

Washington University in St. Louis

## Washington University Open Scholarship

---

All Computer Science and Engineering  
Research

Computer Science and Engineering

---

Report Number: WUCS-96-04

1996-01-01

### Bringing Real-time Scheduling Theory and Practice Closer for Multimedia Computing

R. Gopalakrishnan and Guru M. Parulkar

This paper seeks to bridge the gap between theory and practice of real-time scheduling in the domain of multimedia computer systems. We show that scheduling algorithms that are good in theory, often have practical limitations. However when these algorithms are modified based on practical considerations, existing theoretical results cannot be used as they are. In this paper we motivate the need for new scheduling schemes for multimedia protocol processing, and demonstrate their real-time performance in our prototype implementation. We then explain the observed results by analysis and measurement. More specifically, we show that using strict preemption can introduce overheads... **Read complete abstract on page 2.**

Follow this and additional works at: [https://openscholarship.wustl.edu/cse\\_research](https://openscholarship.wustl.edu/cse_research)



Part of the [Computer Engineering Commons](#), and the [Computer Sciences Commons](#)

---

#### Recommended Citation

Gopalakrishnan, R. and Parulkar, Guru M., "Bringing Real-time Scheduling Theory and Practice Closer for Multimedia Computing" Report Number: WUCS-96-04 (1996). *All Computer Science and Engineering Research*.

[https://openscholarship.wustl.edu/cse\\_research/396](https://openscholarship.wustl.edu/cse_research/396)

Department of Computer Science & Engineering - Washington University in St. Louis  
Campus Box 1045 - St. Louis, MO - 63130 - ph: (314) 935-6160.

## Bringing Real-time Scheduling Theory and Practice Closer for Multimedia Computing

R. Gopalakrishnan and Guru M. Parulkar

### Complete Abstract:

This paper seeks to bridge the gap between theory and practice of real-time scheduling in the domain of multimedia computer systems. We show that scheduling algorithms that are good in theory, often have practical limitations. However when these algorithms are modified based on practical considerations, existing theoretical results cannot be used as they are. In this paper we motivate the need for new scheduling schemes for multimedia protocol processing, and demonstrate their real-time performance in our prototype implementation. We then explain the observed results by analysis and measurement. More specifically, we show that using strict preemption can introduce overheads in protocol processing such as more context switching and extra system calls. We present our scheduling scheme called rate-monotonic with delayed preemption (RMDP) and show how it reduces both these overheads. We then develop the analytical framework to analyze RMDP and other scheduling schemes that lie in the region between strict (immediate) preemption and no preemption. Byproducts of our analysis include simpler schedulability tests for non-preemptive scheduling, and a variant of rate-monotonic scheduling that has fewer preemptions. Finally, we measure the overhead due to context switching on Pentium and Sparc-1 machines and its impact on real-time performance. We show that when scheduling clock interrupts occur every 1 millisecond, RMDP can lessen the overhead of context switching leading to an increase in utilization of as much as 8%.

1

**Bringing Real-time Scheduling Theory and  
Practice Closer for Multimedia Computing**

**R. Gopalakrishnan**

**Guru M. Parulkar**

**WUCS-96-04**

**March 4, 1996**

**Department of Computer Science  
Campus Box 1045  
Washington University  
One Brookings Drive  
St.Louis, MO 63130-4899**



