At the end of the section, we address the issues related to controlling inconsistency when there is a hierarchical specification and when there are queries of different types.

5.1 Import Inconsistency

Let us consider the scenario shown in *Figure 5*. We will focus on the read operation $R_{Q1}(x)$. We defined earlier the *proper* value of an object that is being read by a query ET as the value it would have seen had there been no concurrent update transactions. For $R_{Q1}(x)$, the *proper* value is the value written by $W_{U1}(x)$ as that was the last write on x before $Q1$ started (which we will call $\mathcal{P}1$). In [19] the *proper* value of object x is denoted by $x_{initial}^Q$. $U2$, $U3$ and $U4$ are the update ETs

that run concurrently with $Q1$. If $R_{Q1}(x)$ is successful it will read the value written by $W_{U4}(x)$ which is the *present* value (which we will call $\mathcal{N}4$).

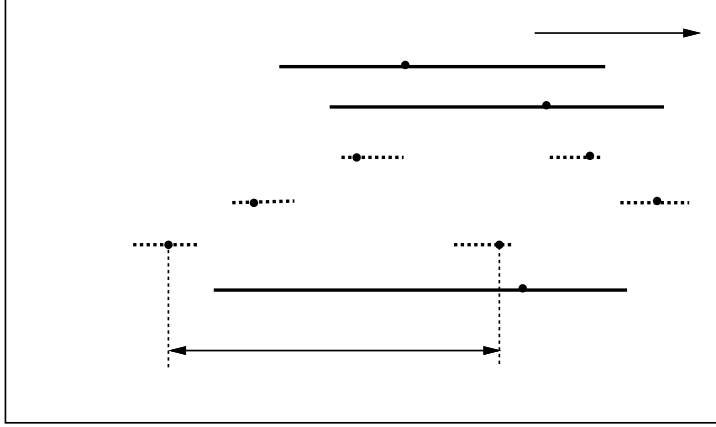


Figure 5: ESR: Check for Read (from Query ET) to go through

In the figure, d^2 ($= \mathcal{N}4 - \mathcal{P}1$) denotes the difference between the *present* and the *proper* values and hence is the *magnitude* of the inconsistency in the data that is being read by $R_{Q1}(x)$ (*Inconsistency_from_read*). Recall that the distance function gives the difference between any two states (between states 1 and 2, 2 and 3, 1 and 3, etc.). If the magnitude of the total inconsistency imported by $Q1$ before $R_{Q1}(x)$ is denoted by \mathcal{I} then for $R_{Q1}(x)$ to succeed the following needs to be satisfied in our two level system:

$$\begin{aligned} d &\leq OIL_x && \text{(object-level)} \\ \mathcal{I} + d &\leq TIL_{Q1} && \text{(transaction-level)} \end{aligned}$$

If these conditions are satisfied then $R_{Q1}(x)$ is successful and \mathcal{I} is incremented by the value d and $Q1$ proceeds, else it fails and has to be aborted.

If changes made by W_{U4} were negative, and $U4$ aborts then we could have a problem. The *current* value would then correspond to the one written by $W_{U3}(x)$ which could violate the inconsistency bound. One solution to this would be to always add the maximum change by an update transaction [19] while determining the inconsistency seen by an operation. The probability of aborts from update transactions in our case is small and hence in our implementation we do not consider this.

In our implementation we store the values of the last 20 writes on each object with the corresponding time stamps. The *proper* value of an object is found by indexing backwards through this list until a older timestamp (than the query) is found. 20 is an empirical figure derived by dividing the measured values of the average duration of query ETs by that of update ETs.

It should be noted however that this is not the same as multi-version timestamp ordering (MVTO). In the MVTO case, timestamped versions are maintained so that if a read operation arrives late, based on the versions, the value written by the last write with a timestamp lesser than this read is returned. However in our case, the value read is the value of the current instance

²In [19] d is defined as $distance(x_{current}x_{initial}^Q)$.

of the object which is the *present* value. The value written by the last write with a timestamp lesser than this read, the *proper* value, is only used in determining the amount of inconsistency viewed by the read operation.

5.2 Export Inconsistency

We now move to *Figure 6*, focusing our attention on $W_{U5}(x)$. $Q1$, $Q2$ and $Q3$ are three query ETs that are running concurrently with update ET $U5$. If $W_{U5}(x)$ is executed, it exports some amount of inconsistency to those concurrent query ETs that have accessed x .

