

 Open access • Proceedings Article • DOI:10.1109/REAL.1990.128734

Dynamic real-time optimistic concurrency control — [Source link](#)

Jayant R. Haritsa, Michael J. Carey, Miron Livny

Institutions: University of Wisconsin-Madison

Published on: 05 Dec 1990 - Real-Time Systems Symposium

Topics: Optimistic concurrency control, Serializability, Non-lock concurrency control, Multiversion concurrency control and Isolation (database systems)

Related papers:

- [Scheduling real-time transactions: a performance evaluation](#)
- [On being optimistic about real-time constraints](#)
- [Concurrency Control and Recovery in Database Systems](#)
- [Experimental Evaluation of Real-Time Optimistic Concurrency Control Schemes](#)
- [Experimental evaluation of real-time transaction processing](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/dynamic-real-time-optimistic-concurrency-control-2vkpun4ziy>

**DYNAMIC REAL-TIME OPTIMISTIC
CONCURRENCY CONTROL**

by

**Jayant R. Haritsa
Michael J. Carey
Miron Livny**

Computer Sciences Technical Report #981

November 1990

Dynamic Real-Time Optimistic Concurrency Control

Jayant R. Haritsa

Michael J. Carey

Miron Livny

Computer Sciences Department

University of Wisconsin

Madison, WI 53706

ABSTRACT — In a recent study, we have shown that in real-time database systems that discard late transactions, optimistic concurrency control outperforms locking. Although the optimistic algorithm used in that study, OPT-BC, did not factor in transaction deadlines in making data conflict resolution decisions, it still outperformed a deadline-cognizant locking algorithm. In this paper, we discuss some alternative methods of adding deadline-based priority information to optimistic algorithms. We present a new real-time optimistic concurrency control algorithm, WAIT-50, that monitors transaction conflict states and gives precedence to urgent transactions in a controlled manner. WAIT-50 is shown to provide significant performance gains over OPT-BC under a variety of operating conditions and workloads.

Dynamic Real-Time Optimistic Concurrency Control

Jayant R. Haritsa

Michael J. Carey

Miron Livny

Computer Sciences Department

University of Wisconsin

Madison, WI 53706

To appear in

Proceedings of the 1990 IEEE 11th Real-Time

Systems Symposium (RTSS)

1. INTRODUCTION

A Real-Time Database System (RTDBS) is a transaction processing system that attempts to satisfy the timing constraints associated with each incoming transaction. Typically, a constraint is expressed in the form of a *deadline*, that is, the user submitting the transaction would like it to be completed before a certain time in the future. Accordingly, greater value is associated with processing transactions before their deadlines as compared to completing them late. Therefore, in contrast to a conventional DBMS where the goal usually is to minimize response times, the emphasis here is on satisfying the timing constraints of transactions.

The problem of scheduling transactions in an RTDBS with the objective of minimizing the percentage of late transactions was first addressed in [Abbo88, Abbo89]. Their work focused on evaluating the performance of various real-time scheduling policies. All these policies enforced data consistency by using a two-phase locking protocol as the underlying concurrency control mechanism. Performance studies of concurrency control methods for conventional DBMSs (e.g., [Agra87, Care88]) have concluded that locking protocols, due to their conservation of resources, tend to perform better than optimistic techniques when resources are limited. In a recent study [Hari90], we investigated the behavior of these concurrency control schemes in a real-time environment. The study showed that for *firm deadline* real-time database systems, where late transactions are immediately discarded, optimistic concurrency control outperforms locking over a wide range of system loading and resource availability. The key reason for this surprising result is that the optimistic approach, due to its validation stage conflict resolution, ensures that eventually discarded transactions do not restart other transactions. The locking approach, on the other hand, allows soon-to-be-discarded transactions to cause other transactions to be either blocked or restarted due to lock conflicts, thereby increasing the number of late transactions.

An important difference between the locking algorithm and the optimistic algorithm that were compared in the above study lies in their use of transaction deadline information. The locking algorithm used this information, which was encoded in the form of transaction priorities, to provide preferential treatment to urgent transactions. The optimistic algorithm, however, was just the conventional broadcast commit optimistic scheme [Mena82, Robi82], and ignored transaction priorities in resolving data contention. The study therefore concluded that, in the firm real-time domain, a "vanilla" optimistic algorithm can perform better than a locking algorithm that is "tuned" to the real-time environment. The following question then naturally arises: How can we use priority information to improve the performance of the optimistic algorithm and thus further decrease the number of late transactions?

A simple answer to this question would be to use priority information in the resolution of data conflicts, that is, to resolve data conflicts always in favor of the higher priority transaction. This solution, however, has two problems: First, giving preferential treatment to high priority transactions may result in an increase in the number of missed deadlines. This can happen, for example, if helping one high priority transaction to make its deadline causes several lesser priority transactions to miss their deadlines. Second, if fluctuations can occur in transaction priorities, repeated conflicts between a pair of transactions may be resolved in some cases in favor of one transaction and in other cases in favor of the other transaction. This would hinder the progress of both transactions and hence degrade performance. Therefore, a priority-cognizant optimistic algorithm must address these two problems in order to perform better than a simple optimistic scheme.

In this paper, we report on our efforts to develop such an algorithm, and present a new real-time optimistic concurrency control algorithm, called **WAIT-50**. The algorithm incorporates a *priority wait* mechanism that makes low priority transactions wait for conflicting high priority transactions to complete, thus enforcing preferential treatment for high priority transactions. To address the first problem raised above, WAIT-50 features a *wait control* mechanism. This mechanism monitors transaction conflict states and, with a simple "50 percent" rule, dynamically controls when and for how long a transaction is made to wait. The second problem is handled by having the priority wait mechanism resolve conflicts in a manner that results in the commit of at least one of the conflicting transactions. Simulation results show that WAIT-50 performs significantly better than OPT-BC, the optimistic algorithm used in our earlier study.

The remainder of this paper is organized in the following fashion: Section 2 reviews our earlier study. In Section 3, we discuss deficiencies of OPT-BC. The new optimistic algorithm, WAIT-50, is presented in Section 4. Then, in Section 5, we describe our RTDBS model and its parameters, while Section 6 highlights the results of the simulation experiments. Finally, Section 7 summarizes the main conclusions of the study and outlines future avenues to explore.

2. BACKGROUND

Our earlier study [Hari90] investigated the relative performance of locking protocols and optimistic techniques in an RTDBS environment. In particular, the performance of a locking protocol, 2PL-HP, was compared with that of an optimistic technique, OPT-BC. These particular instances were chosen because they are of comparable complexity and are general in their applicability since they make no assumptions about knowledge of transaction semantics or resource demands. The details of these algorithms are explained below.

In 2PL-HP, classical two phase locking [Eswa76] is augmented with a *High Priority* [Abbo88] conflict resolution scheme to ensure that high priority transactions are not delayed by low priority transactions. This scheme resolves all data conflicts in favor of the transaction with the higher priority. When a transaction requests a lock on an object held by other transactions in a conflicting lock mode, if the requester's priority is higher than that of all the holders, the holders are restarted and the requester is granted the lock; otherwise, the requester waits for the lock holders to release the object. The High Priority scheme also serves as a deadlock prevention mechanism.¹

In OPT-BC, classical optimistic concurrency control [Kung81] is modified to implement the notion of a *Broadcast Commit* [Mena82, Robi82]. Here, when a transaction commits, it notifies other running transactions that conflict with it and these transactions are immediately restarted. Since there is no need to check for conflicts with already committed transactions, a transaction which has reached the validation stage is guaranteed to commit. The broadcast commit method² detects conflicts earlier than the basic optimistic algorithm, resulting in less wasted resources and earlier restarts; this increases the chances of meeting transaction deadlines. An important point to note is that transaction priorities are *not* used in resolving data conflicts.

¹ This is true only for priority assignment schemes that assign unique priority values and do not change a transaction's priority during the course of its execution.

² Refer to the Appendix for a discussion of implementation details of the broadcast commit method.

The results of our study showed that both the policy for dealing with late transactions and the availability of resources have a significant impact on the relative behavior of the algorithms. In particular, for a *firm deadline* system, where late transactions are discarded without being run to completion, OPT-BC outperformed 2PL-HP over a wide range of system loading and resource availability. Figures 1a and 1b present sample graphs of how the percentage of late(discarded) transactions varied as a function of the transaction arrival rate. These graphs were derived for the baseline model of the study, which characterized an RTDBS system with high data contention, under conditions of limited resources and plentiful resources, respectively.

In the above scenario, 2PL-HP suffered from two major problems: *wasted restarts* and *mutual restarts*. A "wasted restart" occurs when an executing transaction is restarted by another transaction that later misses its deadline. Such restarts are useless and cause performance degradation. In OPT-BC, however, we are guaranteed the commit of any transaction that reaches the validation stage. Since only validating transactions can cause restarts of other transactions, *all* restarts generated by the OPT-BC algorithm are useful.

The problem of "mutual restarts" arises when fluctuations occur in transaction priority profiles. For certain types of dynamic transaction priority assignment schemes (e.g., Least Slack [Jens85]), it is possible for a pair of concurrently running transactions to have opposite priorities relative to each other at different points in time during

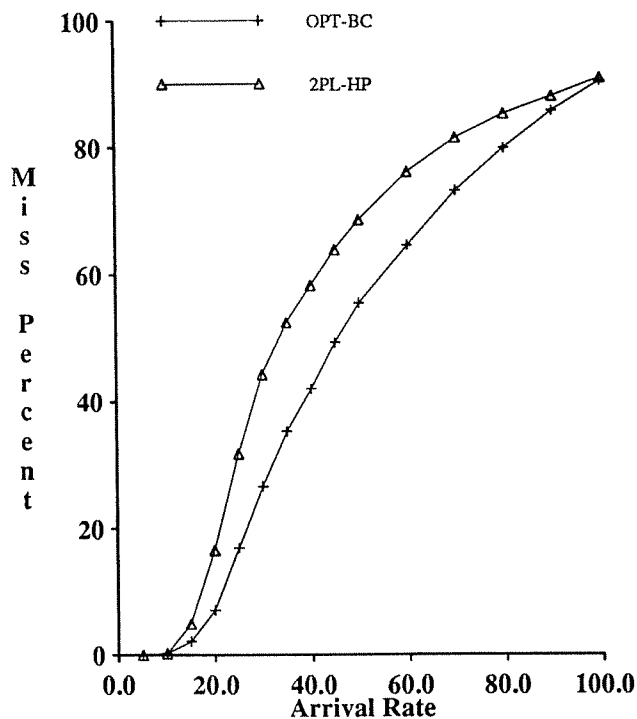


Figure 1a: Baseline Model (Limited Resources)

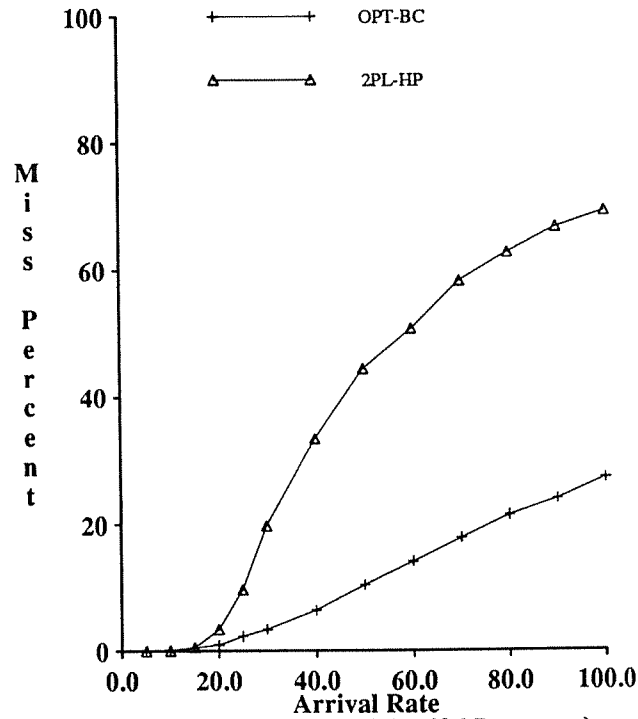


Figure 1b: Baseline Model (Plentiful Resources)

their execution. We will refer to this phenomenon as "priority reversal".³ For algorithms like 2PL-HP, which use transaction priorities to resolve data conflicts, priority reversals may lead to "mutual restarts" – a pair of transactions restart each other, thus hindering the progress of both transactions. Since OPT-BC does not use transaction priorities in resolving data contention, such problems simply *do not arise*.