# Bringing Real-time Scheduling Theory and Practice Closer for Multimedia Computing

R. Gopalakrishnan  
*gopal@dworin.wustl.edu*

Gurudatta M. Parulkar  
*guru@flora.wustl.edu*

Department of Computer Science  
Washington University in St. Louis

## Abstract

This paper seeks to bridge the gap between theory and practice of real-time scheduling in the domain of multimedia computer systems. We show that scheduling algorithms that are good in theory, often have practical limitations. However when these algorithms are modified based on practical considerations, existing theoretical results cannot be used as they are. In this paper we motivate the need for new scheduling schemes for multimedia protocol processing, and demonstrate their real-time performance in our prototype implementation. We then explain the observed results by analysis and measurement. More specifically, we show that using strict preemption can introduce overheads in protocol processing such as more context switching and extra system calls. We present our scheduling scheme called rate-monotonic with delayed preemption (RMDP) and show how it reduces both these overheads. We then develop the analytical framework to analyze RMDP and other scheduling schemes that lie in the region between strict (immediate) preemption and no preemption. Byproducts of our analysis include simpler schedulability tests for non-preemptive scheduling, and a variant of rate-monotonic scheduling that has fewer preemptions. Finally, we measure the overhead due to context switching on Pentium and Sparc-1 machines and its impact on real-time performance. We show that when scheduling clock interrupts occur every 1 millisecond, RMDP can lessen the overhead of context switching leading to an increase in utilization of as much as 8%.

## 1 Introduction

Multimedia applications represent an important class of network applications. These applications are enabled by the availability of high speed networks, and affordable high performance personal computers and workstations. We identify two important requirements of multimedia applications:

- Continuous media communication and processing must occur in real-time. In addition, both the network and the endsystems must provide quality-of-service (QoS) guarantees to shield such applications from unpredictable network and system loads.
- Multimedia applications are compute and data intensive. So they require high performance in terms of both transmission speeds and processing power. Efficient use of network and compute resources is therefore important.

This paper is concerned with operating system (OS) support for multimedia applications. An essential feature required of an OS to support multimedia is real-time scheduling. However, providing guaranteed real-time performance to applications without compromising their efficient operation in a practical implementation is an interesting and elusive problem. This problem we believe points to a gap between theory and practice of real-time scheduling in the following sense—scheduling algorithms that have a good theoretical basis are not the most suitable choice from a practical standpoint, and scheduling schemes that are dictated by practical considerations do not have theoretical results to guarantee their real-time operation.

The purpose of this paper is twofold. We first illustrate the above situation, that is we describe a real-time scheduling scheme that is highly efficient but for which none of the existing theoretical results can be used as they are. Second we present new theoretical results that can be used for our scheduling scheme as well as for a variety of other efficient scheduling algorithms. More specifically, we show that traditional rate-monotonic (RM)

and earliest-deadline-first (EDF) scheduling lead to inefficiency due to their strict preemptive nature, and due to excessive involuntary context switching. These two problems are especially true in the case of protocol processing for multimedia applications. Strict preemption requires that distinct protocol tasks that share data (such as buffer lists and packet queues) lock them before access. In a real-time scheduling environment lock operations must be mediated by the OS to prevent unbounded priority inversion [13]. Thus system calls for locking increase per packet processing costs, counteracting efforts that have gone into speeding up packet processing. Likewise, involuntary context switches due to strict preemption causes CPU cycles to be wasted and reduces useful utilization.

We present our modified RM scheme called RMDP (RM with delayed preemption) as an example that solves both the problems mentioned above. The scheme is specially suitable for protocol processing as well as media processing for multimedia applications. However a survey of existing theoretical results show that they either deal with the case when preemption is immediate *i.e.*, the instant a higher priority task arrives (fully preemptive), or when there is no preemption (non-preemptive). The main contribution of this paper is a general technique to analyze real-time scheduling schemes in the region that lies between the fully preemptive and the non-preemptive cases. A new feature of our theoretical framework is the notion of a *preemption policy*. Our proof method, that we call *idealized scheduler simulation*, allows us to evaluate the impact of a given preemption policy on well known scheduling algorithms such as RM and EDF.

Our study of preemption policies have led to other interesting results. One result consists of schedulability tests for non-preemptive scheduling by treating it as a special case of our general preemption policies. Another interesting result is a variant of RM scheduling (that we call *mixed* scheduling) that uses an EDF preemption policy and is able to achieve almost as few preemptions as EDF scheduling. All the results above are derived assuming that context switching costs are zero. Since our main concern is efficiency issues in protocol processing, it is essential to quantify the cost of context switching that arises due to real-time scheduling. We therefore instrumented our implementation to measure and quantify the effect of context switching on performance. Our results show that in most cases, useful CPU utilization can be increased by using our efficient scheduling schemes.