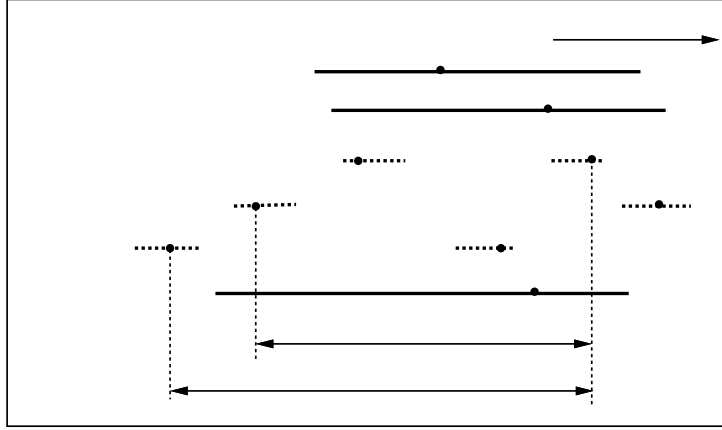


Figure 6: ESR: Check for Write to go through

The *proper* values of x for $R_{Q1}(x)$ is the value written by $W_{U1}(x)$, for $R_{Q2}(x)$ and $R_{Q3}(x)$ it is the value written by $W_{U2}(x)$. For each object x , we maintain a list of uncommitted query ETs which have read its value, along with the respective *proper* values (in practical environments, a better mechanism could be used). Let us denote these values as $\mathcal{P}1$, $\mathcal{P}2$ and $\mathcal{P}3$ for object x corresponding to $Q1$, $Q2$ and $Q3$ respectively. Let $\mathcal{N}5$ denote the *new* value that will be written by $W_{U5}(x)$ if it is allowed to execute. It should be noted that in this particular example, $\mathcal{P}1$ and $\mathcal{P}2$ happen to be the same. Since there are many concurrent query ETs, $W_{U5}(x)$ would export some inconsistency to each of these. Hence the inconsistency that will be exported by $W_{U5}(x)$ denoted by d is chosen as the maximum of the inconsistencies that will be exported to the concurrent query ETs (*Inconsistency_from_write*). In this case, it is the higher of $|d1|$ and $|d2|$ whose values are $(\mathcal{N}5 - \mathcal{P}1)$ and $(\mathcal{N}5 - \mathcal{P}2)$ respectively. Our notion of *Inconsistency_from_write* differs from [21] at this point. In [21] d is determined by adding the inconsistencies exported to all the concurrent query ETs. We chose to do it this way, as we only have a maximum of one read on an object per transaction as mentioned earlier. If the magnitude of the total inconsistency exported by $Q1$ before $W_{U5}(x)$ is denoted by \mathcal{E} then for $W_{U5}(x)$ to succeed, the following should to be satisfied in our two level system:

$$\begin{aligned} d &\leq OEL_x && \text{(object-level)} \\ \mathcal{E} + d &\leq TEL_{U5} && \text{(transaction-level)} \end{aligned}$$

If these conditions are satisfied then $W_{U5}(x)$ is executed and \mathcal{E} is incremented by the value d and $U5$ proceeds, else it fails and $U5$ is aborted.

5.3 Other Issues

5.3.1 Controlling Inconsistency in a Multi-level Hierarchy

To determine the inconsistency from an operation, the mechanisms explained earlier are to be used. However additional checks are required to determine if an operation can execute without violating the hierarchical inconsistency bounds.

As we mentioned earlier, in a multi-level hierarchy, during the control stage the information flows bottom-up *i.e.* the inconsistency viewed by an object at the leaf level percolates to the root level through intermediate levels. As it percolates, the inconsistency is checked at the various nodes it passes through. Revisiting the banking system shown in *Figure 1*, the following set of checks are required to maintain the inconsistency within the bounds (if the inconsistency in the value read for a division k of company j is d).

$$\begin{aligned} & d \leq Limit_{div_k} && \text{(leaf level)} \\ Inconsistency_{com_j} + d & \leq Limit_{com_j} && \text{(level 3)} \\ Inconsistency_{company} + d & \leq Limit_{company} && \text{(level 2)} \\ Inconsistency_{overall} + d & \leq Limit_{overall} && \text{(root level)} \end{aligned}$$

If there is no violation at any of the above checks, then the total inconsistency accumulated at each of the levels have to be incremented (by making another pass through the hierarchy or some other means) by a value that corresponds to the inconsistency of the operation under consideration. If the bounds are violated at any stage, the operation is unsuccessful and the transaction has to be aborted.