3. PROBLEMS WITH OPT-BC

In this section, we will motivate why there is room for improvement on the OPT-BC algorithm. The validation algorithm of OPT-BC can be succinctly written as:

```
restart all conflicting4 transactions;
commit the validating transaction;
```

Although this algorithm provides immunity from priority dynamics due to its unilateral commit, it does not allow for the use of transaction priorities to further decrease the number of missed deadlines. To illustrate this problem, consider the scenario in Figure 2, where the execution profile of two concurrently executing transactions, X and Y , is shown. X has an arrival time A_X and deadline D_X , and Y has an arrival time A_Y and deadline D_Y . Also, assume that transaction X , by virtue of its earlier deadline, has a higher priority than transaction Y . Now, consider the situation where at time $t = V_Y$, when transaction X is close to completion, transaction Y reaches its validation point and detects a conflict with X . Under the OPT-BC algorithm, Y would immediately commit and in the process restart X . Restarting X at this late stage guarantees that it has no chance of meeting its deadline.

If a priority-cognizant algorithm had been used instead, it would have recognized that X 's priority was higher than that of Y . Then, in some fashion, it would have *prevented Y from committing* until X had completed. With this decision, we could possibly gain the completion of *both* transactions X and Y before their deadlines, as shown in

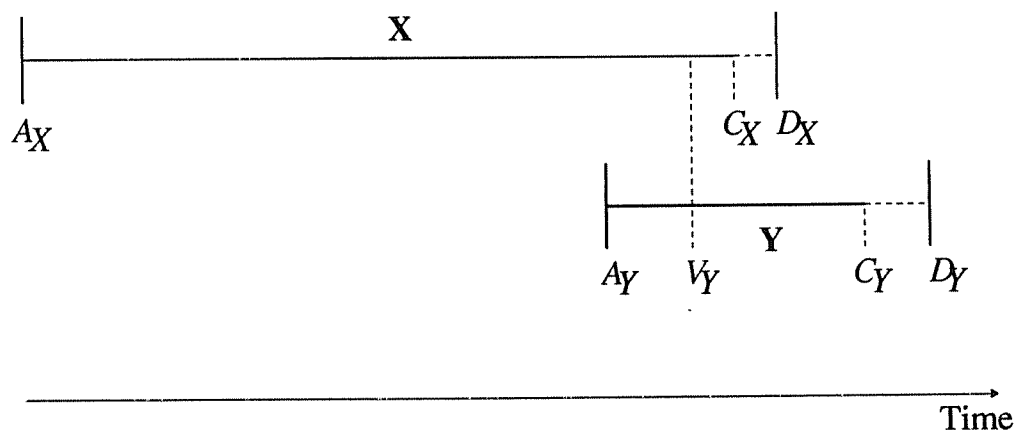


Figure 2: Poor OPT-BC data conflict decision

³ This is different from *priority inversion* [Sha87], which refers to the situation where a transaction is blocked (due to data or resource conflict) by another transaction with a lower priority.

⁴ A conflict exists between the validating transaction V and an executing transaction E if the intersection of the write set of V and the current read set of E is non-empty.

Figure 3 where X completes at time $t = C_X$ and Y completes later at time $t = C_Y$.

It is important to note that preventing transaction Y from committing does not *guarantee* that both transactions will eventually make their deadline. In fact, there are several possible outcomes: (1) transactions X and Y both miss their deadlines; (2) only one of the transactions makes its deadline; (3) both transactions make their deadlines. The above example, however, serves to show how OPT-BC's absolute indifference to transaction priorities can result in degraded performance. Another drawback of OPT-BC is that it has an inherent bias against long transactions, just like the classical optimistic algorithm.⁵ The use of priority information in resolving conflicts can help to counter this bias as well.

4. PRIORITY-COGNIZANT ALGORITHMS

As explained in the previous sections, although the OPT-BC algorithm highlights some major strengths of optimistic concurrency control in firm real-time database systems, there remains potential for improving its performance. We therefore tried to develop new optimistic algorithms that address the problems of OPT-BC without sacrificing the performance-beneficial aspects of the broadcast commit scheme. These algorithms are described in this section. In the subsequent discussion, we will use the term *conflict set* to denote the set of currently running transactions that conflict with a validating transaction. The acronym *CHP* (*Conflicting Higher Priority*) will be used to refer to transactions that are in the conflict set and have a higher priority than the validating transaction. Similarly, the acronym *CLP* (*Conflicting Lower Priority*) will be used to refer to transactions that are in the conflict set and have a lower priority than the validating transaction. In this section, our aim is to motivate the development of the algorithms and discuss, at an intuitive level, their potential strengths and weaknesses.

The example in Section 3, illustrating poor conflict decisions by OPT-BC, showed that we need a scheme to prevent low priority transactions that conflict with higher priority transactions from unilaterally committing. The following two options are available:

- (1) *Restart*: The low priority transaction is restarted.
- (2) *Block*: The low priority transaction is blocked.

Two algorithms, OPT-SACRIFICE and OPT-WAIT, were developed based on these options. WAIT-50 was then developed as an extension of the OPT-WAIT algorithm. These three algorithms are presented below.

4.1. OPT-SACRIFICE

In the OPT-SACRIFICE algorithm, when a transaction reaches its validation stage, it checks for conflicts with currently executing transactions. If conflicts are detected and at least one of the transactions in the conflict set is a CHP transaction, then the validating transaction is restarted – that is, it is *sacrificed* in an effort to help the higher priority transactions make their deadlines. The validation algorithm of OPT-SACRIFICE can therefore be written as:

⁵ It should be obvious that both 2PL and 2PL-HP do not have this bias.

```

if CHP transactions in conflict set then
  restart the validating transaction;
else
  restart transactions in conflict set;
  commit the validating transaction;

```

Referring back to Figure 2, if we were using OPT-SACRIFICE, then at time $t = V_Y$, transaction Y would restart itself due to the conflict with the higher priority transaction X .

OPT-SACRIFICE is priority-cognizant and satisfies the goal of giving preferential treatment to high priority transactions. It suffers, however, from two potential problems. First, there is the problem of *wasted sacrifices*, where a transaction is sacrificed on behalf of another transaction that is later discarded. Such sacrifices are useless and cause performance degradation. Second, the algorithm does not have immunity to priority dynamics. For example, the situation may arise where transaction A is sacrificed for transaction B because B 's priority is currently greater than that of A , and transaction B at a later time is sacrificed for transaction A because A 's priority is now greater than B 's priority. Therefore, priority reversals may lead to *mutual sacrifices*. These two drawbacks are analogous to the "wasted restarts" and "mutual restarts" problems of 2PL-HP.

4.2. OPT-WAIT

The OPT-WAIT algorithm incorporates a *priority wait* mechanism: a transaction that reaches validation and finds CHP transactions in its conflict set is "put on the shelf", that is, it is made to wait and not allowed to commit immediately. This gives the higher priority transactions a chance to make their deadlines first. While a transaction is waiting, it is possible that it will be restarted due to the commit of one of the CHP transactions. The validation algorithm of OPT-WAIT can therefore be written as:

```

while CHP transactions in conflict set do
  wait;
  restart transactions in conflict set;
  commit the validating transaction;

```

Referring back to Figure 2, if we were using OPT-WAIT, then at time $t = V_Y$, transaction Y would wait, without committing, for transaction X to complete first. Of course, X 's completion may cause Y to be restarted.

There are several reasons that suggest that the priority wait mechanism may have a positive impact on performance, and these are outlined below:

- (1) In keeping with the original goal, precedence is given to high-priority transactions.
- (2) The problem of "wasted sacrifices" does not exist because if a CHP transaction is discarded due to missing its deadline, or is restarted by some other transaction, then the waiter is immediately "taken off the shelf" and committed if no other CHP transactions remain.
- (3) Priority reversals are not a problem because, if a CHP transaction being waited for were to become a CLP transaction, the waiting transaction will no longer wait for it, and will immediately commit if no other CHP transactions remain.
- (4) Since transactions wait instead of immediately restarting, a blocking effect is derived – this results in conservation of resources, which can be beneficial to performance [Agra87].

- (5) The fact that a CHP transaction commits does not necessarily imply that the waiting transaction has to be restarted (!).

The last point requires further explanation: The **key** observation here is that if transaction A conflicts with transaction B , this does not necessarily mean that the **converse** is true [Robi82]. This is explained as follows: Under the broadcast commit scheme, a validating transaction A is said to conflict with another transaction B if and only if

$$WriteSet_A \cap ReadSet_B \neq \phi \quad (1)$$

We will denote such a conflict from transaction A to B by $A \rightarrow B$. For transaction B to also conflict with transaction A , i.e. for $B \rightarrow A$, it is necessary that

$$WriteSet_B \cap ReadSet_A \neq \phi \quad (2)$$

As is obvious from Equations (1) and (2), $A \rightarrow B$ does not imply $B \rightarrow A$. Therefore, if in fact $B \rightarrow A$ is not true, then by committing the transactions in the order (B, A) instead of the order (A, B) , both transactions can be committed without restarting either one.

As per the explanation given above, it is possible with our waiting scheme for the CHP transaction and the waiting transaction to commit in that order without either transaction being restarted. Therefore, the priority wait mechanism has a potential to actually *eliminate* some data conflicts. In the Appendix, a simple probabilistic analysis is presented of the extent to which waiting can cause a reduction in data conflicts.

Although the waiting scheme has many positive features, it is not an unmixed blessing. One potential drawback is that if a transaction finally commits after waiting for some time, it causes all of its CLP transactions to be restarted at a later point in time. This decreases the chances of these transactions meeting their deadlines, and also wastes resources. A second drawback is that the validating transaction may develop new conflicts during its waiting period, thus causing an increase in conflict set sizes and leading to more restarts. Another way to view this is to realize that waiting causes objects to be, in a sense, "locked" for longer periods of time. Therefore, while waiting has the capability to reduce the probability of a restart-causing conflict between a given pair of transactions, it can simultaneously increase the probability of having a larger *number* of conflicts per transaction. This increase may be substantial when there are many concurrently executing transactions in the system.

4.3. WAIT-50

The WAIT-50 algorithm is an extension of OPT-WAIT – in addition to the priority wait mechanism, it incorporates a *wait control* mechanism. This mechanism monitors transaction conflict states and dynamically decides when, and for how long, a low priority transaction should be made to wait for its CHP transactions. A transaction's conflict state is assumed to be characterized by the index *HPpercent*, which is the percentage of the transaction's total conflict set size that is formed by CHP transactions. The operation of the wait mechanism is conditioned on the value of this index. In WAIT-50, a simple "50 percent" rule is used – a validating transaction is made to wait only while $HPpercent \geq 50$, that is, while half or more of its conflict set is composed of higher priority transactions. The validation algorithm of WAIT-50 can therefore be written as:

```

while CHP transactions in conflict set and
    HPPercent  $\geq$  50 do
    wait;
restart transactions in conflict set;
commit the validating transaction;

```

The aim of the wait control mechanism is to detect when the beneficial effects of waiting, in terms of giving preference to high priority transactions and decreasing pairwise conflicts, are outweighed by its drawbacks, in terms of later restarts and an increased number of conflicts. The value of a conflict set is measured in terms of its high priority component. Based on this value it is decided whether or not waiting will potentially help increase the number of on-time transactions. Therefore, while OPT-WAIT and OPT-BC represent the extremes with regard to waiting – OPT-WAIT always waits for a CHP transaction, and OPT-BC never waits – WAIT-50 is a *hybrid* algorithm that controls the amount of waiting based on transaction conflict states. In fact, we can view OPT-WAIT, WAIT-50, and OPT-BC as all being special cases of a general algorithm WAIT-X, where X is the cutoff HPPercent level, with X taking on the values 0, 50, and ∞ , respectively, for these algorithms.

We conducted experiments to evaluate the performance of the various optimistic algorithms, and the following sections describe our experimental framework and results.

5. REAL-TIME DBMS MODEL

The real-time database system model employed here is the same as that of our earlier study – in this model, the system consists of a shared-memory multiprocessor DBMS operating on disk resident data.⁶ The database itself is modeled as a collection of pages. Transactions arrive in a Poisson stream and each transaction has an associated deadline time. A transaction consists of a sequence of read and write accesses. A read access involves a concurrency control request to get access permission, followed by a disk I/O to read the page, followed by a period of CPU usage for processing the page. Write requests are handled similarly except for their disk I/O – their disk activity is deferred until the transaction has committed. The basic structure of the model is shown in Figure 3.

The model has five components: a *source* that generates transactions; a *transaction manager* that models the execution of transactions; a *concurrency control (CC) manager* that implements the details of the concurrency control algorithms; a *resource manager* that models the CPU and I/O resources; and a *sink* that gathers statistics on completed transactions. The following two subsections describe the workload generation process and the hardware resource configuration.

5.1. Workload Model

The workload model characterizes transactions in terms of the pages that they access and the number of pages that they update. Table 1 summarizes the key parameters of the workload model. The *ArrivalRate* parameter specifies the rate of transaction arrivals. The *DatabaseSize* parameter gives the number of pages in the database. The number of pages accessed by a transaction varies uniformly between half and one-and-a-half times the value of *PageCount*. Page requests are generated from a uniform distribution spanning the entire database. *WriteProb* gives

⁶ It is assumed, for simplicity, that all data is accessed from disk and buffer pool considerations are therefore ignored.

