

Priority assignment in real-time active databases¹

Rajendran M. Sivasankaran, John A. Stankovic, Don Towsley, Bhaskar Purimetla, Krithi Ramamritham

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

Edited by Henry F. Korth and Amit Sheth. Received November 1994 / Accepted March 20, 1995

Abstract. Active databases and real-time databases have been important areas of research in the recent past. It has been recognized that many benefits can be gained by integrating real-time and active database technologies. However, not much work has been done in the area of transaction processing in real-time active databases. This paper deals with an important aspect of transaction processing in real-time active databases, namely the problem of assigning priorities to transactions. In these systems, time-constrained transactions trigger other transactions during their execution. We present three policies for assigning priorities to parent, immediate and deferred transactions executing on a multiprocessor system and then evaluate the policies through simulation. The policies use different amounts of semantic information about transactions to assign the priorities. The simulator has been validated against the results of earlier published studies. We conducted experiments in three settings: a task setting, a main memory database setting and a disk-resident database setting. Our results demonstrate that dynamically changing the priorities of transactions, depending on their behavior (triggering rules), yields a substantial improvement in the number of triggering transactions that meet their deadline in all three settings.

Key words: Active databases – Coupling mode – Deadlines – ECA – Priority assignment – Real-time databases

1 Introduction

Traditionally, in soft real-time database systems, a transaction is considered a monolithic unit of work with a given deadline. In these systems, priorities are assigned to these transactions and the transactions are scheduled based on their priorities. The priority assignment usually takes into account the deadlines of the transactions because the underlying assumption is that the deadline reflects the urgency of completing the transaction. Scheduling policies such as earliest

deadline first (EDF) and least slack first (LSF) are examples of policies that account for deadlines. The performance implications of time cognizant priority assignment policies have been studied in detail in soft real-time database systems (Abbott and Garcia Molina 1992; Huang et al. 1989). These studies have concluded that deadline-cognizant priority assignments provide significantly higher performance than priority assignments that ignore deadlines. In this paper, our goal is to study and evaluate time-cognizant priority assignment policies in a real-time active database. A real-time active database is a database system where transactions have timing constraints such as deadlines, where transactions may trigger other transactions, and where data may become invalid with the passage of time. There are many applications, such as cooperative distributed navigation systems and intelligent network services, where real-time active database technology is extremely useful (Purimetla et al. 1993; Sivasankaran et al. 1993).

Before presenting a detailed description of the problem we address in this paper, we give a brief introduction to active databases. The building block of an active database system is the event-condition-action (ECA) rule. The semantics of the ECA rule is that, if the specified event occurs and if the condition is true, then the specified action is to be executed. Some examples of events are begin,² commit/abort of a transaction, accessing a data item, or reaching a specific point in time. A condition is usually a predicate on the database state. An action is the transaction that is executed in *reaction* to a specific *situation* which is a combination of events and conditions. The transaction that fires the rules is called the triggering transaction, and the action that is executed because of the rule firing is called the triggered transaction. In this paper we refer to the transactions that trigger other transactions as *active* transactions or *parent* transactions. An active transaction has a set of triggered transactions that are executed either as part of the active transaction or separately, depending on the type of the *coupling mode* between the parent and the triggered transactions (Dayal et al. 1990). There are three types of coupling

² The begin event may, in turn, have been caused by some external environment event such as an obstacle identified by a sensor (Purimetla et al. 1993; Sivasankaran et al. 1993)

¹ A shorter version of this paper appeared in PDIS 1994

modes: *immediate*, *deferred* and *independent*. The transactions triggered in those modes are referred to as immediate, deferred and independent transactions, respectively. Immediate and deferred transactions are executed as part of the parent transaction, whereas independent transactions are executed independently. Since immediate and deferred transactions are part of the parent transaction, we also refer to them as subtransactions. In our model, an immediate transaction is executed as soon as it is triggered and the parent transaction is suspended until it completes. A deferred transaction is executed after the parent completes execution, but before it commits. In this paper, by a parent completing execution we mean that the parent has finished all its work but has not committed its results. Immediate and deferred transactions commit if and only if the parent commits. We do not consider independent transactions in our study.

Due to the rule firings, an active transaction dynamically generates additional work. In order for the time-cognizant scheduling policies to perform well, they should take into account the work that is dynamically generated. This aspect of the problem makes this scheduling problem different from the classical hard real-time scheduling problem, where execution times are assumed to be known in advance (Klein et al. 1993; Lawler 1983; Xu and Parnas 1990). The priority-driven nature of real-time transaction processing gives rise to the question of how to assign priorities to the parent and to all triggered actions (Purimetla et al. 1993; Sivasankaran et al. 1993). We believe that the strategy to assign priorities can have a significant impact on the performance of the system, as triggered actions must contend with ongoing transactions for resources. In this paper, we address the problem of assigning priorities to triggering and triggered transactions. We introduce and evaluate three policies for assigning priorities to *immediate* and *deferred* transactions.

Our main contributions are:

- The development of two priority assignment policies that account for the work dynamically generated by active transactions
- A comparison of the two priority assignment policies with a baseline policy using a real-time active database simulator in three settings, in a real-time task setting, a main memory database setting and a disk-resident database setting
- A demonstration that for mixed workloads consisting of triggering and non-triggering transactions, priority assignment policies that take into account the dynamic work generated reduce the deadline miss ratio of the triggering transactions significantly at the cost of a very small increase in the deadline miss ratio of non-triggering transactions when compared to the baseline policy
- The identification of the trade-offs between the performance of triggering and non-triggering transactions offered by the different policies, thereby enabling an implementor to select from various policies depending on the relative importance of the triggering and non-triggering transactions in the system
- The identification of differences and probable reasons for the differences in the relative performance of the policies in the three settings

We discuss related work in Sect. 2. In Sect. 3, we explain our transaction model. Section 4 gives a detailed explanation of the priority assignment policies. We discuss the simulator, experiments and results in Sect. 5. We summarize our main results and discuss future work in Sect. 6.

2 Related work

Over the past few years active databases and real-time databases have become important areas of research. Experimental studies reported (Abbott and Garcia Molina 1992; Huang et al. 1989, 1991a, 1991b; Son and Park 1994) are very comprehensive and cover most aspects of real-time transaction processing, but have not considered active workloads (workloads that generate additional transactions dynamically) and have not addressed the problem of subtransaction priority assignment. There have been both theoretical and experimental studies in active databases (Carey et al. 1991; Dayal et al. 1988, 1990; McCarthy and Dayal 1989). Most of the studies have concentrated on the specification of ECA rules (Anon 1992). Experimental work on active databases (Carey et al. 1991) has been performed in a non-real-time setting. Experiments (Carey et al. 1991) show the impact of transaction boundaries and data sharing on the performance of active databases. We address a different problem; that of assigning priorities to triggered transactions in real-time active databases.

The relationship between real-time databases and active databases was discussed briefly by Ramamritham (1993) where it was noted that a missing ingredient in active databases is the active pursuit of timely processing of actions to do real-time processing. In short, past studies on real-time transaction processing have not dealt with active workloads, and studies on active databases have not dealt with real-time transactions.

The problem of assigning deadlines to the parallel and serial subtasks of complex distributed tasks in a real-time system has been studied (Kao and Garcia Molina 1993) through simulation. This work differs from the present work in two ways. First, in Kao and Garcia Molina (1993) the structures of all of the complex tasks are assumed to be known in advance. Second, the experimental system considered in Kao and Garcia Molina (1993) is not a database system, nor does it include an active component. The system we consider for our experiments is a real-time active database with unpredictable data accesses and rule firings.

3 Transaction model

In this section we describe the transaction model we have considered and the simplifying assumptions we have made to study priority assignment in a real-time active database. The performance metric we use in our study is missed deadline percentage (MDP), i.e., the percentage of transactions that miss their deadlines, which is a traditional metric used to evaluate performance in real-time systems.

Data in our system is modeled as objects that have methods associated with them. A transaction in our system is a series of method executions on objects. We consider two

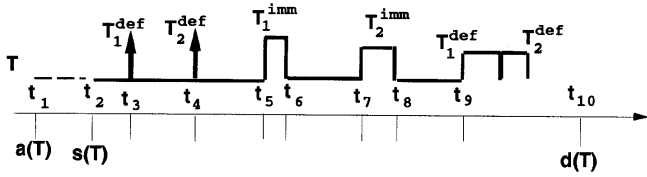


Fig. 1. Life of a complex active transaction

classes of transactions in the system: *non-triggering* (class **NT**) and *triggering* (class **T**). An **NT** transaction is a simple transaction which does not cause any rules to fire. A **T** transaction, on the other hand, can trigger subtransactions upon the occurrence of an event. We consider class **T** transactions that trigger only immediate and deferred transactions. At most one subtransaction is triggered upon the occurrence of an event. We do not consider cascading rule firings, i.e., in our model triggered transactions do not trigger further transactions. We address only *firm real-time* systems where the *value* of the transaction drops to zero once the deadline expires. We abort the transaction that misses its deadline, along with all its related subtransactions. In our current model, only the object events, i.e., an execution of a method on an object, activate rules and cause the triggering of subtransactions. Transaction and temporal events do not activate rules. We make this assumption to simplify the experimentation, since the only effect of transaction and temporal events triggering rules in our study would be to increase the number of subtransactions triggered. This effect can be achieved by increasing the probability that an object event can trigger a rule.

4 Priority assignment for triggered transaction

In real-time active databases, a transaction with time constraints can trigger subtransactions. Traditionally, priority-driven scheduling has been used in real-time systems. Hence, the problem is how to assign a priority to a subtransaction given the priority of the parent transaction. Another problem is that, as a transaction triggers subtransactions either in immediate or deferred modes, the amount of work to be done on behalf of the transaction before it commits also increases. Consequently, the transaction is less likely to complete successfully than another transaction with the same deadline which does not trigger any subtransactions, since it faces more contention for resources during its lifetime.