The outline of the paper is as follows. Section 2 describes relevant aspects of our target environment that includes the real-time upcall (RTU) facility, and its scheduling mechanism that we call RMDP (rate-monotonic with delayed preemption). Section 3 presents existing results and shows why we cannot use them as they are to analyze RMDP and similar schemes. Section 4 precisely defines what aspects of real-time scheduling we will examine in this paper, and gives an overview of our experimental and theoretical results. The ISS method is presented in Section 5 to analyze different preemption policies. Proofs relating to the ISS method are in Appendix A. The next three sections apply the ISS methodology. Section 6 accounts for the effect of preemption policies on the schedulability tests for RM and EDF. Sections 7 and 8 discuss extensions of these results to non-preemptive and mixed scheduling respectively. Section 9 describes our experiments with the RTU facility on the Pentium and Sparc platforms to quantify context switching costs. We finally present our conclusions and future work.

## 2 Target Environment

Our work is targeted towards protocol processing support for networked multimedia applications. It is therefore important to understand some important details of this environment as shown in Figure 1. Processes 1, 2, and 3 are sample multimedia applications. Media processing and protocol processing for each connection are modeled as periodic activity with certain compute and communication resources allotted each period. The CPU and the network adaptor are the two resources shown. In our context, periodic activity in each process is implemented as a real-time upcall (RTU). A downward pointing RTU takes user data, processes it to form packets, and enqueues it at the output queues. An upward pointing RTU takes packets from its input queue, processes them, and possibly causes the next outgoing packet to be processed. Each RTU has a period and a computation time. It is also associated with a handler routine in the user program that gets upcalled by the kernel when the RTU is run. The OS performs a schedulability test to determine if all RTUs can complete before their deadlines *i.e.*, before their next period. The kernel makes an RTU runnable at the start of its period by placing it in the run queue. From the run queue the RTU scheduler picks an RTU to run based on its priority. The priority is determined by the scheduling policy used. The RTU scheduler uses RM priorities which means that smaller the period of an RTU, higher is its priority. The key feature of the scheduling scheme is that it takes advantage of the iterative nature of protocol processing to preempt a running RTU only at the end of an iteration. An iteration in this case corresponds to processing one (or more) protocol data units (PDUs). To implement this preemption scheme, each RTU and the kernel communicate via a data structure in shared memory. The kernel sets a flag in this shared