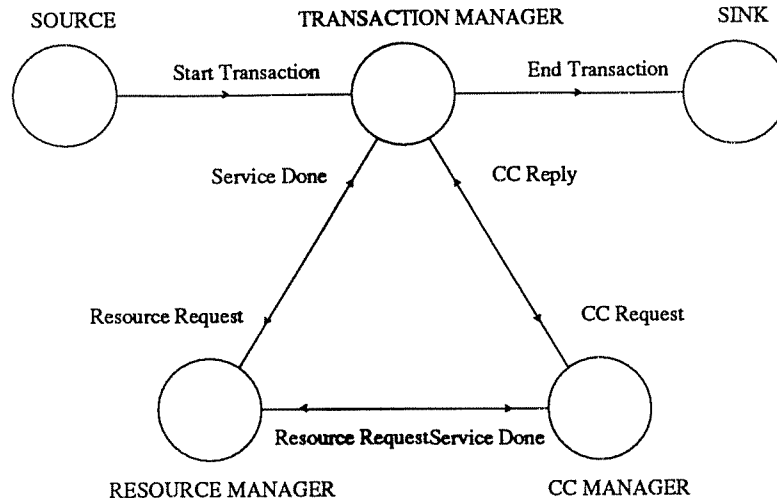


Figure 3: RTDBS Model Structure

the probability that a page which is read will also be updated.

We use three transaction deadline formulas in this study. The first formula, which is the same as the one used in our previous study, is:

$$D_T = A_T + SF * R_T \quad (DF1)$$

where D_T , A_T , and R_T are the deadline, arrival time and resource time, respectively, of transaction T , while SF is a slack factor. The *resource time* is the total service time at the resources that the transaction requires for its data processing. The *slack factor* is a constant that provides control over the tightness/slackness of deadlines. With this formula, all transactions, independent of their service requirement, have the same *slack ratio* – this is defined to be the ratio $\frac{D_T - A_T}{R_T}$. Therefore, all transactions have SF as their slack ratio.

Parameter	Meaning
<i>ArrivalRate</i>	Transaction arrival rate
<i>DatabaseSize</i>	Number of pages in database
<i>PageCount</i>	Avg. pages accessed/transaction
<i>WriteProb</i>	Write probability/accessed page
<i>DeadlineFormula</i>	DF1 or DF2 or DF3
<i>LSF</i>	Low Slack Factor
<i>HSF</i>	High Slack Factor
<i>PriorityPolicy</i>	Transaction priority policy

Table 1: Workload Model Parameters

In order to evaluate the effects of variability in transaction slack ratios, two additional deadline assignment formulas are used in the present study. The first formula is:

$$D_T = \begin{cases} A_T + LSF * R_T \\ A_T + HSF * R_T \end{cases} \quad (DF2)$$

With this formula, transactions will have either a slack factor of *LSF* or *HSF*, with both choices being equally likely. Therefore, the slack ratio for a transaction will be either *LSF* or *HSF*. The other new formula is:

$$D_T = \begin{cases} A_T + LSF * R_{max} \\ A_T + HSF * R_{max} \end{cases} \quad (DF3)$$

This formula is the same as (DF2) except that the resource time used here, R_{max} , is that of the largest possible transaction, not that of the particular transaction for which the deadline assignment is being made. Thus, (DF3) introduces a different type of variation in the slack ratio. With (DF3), transaction slack ratios are spread out over a range of values, based on the ratio of R_{max} to the R_T 's, rather than being restricted to a single value as in (DF1) or two values as in (DF2). The *LSF* and *HSF* workload parameters set the slack factors to be used in each of the deadline formulas. (For DF1, these two parameters have the same value, *SF*).

The transaction priority assignment scheme used in virtually all of the experiments was *Earliest Deadline* – where transactions with earlier deadlines have higher priority than transactions with later deadlines. For one experiment, however, which was designed to show that priority dynamics do not affect the wait-based algorithms, the assignment scheme was *static LTD (Least Time to Deadline)*. In the basic LTD scheme, the priority of a transaction is evaluated as $(D_T - clock)$, and the smaller this value, the higher the priority of the transaction. For the static variant, the priority of a transaction is evaluated when it arrives and reevaluated only whenever it is restarted. In between reevaluations, the priority remains at the previously computed value. The *PriorityPolicy* workload parameter fixes the transaction priority scheme.

A feature common to all of the experiments is that the system operates under firm deadlines, and therefore discards late transactions. It is important to note that while the workload generator uses transaction resource requirements in assigning deadlines, we assume that the system itself lacks any knowledge of these requirements. This implies that a transaction is detected as being late only when it actually misses its deadline.

5.2. Resource Model

The physical resources in our model consist of multiple CPUs and multiple disks. There is a single queue for the CPUs and the service discipline is preemptive-resume, with the preemption being based on transaction priorities. Each of the disks has its own queue and is scheduled according to the priority-based variant of the elevator disk scheduling algorithm described in [Care89]. Table 2 summarizes the key parameters of the resource model. Requests at each disk are grouped into *NumDiskPrio* priority levels and the elevator algorithm is applied within each priority level. (The algorithm for mapping requests to disk priority levels is described in the Appendix). Requests at a priority level are served only when there are no pending requests at higher priority levels, and the mapping of requests to priority levels is periodically reevaluated to reflect the change in transaction priority values over time. The data itself is modeled as being uniformly distributed across all of the disks and across all tracks

within a disk, with the number of tracks on each disk being set by the *NumTracks* parameter. The *PageCpu*, *DiskDelay* and *SeekFactor* parameters capture CPU and disk processing times per data page, and are detailed in [Care89].

6. EXPERIMENTS and RESULTS

In this section, we present performance results for our experiments comparing the various optimistic algorithms in a real-time database system environment. The simulator used to obtain the results is written in the Modula-2-based DeNet simulation language [Livn88]. The performance metric is *MissPercent*, which is the percentage of transactions that do not complete before their deadline.⁷ *MissPercent* values in the range of 0 to 20 percent are taken to represent system performance under "normal" loadings, while *MissPercent* values in the range of 20 to 100 percent represent system performance under "heavy" loading.⁸ The simulations also generated a host of other statistical information, including CPU and disk utilizations, the number of data conflicts, system population, etc. These secondary measures help to explain the behavior of the algorithms under various loading conditions. The resource parameter settings are such that the CPU time to process a page is 10 milliseconds while disk access times are between 15 and 30 milliseconds, depending on the level of disk utilization. Disk access times depend on disk utilization due to the elevator scheduling policy.

For experiments that were intended to factor in the effect of resource contention on the performance of the algorithms, the number of processors and number of disks were set to 10 and 20, respectively. For experiments intended to isolate the effect of data contention, we approximately simulated an "infinite" resource situation [Agra87], that is, where there is no queuing for resources. This was done by increasing by twenty-fold the number of processors and the number of disks, from their baseline values of 10 and 20 to 200 and 400, respectively. A point to note here is that while abundant resources are usually not to be expected in conventional database systems, they may be more common in RTDBS environments since many real-time systems are sized to handle transient heavy loading. This directly relates to the application domain of RTDBSs, where functionality, rather than cost, is usually

Parameter	Meaning
<i>NumCPUs</i>	Number of processors
<i>NumDisks</i>	Number of disks
<i>NumDiskPrio</i>	Number of disk priority levels
<i>NumTracks</i>	Number of tracks per disk
<i>PageCpu</i>	CPU time for processing a data page
<i>DiskDelay</i>	Disk rotational + transfer delays
<i>SeekFactor</i>	Factor relating seek time to distance

Table 2: Resource Model Parameters

⁷ All *MissPercent* graphs in this paper show mean values that have relative half-widths about the mean of less than 10% at the 90% confidence interval, with each experiment having been run until at least 5000 transactions were processed by the system. Only statistically significant differences are discussed here.

⁸ Any long-term operating region where the miss percent is large is obviously unrealistic for a viable RTDBS. Exercising the system to high miss levels, however, provides valuable information on the response of the algorithms to brief periods of stress loading.

the driving consideration.

We began our experiments by evaluating the algorithms for the baseline model of our earlier study. This was done in order to provide continuity from that study to the present work. Subsequently, for reasons explained in the following discussion, we moved to a new baseline model. After initial experiments with this model, further experiments were constructed around it by varying a few parameters at a time. These experiments evaluated the impact of data contention, resource contention, deadline slack variation, transaction write probabilities, priority dynamics and the wait control mechanism parameters. We will hereafter refer to the old baseline model as FIX-SR (Fixed Slack Ratio), and the new baseline model as VAR-SR (Variable Slack Ratio).

6.1. FIX-SR Baseline Model

The settings of the workload parameters and resource parameters for the FIX-SR baseline model are listed in Tables 3 and 4, respectively. These settings generate an appreciable level of both data contention and resource contention. For this model, Figures 4a and 4b show MissPercent behavior under normal load and heavy load conditions, respectively. When the same experiment is carried out under infinite resources, Figures 5a and 5b are obtained. From this set of graphs, we can make the following observations:

- (1) OPT-SACRIFICE performs significantly worse than the wait-based algorithms over the entire operating region, and for the most part, also performs worse than OPT-BC. The poor performance of this algorithm is primarily due to the problem of "wasted sacrifices" discussed in Section 4. Also, in the infinite resource case, the sacrifice policy generates a steep rise in the number of data conflicts, when compared to that of

Parameter	Value
<i>DatabaseSize</i>	1000 pages
<i>PageCount</i>	16 pages
<i>WriteProb</i>	0.25
<i>DeadlineFormula</i>	DF1
<i>LSF</i>	4.0
<i>HSF</i>	4.0

Table 3: FIX-SR Baseline Model Workload Settings

Parameter	Value
<i>NumCPUs</i>	10
<i>NumDisks</i>	20
<i>NumDiskPrio</i>	5
<i>NumTracks</i>	1000
<i>PageCpu</i>	10 ms
<i>DiskDelay</i>	15 ms
<i>SeekFactor</i>	0.5 ms

Table 4: FIX-SR Baseline Model Resource Settings

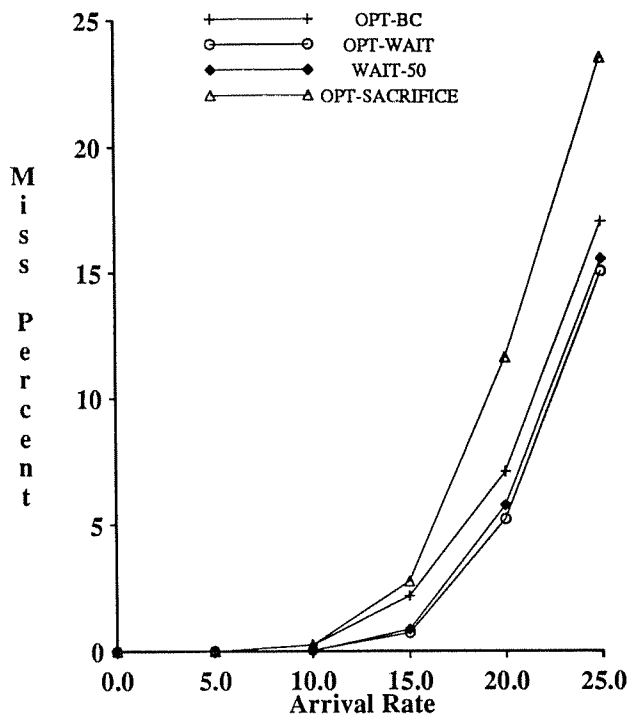


Figure 4a: FIX-SR Baseline Model (Normal Load)

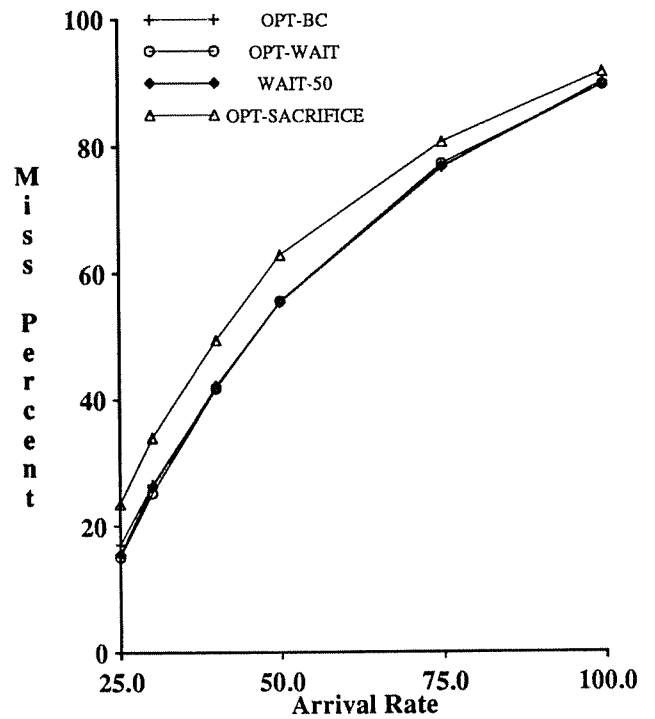


Figure 4b: FIX-SR Baseline Model (Heavy Load)