To better illustrate the problem being addressed, we provide an example of the structure of a complex active transaction executing on a uniprocessor. Figure 1 shows the life of a complex active transaction T that triggers transactions in immediate and deferred modes. Transaction T arrives at time t_1 ($a(T)$) with t_{10} as its deadline ($d(T)$) and is started at time t_2 ($s(T)$). Deferred transactions T_1^{def} and T_2^{def} are triggered at times t_3 and t_4 , respectively. Immediate transactions T_1^{imm} and T_2^{imm} are triggered at times t_5 and t_7 , respectively. Fig. 1 shows that the transactions T_1^{imm} and T_2^{imm} execute immediately while T is suspended. Finally, the figure shows that, once the parent T completes at t_9 , the deferred transactions T_1^{def} and T_2^{def} execute. The problem we address is

that of assigning priorities to subtransactions T_1^{imm} , T_2^{imm} , T_1^{def} and T_2^{def} , and the problem of dynamically reassigning the priority of the parent transaction T during its lifetime. It should be noted that the deferred transactions can execute in parallel if operating on a multiprocessor or in a distributed system.

Later in this section, we describe three priority assignment policies, PD, DIV and SL³, for assigning priorities to the parent, immediate and deferred subtransactions in an active real-time database system. PD is a *static* baseline policy where the priority of the parent does not change with time. DIV and SL are *dynamic* policies that change the priority of the parent, depending on the amount of dynamic work it has generated. In all cases, the priorities are assigned to the triggered transactions when they start execution and priorities are not changed subsequently during their execution. It should be noted that the three algorithms use increasingly more knowledge (estimates) about the transactions in their assignment of priorities. PD is a pure EDF policy for both transactions and subtransactions and uses no other information about transactions. DIV makes a simple adjustment to EDF where priorities (based on deadlines) are dynamically changed when new subtransactions are triggered. Our hypothesis is that very simple modifications will give significant improvements. Note that DIV uses estimates of execution times of the parent transaction and subtransactions that it has triggered. SL is an LSF policy that is slightly more sophisticated than DIV; it uses estimates of execution times of the parent transaction and triggered subtransactions and estimates regarding the subtransactions that might be triggered in the future. Our hypothesis is that, when such information is available, it can provide additional performance benefits.

Below we describe how the policies PD, DIV and SL assign priorities to the parent transaction and the triggered subtransactions for two cases: when the triggered subtransaction is of type immediate and when the triggered subtransaction is of type deferred. Before describing the policies, we introduce some attributes of a transaction and other related terms that we require in order to explain our priority assignment policies.

4.1 Attributes of a transaction

Let us consider a transaction T for which we define the following attributes. Some attributes of the transaction may change during the execution of the transaction and these attributes are subscripted with t , where t is the time at which the attribute's value is represented.

- $a(T)$: the arrival time of T
- $s(T)$: the start time of T
- $d(T)$: the deadline of T
- $n_t^{def}(T)$: the number of deferred transactions triggered
- $P_t(T)$: the priority of T at time t

The policies we have developed use estimates of certain quantities associated with a transaction such as its length and execution time. The following are some of the estimates that we use in our policies. An example of how these estimates

³ The name PD derives from Parent Deadline; DIV from DIViding the parent's effective slack; SL from adjusting the average case SLack

are obtained can be found in Sect. 5 where we describe the parameters for the experiments.

- $X_t(T)$: the estimated remaining execution time for T at time t
- $C_t(T)$: the estimated completion time of T at time t ;
 $C_t(T) = t + X_t(T)$
- $S_t(T)$: the estimated slack of T at time t ;
 $S_t(T) = d(T) - C_t(T)$
- $m_t^{imm}(T)$: the estimated number of immediate transactions triggered by T after time t
- $m_t^{def}(T)$: the estimated number of deferred transactions triggered by T after time t
- $\bar{X}^{imm}(T)$: the estimated average execution time of an immediate subtransaction triggered by T
- $\bar{X}^{def}(T)$: the estimated average execution time of a deferred subtransaction triggered by T
- $X_t^{avg}(T)$: the estimated average case remaining execution time of T including its subtransactions at time t ;
 $X_t^{avg}(T) = X_t(T) + n_t^{def}(T) * \bar{X}^{def}(T) + m_t^{def}(T) * \bar{X}^{def}(T) + m_t^{imm}(T) * \bar{X}^{imm}(T)$
- $S_t^{avg}(T)$: the estimated average case slack of T including its subtransactions at time t ;
 $S_t^{avg}(T) = d(T) - t - X_t^{avg}(T)$

$X_t(T)$ is the estimate (at time t) of the remaining time it will take a transaction or subtransaction to execute. This quantity, when estimated at the start of a transaction or when a subtransaction is triggered, is the estimated execution time of that (sub)transaction. The estimated average case execution time of transaction T , $X_t^{avg}(T)$ includes the execution time left at time t for T , execution times of all of the deferred transactions T has triggered prior to time t , and estimates of execution times of the subtransactions which the transaction T may trigger during its remaining lifetime. We believe that by analyzing the characteristics of a real-time database application one might be able to obtain information such as the estimated execution time of transactions, the kind of actions they trigger and the triggering probability.

It should be noted that a lower value of $P_t(T)$ indicates higher priority. For instance, if there are two transactions T_1 and T_2 , and $P_t(T_1) < P_t(T_2)$, then T_1 gets priority over T_2 . Also, in all the policies, the deadline of the subtransaction is set to the parent's deadline. For instance, if a transaction T triggers a subtransaction T^{sub} transaction at time t , then

$$d(T^{sub}) = d(T)$$

4.2 Priority assignment when immediate subtransactions are triggered

1. *PD*: Immediate subtransactions are assigned a priority equal to the deadline of the parent. Further, the priority of the parent transaction which is based on its deadline does not change with the triggering of subtransactions. This is a very simple baseline algorithm. All the actions done on behalf of a transaction get the same priority as the transaction itself at any point during its lifetime. Let us consider a transaction T triggering its i th immediate transaction T_i^{imm} at time t . Then

$$P_t(T_i^{imm}) = d(T)$$

2. *DIV*: The parent's estimated effective slack is divided equally among the *current* immediate subtransaction triggered, the deferred subtransactions triggered prior to the triggering point (current time) and the parent. This quantity (obtained by dividing the parent's estimated effective slack) is added to the estimated completed time of the immediate subtransaction to give the priority of the subtransaction. The parent's deadline is also adjusted dynamically to reflect the work that has been triggered dynamically (Eq. 2).

Let us consider a transaction T triggering its i th immediate transaction T_i^{imm} at time t . The estimated effective slack is calculated by subtracting the sum of estimated execution times of deferred transactions that have been triggered until time t and estimated execution time of the i th immediate transaction from the parent's estimated slack as shown in Eq. 1.

$$P_t(T_i^{imm}) = C_t(T_i^{imm}) + \frac{S_t(T) - \left(X_t(T_i^{imm}) + \sum_{j=1}^{n_t^{def}(T)} X_t(T_j^{def}) \right)}{n_t^{def}(T) + 2} \quad (1)$$

Let us assume T_i^{imm} finishes at time $w > t$. The parent's priority is reassigned at time w as follows

$$P_w(T) = P_t(T) - (w - t)$$

For instance, in Fig. 1, for T_1^{imm} at the triggering point t_5 , its priority is assigned as follows:

$$P_{t_5}(T_1^{imm}) = C_{t_5}(T_1^{imm}) + \frac{S_{t_5}(T) - \left(X_{t_5}(T_1^{imm}) + X_{t_5}(T_1^{def}) + X_{t_5}(T_2^{def}) \right)}{4}$$

The priority of the parent transaction will be modified when T_1^{imm} completes at t_6 as follows:

$$P_{t_6}(T) = P_{t_5}(T) - (t_6 - t_5) \quad (2)$$

The main idea behind this policy is that of giving higher priorities to class **T** transactions, which have more work to do before completion. This should increase the likelihood of the class **T** transactions meeting their deadlines. This policy only uses the estimates of execution times of subtransactions that have already been triggered. It does not use any knowledge about future triggering of transactions. The priority assigned to the subtransaction can be thought of as a virtual deadline which is the sum of its estimated completion time and some slack that it gets from the parent. This virtual deadline is then used to schedule the subtransaction using the EDF algorithm.

3. *SL*: The average case slack [$S_t^{avg}(T)$] of the parent is adjusted at each potential triggering point and used as the priority for both the triggered and triggering transactions. The initial value of slack is assigned based on estimates of the remaining execution time for a transaction and its subtransactions. The slack is then adjusted at each object event based on whether the parent transaction triggers a subtransaction or not. The triggered transactions are assigned the same slack as the parent, i.e., they are executed at the same

priority. Let us assume a transaction T starts at time t_0 and time t is a potential triggering point where T could trigger its i th immediate transaction T_i^{imm} . Initially the priority of T (average case slack) at time t_0 is set as follows:

$$P_{t_0}(T) = S_{t_0}^{avg}(T)$$

Now at time t , if subtransaction T_i^{imm} is triggered, the priority is adjusted as follows:⁴

$$P_t(T) = d(T) - t - X_t^{avg}(T) - X_t(T_i^{imm})$$

$$P_t(T_i^{imm}) = P_t(T)$$

If a subtransaction is not triggered at time t then the slack adjustment is as follows:

$$P_t(T) = d(T) - t - X_t^{avg}(T) \quad (3)$$

In Fig. 1, the priorities for T and T_1^{imm} will be

$$P_{t_5}(T) = d(T) - t_5 - X_{t_5}^{avg}(T) - X_{t_5}(T_1^{imm})$$

4.3 Priority assignment when deferred subtransactions are triggered

The priority assignment for deferred transactions is very similar to that of immediate transactions. Under the DIV and SL policies the priority of the parent changes when it triggers a deferred transaction. The deferred transaction is executed after the parent transaction finishes execution. Let us consider a transaction T that is triggering its i th deferred transaction at time t . The priority assignment for the parent transaction T under the three policies is:

1. PD protocol:

$$P_t(T) = d(T)$$

2. DIV protocol:

$$P_t(T) = P_t(T) - X_t(T_i^{def})$$

3. SL protocol:

$$P_t(T) = d(T) - t - X_t^{avg}(T) - X_t(T_i^{def})$$

For the SL protocol, if the transaction T does not trigger a deferred transaction then the priority assignment is the same as Eq. 3 in the last section. The following equations illustrate the deadlines/slack assignments for the deferred transactions. Let us assume that transaction T has triggered m deferred transactions before completing execution. In the following equations, we consider the assignment of priority for T_j^{def} , $j \leq m$, at time w , which is the time when the deferred transactions start their execution (after the parent transaction completes execution).