5.3.2 Controlling Inconsistency in Different Types of Queries

So far the mechanism we have explained is sufficient if all the queries perform *sum* operation on the object read. This is not going to work however if the queries required for example is the *average* of the values read. In the case of *average*, the inconsistency in the result will actually depend upon the maximum and minimum values viewed by the reads of the transaction. Though it is not easy to specify general mechanisms for all kinds of queries we show how this can be done for a specific case. In this discussion we will be concerned only with query inconsistency limits and not with object limits as the criterion for object inconsistency is going to remain unchanged. We also assume that an object could be read any number of times during a transaction. We have not implemented this in our system as we were only working with the *sum* operation and hence the performance results hold.

The mechanism would work as follows. For every transaction, we maintain the maximum and minimum values viewed, for each object it accesses. When an aggregate operation is encountered, the actual inconsistency is calculated based on the object's maximum and minimum values viewed. *i.e.* if the operation is $avg(o_1, \dots, o_n)$ then the *min_result* is determined by summing up the minimum values of $o_1 \dots o_n$ and dividing it by n and the *max_result* is obtained similarly.

Then the *result_inconsistency*³ is half the difference between *max_result* and *min_result*. The *result_inconsistency* is compared with the TIL. It is at this point a decision is made whether to reject or execute the operation as opposed to dynamically doing it when objects are read, as was done in the previous discussion. This is a viable solution to this problem as predeclaration of objects to be accessed or number of operations in a query, is not practicable.

As we saw above, there are some overheads in maintaining the various data structures and moving information around for hierarchical inconsistency control. However paying a small cost, the user gets some flexibility in deciding the precision with which the data items are to be estimated.

6 Details of the Prototype System

We used the client server model for our implementation. Multiple transaction clients submit transactions to a central transaction server. The server and each of the clients run on different DECstation 5000's running Unix, all connected in the same LAN. Since the environment supported threads and RPC library routines, the implementation of the prototype system was much simplified. We have about 10 workstations in the LAN and because of the synchronous nature of our RPC, the maximum multiprogramming level is limited to 10.

A null RPC call takes about 11 milliseconds to return while the average RPC call takes somewhere between 17 and 20 milliseconds. This along with the fact that the server is multithreaded, allowed the prototype system to process about 50 to 60 transactions⁴ per second, with each transaction having an average of 10 operations.

In implementing a time stamp ordered mechanism, one of the important functions is the generation of timestamps. As there was a two minute range of variation between the local system clocks of the different client sites, to ensure that the timestamps from all the sites are given a fair treatment, a correction factor was applied to the local time to achieve virtual clock synchronization. Also to ensure that the timestamps were unique, we used the standard technique of appending the site-id's to the timestamp.

The system supports the five basic operations *Read*, *Write*, *Begin*, *Commit* and *Abort*. *Read id* reads the value of object *id* and *Write id, val* changes the value of object *id* to *val*. The recovery part of the system is very simple and is just sufficient enough to recover from aborted transactions. This is so because (i) we have strict ordering of operations and hence do not have any commit/abort dependencies (ii) we use the shadow paging technique while updating the value of the data items, so that if a transaction aborts, instead of rollbacks, all the data items are restored to their previous values and the transaction is started with a new timestamp. However we indirectly pay some price in the form of some delay in the strict ordering, but does not affect the results of our performance tests.

The clients are supplied with data files consisting of a number of transactions that are randomly generated, to serve as the load of transactions. The clients read transactions from the file

³This is explained in more detail in [19].

⁴The number of conflicting operations during these tests were small.

and submit operations to the server successively until all transactions have been processed. If a transaction is aborted the client resubmits it with a new timestamp, and does so, until it is successfully completed. The server primarily consists of a scheduler, a transaction manager and a data manager. The scheduler which acts as the front-end, receives transaction requests from the clients and schedules the operations based on timestamp ordering by submitting it to the transaction manager. The transaction manager passes those operation to the data manager and based on the return value the transaction manager takes the appropriate action. The transaction manager uses many other data structures, like counters that are required for maintaining inconsistency accumulated by transactions etc. The data manager deals with the maintenance of object inconsistency. Objects are defined in a simple way, each has an *id*, a *value* associated with it, and the respective *OIL* and *OEL*. The database is maintained in the main memory on the server side and hence writing an object is simulated by changing its value in memory. This has been done to avoid overheads and complications in recovery if data is stored in files. When the server is invoked, it initializes all the objects by reading the start-up data file. The object limits are actually defined at the server side. These values can be set by the application designer/programmer. The values of *OIL* and *OEL* are randomly generated within a specified range, which is varied while the performance tests on object inconsistency limits are carried out. All the data required for evaluating the performance are available at the server.