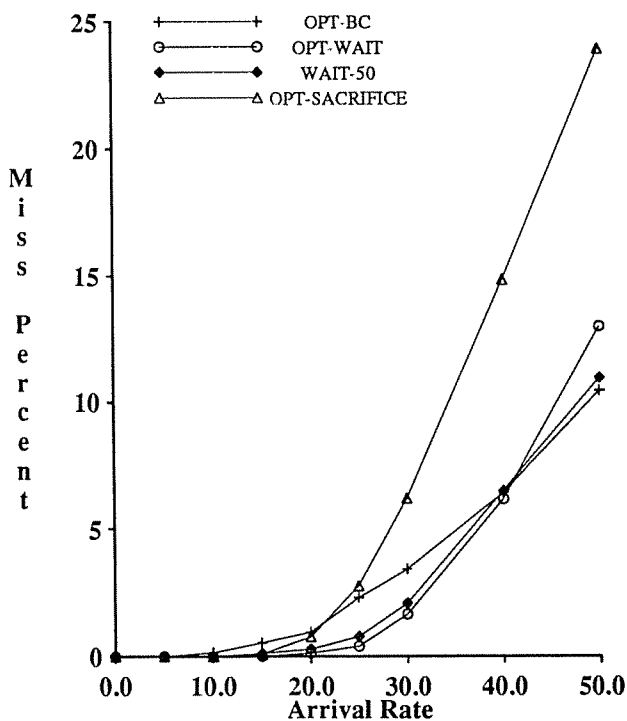


Figure 5a: Infinite Resources (Normal Load)

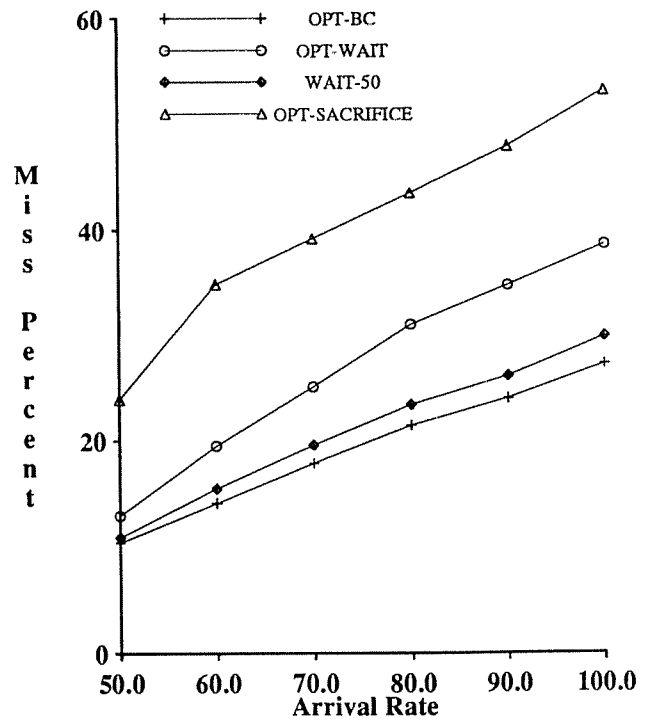


Figure 5b: Infinite Resources (Heavy Load)

OPT-BC, by causing a significant increase in the average number of transactions in the system.⁹ This is brought out quantitatively in Figure 5c, which plots the average number of conflicts per input transaction for the infinite resource scenario.

- (2) OPT-WAIT, due to the beneficial effects of its priority cognizance, performs very well at low levels of data contention (Figs. 4a, 5a). As data contention increases, however, its performance steadily degrades. Finally, at high contention levels under infinite resources (Fig. 5b), it performs significantly worse than OPT-BC. The reason for OPT-WAIT's poor performance in this region is that its priority wait mechanism, just like the sacrifice policy, causes an increase in the average number of transactions in the system. This population increase generates a corresponding rise in the number of data conflicts (see Fig. 5c), resulting in higher miss percentages.
- (3) WAIT-50 provides the *best overall* performance. At low data contention levels, it behaves like OPT-WAIT, and at high contention levels it behaves like OPT-BC. The explanation for this behavior is given in the next section.
- (4) Under high resource contention (Fig.4b), WAIT-50 and OPT-WAIT behave identically to OPT-BC. This is because, with heavy resource contention, it is uncommon for a low priority transaction to reach its validation stage much before its deadline, and therefore the wait-times of transactions are mostly small. Accordingly, the priority wait mechanism has very limited impact, and WAIT-50, OPT-WAIT, and OPT-BC become essentially the same algorithm.

The above set of experiments were encouraging because they showed that there were performance benefits to be gained by using priority-cognizant algorithms. It was all the more encouraging that these performance improvements were obtained despite all transactions having had the same slack ratio (the model used deadline formula DF1). A fixed transaction slack ratio reduces the likelihood of a validating transaction finding a higher priority transaction in its set of conflicting transactions. This creates favorable circumstances for OPT-BC since the detrimental effects of its priority insensitivity are reduced.

6.2. VAR-SR Baseline Model

In order to generate a workload with variation in transaction slack ratios, the VAR-SR baseline model was developed for the current study. This model uses deadline assignment formula DF2 to generate variation in transaction slack ratios. The workload parameters *LSF* and *HSF* are set at 2.0 and 6.0, respectively.¹⁰ The remaining workload parameter settings and resource parameter settings are the same as those for the FIX-SR baseline model (see Tables 2 and 3). In the subsequent discussions, we will compare the performance of only the OPT-BC, OPT-WAIT and WAIT-50 algorithms since OPT-SACRIFICE invariably performed worse than the wait-based algorithms.

For the VAR-SR baseline model, Figures 6a and 6b show the behavior of the algorithms under normal load and heavy load, respectively. When the same experiment was carried out under infinite resources, Figures 7a and

⁹ The data conflict counter is incremented by the conflict set size whenever a transaction reaches its validation stage.

¹⁰ These parameter selections ensure that the *mean* slack ratio is the same as that of the FIX-SR baseline model, namely 4.0.

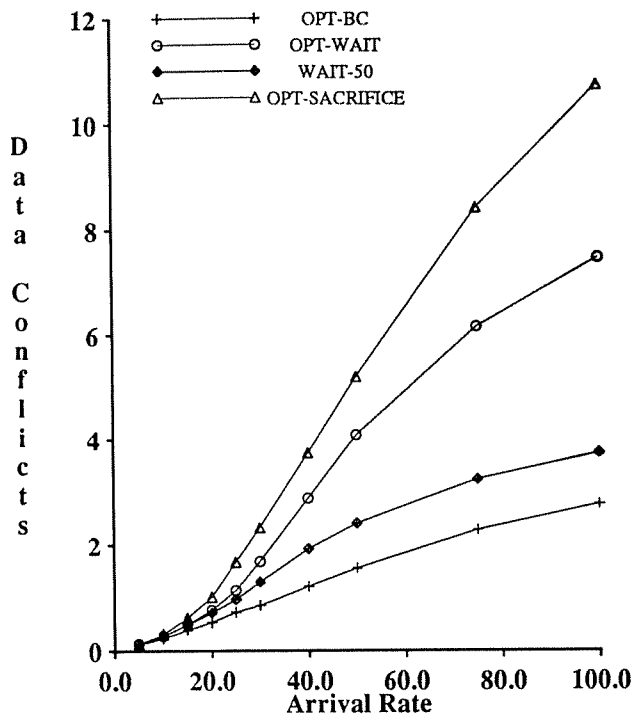


Figure 5c: Conflicts (Infinite Resources)

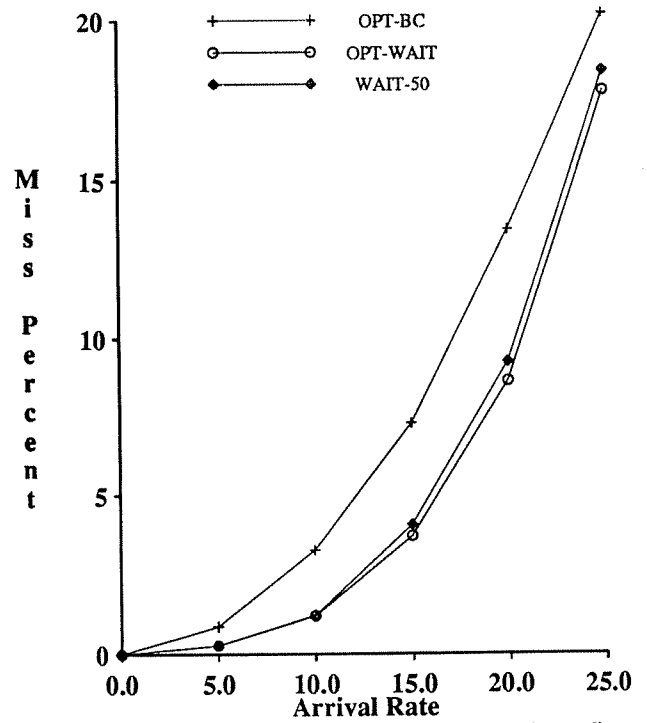


Figure 6a: VAR-SR Baseline Model (Normal Load)

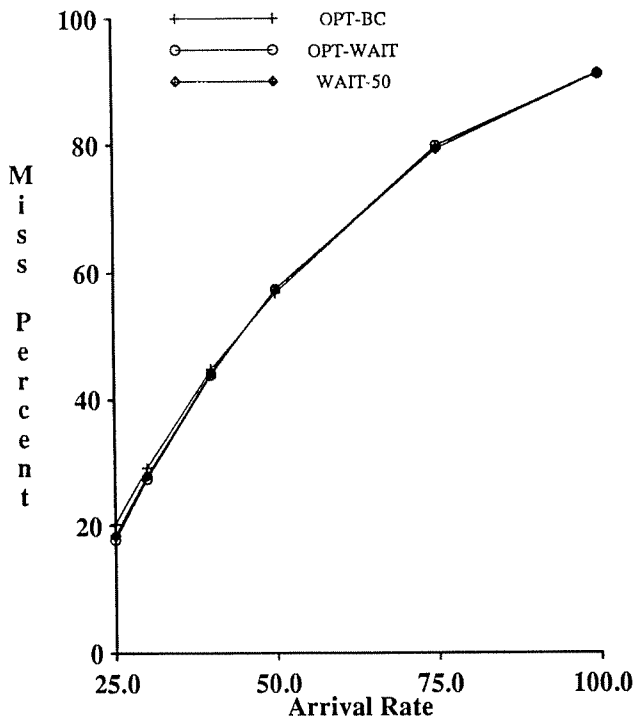


Figure 6b: VAR-SR Baseline Model (Heavy Load)

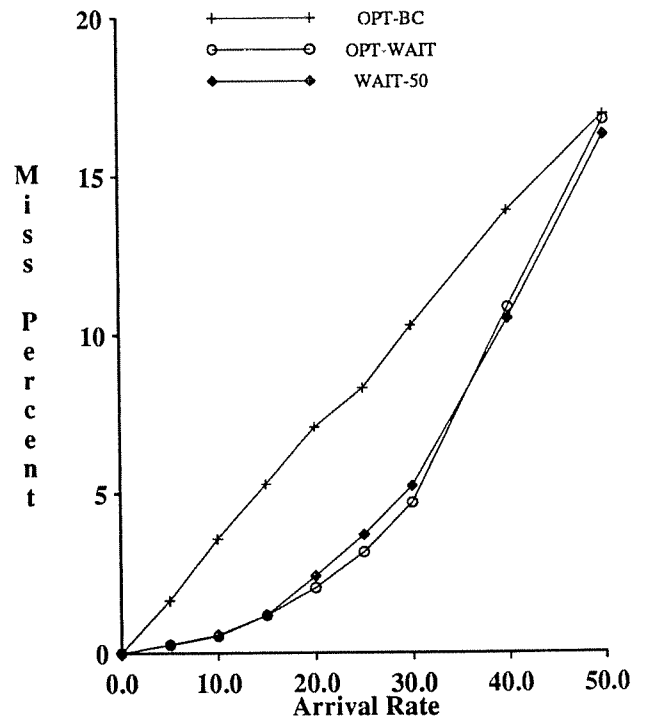


Figure 7a: Infinite Resources (Normal Load)