1. PD protocol:

$$P_w(T_j^{def}) = d(T)$$

2. DIV protocol:

$$P_w(T_j^{def}) = C_w(T_j^{def}) + \frac{S_w(T) - \sum_{k=1}^m (X_w(T_k^{def}))}{m}$$

⁴ Note that we update priorities only at object events

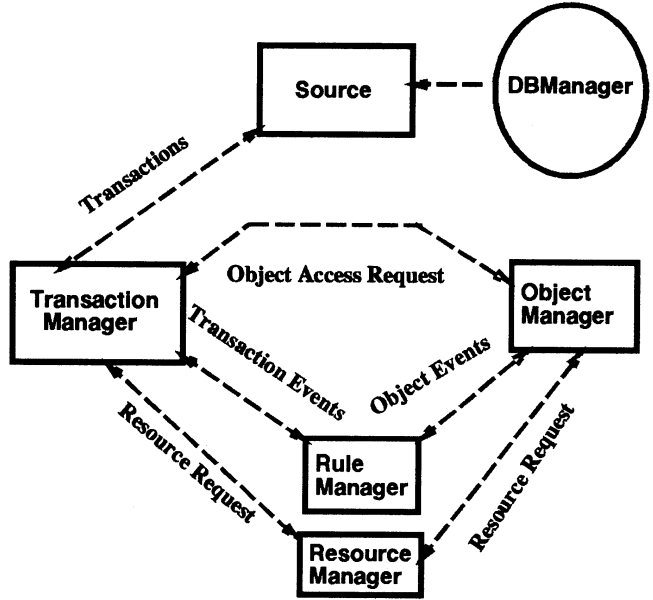


Fig. 2. Simulator architecture

3. SL protocol:

$$P_w(T_j^{def}) = d(T) - w - \sum_{k=1}^m (X_w(T_k^{def}))$$

5 Experimental Results

We begin this section with a brief description of *RADEx* (Real-time Active Database Experimental simulator) and its validation. We then discuss the experimental setup, along with the assumptions made in our experiments. We also present a table of important parameters and their values. Finally, we describe each set of experiments and an analysis of the results. In the experiments, 95% confidence intervals have been obtained whose widths are less than $\pm 5\%$ of the point estimate for the MDP.

5.1 Simulation model

Our performance model of an active real-time database was implemented using the DeNet Simulation Language (Livny 1990). *RADEx* is made up of five active modules – *source*, *transaction manager*, *object manager*, *resource manager*, *rule manager*, and a passive module *DB manager*. Fig. 2 illustrates the architecture of the simulator. The following is a detailed description of the modules:

- *DB manager*: This is the passive module that models the data. The data is modeled as having a certain number of object classes and each object class has a certain number of instances. Each object class has a certain number of methods defined which are used to access the object. Each object instance in the database is mapped to a page or number of pages in secondary storage.

- *Source*: The source (transaction generator) generates the incoming transactions into the system. One can view the source as the application or the environment in which the real-time active database is used. It generates transactions with timing constraints with a specified arrival distribution. In our study we consider only aperiodic transaction streams with firm deadlines.
- *Transaction manager*: The transaction manager is responsible for scheduling and execution of the transactions it receives from the source. It executes the submitted transactions by requesting the object manager to execute the specific methods on specific objects. It handles the various transaction events: begin, commit and abort. It informs the rule manager of the transaction events.
- *Object manager*: The object manager is responsible for concurrency control and sending messages to the rule manager when an object event occurs. The transaction manager sends requests to the object manager for access to objects. If the request for an object can be satisfied, the object manager sends the *object granted* message back to the transaction manager. If the request cannot be satisfied it sends an abort message back to the transaction manager. Conflict resolution in this module is based on priorities, where the lower priority conflicting transaction waits or gets aborted depending on whether it is the requester or holder of locks.
- *Rule manager*: The rule manager models the active workload in the system. The rule firings are modeled probabilistically, i.e., a rule is fired with a certain probability. The rule manager checks to see if any rules are triggered whenever it gets an event notice from the transaction manager or the object manager. It models the condition evaluation, and, finally, generates the transactions corresponding to the actions of the rules triggered if their conditions are satisfied and submits them to the transaction manager.
- *Resource manager*: The resource manager simulates the CPUs, disks and the main memory buffer. The object manager makes requests to the resource manager for the necessary pages or for CPU time to execute the methods. The transaction manager requests CPU time and buffer space to load transactions from the resource manager. The CPU, disk, and memory resource scheduling are priority driven. Our resource model is a multiprocessor, multidisk, shared-memory system. The incoming resource requests are queued in a common CPU queue or in a randomly selected disk queue depending on the kind of request.

We do priority-driven preemptive scheduling. The scheduling and conflict resolution decisions made in different modules of the simulator are independent of each other. No global scheduling decisions are made. For instance, the scheduling decision made in the resource manager is independent of the one made in the transaction manager. Overload management in our study is based on the *not tardy* policy, i.e., a transaction is aborted as soon as its deadline expires. This corresponds to firm real-time transactions. This policy assumes that finishing a transaction after its deadline expires does not impart any value to the system.

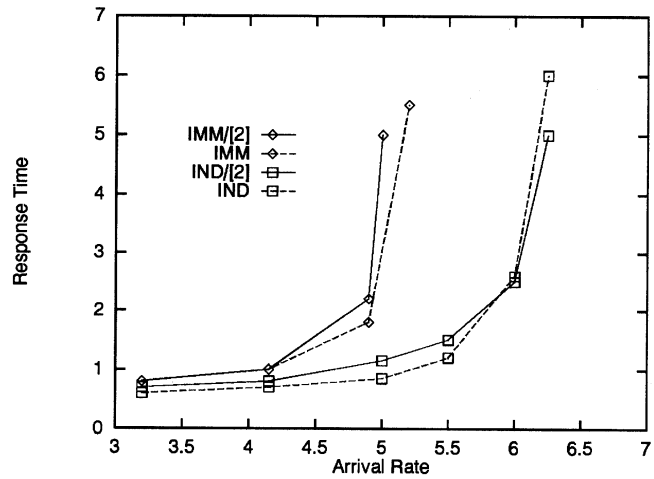


Fig. 3. Validation of active part of simulator

5.2 Validation of the simulator

Experiments were conducted to validate the simulator. The validation was accomplished in three steps.

1. We validated the *active* part of the simulator against the results in Carey et al. (1991). The results are illustrated in Fig. 3. We mapped the model in Carey et al. (1991) onto ours as closely as possible. We were not able to obtain the exact results because of the following differences in the two models. The buffer is explicitly modeled in Carey et al. (1991), whereas we model our buffer using a parameter h , which is the probability that a page is resident in the buffer. h is set to 0.9 for all the experiments in this validation. The explanation of the parameters and the experiments can be found in Carey et al. (1991). In Fig. 3, we observe that our results are within 10% of the original results.

2. We validated the *real-time* part of the simulator by trying to duplicate the results in Abbott and Garcia Molina (1992). Under the NT⁵ (*not tardy*) overload management policy, a transaction is aborted as soon as it becomes tardy. Under the AE (*all eligible*) policy, a transaction is run until it finishes. The results are illustrated in Fig. 4. Again our results were very similar to previously published results. The slight performance improvement obtained by our policies in the *not tardy* case can be explained by the fact that checking for tardiness is done more often in our model.

3. Finally, we validated our simulator against the results in Kao and Garcia Molina (1993), where policies for assigning deadlines to parallel and serial subtasks of complex distributed tasks in a real-time system have been evaluated. We do not present the result graph for this experiment to save space. There was an inherent difference in the two models because the one in Kao and Garcia Molina (1993) is for a distributed system, whereas ours is for a single site multiprocessor system. The system in Kao and Garcia Molina (1993) has multiple servers with a queue for each server, whereas ours has multiple servers with a single queue. We experimented with *UD* and *DIV-1* policies mentioned in Kao and

⁵ Please note this is not the same as class NT which denotes the class of non-triggering transactions

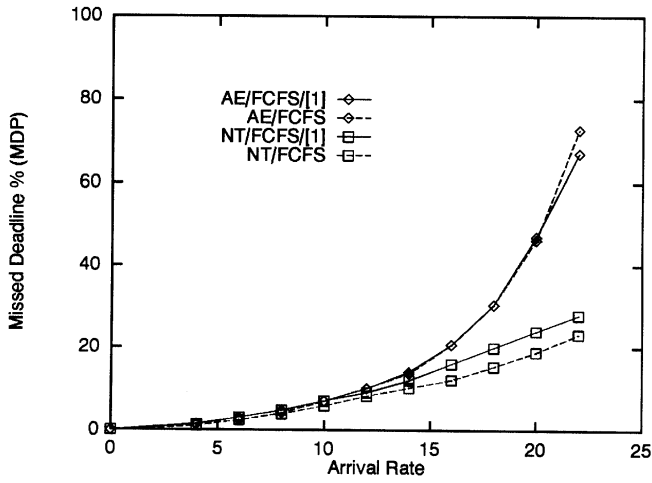


Fig. 4. Validation of real-time part of simulator

Garcia Molina (1993) and our missed deadline percentages were lower by no more than 5–10%.

In our validations we did not perform the complete set of experiments that are found in Carey et al. (1991), Abbott and Garcia Molina (1992) or Kao and Garcia Molina (1993), but just certain baseline experiments.