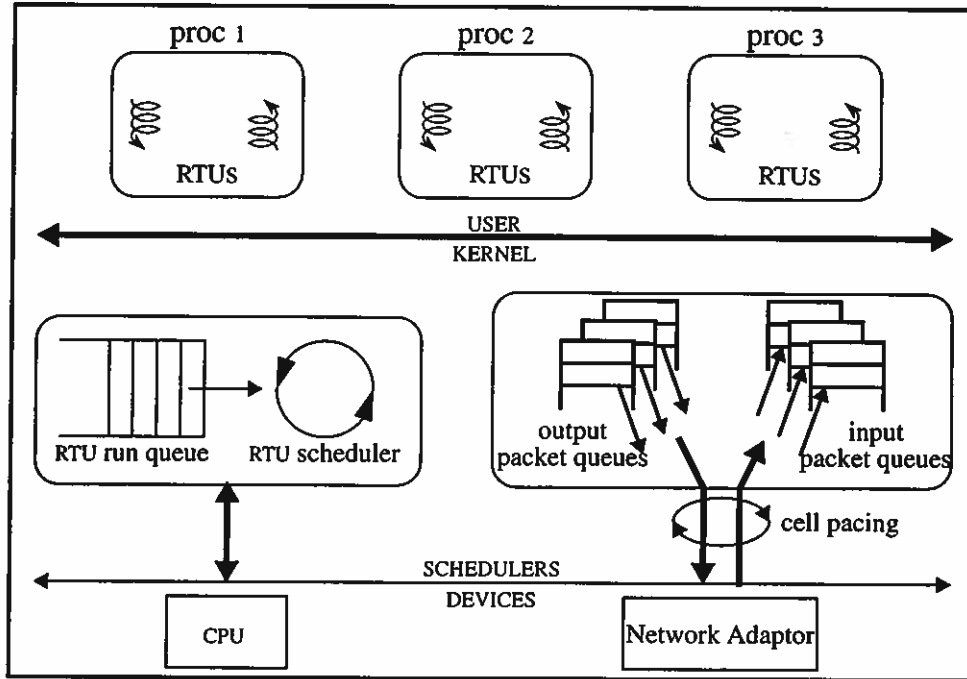


Figure 1: System Model

structure to inform a running RTU that a higher priority RTU has become ready to run. The handler checks this flag after each iteration and if it is set, it yields the CPU by returning from the handler. If the yielding RTU has not completed its work for the current period, the kernel puts it back in the run queue. The RTU is invoked later when it regains the highest priority status in the run queue. Because a running RTU can delay its preemption by one iteration, we call this scheme RM with delayed preemption (RMDP). Since RMDP does not preempt a handler asynchronously, RTU handlers do not need to lock shared data that is accessed only inside an iteration. Also we show by experiment that RMDP reduces the number of involuntary context switches thus increasing useful CPU utilization. Full details of the RTU implementation including security issues, system call interface, other efficiency benefits, and protocol performance measurements can be found in [2, 4]. Other research groups have used periodic real-time threads [15, 9] to implement protocol processing. However these solutions do not cater to the special needs of protocol processing, leading to overheads caused by locking and context switching.

From an efficiency standpoint, two aspects are important in an actual implementation of real-time threads (or in our case RTUs).

- *Involuntary context switches* : An involuntary context switch occurs when a thread with priority higher than the running thread is inserted in the run queue. In UNIX this happens in the clock interrupt that occurs periodically. In conventional time shared UNIX involuntary switches occur only at 100 millisecond intervals, and the clock ticks 100 times per second. To support multimedia (such as audio), millisecond granularity is required, and so the RTU facility uses a clock interrupt rate of 1000 times per second. In an involuntary switch the state of the process containing the running thread must be saved, and restored back when it is resumed. In a voluntary context switch that occurs when the running thread yields after completion, save and restore operations are not required. Since involuntary switches can occur only during the clock interrupt, a higher clock interrupt frequency should increase the number of such switches as we show later. For sake of efficiency it is all the more important for the RTU facility to reduce the number of involuntary context switches.
- *System calls* : A system call is made when a process performs an operation that requires OS intervention. A system call involves crossing the user–kernel boundary. For our Pentium machine running NetBSD, this overhead is about 770 processor cycles (or 280 instructions) [1]. As mentioned before, a lock operation in

	Utilization Bound	
	Rate Monotonic	Earliest-deadline-first
Fully Preemptive	$n(2^{1/n} - 1)$	1.0
No Preemption	N/A	pseudopolynomial time test
With Blocking	$n(2^{1/n} - 1) - \max(B_i/T_i)$	$1 - \max(B_i/T_i)$

Figure 2: Related Results

the real-time scheduling context entails a system call. Since PDU processing itself takes only a few hundred instructions, getting and releasing a lock per PDU increases packet processing time manyfold.

Thus RMDP increases efficiency in at least two ways: it eliminates per packet system calls for locking, and reduces involuntary context switching. To provide guarantees we also need to have a schedulability test that takes into account the preemption delay introduced by RMDP. In the next section we survey some of the scheduling algorithms that have been studied in previous work and show why they cannot be used as they are to analyze RMDP.

### 3 Related Work and their Limitations

From now on we refer to periodic activity (an RTU or real-time thread) as a *task*. A task has a period  $T$  and a computation requirement  $C$ . The utilization of a task is  $C/T$ , and the total utilization  $U$  of a set of tasks is the sum of their individual utilizations. We start with the well known RM and EDF scheduling algorithms introduced in [11]. RM gives higher priorities to tasks with smaller periods, and EDF gives higher priorities to tasks with earlier deadlines. The analysis in [11] assumed a uniprocessor system where tasks are independent, and deadlines are equal to the periods. It also assumes that preemption is immediate *i.e.*, a running task is preempted as soon as a higher priority task becomes runnable and the context switch time is zero. We refer to such a scheduler as an *idealized scheduler*.