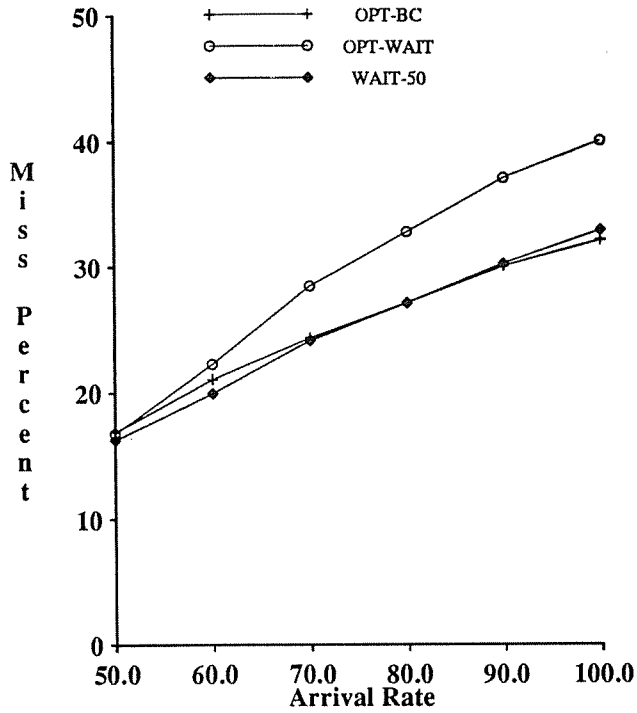


Figure 7b: Infinite Resources (Heavy Load)

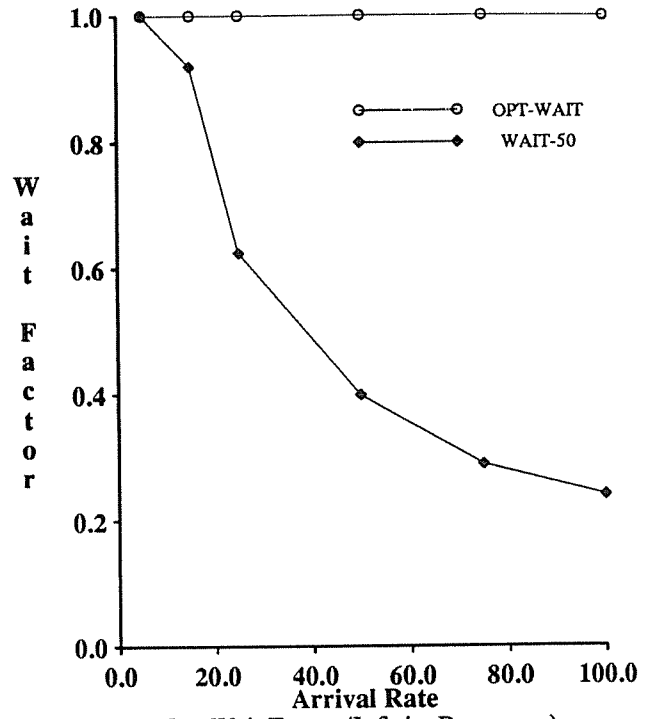


Figure 7c: Wait Factor (Infinite Resources)

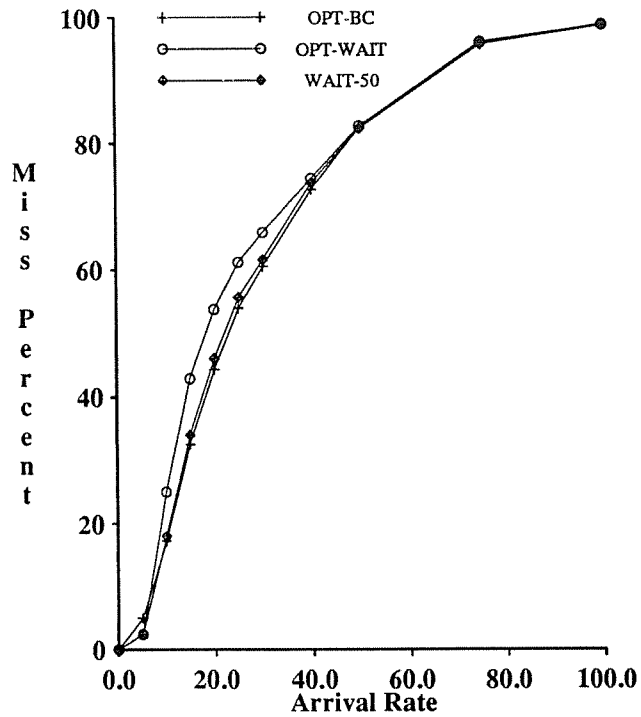


Figure 8: Write Prob. = 1.0 (Finite Resources)

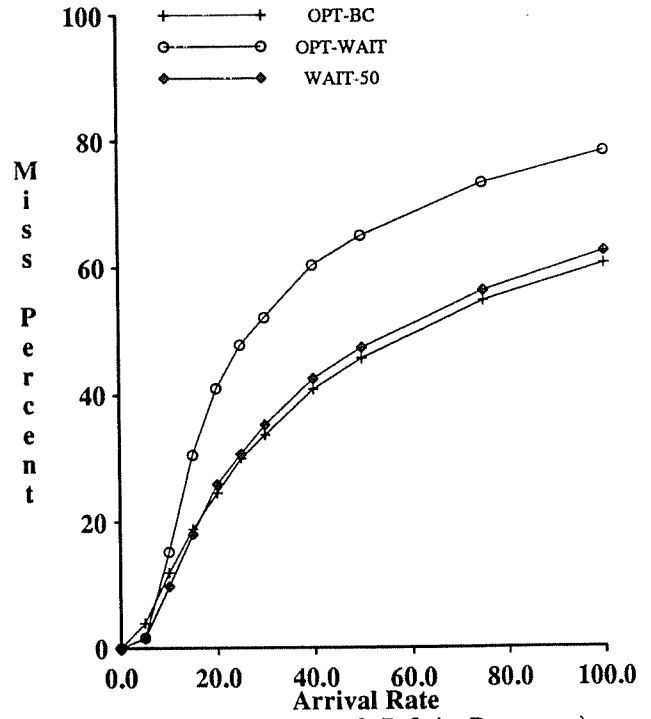


Figure 9a: Write Prob. = 1.0 (Infinite Resources)

7b were obtained. From this set of graphs we can make the following observations:

- (1) The priority-cognizant algorithms, WAIT-50 and OPT-WAIT, now perform *significantly better* than OPT-BC under normal loads.
- (2) WAIT-50 again turns in the best overall performance by behaving like OPT-WAIT at low data contention levels and like OPT-BC at high data contention levels.

As can be seen from this experiment, and will be further confirmed in subsequent experiments, WAIT-50 provides performance close to either OPT-BC or OPT-WAIT in operating regions where they behave well, and provides the same or slightly better performance at intermediate points. Therefore, in an overall sense, *WAIT-50 effectively integrates priority and waiting* into the optimistic concurrency control framework. The control mechanism is clearly quite competent at deciding when the benefits of waiting, in terms of helping high priority transactions to make their deadlines, are outweighed by the drawbacks of causing an increased number of conflicts. In Figure 7c, we plot the "wait factor" of WAIT-50 with respect to that of OPT-WAIT, which measures the total time spent in priority-waiting using WAIT-50, normalized by the waiting time of OPT-WAIT.¹¹ As can be seen from this figure, WAIT-50's wait factor is close to that of OPT-WAIT at low contention levels but decreases steadily as the data contention level is increased. Therefore, while OPT-WAIT and OPT-BC represent the extremes with regard to waiting, WAIT-50 gracefully controls the waiting to match the data contention level in the system.

6.3. Write Probability

All the previously described experiments were carried out for a write probability of 0.25. The next set of experiments look into the performance effects of varying transaction write probabilities. In the first experiment, the write probability was increased to 1.0, keeping the other parameters the same as those of the baseline model. This experiment was conducted for both finite resource and infinite resource scenarios, and the results are shown in Figures 8 and 9a. From this set of figures, we can make the following observations:

- (1) OPT-WAIT suffers a substantial performance degradation and does worse than OPT-BC over almost the entire operating region. There are two reasons for this: First, the increased write probability generates higher levels of data contention which, in combination with the population increase effect of the priority wait mechanism, results in a steep increase in the number of conflicts. Second, the conflict-elimination capability of OPT-WAIT vanishes since *all* conflicts are now *bi-directional*. These effects are captured dramatically in Figure 9b, which profiles the average number of conflicts per input transaction under infinite resources.
- (2) Although WAIT-50 also employs the priority wait mechanism, it does not suffer OPT-WAIT's performance degradation. This is due to its control mechanism, which ensures OPT-BC-like behavior when high data contention levels are reached by sharply reducing its wait factor. Figure 9c, which plots the wait factor of the algorithms for the infinite resources case, shows this effect quantitatively.

¹¹ The wait factor of OPT-BC is trivially zero as the algorithm has no wait component.

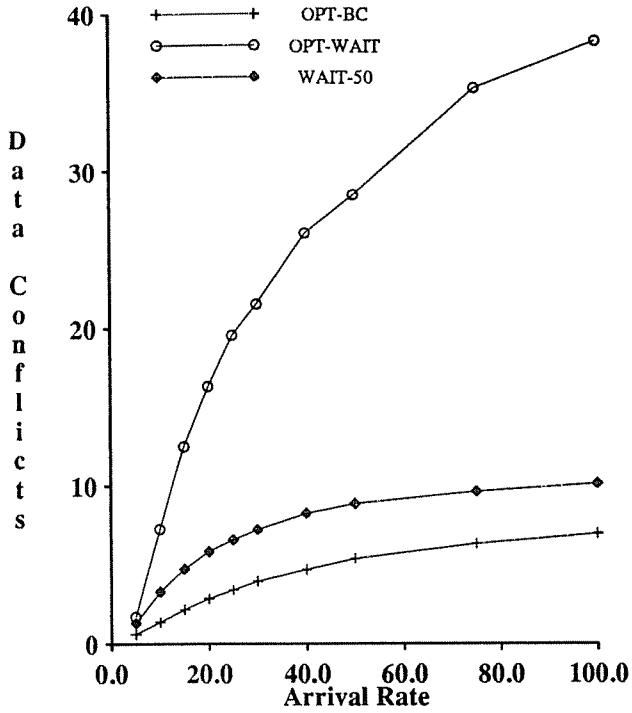


Figure 9b: Conflicts (Infinite Resources)

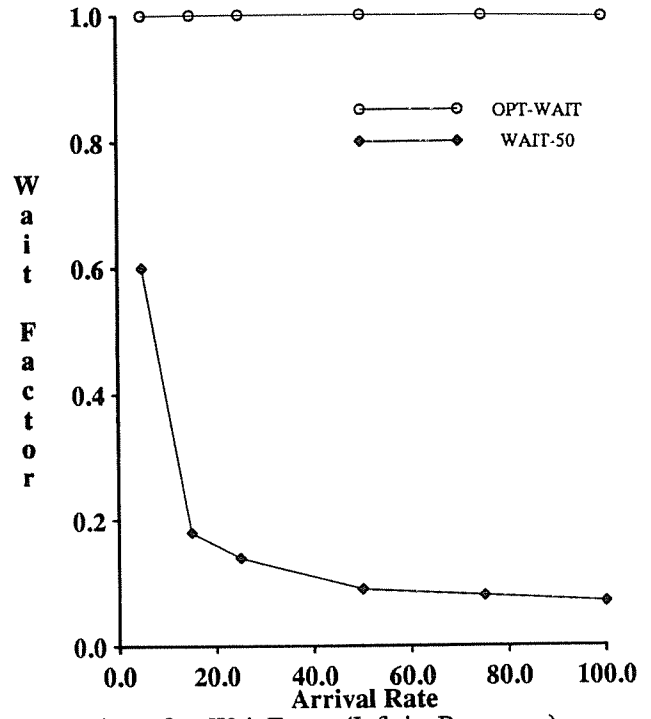


Figure 9c: Wait Factor (Infinite Resources)

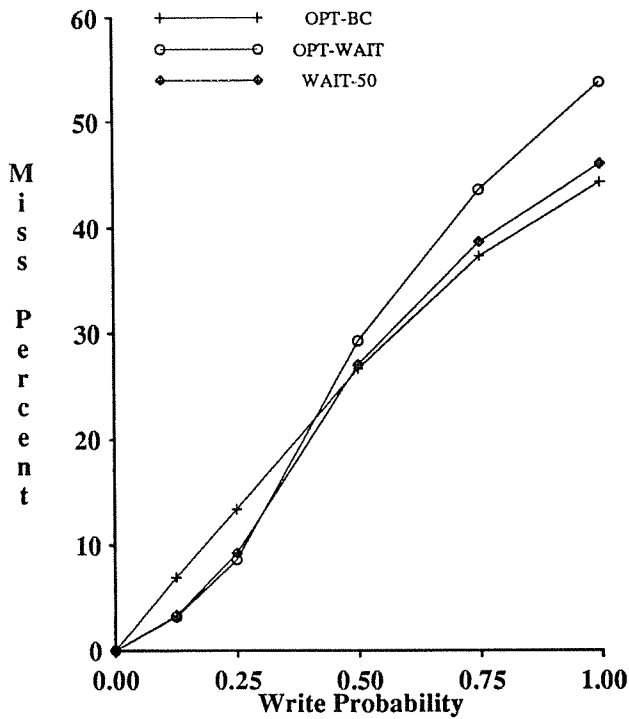


Figure 10: Finite Resources (Arrival Rate = 20)