5.3 Baseline parameters

In this section we describe the workload model and the calculation of estimates required by DIV and SL. Let us assume a transaction T arrives at time t_0 and triggers a subtransaction at time t_1 . Let $U(i, j)$ denote a uniformly distributed integer valued random variable in the range $[i, j]$. The lengths, $L(T)$, of a transaction T and, $L(T^{sub})$, of a subtransaction T^{sub} (immediate or deferred) that might be triggered by T are given in number of method invocations as follows.

$$\begin{aligned} L(T) &= U(5, 7), T \in \text{class } \mathbf{T} \\ &= U(4, 6), T \in \text{class } \mathbf{NT} \\ L(T^{sub}) &= U(4, 6), T \in \text{class } \mathbf{T}, T^{sub} \text{ subtransaction of } T \end{aligned}$$

Each method takes one unit of time to execute. At the start of every method execution, a transaction T belonging to class \mathbf{T} triggers a subtransaction with probability p . The probability that the triggered subtransaction is of type immediate is q and the probability that it is of type deferred is $1 - q$. For any transaction T , we can compute the length of the transaction when it arrives ($L_{a(T)}(T) \equiv L(T)$) and hence, we know the remaining length of the transaction at any time t . Let $L_t(T)$ denote the remaining length of transaction T at time t . The estimates that we use in the priority assignment policies are calculated as follows:

$$\begin{aligned} X_t(T) &= L_t(T) \\ m_t^{imm}(T) &= L_t(T) * p * q \\ m_t^{def}(T) &= L_t(T) * p * (1 - q) \\ \bar{X}^{imm}(T) &= L(T^{sub}) \\ \bar{X}^{def}(T) &= L(T^{sub}) \end{aligned}$$

Table 1. System parameters

Parameter	Setting
N_{CPU}	6
Time taken to execute a method	1 unit
Number of object instances	2000
Overload management policy	not tardy
$frac_T$	0.15

It should be noted that all other estimates discussed in Sect. 4 can be calculated from the above estimates.

The deadline of a transaction T is set using the following formula:

$$d(T) = a(T) + (1 + \beta) * X_{at(T)}(T)$$

where β is a uniformly distributed random variable within a specified range. We consider three types of workloads: one where the class \mathbf{T} transactions trigger only immediate subtransactions ($q = 1$), one where they trigger only deferred transactions ($q = 0$), and one where they trigger both immediate and deferred subtransactions with equal probability ($q = 0.5$). We normalize the slacks for the three types of class \mathbf{T} transactions taking into account the fact that the deferred transactions can be executed in parallel, and the immediate transactions are executed in sequence. Therefore, transactions that trigger only immediate subtransactions get more slack than those that trigger deferred subtransactions.

We use a parameter *load* in our experiments which is very similar to the one in Kao and Garcia Molina (1993). In order to define *load* we specify the arrival rates and service rates of class \mathbf{T} and \mathbf{NT} transactions. The arrivals of class \mathbf{T} and class \mathbf{NT} transactions are generated according to Poisson processes with mean interarrival times of $1/\lambda_T$ and $1/\lambda_{NT}$ time units, respectively. The arrival rates are calculated using the following two equations, where all other quantities except the arrival rates are assumed to be known. In the first equation, we define the *load* to be the ratio of work generated to the total processing capacity of the system. Let $1/\mu_T$ and $1/\mu_{NT}$ denote the average total execution time of class \mathbf{T} and \mathbf{NT} transactions, respectively, and N_{CPU} the number of CPUs in the system. In the second equation, $frac_T$ is the fraction of *load* that is contributed by the class \mathbf{T} transactions.

$$\begin{aligned} load &= \frac{\frac{\lambda_T}{\mu_T} + \frac{\lambda_{NT}}{\mu_{NT}}}{N_{CPU}} \\ frac_T &= \frac{\frac{\lambda_T}{\mu_T}}{\frac{\lambda_T}{\mu_T} + \frac{\lambda_{NT}}{\mu_{NT}}} \end{aligned}$$

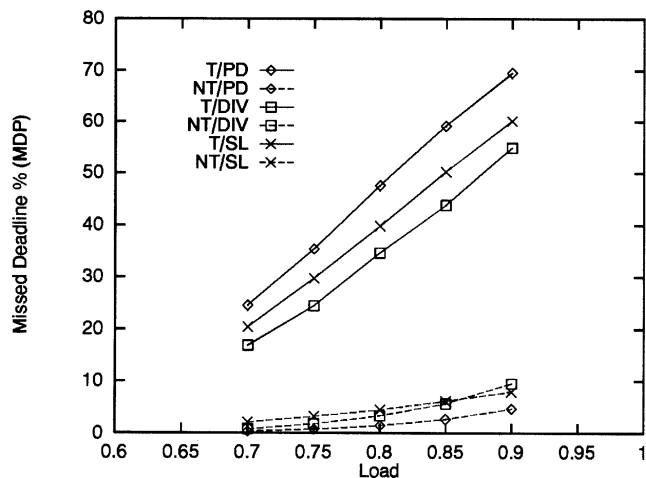
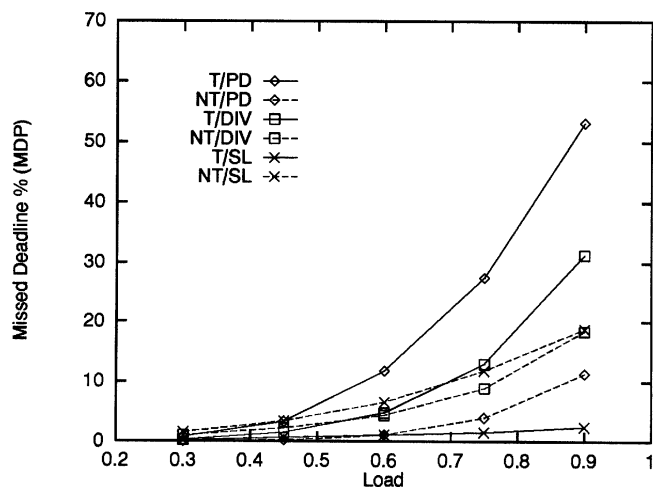
From the above equations, for given values of *load*, $frac_T$, $1/\mu_T$ and $1/\mu_{NT}$, we can compute $1/\lambda_T$ and $1/\lambda_{NT}$. Tables 1 and 2 show the system and transaction parameter settings, respectively, for our baseline experiments.

5.4 Real-time tasks

In the first set of experiments we deal with real-time active tasks executing in a multiprocessor environment. The purpose of these experiments is to isolate and study the effect

Table 2. Transaction parameters

Class	Parameter	Setting
class NT	Length in methods ($L(T)$)	$U(4, 6)$
	Slack parameter (β)	$U(0.5, 1.25)$
class T	Length of the parent in methods ($L(T)$)	$U(5, 7)$
	Probability of triggering by object event (p)	0.8
$q = 1.0$ (all imm.)	Length of the subtransaction ($L(T^{sub})$)	$U(4, 6)$
	Slack parameter (β)	$U(6.0, 6.5)$
$q = 0.0$ (all def.)	Length of the subtransaction ($L(T^{sub})$)	$U(4, 6)$
	Slack parameter (β)	$U(2.0, 2.5)$
$q = 0.5$ (both imm. and def.)	Length of the subtransaction ($L(T^{sub})$)	$U(4, 6)$
	Slack parameter (β)	$U(4.0, 4.5)$

**Fig. 5.** Only immediate/task model**Fig. 6.** Only deferred/task model

of scheduling on performance. We simulate a main memory database system where there are no data conflicts, i.e., every data access is a shared access. The experiments presented here are:

- Load versus MDP for a fixed slack distribution
- Average slack [$E(\beta)$] versus MDP for a fixed load
- Analysis of trade-offs between DIV and SL for a fixed slack distribution and load

In the following discussion we use slack to denote the average slack parameter that is used to calculate the deadline (the initial slack) of the transaction belonging to that class. Note that slack corresponds to $E(\beta)$. We use $estSlack^6$ to denote the remaining slack time that is estimated by the policies.

5.4.1 Load versus MDP

The first set of performance results are presented in Figs. 5–7, respectively. Figure 5 deals with the case where all the subtransactions are triggered in immediate mode ($q = 1$). We observe from the graph that, when compared to PD, both the DIV and SL decrease the MDP of class T transactions at higher loads by as much as 10–15%, at the cost of a small increase of around 4% in the MDP of class NT transactions. DIV reduces the MDP of class T by a greater amount than

SL, accompanied by a smaller increase in the MDP of class NT transactions. But, as the load increases, the difference between the MDPs of DIV and SL for class NT transactions reduces. When the load is 0.9, DIV performs better than SL for class NT transactions. The first evidence for our hypothesis, that extra information like the estimates of execution time and accounting for the dynamic work generated improves the performance of triggering transactions, is seen here. This performance improvement, however, comes at the cost of decreased performance of non-triggering transactions.

In Fig. 6, we present the results for the case where all the subtransactions triggered are executed in deferred mode ($q = 0$). Again, SL and DIV decrease the MDP of class T transactions over PD. In this case, SL provides substantially better performance to class T transactions than NT, whereas the reverse is true for PD and DIV. SL performs better than DIV at higher loads for class T. However, in the case of class NT transactions, DIV performs better than SL. SL reduces the MDP of class T by 30% at high loads, compared to DIV with a slight increase in the MDP of class NT. Essentially SL gives higher preference to class T transactions over class NT transactions than DIV. This difference in performance from the previous results in Fig. 5 can be explained by the fact that deferred subtransactions can be executed in parallel, on a multiprocessor system, whereas the immedi-

⁶ $estSlack$ at time t for a transaction T is $S_t(T)$

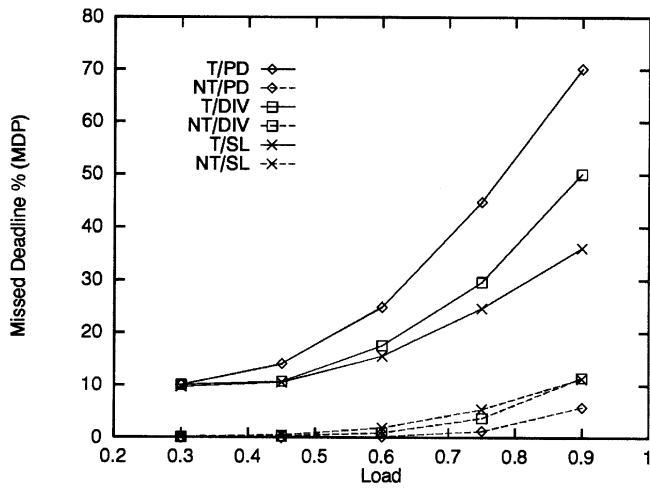


Fig. 7. Immediate-deferred/task model

ate subtransactions are executed sequentially⁷. Hence, SL gives a very high preference to class **T** transactions when deferred subtransactions are present. This explains the fact that SL keeps the MDP of class **T** transactions nearly constant, while the MDP of class **NT** transactions increases. It also explains the fact that SL gives lower MDP for class **T** than the class **NT**.