Table 2 summarizes the relevant schedulability results for RM and EDF scheduling. The first row corresponds to the fully preemptive case and is the most widely used. Each table entry gives the maximum task utilization that can be supported by the scheduling scheme to guarantee schedulability (*i.e.*, all tasks meet their deadlines). As can be seen, EDF can support 100% CPU utilization. RM on the other hand limits utilization to  $n(2^{1/n} - 1)$ , which is about 69% for reasonable values of  $n$  (the number of tasks). The second row corresponds to the case when there is no preemption. There are no simple schedulability tests for this case as there are in the fully preemptive case. We discuss non-preemptive EDF scheduling in Section 7. The third row presents schedulability tests for RM and EDF that take into account the effect of *blocking*. Blocking occurs when a task is prevented from running by a lower priority task. Blocking that occurs when a lower priority task is in a critical section of code is considered in [13]. It can be seen that blocking reduces the utilization bound compared to the idealized case. The reduction is equal to the maximum value of  $B_i/T_i$  over all tasks, where  $B_i$  is the time for which task  $i$  can be blocked, and  $T_i$  is its period. Blocking due to other system events such as periodic timer interrupts is considered in [7].

#### 3.1 Inadequacy of Existing Results

None of the results in Table 2 can be directly used to model delayed preemption. The fully preemptive test does not hold because it assumes tasks cannot be blocked. Using the non-preemptive test is a conservative approach because it may reject task sets that are schedulable with delayed preemption. The blocking case at first glance seems to capture the effect of delayed preemption. However its derivation assumes that during the duration that a task is blocked, no task makes any progress. In other words, the blocking time is regarded as a system overhead (somewhat like a high priority system activity). Thus it underestimates the utilization that can actually be supported. Theoretically the amount by which this test can underestimate utilization in the worst case can be upto 49% although in practice it will be lower. An example of the worst case is 2 tasks with periods 100 and 99, and with the same computation time of 49. If the lower priority task blocks the higher priority task, then the loss due to blocking predicted by the test is around 49% whereas the two tasks can indeed be scheduled using non-preemptive EDF. This shows that it is important to distinguish between blocking caused due to overhead and blocking due to other tasks. Borrowing terminology in [8] we characterize existing results as analyses of *mechanism blocking i.e.*, concerned with overheads caused due to system activities. Our use of



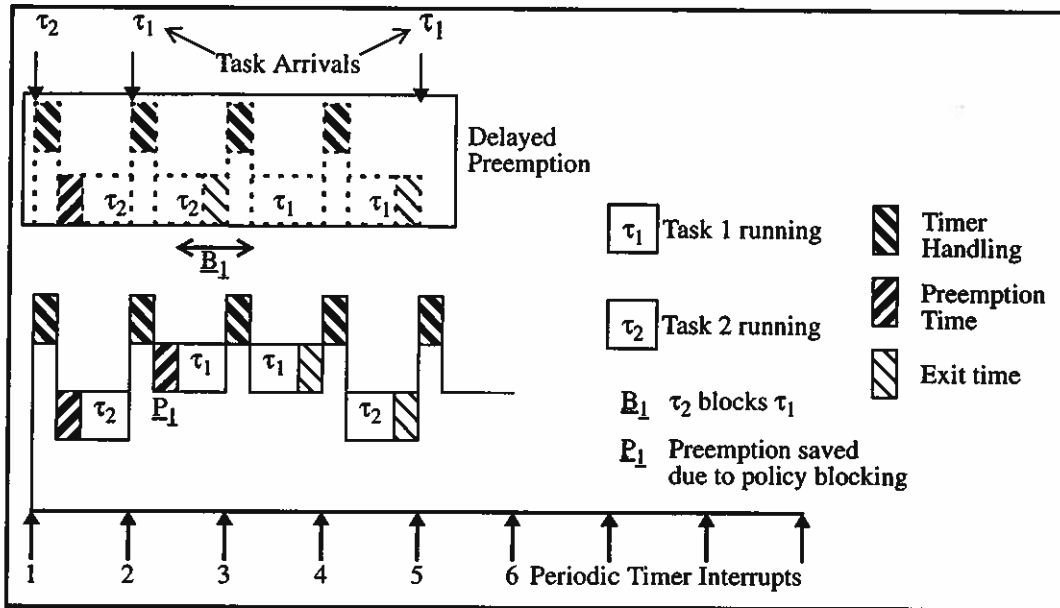


Figure 3: Mechanism Vs. Policy Blocking

preemption policies on the other hand is an example of *policy blocking* and the distinction between the two is made clearer in Section 4.