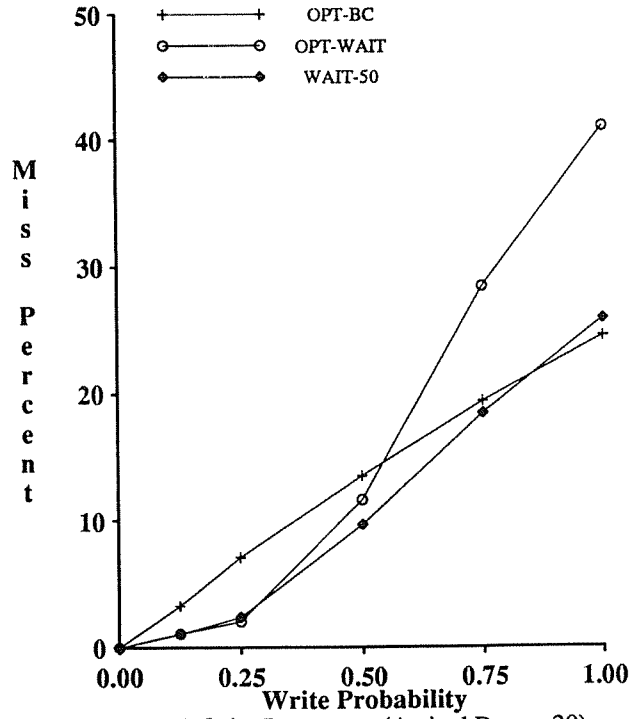


Figure 11: Infinite Resources (Arrival Rate = 20)

In the second experiment, the write probability was varied from 0.0 to 1.0, keeping the arrival rate constant at 20 transactions/sec. Figures 10 and 11 show how the algorithms behave under conditions of finite and infinite resources, respectively. These graphs clearly show that while OPT-WAIT performs well at low conflict levels, OPT-BC does much better at high conflict levels. We also observe that WAIT-50 again provides good performance over the entire range.

6.4. Wait Control Mechanism

The next experiment presented here examines the effect of the choice of 50 percent as the cutoff value for the HPpercent control index. Keeping all parameters the same as those of the baseline model, we measured the performance of WAIT-25 and WAIT-75 under conditions of finite and infinite resources. Figures 12a and 12b give the results of the finite resources experiment under normal load and heavy load, respectively, while Figures 13a and 13b give the corresponding results under infinite resources. From these graphs, we can make the following observations:

- (1) Lowering the cutoff value to 25 percent results in a slightly better normal load performance, but worse heavy load performance. This behavior is due to the increased wait factor that is delivered by the lowered cutoff value.
- (2) Raising the cutoff value to 75 percent has the opposite effect: the normal load performance becomes worse, while the heavy load performance is slightly better. This behavior is due to the decreased priority cognizance that is delivered by the increased cutoff value.

A 50 percent cutoff, therefore, appears to establish a reasonable tradeoff between these opposing forces, providing good performance across the entire range of loading. The basic philosophy is that under light loads, when data contention levels are low, priority-based waiting is always beneficial. Under heavy loads, however, when data contention levels are high, waiting is the wrong thing to do. WAIT-50 is effective in dynamically making this transition.

6.5. Transaction Slack Ratios

The next experiment examined the effect of having a range of transaction slack ratios on the relative performance of the algorithms. For this experiment, deadline formula DF3 was used to generate a range of slack ratios. *LSF* and *HSF* were set at 2.0 and 4.0, respectively, keeping all other parameters the same as those of the baseline model. When the experiment was conducted under conditions of finite resources, Figures 14a and 14b were obtained. Figures 15a and 15b show the same experiment under conditions of infinite resources. From this set of figures, we draw the following conclusions:

- (1) The priority wait mechanism again exhibits a significant performance improvement potential under normal loads, as WAIT-50 and OPT-WAIT perform much better than OPT-BC there.
- (2) WAIT-50 again provides reasonably good performance across the full range of data contention by switching over from OPT-WAIT-like behavior at low contention levels to OPT-BC-like behavior at high contention levels. It should be noted, however, that the control mechanism slightly underestimates the benefits of waiting at normal loads under infinite resources, and WAIT-50 therefore does not track OPT-WAIT as closely

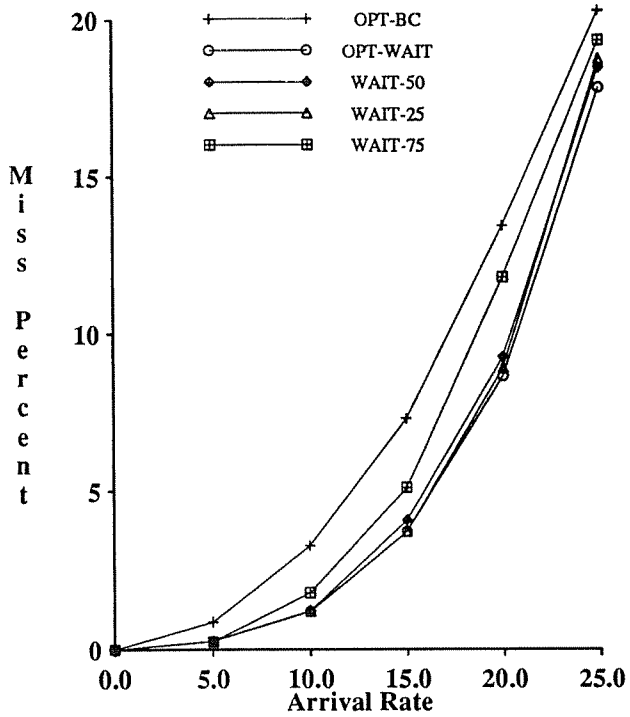


Figure 12a: Control Mechanism (Normal Load)

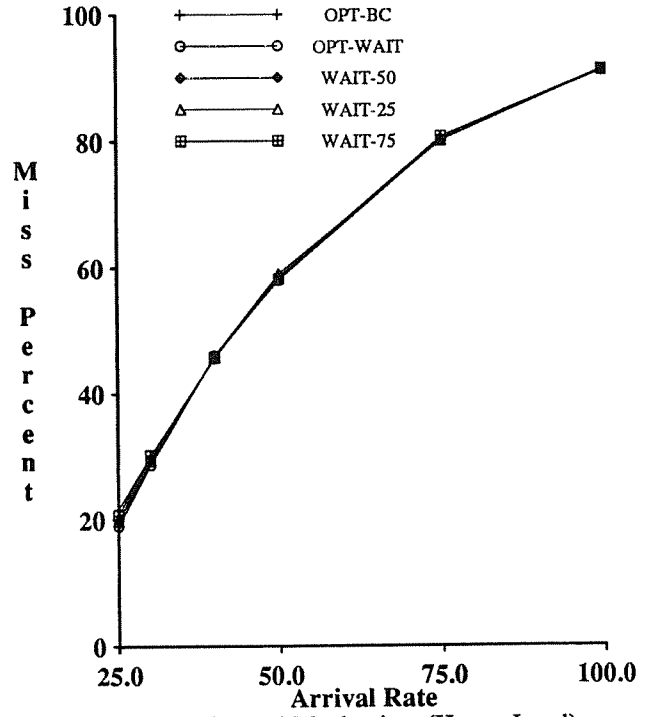


Figure 12b: Control Mechanism (Heavy Load)

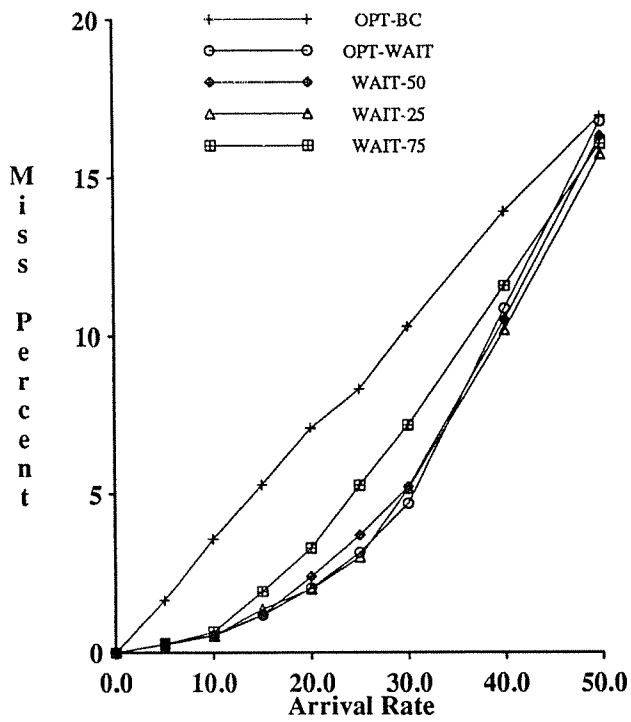


Figure 13a: Infinite Resources (Normal Load)

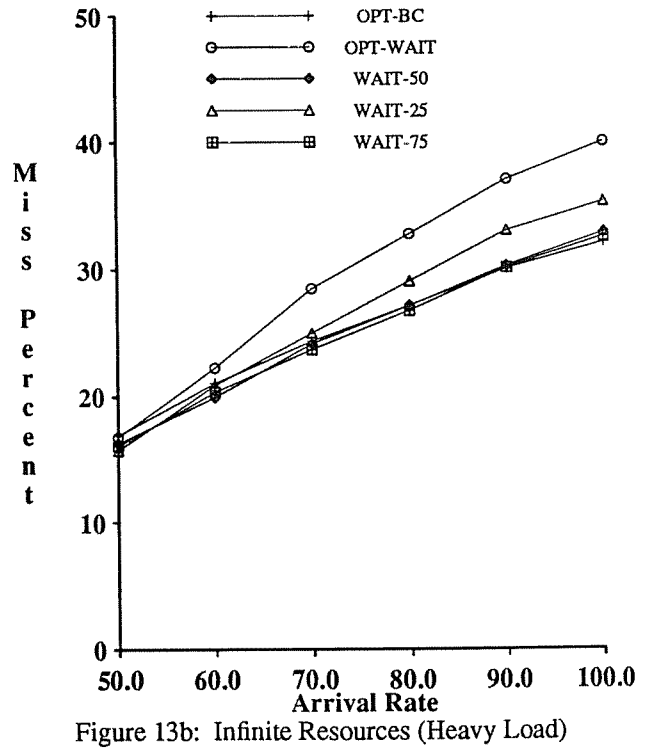


Figure 13b: Infinite Resources (Heavy Load)

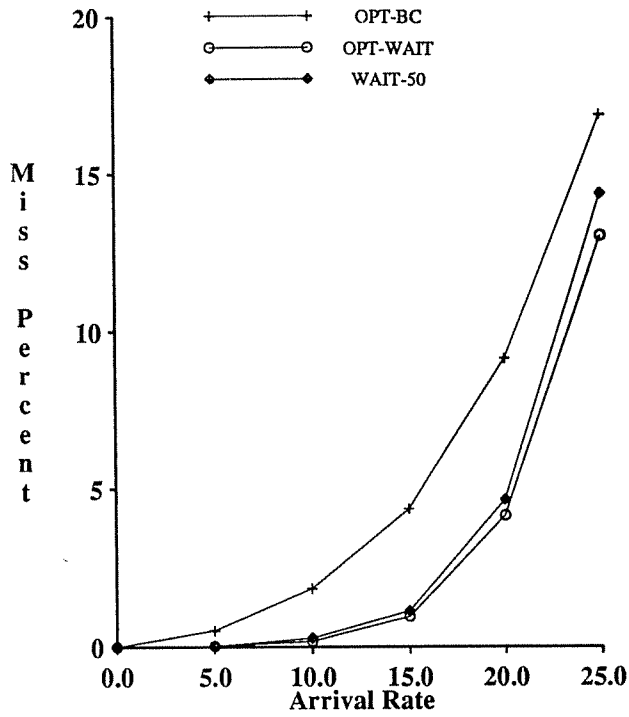


Figure 14a: Slack Ratio (Normal Load)

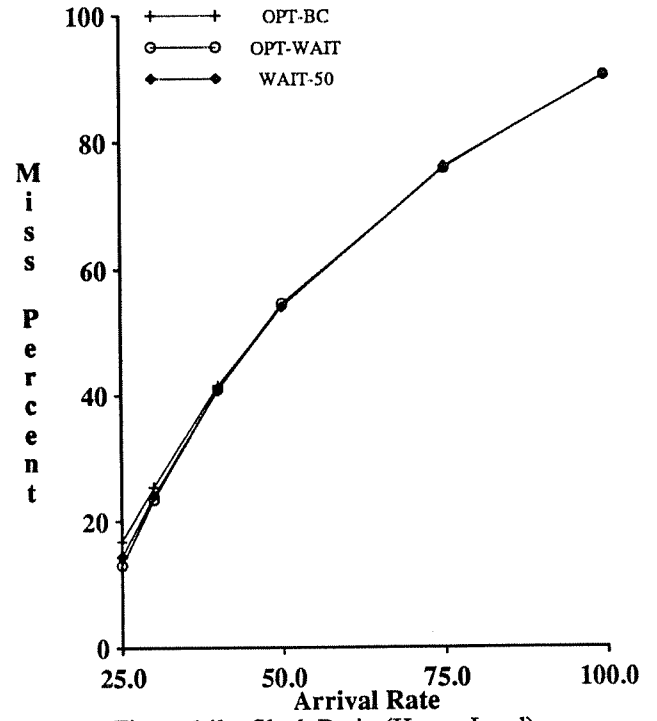


Figure 14b: Slack Ratio (Heavy Load)