The results for the case where the triggered subtransactions execute either in deferred or immediate mode with equal probability ($q = 0.5$) are illustrated in Fig. 7. The relative performance behaviors of the three policies is similar to those in the case where all subtransactions are of the deferred type. However, the difference between the MDPs of class **T** transactions for the three policies is lower than the only deferred case because of the presence of immediate subtransactions. The qualitative performance of SL with respect to the triggering and non-triggering transactions is similar to that of the case where all the subtransactions are of the immediate type.

We observe from these results that, while the MDP of class **T** transactions is reduced by the SL and DIV policies, the MDP of the class **NT** transactions increases. This is desirable when class **T** transactions are more valuable to the application than class **NT** transactions. In order to examine the performance of the policies when transactions of both classes have the same value, we evaluated the combined MDP of all the transactions. Again, 85% of the workload comes from class **NT** transactions. The results are illustrated in Figs. 8–10. We observe that PD always performs best if we give equal value to both class **T** and class **NT** transactions. DIV is the next best and SL is the worst. PD works best here because DIV and SL are biased toward class **T** transactions which require more CPU time than class **NT** transactions. Although DIV and SL reduce the MDP of class **T** transactions, the overall MDP increases in this case, since the majority of transactions belong to class **NT** (85%). It

⁷ We disallow multiple immediate subtransactions being fired at the same time. The parent transaction is suspended when an immediate transaction is executed. In effect, the immediate subtransactions are executed *sequentially*, as opposed to deferred subtransactions that are executed in parallel at the end

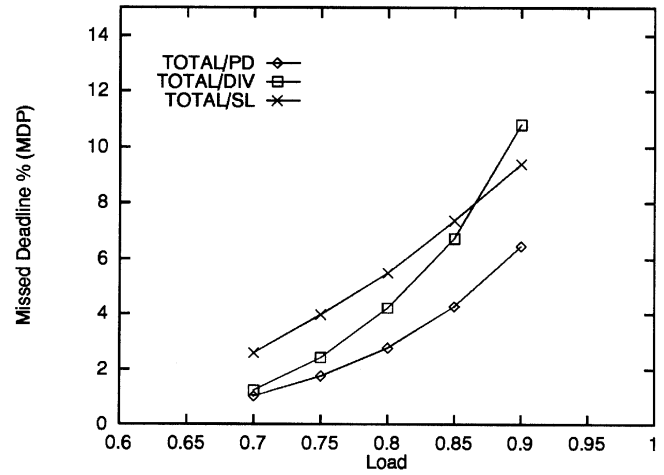


Fig. 8. Only immediate/task model/total

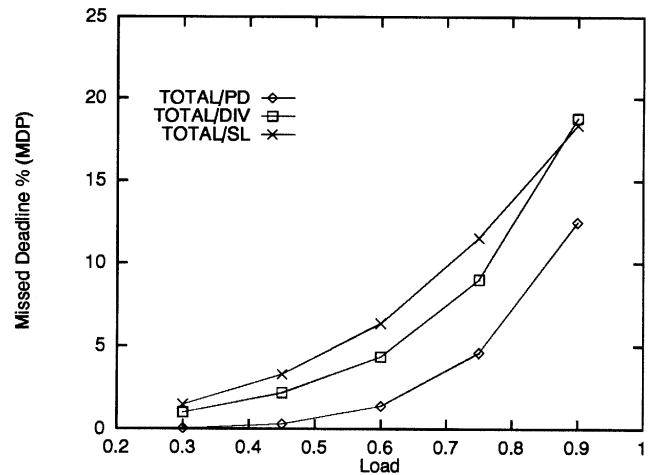


Fig. 9. Only deferred/task model/total

should be noted that the total MDP can be a biased performance measure (Pang et al. 1992). Total MDP will be low for policies that favor short transactions which is class **NT** transactions in our study. PD, by favoring class **NT** transactions, gives a lower total MDP than both SL and DIV which favor class **T** transactions.

5.4.2 Average slack versus MDP

It is understandable that accounting for the dynamic work generated could result in better performance for class **T** transactions. However, the explanation of the difference in the performances of DIV and SL is not clear from the above experiments. One of the main differences between DIV and SL is that, while DIV always favors class **T** transactions, SL, depending on the slack parameters, could favor either class. DIV changes (increases) the priority of class **T** transactions as it triggers other transactions, but keeps the priority of class **NT** unchanged. SL, on the other hand, modifies the priority of both the classes as they are executing. Hence, it is possible that depending on the slack parameters, the relative performance of DIV and SL could change. Experiments

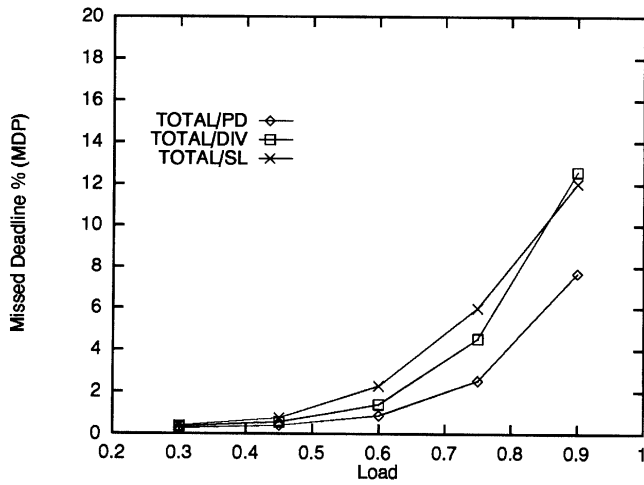


Fig. 10. Immediate-deferred/task model/total

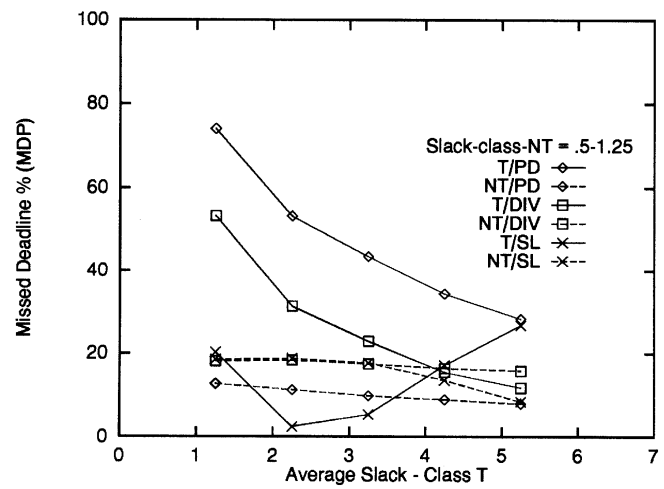


Fig. 12. Only deferred/task model/vary class T slack

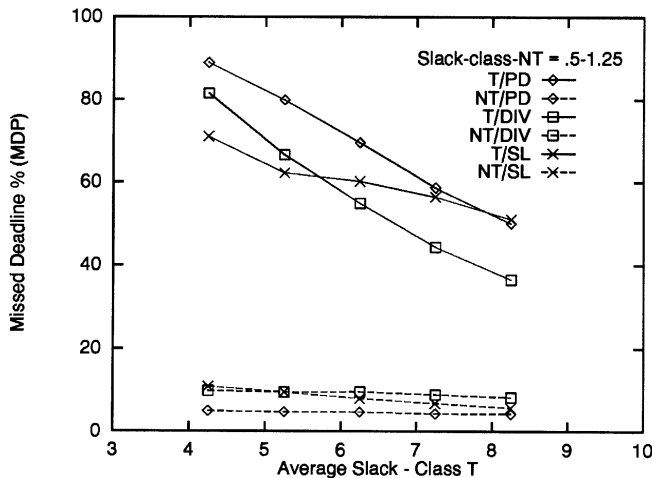


Fig. 11. Only immediate/task model/vary class T slack

were conducted where we varied the slack of class **T** transactions, holding the slack of class **NT** transactions constant for a load of 0.9. Figure 11 shows the result of this experiment where the class **T** transactions trigger only immediate subtransactions. The MDP of each transaction class is given as a function of average slack of class **T** transactions for all three policies. The average slack of class **T** is plotted on the x -axis and MDP is plotted on the y -axis. We observe that when the slack of class **T** is low, SL provides a lower MDP for class **T** transactions than DIV, but performs worse for class **NT**. At higher slack values of class **T**, DIV starts performing better than SL for both classes. Figure 5 is an instance of this latter behavior at different loads. Also, increasing the slack of class **T** means increase in the *estSlack* (remaining slack time) of class **T** transactions which results in relatively lower priorities for class **T** transactions compared to class **NT** transactions. This effect can be observed from the fact that SL starts performing better for class **NT** at high slack values of class **T** as seen in Fig. 11. We also conducted experiments where we held the slack of class **T** constant, while varying the slack of class **NT**. We do not present the graphs of these experiments, but it suffices to say that the same kind of phenomenon was observed.

Figure 12 shows the result of an experiment where class **T** transactions trigger only deferred subtransactions. The MDP of the different transaction classes is given as a function of average slack time for class **T** transactions. The average slack of class **T** transactions is plotted along the x -axis and the MDP is plotted along the y -axis. SL performs better than DIV for class **T** with low slack values and DIV starts performing better than SL when the slacks are really high. As mentioned earlier, very large slacks for class **T** implies large *estSlack* for class **T** transactions and hence lower priorities for class **T** relative to class **NT** transactions. Hence, SL starts performing better for class **NT** at high slack values of class **T** as seen in Fig. 12. We observe that under SL, MDP for class **T** reaches a minimum and then slowly starts to increase. This is the slack value of class **T** which performs the best for SL for a given slack value of class **NT**. Other slack values to the left are too tight and values to the right are too loose. The difference between the case where class **T** transactions trigger only immediate transactions and where they trigger only deferred transactions is that the MDP of class **T**, where DIV starts performing better than SL, is much higher in the first case. In other words, DIV starts performing better than SL for class **T** transactions that trigger only immediate subtransactions at *relatively* low slack values, when compared to case where class **T** transactions trigger only deferred subtransactions. In all the cases seen so far, we observe that PD does worse than DIV and SL for class **T** and does better than DIV and SL for class **NT** transactions.