7 Objectives of the Performance Tests

The main focus of the first set of tests was on measuring the actual increase in concurrency, resulting from the use of ESR over SR in terms of the throughputs of ESR and SR for various MPL (multiprogramming levels) and different ranges of inconsistency bounds. We also studied some related metrics like the number of retries(aborts), the total number of operations performed (number of reads and writes) and the number of inconsistent operations that succeed (even after viewing some inconsistency) at various MPL and inconsistency bounds. The tests were performed for *high*, *medium*, *low* and *zero* values for the bounds. When we say high we mean high values for both *TIL* and *TEL*. Note that the *zero* bound corresponds to the SR case. During these tests we kept *OIL* and *OEL* constant at high values so that they do not affect the results. As the bounds on inconsistency increase, we would expect the throughput to increase.

The thrashing point (the MPL where the throughput begins to drop) is highly dependent on the conflict ratio. For example, for the conflict ratios considered in [2], thrashing occurs when the MPL is around 30 or 40. However in our case we had to use a higher conflict ratio so that we could observe thrashing at a lower MPL (within 10 to be precise). This has a side effect of producing reduced overall throughputs. Though we had about 1000 objects defined in our database, most of our transactions accessed only about 20 objects to create a high conflict ratio. In our implementation the object values range from 1000 to 9999. The following is the approximate magnitude of the various bounds on inconsistency we used for the first set of tests.

Level	TIL	TEL
high-epsilon	100,000	10,000
medium-epsilon	50,000	5,000
low-epsilon	10,000	1,000

Typical query ETs have about 20 operations per transaction while update ETs have around 6 operations per transaction. Hence the TEL values are on the lower side compared to the TIL values.

We have also studied the affect of different TIL and OIL values, at a given MPL, on system performance. In this study, first the throughput was measured by varying TIL while maintaining the TEL at various constant levels. Next we studied the role of OIL in determining the throughput of the system.

8 Results and Discussion

The performance curves for the various tests are included below. The tests were repeated a few times to eliminate any disturbances caused by the workstation load and network traffic in the LAN. Though we do not include any information about confidence intervals in our graphs, the 90 percent confidence intervals lie within $\pm 3\%$ percentage points of the mean value of the performance metrics shown in the various graphs. The variation here maybe a bit higher compared to the simulation models as there are factors like network load etc. that affect the performance of the prototype system. In the graphs, *epsilon* refers to the inconsistency bounds.

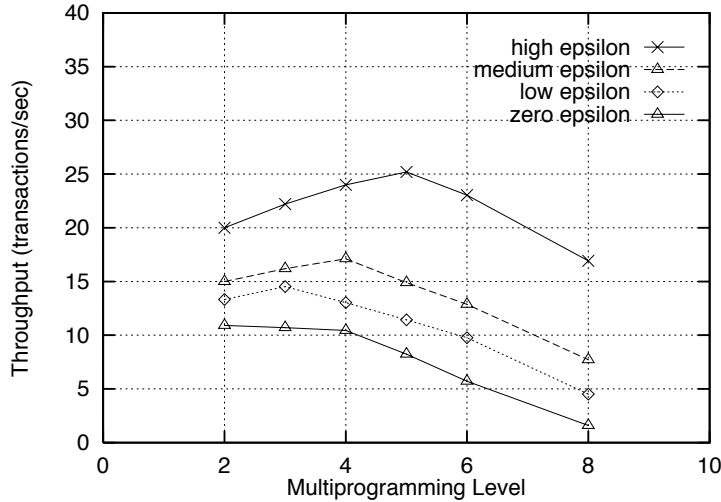


Figure 7: Throughput vs Multiprogramming Level

As expected, in *Figure 7* we see that at higher bounds for inconsistency, the throughput with ESR is much higher than that of SR. As the bounds decrease in value, ESR starts approaching SR. However the key observation that we can make is the shift in thrashing point to higher MPL with increase in inconsistency bounds. Thrashing occurs at a MPL of 3 in the case of lower bounds and it shifts to a MPL of 5 at higher bounds. In all these tests, OIL and OEL were maintained at high values so that they do not affect the throughputs. This significant shift of the thrashing point can

be attributed to some of the other performance metrics noted in the next three graphs. These give us a picture of the dynamics that take place internally when the transactions are processed.