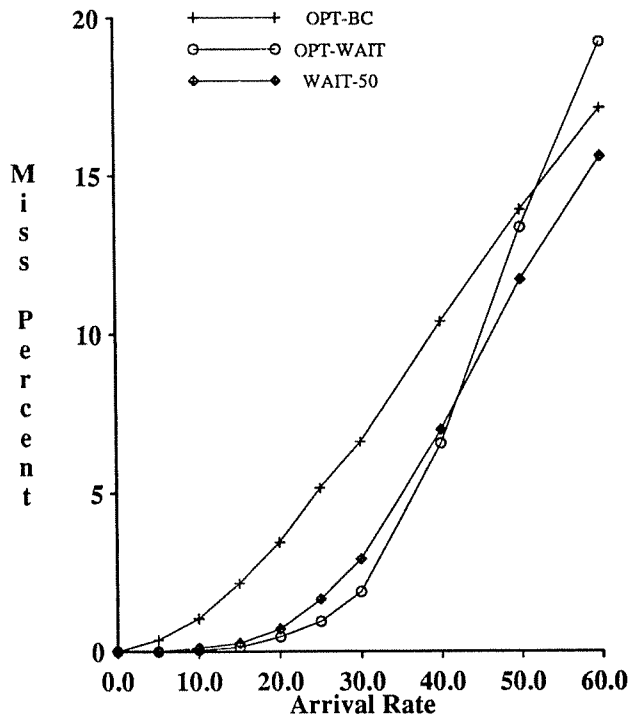


Figure 15a: Infinite Resources (Normal Load)

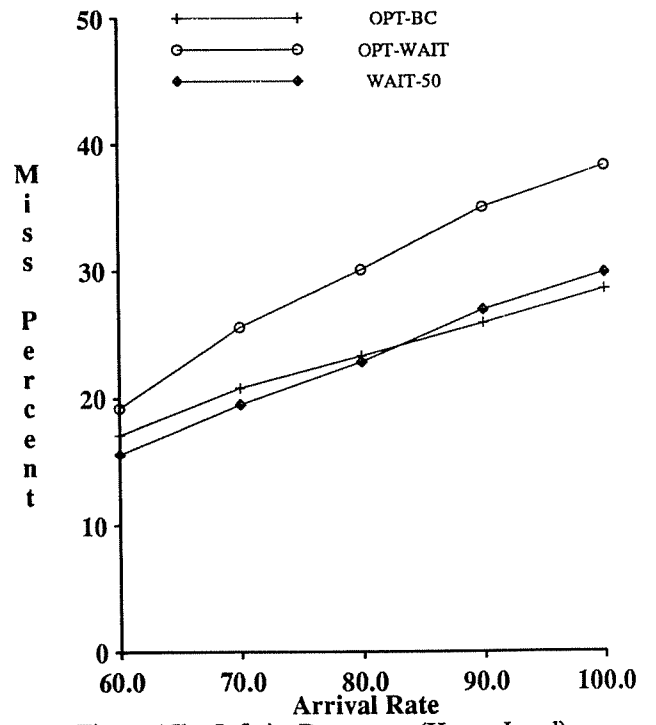


Figure 15b: Infinite Resources (Heavy Load)

there as in some of the other experiments. This is related to the fact that the control mechanism uses just a single, simple parameter to characterize transaction conflict states and therefore cannot be expected to completely capture the performance tradeoffs of waiting versus not waiting. For the most part, however, the parameter is quite effective in delivering good performance.

6.6. Priority Dynamics

We conducted one experiment using *static LTD* as the priority assignment mechanism, keeping the other parameters the same as those of the baseline model. This priority assignment scheme causes priority fluctuations that may lead to priority reversals. The objective of this experiment was to experimentally confirm that WAIT-50 and OPT-WAIT, while being priority cognizant, were, like OPT-BC, immune to priority reversals.¹² The results of the experiment are given in Figures 16a and 16b, which were obtained under finite resources. When the same experiment was carried out under infinite resources, Figures 17a and 17b were obtained. These graphs are qualitatively similar to those obtained in the absence of priority reversals. Therefore, while priority reversals can degrade the performance of priority-cognizant algorithms like 2PL-HP [Hari90] and OPT-SACRIFICE, they do not have an impact on OPT-WAIT and WAIT-50.

7. CONCLUSIONS

In this paper, we have addressed the problem of incorporating transaction deadline information into optimistic concurrency control algorithms. We presented a new real-time optimistic concurrency control algorithm, called WAIT-50, that uses deadline-based transaction priorities to improve data conflict resolution decisions. The algorithm features a *priority wait* mechanism that gives precedence to urgent transactions. This mechanism forces low priority transactions to wait for conflicting high priority transactions to complete, thus enforcing preferential treatment for high priority transactions. We showed that the mechanism has a capacity to eliminate some data conflicts due to its wait component, which causes changes to be made to the commit order of transactions. The priority-wait mechanism provides immunity to priority fluctuations by resolving conflicts in a manner that results in the commit of at least one of the conflicting transactions.

While the priority wait mechanism works well at low data contention levels, it can cause significant performance degradation at high contention levels by generating a steep increase in the number of data conflicts. A simple *wait control* mechanism consisting of a "50 percent" rule is used in the WAIT-50 algorithm to address this problem. The "50 percent" rule is the following: If half or more of the transactions conflicting with a validating transaction are of higher priority, the transaction is made to wait; otherwise, it is allowed to commit.

Using a simulation model of a RTDBS, we studied the performance of the WAIT-50 algorithm over a range of workloads and operating conditions. WAIT-50 was shown to provide significant performance gains at normal loads over OPT-BC, a priority-insensitive optimistic algorithm. The wait control mechanism of WAIT-50 was

¹² With the static LTD priority assignment scheme, it is possible for a waiting transaction to have higher priority transactions in its conflict set even when its deadline is reached. OPT-WAIT and WAIT-50 were therefore modified to ensure that a waiter would commit at its deadline, independent of the composition of its conflict set. Note that this problem does not arise with the Earliest Deadline priority assignment policy, which guarantees that a transaction at its deadline has the highest priority in the system.

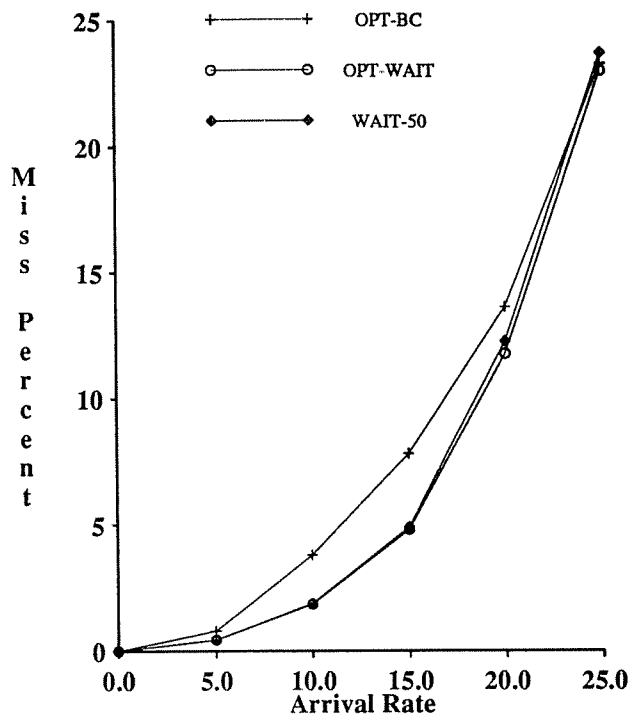


Figure 16a: Priority Reversals (Normal Load)

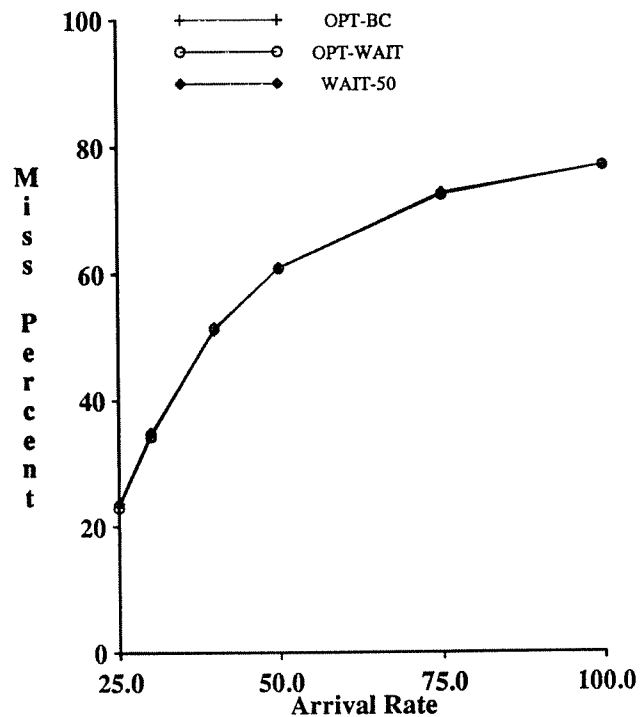


Figure 16b: Priority Reversals (Heavy Load)

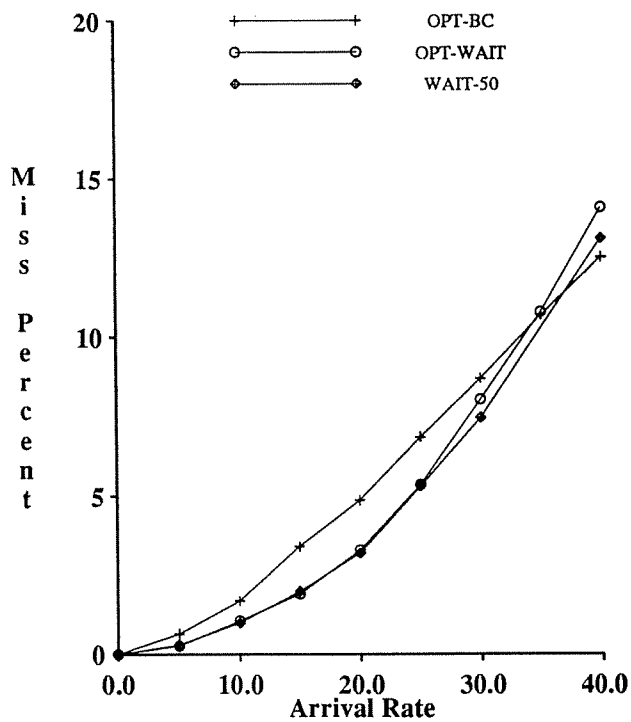


Figure 17a: Infinite Resources (Normal Load)

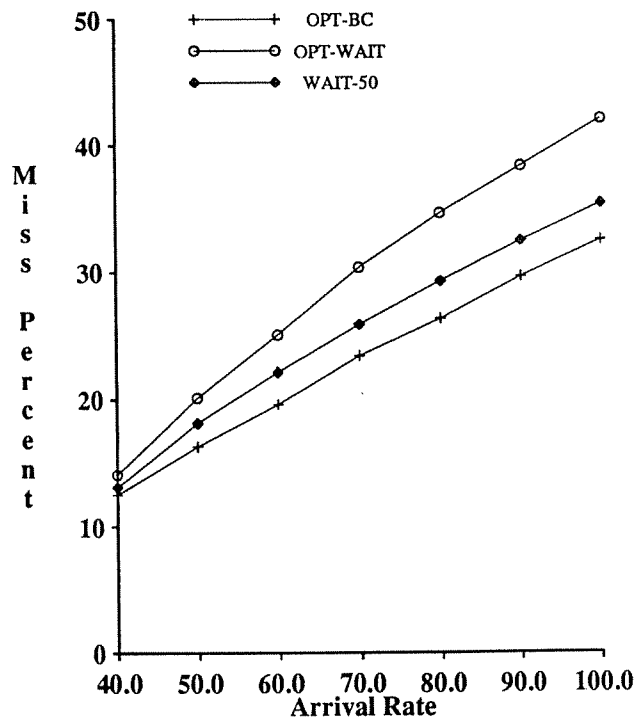


Figure 17b: Infinite Resources (Heavy Load)

found to be effective in maintaining good performance over the entire range of data and resource contention levels. In summary, we conclude that the WAIT-50 algorithm successfully utilizes transaction priority information to provide improved performance in a stable manner.