5.4.3 Analysis of trade-offs between DIV and SL policies

We observed in the previous experiments that DIV and SL give preference to **T** class transactions over the class **NT** transactions. In this section we describe and evaluate two algorithms that provide the capability to trade off the MDPs for classes **T** and **NT**. The algorithms are constructed from DIV and SL in the following way. We introduce a parameter α which controls the priority assignment in DIV and SL policies. We call these parameterized policies ALPHA-DIV

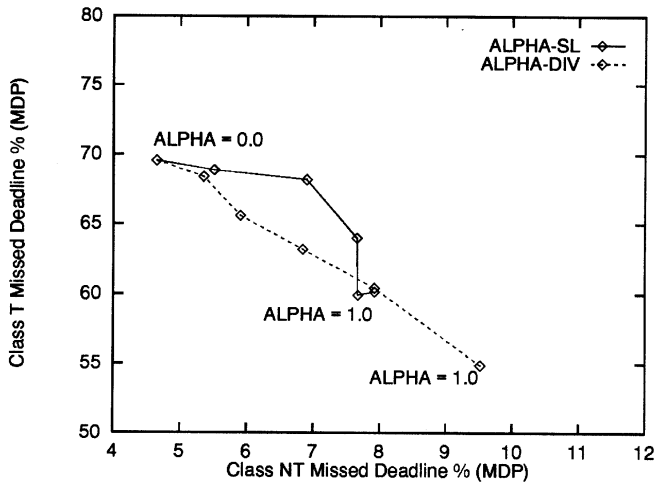


Fig. 13. Only immediate/task model/trade-offs

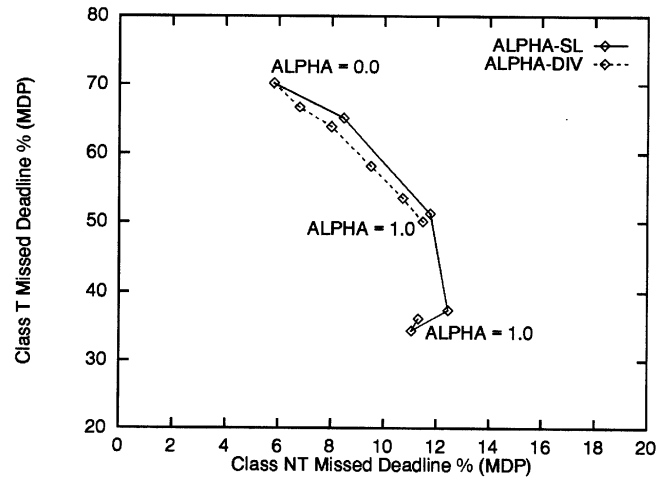


Fig. 15. Immediate-deferred/task model/trade-offs

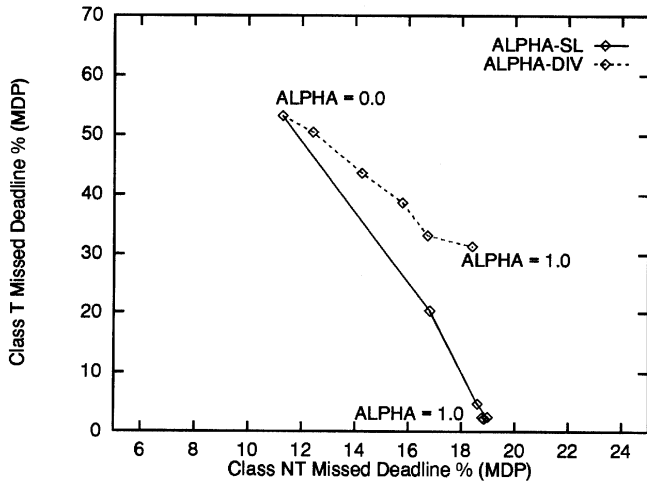


Fig. 14. Only deferred/task model/trade-offs

and ALPHA-SL. It should be noted that the experiment is conducted for a particular load and slack parameter values.

$$\text{ALPHA-DIV} : P_t(T) = \alpha * P_t^{DIV}(T) + (1 - \alpha) * P_t^{PD}(T)$$

where P_t^{DIV} is the priority assigned to the transaction T by the DIV protocol and P_t^{PD} is the deadline of the parent transaction.

$$\text{ALPHA-SL} : P_t(T) = \alpha * P_t^{SL}(T) + (1 - \alpha) * P_t^{PD}(T)$$

where P_t^{SL} is the priority assigned to the transaction by the SL protocol.

When α is zero both the policies reduce to PD. When α is one they reduce to DIV and SL, respectively. For this set of experiments the load was kept constant at 0.9. We studied the performance of the ALPHA-DIV and ALPHA-SL as α is varied from zero to one. The results are found in Figs. 13–15 where the MDPs of classes NT and T transactions are plotted. The points on the curves correspond to α values varied from 0 to 1 with an increment of 0.2. We also plotted some points with intermediate values of α to clearly distinguish between the different policies. We observe that ALPHA-DIV works better than ALPHA-SL throughout the

range for the immediate only case. It reduces the MDP of class T transactions by 15% with an increase of 4% in the MDP of class NT transactions, compared to the ALPHA-SL protocol which achieves a lesser reduction at the cost of a 3% increase. The same experiment was run with the load kept constant at 0.75. In this case, ALPHA-DIV reduced the MDP of class T by a larger amount than that of ALPHA-SL and with a smaller increase in the MDP of class T compared to ALPHA-SL. Hence, the ALPHA-DIV protocol performs better than the ALPHA-SL protocol for the case where all subtransactions are triggered in immediate mode.

However, when the workload consists only of deferred transactions, the ALPHA-SL protocol reduces the MDP of class T transactions to a very negligible value with a rise in the MDP of class NT transactions. ALPHA-DIV is not able to reduce the MDP of class T transactions so significantly. So if class T transactions have a higher value than class NT transactions, then ALPHA-SL is the protocol of choice in the deferred only case. ALPHA-SL consistently performs better than ALPHA-DIV for all α values and gives more flexibility to trade off performance between the two classes. The results for the case where 50% of the subtransactions are deferred and 50% immediate, are illustrated in Fig. 15. Here, like in this previous case, ALPHA-SL gives more flexibility to trade off performance between the two classes. In general, where most subtransactions are deferred, ALPHA-SL provides more flexibility to achieve a higher reduction in the MDP of class T transactions with an increase in the MDP of the class NT transactions than the ALPHA-DIV policy.

5.4.4 Summary of performance for the real-time task model

In summary, for a real-time task model, the dynamic priority assignment policies DIV and SL reduce the MDP of class T transactions, while increasing the MDP of class NT transactions. The choice between DIV and SL depends upon the slacks of the transactions and the load in the system. SL is a pure slack-based policy that performs well for a certain class of transactions depending on the slack of that class, its absolute value and its value relative to slack of the other class. DIV always favors class T transactions. We

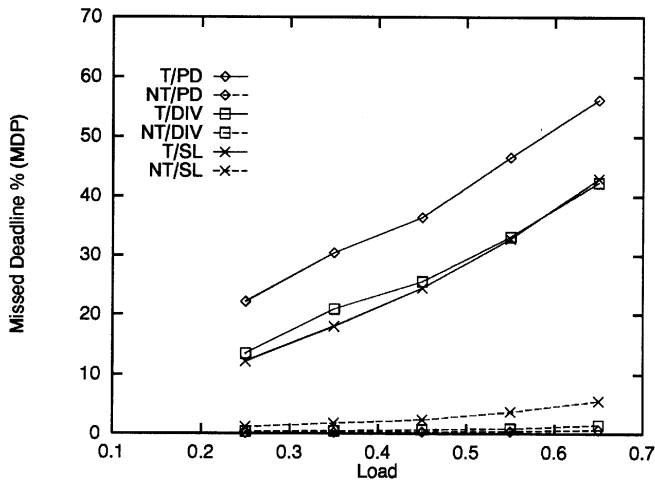


Fig. 16. Only immediate/main memory DB/2000 objects

also looked at two policies, ALPHA-DIV and ALPHA-SL, that enable us to trade off performance between different classes. Our hypothesis, that accounting for the work generated dynamically and using knowledge about transactions will benefit the triggering transactions, is substantiated by the results in this section.

5.5 Main memory database

We now consider the effects that data contention can have on the behavior of the policies. Since we are changing just one parameter, any difference in the behavior is due to data contention. We introduce data contention by requiring every data access to be an exclusive access. It should be noted that the way we calculate the estimates in this case is the same as in the real-time task case. The concurrency control algorithm we use is the high-priority (HP) algorithm (Abbott and Garcia Molina 1992), modified to deal with subtransactions. According to our concurrency control mechanism, all of the subtransactions triggered by a transaction (whether in deferred or immediate mode) are considered part of the transaction and share the locks. Similarly, two subtransactions of the same parent transaction share the locks. Hence, in our model a parent transaction and its subtransactions represent a set of cooperating transactions to complete a single task. All of the subtransactions and the parent transaction release the locks at their commit time which occurs after the parent transaction and the deferred subtransactions have finished. Deadlocks could occur between transactions in spite of using the HP protocol, because of dynamically changing priorities. Deadlocks are prevented by checking for deadlocks each time a transaction waits for another transaction and aborting the transaction that causes the deadlock. It has been shown in previous studies that the choice of the transaction to be aborted to resolve a deadlock does not have a significant impact on the performance (Huang et al. 1989).