#### 4 Problem Definition and Overview of Results

This section precisely defines the nature of the policy blocking that we wish to study, and presents experimental results for RMDP to demonstrate the effect of policy blocking on schedulability. We then give a brief overview of our results that will be used to explain our observations. Figure 3 shows essential features of a real-time scheduler. Periodic timer interrupts cause tasks to arrive into the run queue. For example task  $\tau_2$  arrives at  $t = 1$ , and task  $\tau_1$  arrives at  $t = 2$  and  $t = 5$ . In the fully preemptive case (solid lines),  $\tau_1$  preempts  $\tau_2$  at  $t = 2$ . When  $\tau_1$  finishes at  $t = 4$ ,  $\tau_2$  resumes and finishes at  $t = 5$ . In the case of delayed preemption (shown on top with dotted lines)  $\tau_1$  does not preempt  $\tau_2$  at  $t = 2$ . Instead it runs until  $t = 3$  and then finishes, after which  $\tau_1$  runs. Thus  $\tau_1$  experiences *policy* blocking for a duration shown as  $B_1$ . During time  $B_1$ ,  $\tau_2$  makes progress and actually finishes. This is quite different from the *mechanism* blocking that is caused at every tick due to timer handling. Other examples of mechanism blocking shown in Figure 3 are the time to cleanup after a task completes (*exit* time), and the time to do a context switch (*preemption* time). An advantage of delayed preemption is that it can save involuntary context switches an example of which is  $P_1$ . The results in [7] only account for timer handling and exit times, and makes a simplified assumption that every arrival causes a preemption. On the other hand, our work studies policy blocking in a systematic manner. It also measures preemption costs with greater precision and shows how it can be reduced.

In Section 5 we show policy blocking caused by preemption policies reduces the utilization bound, just as mechanism blocking does. However the savings in context switching due to the use of preemption policies increases the utilization bound compared to the case when immediate preemption is used. The net effect is that there can often be an increase in the useful utilization. We demonstrate this with an experiment in which we created 16 processes, each with 6 RTUs with periods ranging from 40 milliseconds (msec) to 90 msec in increments of 10 msec. We compared two cases, one in which preemption was immediate, and the other in which preemption is delayed until the running handler completes. The running time of each handler was 0.54 msec giving a total utilization of 0.86. The experiment was run long enough (least common multiple of RTU periods) on a 100 Mhz Pentium running the NetBSD OS with a clock interrupt frequency (hz) of 1000. Figure 4 plots the smallest difference measured between the deadline and completion time for each RTU. It can be seen that for RM with immediate preemption, all the 90 msec have negative ordinates, which means that each of them finished after their deadline

in at least one invocation. In the case of RMDP, no RTU was found to miss its deadline. This result can be interpreted as saying that the increase in utilization due to RMDP is equal to the utilization of the 90 msec RTUs. This gain comes to around 9%. The rest of the paper essentially explains this result by analysis and measurement. Section 5 develops the tools to quantify the effect of policy blocking. Section 9 quantifies the effect of reducing the number of involuntary context switches by measurements.