An interesting future challenge is to develop an analytical model of the WAIT-50 algorithm and to theoretically account for the effectiveness of the 50 percent control rule. Also, our control mechanism, which monitors transaction conflict states, operates on a per-transaction basis. It would be instructive to compare its performance with that of a control mechanism that monitors global data contention levels (e.g., [Care90]).

REFERENCES

- [Abbo88] Abbott, R., and Garcia-Molina, H., "Scheduling Real-Time Transactions: a Performance Evaluation," *Proc. of the 14th International Conference on Very Large Database Systems*, Aug. 1988.
- [Abbo89] Abbott, R., and Garcia-Molina, H., "Scheduling Real-Time Transactions with Disk Resident Data," *Proc. of the 15th International Conference on Very Large Database Systems*, Aug. 1989.
- [Agra87] Agrawal, R., Carey, M., and Livny, M., "Concurrency Control Performance Modeling: Alternatives and Implications," *ACM Trans. on Database Systems*, Dec. 1987.
- [Bern87] Bernstein, P.A., Hadzilacos, V., and Goodman, N., *Concurrency Control and Recovery in Database Systems*, Addison-Wesley, 1987.
- [Care88] Carey, M., and Livny, M., "Distributed Concurrency Control Performance: A Study of Algorithms, Distribution, and Replication," *Proc. of the 14th International Conference on Very Large Database Systems*, Aug. 1988.
- [Care89] Carey, M., Jauhari, R., and Livny, M., "Priority in DBMS Resource Scheduling," *Proc. of the 15th International Conference on Very Large Database Systems*, Aug. 1989.
- [Care90] Carey, M., Krishnamurthi, S., and Livny, M., "Load Control for Locking: The Half-and-Half Algorithm," *Proc. of the 1990 ACM PODS Symposium*, April 1990.
- [Eswa76] Eswaran, K., Gray, J., Lorie, R., and Traiger, I., "The Notions of Consistency and Predicate Locks in a Database System," *Communications of the ACM*, Nov. 1976.
- [Gray79] Gray, J., "Notes on Database Operating Systems," in *Operating Systems: An Advanced Course*, R. Bayer, R. Graham, and G. Seegmuller, eds., Springer-Verlag, 1979.
- [Hari90] Haritsa, J., Carey, M., and Livny, M., "On Being Optimistic about Real-Time Constraints," *Proc. of the 1990 ACM PODS Symposium*, April 1990.
- [Jens85] Jensen, E., Locke, C., and Tokuda, H., "A Time-Driven Scheduling Model for Real-Time Operating Systems," *Proc. IEEE Real-Time System Symposium*, Dec. 1985.
- [Kung81] Kung, H., and Robinson, J., "On Optimistic Methods for Concurrency Control," *ACM Trans. on Database Systems*, June 1981.
- [Lin88] Lin, K., and Lin, M., "Enhancing Availability in Distributed Real-Time Databases," *ACM SIGMOD Record*, March 1988.
- [Livn88] Livny, M., *DeNet User's Guide*, Version 1.0, Comp. Sci. Dept., Univ. of Wisconsin, Madison, 1988.
- [Mena82] Menasce, D., and Nakanishi, T., "Optimistic versus Pessimistic Concurrency Control Mechanisms in Database Management Systems," *Information Systems*, vol. 7-1, 1982.
- [Robi82] Robinson, J., "Design of Concurrency Controls for Transaction Processing Systems," *Ph.D. Thesis*, Carnegie Mellon University, 1982.
- [Sha87] Sha, L., Rajkumar, R., and Lehoczky, J., "Priority Inheritance Protocols: An Approach to Real-Time Synchronization," *Tech. Report No. CMU-CS-87-181*, Carnegie Mellon University, Dec. 1987.

-

APPENDIX

A.1 Implementation of OPT-BC

We present here a lock-based mechanism for implementing the broadcast commit optimistic concurrency control algorithm. Our method is as follows: Whenever a transaction wants to read an object, it sets a read lock on it. If an object is to be written, the update is made to a private copy. After the transaction comes to validation and decides to commit, which requires all of its private copies to be made public, write locks are set on all the objects that are to be updated, and the updates are then performed.

Read locks are individually requested with the synchronous lock call (using System R lock manager notation [Gray79]) LOCK (*ReadSet_i*, SHARED, WAIT), where *ReadSet_i* refers to the specific data object on which the read lock is being requested. All write locks are requested simultaneously at commit time with the non-blocking lock call LOCK (*WriteSet*, EXCLUSIVE, TEST). Write locks preempt read locks and therefore the writeset lock request will always succeed¹³, except as noted below.

In the process of giving locks on the writeset, the lock manager informs the recovery manager of the list of preempted conflicting transactions to be aborted. The transaction abort processing can be done asynchronously by having the lock manager maintain a table of the current status of all executing transactions. The status of a transaction can be either *Running*, *AbortScheduled* or *Committing*. Any further lock request from a transaction scheduled for abort is refused by the lock manager. In particular, a validating transaction will be refused its writeset lock request if it has been scheduled for abort by a previously committed (or currently committing) transaction. A transaction whose lock request is thus denied is put to sleep and aborted at a later time.

Once a transaction has obtained its write locks, it releases all its read locks with the following call to the lock manager UNLOCK(*ReadSet* – *WriteSet*). It then makes its updates and commits, releasing all write locks at the end with the call: UNLOCK(*WriteSet*).

A.2 Implementation of OPT-WAIT

In the OPT-BC algorithm, when a validating transaction makes the request for simultaneous locks on its entire writeset, the request is guaranteed to succeed unless the transaction has been scheduled for abort. For the OPT-WAIT algorithm, however, the success of the request is also dependent on the current composition of the conflict set of the validating transaction. If there are higher priority transactions in the conflict set, the writeset lock request is refused and the transaction is made to wait. A waiting transaction needs to periodically re-issue its writeset lock request since the composition of its conflict set is a function of time. A reasonably efficient method to do this is for the lock manager to maintain a list of waiters for each running transaction together with a count of conflicting higher priority transactions for each waiting transaction. This count is evaluated at the time of making the writeset lock request, and if the count is non-zero, the requesting transaction is put to sleep. Whenever a transaction commits, restarts, or is discarded, the high priority count of each of its waiters is decremented. The terminated transaction is

¹³ A special case where the writeset lock request may fail temporarily is for transactions with "blind writes" [Bem87]. A transaction with blind writes will have its writeset lock request granted once the write locks currently existing on objects in its "blind write set" are released.

also taken off any list of waiters in which it is a member. If a waiter's high priority count goes to zero, it is woken up. The awakened waiter then reissues its writeset lock request. This process continues until the waiter is either restarted or is given the writeset lock.

The WAIT-50 algorithm can be implemented, in a similar fashion, by making simple extensions to the scheme described above.

A.3 Conflict Elimination Probability

In this section, we carry out a simple probabilistic analysis of the extent to which the waiting scheme can reduce data conflicts. To do this, we ask the following question: Given that transaction A conflicts with transaction B ($A \rightarrow B$), what is the probability that transaction B does not conflict with A ($B \nrightarrow A$)? Assuming a uniform database access pattern and that $WriteSet \subseteq ReadSet$ for all transactions, and using *overlap* to refer to the cardinality of the set $ReadSet_A \cap ReadSet_B$, the probability $\Pi_{B \nrightarrow A | A \rightarrow B}$ is given by the following expression

$$\Pi_{B \nrightarrow A | A \rightarrow B} = \sum_{k=1}^{\max(overlap)} Pr(B \nrightarrow A | k \text{ overlap}) Pr(k \text{ overlap} | A \rightarrow B)$$

which can be re-written as

$$\Pi_{B \nrightarrow A | A \rightarrow B} = \frac{1}{Pr(A \rightarrow B)} \sum_{k=1}^{\max(overlap)} Pr(B \nrightarrow A | k \text{ overlap}) Pr(k \text{ overlap}) Pr(A \rightarrow B | k \text{ overlap}) \quad (A1)$$

The terms in Equation (A1) can be evaluated as shown in Table A1, where the symbols have the following meaning:

N	=	number of data items in the database
R_A	=	size of readset of transaction A
R_B	=	size of readset of transaction B
W_A	=	size of writeset of transaction A
W_B	=	size of writeset of transaction B

For the sake of simplicity, let us assume that every transaction has the same ratio between the sizes of its write set and read set, and denote it by f ($0 \leq f \leq 1$). A sample plot of $\Pi_{B \nrightarrow A | A \rightarrow B}$ as a function of f is shown in Figure 18, for $N = 100$ and $N = 1000$, keeping the read set sizes of both A and B constant at 16. Note that for $N \gg (R_A, R_B)$, which is usually true since database sizes are typically much larger than transaction sizes, Equation (A1) reduces to

$$\Pi_{B \nrightarrow A | A \rightarrow B} \approx 1 - f \quad (A2)$$

and this is evident in the $\Pi_{B \nrightarrow A | A \rightarrow B}$ curve for $N = 1000$ in Figure 18. This means that if the write fraction is 0.25, for example, then by waiting, we can resolve about 75% of the conflicts without restarting either transaction.

It should be noted that neither Equation (A1) nor Equation (A2) gives the *actual probability* with which conflicts will be reduced by waiting in our system – this is because some of the assumptions made in deriving the equations do not hold in our scenario. These assumptions include the facts that only conflicts with a single

$Pr(A \rightarrow B)$	$=$	$1 - \frac{\binom{N - W_A}{R_B}}{\binom{N}{R_B}}$
$Pr(B \rightarrow A \mid k \text{ overlap})$	$=$	$\frac{\binom{R_B - k}{W_B}}{\binom{R_B}{W_B}}$
$Pr(k \text{ overlap})$	$=$	$\frac{\binom{R_A}{k} \binom{N - R_A}{R_B - k}}{\binom{N}{R_B}}$
$Pr(A \rightarrow B \mid k \text{ overlap})$	$=$	$1 - \frac{\binom{R_A - k}{W_A}}{\binom{R_A}{W_A}}$
$\max(\text{overlap})$	$=$	$\min(R_A, R_B)$

Table A1: Term Evaluations

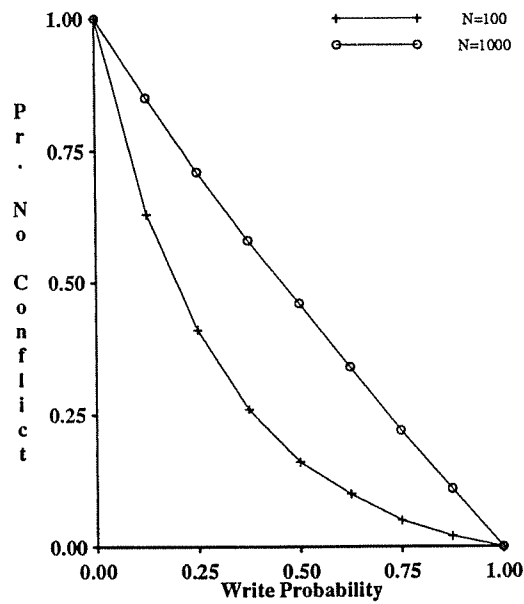


Figure 18: Uni-directional conflict probabilities.

transaction are considered, and that transaction priorities are not taken into account. The equations provide, however, a rough idea of the extent to which waiting can be beneficial with respect to reducing data conflicts.

A.4 Disk Priority Mechanism

When transaction priorities are taken from a continuum, as in the case of the Earliest Deadline priority policy, for example, the assignment of transaction requests to disk priority levels becomes a non-trivial problem. The problem is rendered difficult by the fact that the granularity to be used in demarcating priority levels is a function of the workload characteristics. For example, consider the following scenario: Transactions T_1 , T_2 , and T_3 , have deadlines 10, 20 and 30, respectively. They each make a request to the same disk at clock time = 5. Now, should the system assign all the requests to the same priority level, or instead, perhaps, to three consecutive levels? The decision is really dependent on what constitutes the system's notion of "close" priorities and therefore must be based on workload characteristics.

Since we have full knowledge of the workload generator characteristics in this study, and since the workload is "well-behaved", it was possible for us to develop a mapping formula that resulted in the range of transaction priorities being evenly spread over the disk levels. The formula used was:

$$DiskLevel_T = \left\lceil NumDiskPrio * \left[1 - \frac{DiskValue_T}{MaxDiskValue} \right] \right\rceil$$

In this equation, $DiskValue_T$ is transaction T 's disk value, computed as $(P_T - clock)$ for Earliest Deadline and P_T for static LTD, where P_T is the priority of transaction T . $MaxDiskValue$ is a normalizing constant set equal to the largest possible disk value, $HSF * R_{max}$, where HSF and R_{max} are as defined in Section 5. (Note that the greater a request's disk value, the lower is its priority at the disk). With this mapping, the range of transaction priorities is spread evenly across the NumDiskPrio priority levels at each disk.