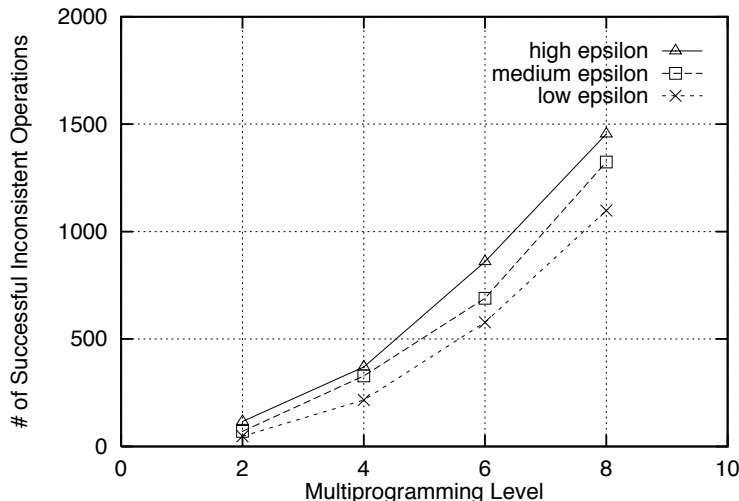


Figure 8: Successful Inconsistent Operations vs Multiprogramming Level

Referring to Figure 8, we see that the number of inconsistent operations that are successful increases with an increase in the inconsistency bounds and the MPL. It should be noted that we do not have the case of *zero epsilon* here as this corresponds to the SR case where inconsistent operations are not permitted. We see a steady increase at each level of inconsistency bound, and this leads us to the next graph.

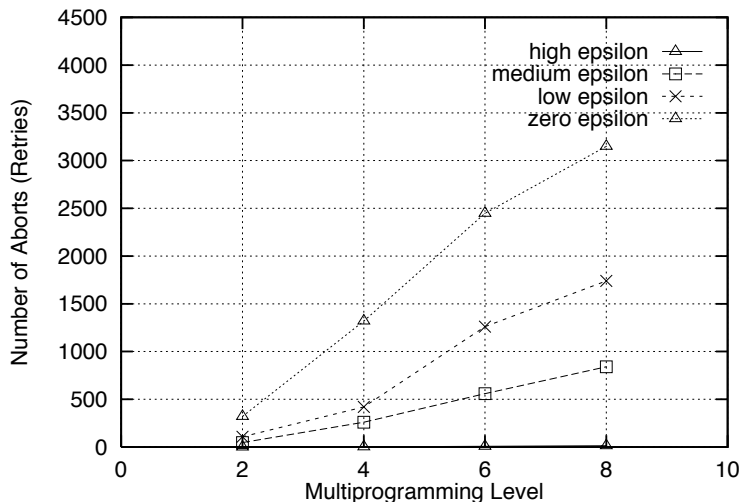


Figure 9: Number of Aborts vs Multiprogramming Level

The number of aborts (retries) that take place in the systems at various levels is shown in Figure 9. A noteworthy observation here is that the number of aborts at high inconsistency bounds is almost *zero*. As expected we see that at lower bounds the number of aborts shoots up rapidly and for the case of *zero epsilon* (SR) the number is very high.

Another performance metric that goes still deeper into the system performance is the total number of operations performed, which is the sum of the read and write operations. This is very useful because it gives an indication of the effort that is wasted by the system, when transactions abort, thereby affecting the throughput. Figure 10 gives the total number of operations executed

at various MPL values. Since we noted in the previous graph that the number of aborts is practically *zero* for the case of high inconsistency bounds, the total number of operations performed in this case should be the same as the actual number of operations required by the transactions.

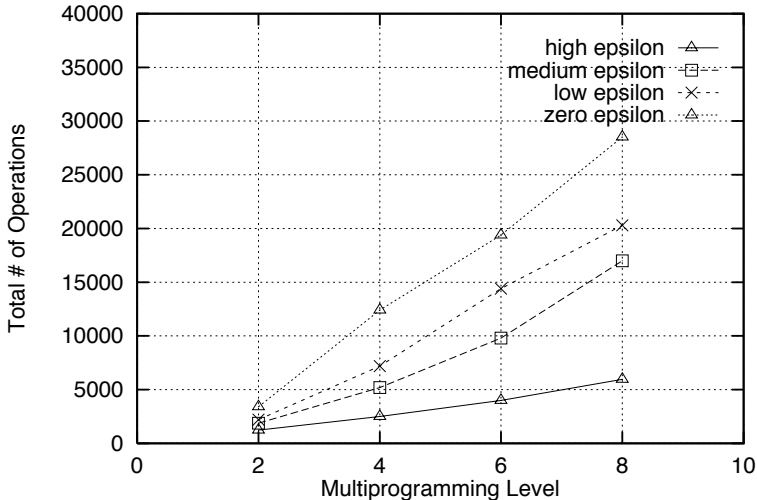


Figure 10: Number of Operations (R+W) vs Multiprogramming Level

Anything that is above this figure gives a measure of the number of useless operations executed by the system, causing a decrease in the throughput. Hence all these graphs give us some insight into the quantitative performance improvement that can result from ESR.