The results are illustrated in Figs. 16 and 17 for workloads consisting of only immediate and only deferred subtransactions, respectively. In the case that all subtransactions are immediate, DIV and SL provide better performance to class T while providing worse performance to class NT as

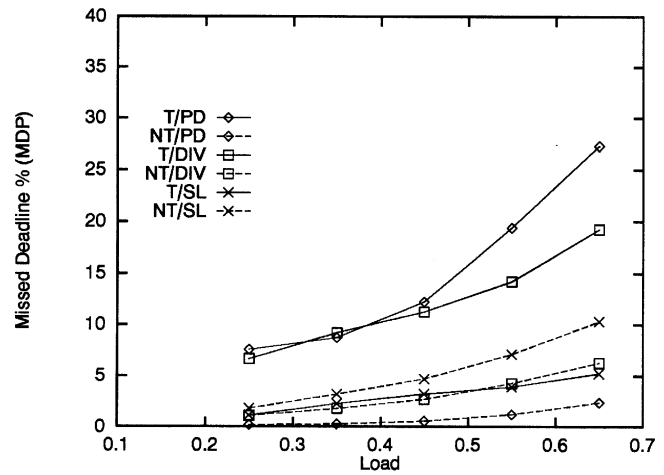


Fig. 17. Only deferred/main memory DB/2000 objects

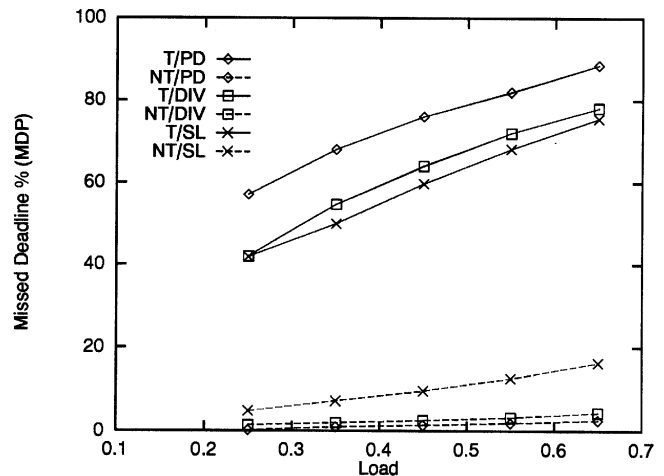


Fig. 18. Only immediate/main memory DB/500 objects

was observed in the case with no data contention. Figure 16 shows the performance for varying loads at a fixed slack distribution. The relative performance of DIV and SL change with the varying load for class T transactions, but for class NT transactions DIV performs better than SL.

The level of data contention could affect the performance of the priority assignment policies. In order to study the effect of data contention we changed the number of objects in the system. For a given load, reducing the number of objects will amount to a higher data contention. Figures 18 and 19 show the immediate and deferred cases of workload where the number of objects is 500. For the case where all the subtransactions are of type immediate, SL performs better than DIV for class T. This is not the case when the number of objects is 2000.

In the case where all the subtransactions are triggered in immediate mode, for the main memory database (Fig. 16) the MDP of class NT keeps increasing for SL when compared to DIV. This is not the case for the immediate only case of the real-time task model (Fig. 5). The difference between the two cases is that there is data contention in the main memory case. Moreover, when the level of data contention is higher (Fig. 18) the difference between the MDPs of class

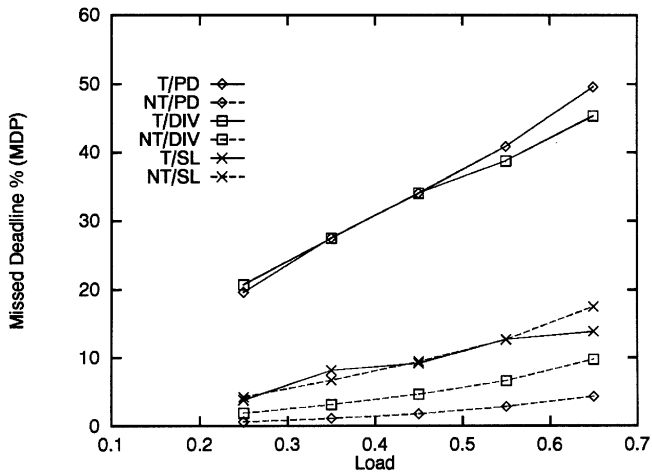


Fig. 19. Only deferred/main memory DB/500 objects

NT between DIV and SL increases. One could attribute this to the slacks being too tight. But that should affect DIV and PD too. In our model, SL estimates the number of steps but does not take into account the blocking delays due to data contention. This could lead SL to underestimate the time taken to execute the remaining steps (i.e., overestimate the *estSlack*) and miss the deadlines. But the *estSlacks* of all the transactions are being uniformly overestimated. If we assume that on an average all transactions undergo equal amount of blocking delays, then SL can overestimate *estSlacks* of transactions that have not been blocked until a certain point in time, and the transaction can miss its deadline.

We observed that the number of concurrency control aborts in the case of SL, was two to three times the number of aborts in PD or DIV. The number of aborts was highest both for class **T** and class **NT** transactions for SL. DIV had more aborts for class **NT** than PD and PD had more aborts for class **T**. This is explicable because DIV favors class **T** transactions over class **NT** and PD favors class **NT** over class **T**. For SL, the high number of aborts for class **NT** was definitely affecting its performance for this class. Under a policy like SL, the relative priority ordering between two transactions can keep changing throughout the life of these transactions depending on how many steps are executed by these two transactions, whereas in a deadline based priority scheme the relative ordering is fixed. When there is data contention, the number of steps executed by transactions that start around the same time can differ, which would affect the *estSlack* of these transactions. Transactions that do not experience contention could end up with more *estSlack* than transactions that have experienced some data contention.

As an example, let us consider two transactions T_1 and T_2 that arrive at times 1 and 3, respectively. Let their lengths be five units each and deadlines be 12 ($d(T_1)$) and 14 ($d(T_2)$), respectively. At time 5 let us say T_1 has completed four steps and T_2 has completed one step because it faced data contention. Therefore, $X_5(T_1)$ is 1 and $X_5(T_2)$ is 4. Now *estSlack* of T_1 at time 5 is $S_5(T_1)$, i.e., $12 - 5 - 1 = 6$ and the *estSlack* of T_2 at time 5 is $S_5(T_2)$, i.e., $14 - 5 - 4 = 5$. *estSlack* of T_2 at time 5 is less than *estSlack* of T_1 at time

5. In such a scenario, if T_2 needs a data item locked by T_1 , then T_2 will abort T_1 , and T_1 will be restarted with a very low slack (a high priority) and abort other transactions that are nearing their end. It is more likely that a transaction that came earlier holds the lock on a data item that is required by a transaction that came later than vice versa. We conjecture that this kind of a scenario occurs often under SL resulting in a large number of concurrency control aborts.

In these experiments, as expected, the absolute MDP values are higher than the MDP values of the previous set of experiments with no data contention; that is, more transactions are aborted due to deadline misses. Due to data conflicts, more transactions experience waits or get aborted. This shows the effect of data contention on the performance. One last observation is that PD does better than SL for the immediate only case at low loads with tight slacks for class **T** transactions (not shown in graphs). This is due to the high number of concurrency control aborts in the case of SL.

5.5.1 Summary of performance in main memory databases

The observations that DIV and SL perform better than PD for class **T** and PD performs better than DIV and SL for class **NT** also hold in the case of main memory databases. But there is a difference in the performance between these two cases mainly due to the high number of concurrency control aborts of SL. At low loads, high concurrency control aborts lead SL to behave worse than PD for class **T** transactions that trigger only immediate subtransactions. The very high number of concurrency control aborts of SL for class **NT** transactions causes SL to behave worse than DIV for class **NT**. This problem can be solved by considering the amount of work performed by a transaction before aborting it, and checking for *feasibility* (i.e., if the transaction has a chance to complete before the deadline). Also, accounting for concurrency control delays in estimating the slack might help.

5.6 Disk resident database

The third scenario considered is one where the objects are not always in main memory, in which case it is necessary to retrieve them from the disk. In our simulator, we do not explicitly manage the buffer. The database buffer is modeled using a parameter h , which is the probability that a page is resident in the buffer. Hence, with a probability $1-h$ the disk subsystem is accessed. In our calculation of estimates we take the i/o time into account. The estimates are calculated as follows:

$t^{i/o}$: average time taken to do an i/o

$$X_t(T) = L_t(T) * (1 + ((1 - h) * t^{i/o}))$$

$$m_t^{imm}(T) = L_t(T) * p * q$$

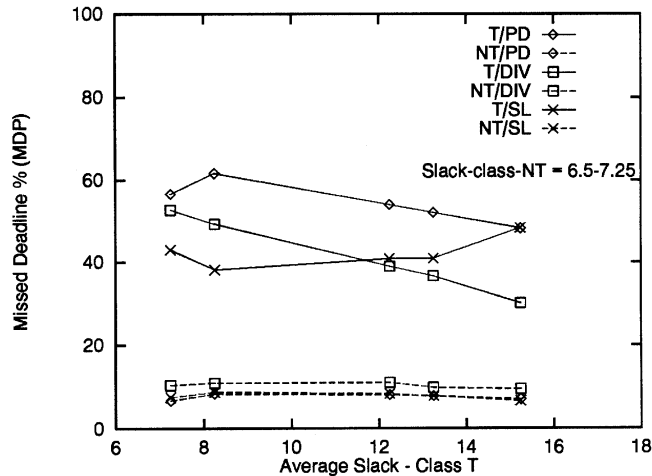
$$m_t^{def}(T) = L_t(T) * p * (1 - q)$$

$$\bar{X}^{imm}(T) = L(T^{sub}) * (1 + ((1 - h) * t^{i/o}))$$

$$\bar{X}^{def}(T) = L(T^{sub}) * (1 + ((1 - h) * t^{i/o}))$$

Table 3. Setting for disk-resident databases

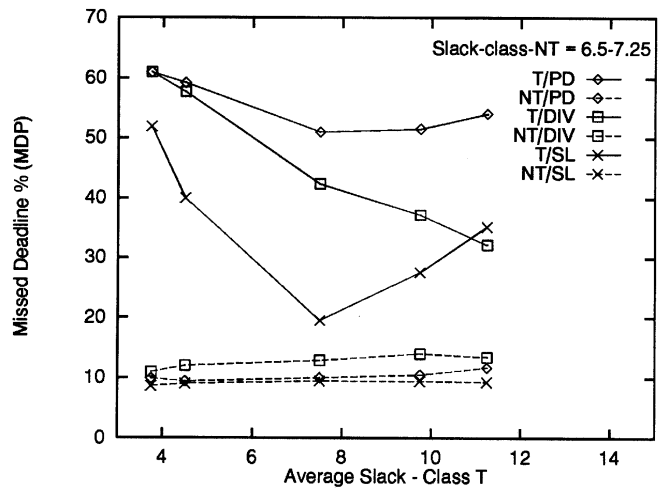
Parameter	Setting
N_{Disk} (number of disks)	2
$t^{i/o}$	30 units
h	0.9

**Fig. 20.** Only immediate/disk-resident DB/2000 objects

It should be noted that all other estimates (see Sect. 4.1) can be derived from the above estimates. Table 3 shows important parameter values that we used in our experiments.

We experimented by varying the slack of class **T** while maintaining the same slack distribution for **NT** for a particular load value (0.35). Figs. 20 and 21 show the results of this experiment for the cases where all the subtransactions are triggered in immediate mode and deferred mode, respectively. The MDPs of the different transaction classes are given as a function of average slack parameter for class **T** transactions. The results are similar to the real-time task case for class **T**, i.e., the MDP of class **T** is reduced by DIV and SL policies. When the slack values are lower, SL performs better than DIV and at higher slack values DIV performs better than SL. Also, as in the real-time task case there is a particular slack value of class **T** (relative to class **NT** slack) where SL performs the best. But for class **NT** SL performs at least as well as DIV. The reason is that SL accounts for *i/o* times in its calculation of slack for class **NT** transactions. The lack of estimate for blocking delays due to concurrency control could get subsumed by the *i/o* estimate since the latter is so much larger. On the other hand, DIV does not account for the *i/o* times in the case of class **NT** transactions and hence does not perform as well as SL. One can think of *i/o* as the dynamic work that is generated by the class **NT** transactions. Since SL takes into account this dynamic work that is being generated by class **NT**, it performs better than DIV and PD for this class.

One anomaly observed in the case where all the subtransactions are triggered in deferred mode (Fig. 21) is that, in spite of increasing the slack, the MDP of class **T** transactions for PD does not decrease at high slack values (9.75, 11.25). This is because very high slacks result in long deadlines (low priorities) that were causing a lot of concurrency

**Fig. 21.** Only deferred/disk-resident DB/2000 objects

control aborts and restarts resulting in a lot of wasted resources. This problem can be solved by considering the amount of work performed by a transaction before aborting it, checking for feasibility (i.e., if the transaction has a chance to complete before the deadline) when restarted and giving higher priorities to transactions that are restarted.

We also experimented by varying the load of the system while maintaining fixed slack distributions. Here SL performs the best for both class **T** and class **NT** for both the cases where only immediate subtransactions are triggered and where only deferred subtransactions are triggered. In the case where all the subtransactions are triggered in deferred mode, SL reduces the MDP of both the classes almost by half in overloaded situations.

5.6.1 Summary of performance in disk-resident database

In the disk-resident database case, DIV and SL perform better than PD for class **T**. DIV performs the worst for class **NT**. SL performs as good or better than PD for class **NT**. At high loads, SL performs the best for both classes **T** and **NT**. Since SL accounts for the *i/o* in its estimates for class **NT**, it tends to perform better than DIV consistently for this class. Our claim, that accounting for dynamic work generated and using the knowledge about the *i/o* requirements and triggering characteristics of transactions improves the performance, is substantiated by the above results.

6 Conclusion

We have studied the problem of assigning priorities to triggered transactions and reassigning priorities of triggering transactions in a firm real-time active database. We introduced a simple baseline policy PD and two other policies, namely DIV and SL, which assign priorities taking into consideration the active work generated by a transaction. The policies use different amounts of information about the transactions. PD uses the deadline to schedule the transactions, DIV uses estimates of execution times of the parent transac-

tion and subtransactions that it has triggered, and SL uses estimates of execution times of the parent transaction, triggered subtransactions and estimates about the subtransactions that might be triggered in the future.

The main conclusions are:

1. There are three *key* parameters – load, slack parameter of class **T** and slack parameter of class **NT** of the workload – which determine the relative performances of the three algorithms.
 - For a real-time task case in all of the load/slack space that we have experimented with, DIV and SL perform better than PD for class **T** and PD performs best for class **NT** transactions. It might be possible for SL to perform better than PD for class **NT** for extremely low slack values. But it would be almost impossible for DIV to perform better than PD for class **NT** because both schedule class **NT**, depending on the deadline, and DIV favors class **T**. Hence, DIV can do only as well as PD for class **NT**. Similarly, it is difficult to conceive cases where PD could do better than DIV for class **T**.
 - In a main memory database setting with data contention, concurrency control aborts could change the performance of the algorithms. For instance, at very low load values for the immediate only case, SL could do worse than PD for class **T** because of the high number of concurrency control aborts in the case of SL. Similarly, for the deferred only case, DIV could do worse than PD for class **T**. SL performs consistently worse than DIV for class **NT** at low slack values because of the high number of concurrency control aborts and restarts of transactions with very high priorities.
 - In a disk-resident database setting, DIV and SL perform better than PD for class **T**. SL's performance is comparable or better than that of PD for class **NT**. This improvement in the performance of SL for class **NT** is because of the fact that it accounts for the i/o time in its calculation of estimates.
2. For class **T** the sources of unpredictability are triggering of other subtransactions, concurrency control delays and i/o, and for class **NT** they are concurrency control delays and i/o. In the real-time task model, SL and DIV account for active work and hence perform better than PD for class **T**. For class **NT** there is no unpredictability. But, since SL and PD perform better for class **T**, they penalize class **NT**. DIV behaves qualitatively the same for a main memory database setting as it does for real-time task setting. But, due to the high number of concurrency control aborts, SL performs relatively poorly for class **NT** compared to the real-time task model. SL, in the case of a disk-resident database, accounts for all sources of unpredictability in its estimate of *estSlack* for both class **T** (triggering and i/o), and class **NT** (i/o) except blocking delay, which might not be significant compared to i/o. DIV accounts for the sources of unpredictability as SL does for class **T**, but does not account for the i/o for class **NT**. Hence, SL performs better than DIV for class **NT**.
 - Using the extra information such as estimates of execution times enhances the performance of class **T**

transactions, as can be seen by DIV and SL performing better than PD in all three settings. This benefit comes at the cost of increased MDP for class **NT**. Using the extra information about i/o makes SL perform better than DIV and PD for class **NT** in the disk-resident database setting.

Some of the extensions we want to address are:

- Experiment with variations of DIV and SL policies that use more information about a transaction; for instance, the exact number of subtransactions that it is going to trigger, and the type of subtransactions it is going to trigger, i.e., immediate or deferred, to see what advantage is obtained by exploiting this information.
- Study variations of DIV and SL that will assign priorities at every scheduling instance, instead of at every object event.
- Study the effect of errors in the knowledge about transactions like $L_t(T)$, $X_t(T)$ and p on the performance of the DIV and SL policies.
- Evaluate algorithms that change the priority of triggered subtransactions during their execution. This will be particularly useful to study systems where the triggered subtransactions can trigger further subtransactions.
- Consider the concurrency control blocking delays in the estimate of slacks.
- Consider the amount of work a transaction has completed before aborting due to a concurrency control conflict.
- Consider the execution model for immediate subtransactions where the parent transaction is not suspended during the execution of the immediate subtransaction.
- Consider different lock-sharing semantics between the triggering and triggered transactions and between triggered transactions.

Acknowledgement. This work was supported in part by NSF under grants IRI-9114197 and IRI-9208920.

References

- Abbott RK, Garcia-Molina H (1992) Scheduling real-time transactions: a performance evaluation. *ACM Trans Database Syst* 17:513–560
- Anon (1992) Special issue on active databases. *Bulletin, Technical Committee on Data Engineering*, vol. 15 (1–4), December
- Carey MJ, Jauhari R, Livny M (1991) On transaction boundaries in active databases: a performance perspective. *IEEE Trans Knowl Data Eng* 3:1–37
- Dayal U, et al (1988) The HIPAC project: combining active databases and timing constraints. *SIGMOD Rec* 17
- Dayal U, Hsu M, Ladin R (1990) Organizing long-running activities with triggers and transactions. *ACM SIGMOD*, Atlantic City, New Jersey
- Huang J, Stankovic J, Towsley D, Ramamritham K (1989) Experimental evaluation of real-time transaction processing. *Proceedings of the Real-Time Systems Symposium*, December, pp 144–153
- Huang J, Stankovic JA, Ramamritham K, Towsley D (1991) Experimental evaluation of real-time optimistic concurrency control schemes. *Proceedings of the 17th Conference on Very Large Databases*, September, Barcelona, Spain
- Huang J, Stankovic JA, Ramamritham K, Towsley D (1991) On using priority inheritance in real-time databases. *Proceedings of the Real-Time Systems Symposium*, December, San Antonio, Texas
- Kao B, Garcia-Molina H (1993) Subtask deadline assignment for complex distributed soft real-time tasks. *Technical Report STAN-CS-93-1491*, Stanford University

- Klein MH, Ralya T, Pollak B, Obenza R, Harbour MG (1993) A practitioners handbook for real-time systems – Guide to rate monotonic analysis for real-time systems. Kluwer Academic, Dordrecht
- Lawler EL (1983) Recent results in theory of machine scheduling. In: Bachem A et al (eds) Mathematical programming: the state of the art. Springer, Berlin Heidelberg New York
- Livny M (1990) DeNet users guide (vers 1.5), Department of Computer Science, University of Wisconsin, Madison
- McCarthy D, Dayal U (1989) The architecture of an active data base management system. ACM SIGMOD, Portland, Oregon
- Pang H, Livny M, Carey MJ (1992) Transaction scheduling in multiclass real-time database system. Proceedings of the IEEE Real-time Systems Symposium, Phoenix, Arizona
- Purimetla B, Sivasankaran RM, Stankovic J (1993) A study of distributed real-time active database applications. IEEE Workshop on Parallel and Distributed Real-time Systems, April, Newport Beach, California
- Ramamritham K (1993) Real-time databases. *Int J Distrib Parallel Databases* 1:199–226
- Sivasankaran RM, Purimetla B, Stankovic J, Ramamritham K (1993), Network services databases – A distributed active real-time database (DARTDB) application. IEEE Workshop on Real-Time Applications, May, New York, New York
- Son SH, Park S (1994) Scheduling transactions for distributed time-critical applications. In: Casavant TL, Singhal M (eds) Readings in Distributed Computing Systems, IEEE Computer Society Press, New York
- Xu J, Parnas DL (1990) Scheduling processes with release times, deadlines, precedence and exclusion relations. *IEEE Trans Software Eng* 16:360–369