Now we come to our next set of tests where we see how various inconsistency bounds individually affect the performance of ESR. We have studied the affects of the import inconsistency bounds *i.e.* TIL and OIL. All these tests have been performed at a constant MPL of 4.

First we look at the affects of TIL.

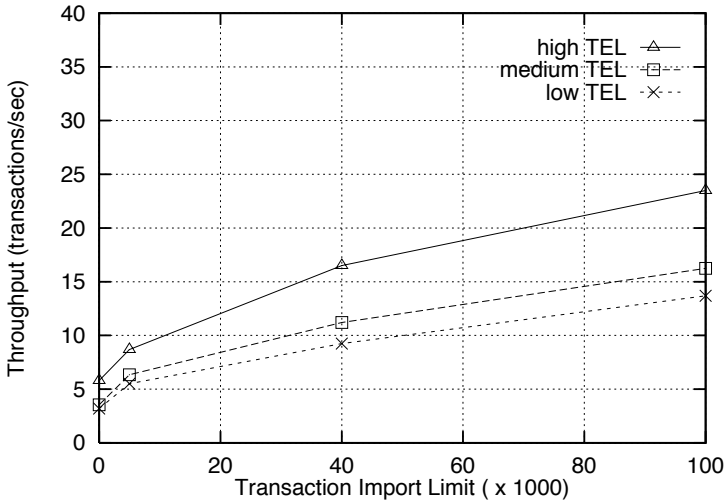


Figure 11: Throughput vs Transaction Import Limit (TEL varies)

Figure 11 shows the variation of throughput with increase in TIL when TEL is held at various constant levels. As expected, the throughput increases with increase in TIL. The slope is the highest at smaller to medium values, since a lot of transactions fall into these categories. But there are some other transactions which are a few in number but need high inconsistency bounds to succeed. They are covered by the higher values of the inconsistency bounds.

The next set of tests study the affect of OIL on the throughput and we have some interesting observations here. Apart from object and transaction inconsistency limits, the average change in value due to a write also affects the throughput. Let the average change in value due to a write be denoted by w . Hence in these graphs, instead of mentioning the absolute value of the bounds, we parameterize it, and specify it in terms of w (we have only done an approximate evaluation of this, as it is difficult to determine this precisely).

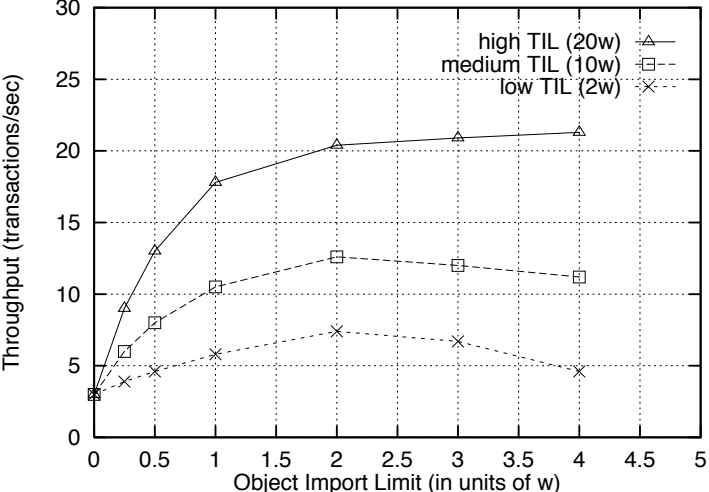


Figure 12: Throughput vs Object Import Limit (TIL varies)

Figure 12 shows the variation of throughput with respect to OIL. The interesting observation here is that in the case of low to medium TIL, the throughput is low at both low and high values of OIL but higher at medium values of OIL. At low OIL, very little inconsistency is tolerated and the throughput is low. At zero OIL, no inconsistency is tolerated on any operation and hence the throughput correspond to the SR case.

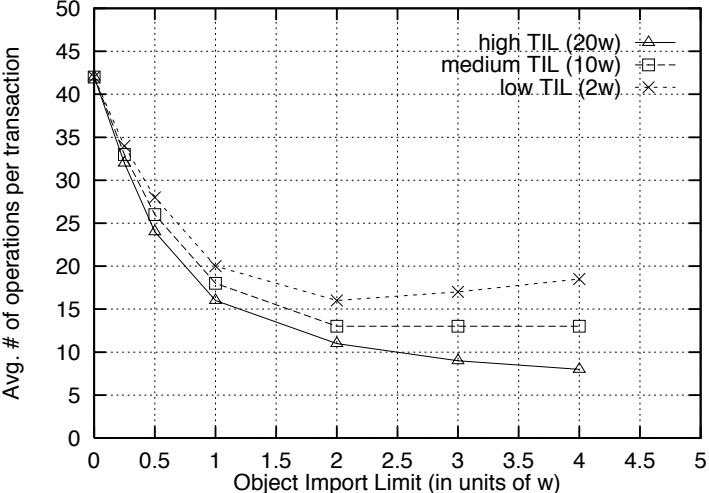


Figure 13: Average number of operations per transaction (TIL varies)

Figure 13 goes into the details of the system that can possibly explain this observation. It gives the average number of operations that are executed for a transaction to complete and this includes the number of operations executed in aborted transactions (which is essentially a waste). For a moment if we ignore TIL, as OIL increases, the number of operations per transaction should keep decreasing and that is what we see in the case of high TIL. However if we consider TIL,

we see its effect slowly creeping in as OIL increases. In fact, for the low TIL case after a certain value of OIL the number of operations per transaction increases. The drop in the throughput at low TIL when the OIL is high can be attributed to the fact that more number of operations are executed by an aborted transaction in this case before it is aborted. The reason is as follows - some operations that cause high inconsistency are not allowed by medium OIL. As a result, the operation fails and the transaction is aborted and resubmitted again. However in the case of high OIL these operations which view high inconsistency are allowed to pass through. In most cases, due to such high inconsistency operations, the total inconsistency viewed by the transaction is drastically increased and because of the low TIL the transaction inconsistency bound is violated and the transaction has to be aborted. Thus the throughput is lower because of the wasted efforts in executing more operations per aborted transaction. Even though the curves indicate that this happens only at low values of TIL, we conjecture that for high TIL also the same would have been observed had we tested with even higher values of OIL. In practice, the object inconsistency limits have to be selected by the application designer/programmer based on the range of values the objects can take and the desired TILs.

We see from the various performance curves that the performance of ESR improves with increases in the transaction inconsistency bounds. Most performance tests yielded the expected results. Interesting observations include the shift in the thrashing point and affects of object inconsistency bounds on performance at different TILs. The actual quantitative performance improvement in an application environment would depend upon the nature of the applications, the typical conflict ratio in those environments etc. In practice, there has to be a compromise between the desired accuracy of the results and the throughput.

9 Conclusion

Previous literature on ESR specified inconsistency bounds with respect to transactions [21] or objects [19]. We have introduced the notion of hierarchical inconsistency bounds that allows inconsistency to be specified at different levels. Using a practical example we demonstrated that with this, a user could gain more flexibility in specifying the inconsistency viewed by ETs. We have also presented detailed mechanisms to control the inconsistency seen by ETs within the specified bounds at different levels of a hierarchy. Wu et. al [21] have suggested some general ways by which the amount of inconsistency can be controlled in ETs. However our notion of hierarchical inconsistency control attempts to provide a finer-grained approach to controlling inconsistencies in ETs.

Several notions of correctness weaker than SR have been proposed previously and evaluated. Ramamritham and Chrysanthis, present a review of these in [18]. Since no work has been done to determine the performance gains from ESR, we have presented an evaluation of the quantitative performance improvements resulting from ESR, obtained through a series of performance tests on our prototype system. For simplicity, we used (i) a two-level hierarchy for inconsistency specification and control *i.e.* the transaction level and the object level and (ii) queries that calculate the *sum* of the values read. Apart from studying the throughput, we studied many other performance metrics (like number of aborts, number of successful inconsistent operations etc.)

that gave us a picture of the dynamics that take place internally when the transactions are processed. While studying the performance of ESR, we made a couple of important observations. First, while studying the throughput of the system at various multiprogramming levels, we noted that the thrashing point shifts to a higher multiprogramming level when inconsistency bounds are raised. Another interesting observation is that, at a particular multiprogramming level, for a given transaction inconsistency bound, the peak throughput occurs at some intermediate value of the object inconsistency bound rather than at low or high values.

This shows that in practice, there has to be a compromise between the desired accuracy of the results with respect to individual objects and the throughput. So far, we have focussed our attention on a centralized server for the performance tests. It will be worthwhile to evaluate ESR in the case of a distributed system with data replication.

Acknowledgements

The Digital Equipment Corporation sponsored Project Pilgrim at the University of Massachusetts provided us the environment for conducting the performance tests.

References

1. R. Agrawal, M.J. Carey and M. Livny. Concurrency control performance modeling: Alternatives and implications. *ACM Transactions on Database Systems*, 12(4): pp 605-654, December 1987.
2. B. R. Badrinath and K. Ramamritham. Performance Evaluation of Semantics-based Multilevel Concurrency Control Protocols. In *Proceedings of ACM SIGMOD Conference on Management of Data*, pp 163-172, May 1990.
3. D. Barbara and H. Garcia Molina. The case for controlled inconsistency in replicated data. In *Proceedings of the Workshop on Management of Replicated Data*, pp 35-42, Houston, November 1990.
4. D. Barbara and H. Garcia Molina. The demarcation protocol: A technique for maintaining arithmetic constraints in distributed database systems. *Technical Report CS-TR-320-91*, Computer Science Department, Princeton University, April 1991.
5. P.A. Bernstein, V. Hadzilacos and N. Goodman. *Concurrency Control and Recovery in Database Systems*. Addison Wesley Publishing Company, first edition, 1987.
6. R. Cordon and H. Garcia-Molina. The Performance of a Concurrency Control mechanism that exploits semantic knowledge. In *Proceedings of the 5th International Conference on Distributed Computing Systems*, pp 350-358, March 1985.
7. W. Du and A. Elmagarmid. Quasi serializability: a correctness criterion for global concurrency control in Interbase. In *Proceedings of the International Conference on Very Large Data Bases*, pp 347-355, Amsterdam, The Netherlands, August 1989.
8. P. Franaszek and J.T. Robinson. Limitations of concurrency in transaction processing. *ACM Transactions on Database Systems*, 10(1): pp 1-28, March 1985.
9. H. Garcia-Molina. Using semantic knowledge for transaction processing in a distributed database. *ACM Transactions on Database Systems*, 8(2): pp 186-213, June 1983.

10. H. Garcia-Molina and K. Salem. Sagas. In *Proceedings of ACM SIGMOD Conference on Management of Data*, pp 249-259, May 1987.
11. H. Garcia-Molina and G. Wiederhold. Read-only transactions in a distributed database. *ACM Transactions on Database Systems*, 7(2): pp 209-234, June 1982.
12. T. Haerder and A. Reuter. Principles of transaction-oriented database recovery. *ACM Computing Surveys*, 15(4); pp 287-317, December 1983.
13. H. Korth, E. Levy and A. Silberschatz. A formal approach to recovery by compensating transactions. In *Proceedings of the 16th International Conference on Very Large Data Bases*, Brisbane, Australia, August 1990.
14. H. Korth and G.D. Speegle. Formal model of correctness without serializability. In *Proceedings of 1988 ACM SIGMOD Conference on Management of Data*, pp 379-386, May 1988.
15. E. Levy, H. Korth and A. Silberschatz. A theory of relaxed atomicity. In *Proceedings of the 1991 ACM Symposium on Principles of Distributed Computing*, August 1991.
16. C. Pu and A. Leff. Replica Control in distributed systems: An asynchronous approach. In *Proceedings of the 1991 ACM SIGMOD International Conference on Management of Data*, pp 377-386, Denver, May 1991.
17. C. Pu and A. Leff. Autonomous transaction execution with epsilon-serializability. In *Proceedings of RIDE Workshop on Transaction and Query Processing*, Phoenix, February 1992.
18. K. Ramamritham and P. K. Chrysanthis. In Search of Acceptability Criteria: Database Consistency Requirements and Transaction Correctness Properties. (to appear in) *Distributed Object Management*, edited by Özsu, Dayal, and Valduriez, Morgan Kaufmann, 1992.
19. K. Ramamritham and C. Pu. A Formal Characterization of Epsilon Serializability. *Technical Report 91-92*, Department of Computer Science, University of Massachusetts, Dec. 1991.
20. A. Sheth and M. Rusinkiewicz. Management of interdependent data: Specifying dependency and consistency requirements. In *Proceedings of the Workshop on Management of Replicated Data*, pp 133-136, Houston, November 1990.
21. K.L. Wu, P.S. Yu and C. Pu. Divergence Control for Epsilon Serializability. In *Proceedings of Eighth International Conference on Data Engineering*, pp 506-515, Phoenix, February 1992.