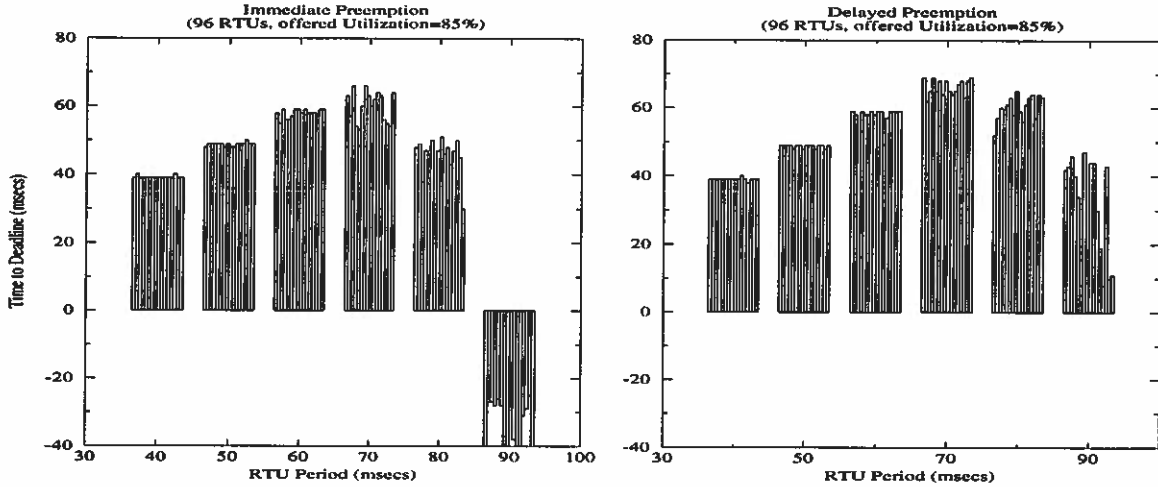


Figure 4: Increase in utilization by reducing context switching

#### 4.1 Preemption Policies

Before we present preemption policies, we introduce some notation and describe how the scheduler operates when a particular preemption policy is in effect. Tasks are represented as  $J_1, \dots, J_n$ . A task  $J_i$  has a period  $T_i$  and a CPU requirement (or usage requirement) of  $C_i$ . We assume that tasks with smaller subscripts have smaller periods. The utilization  $u_i$  of task  $J_i$  is given by  $C_i/T_i$ , and the total utilization  $U$  is the sum of the utilization for all tasks. The scheduler operates as follows. When a task arrives, it is assigned a priority according to the particular scheduling policy (RM or EDF) used. The task is then placed in the run queue which we assume is ordered on the basis of priority. If the new task lands at the head of the queue (i.e it has the highest priority among all runnable tasks in the system), the preemption policy is used to decide whether the running task should be preempted or not. If a preemption decision is made the appropriate action is taken to move the running task back to the run queue. In the case of RTUs, this would involve setting a flag in the shared data structure. The CPU becomes available for the next eligible task some time after the preemption action was initiated. This time delay is called the *preemption interval*. We assume that the preemption mechanism itself does not introduce any delays and only the preemption policy can contribute to the preemption interval. Once the running task yields the processor, the highest priority task in the run queue is assigned the processor. It must be noted that several tasks can arrive during the preemption interval. These tasks are also placed in the run queue and are considered when the choice for the next task to run is made.

We denote the running task by  $J_r$  and an arriving task by  $J_a$ . The threshold and delayed preemption policies are formalized as follows.

**Delayed Preemption** In this scheme, the computation time  $C_i$  for a task is divided into equal sized quanta  $c_i$ .

We allow  $J_r$  to be preempted only at the end of a quantum. If the priority of  $J_a$  is higher than that of  $J_r$ , action is initiated to preempt  $J_r$  at the end of its current quantum. The fully preemptive case occurs when the quantum size for all tasks is set to 0. The non preemptive case occurs when  $c_i$  is set to  $C_i$  for all tasks.

**Threshold Preemption** In this scheme,  $J_a$  preempts  $J_r$  if it has the highest priority among the runnable tasks (including  $J_r$ ), and the periods  $T_a$  and  $T_r$  are related by the inequality  $T_a/T_r < K_r$  where the threshold  $K_r$  is given by  $0 \leq K_r \leq 1$ . It is easy to see that the fully preemptive case occurs when the value of the threshold  $K_r$  is set to 1. Similarly the non preemptive case occurs when  $K_r$  is set to  $T_1/T_r$ . The idea behind

Preemption Policy	Utilization Bound	
	Rate Monotonic	Earliest-deadline-first
Threshold Preemption	$n(2^{1/n} - 1) - \max_{1 < r < n} u_r(1/K_r - 1)$	$1 - \max_{1 < r < n} u_r(1/K_r - 1)$
Delayed Preemption	$n(2^{1/n} - 1) - \max_{1 < r < n} c_r(1/T_1 - 1/T_r)$	$1 - \max_{1 < r < n} c_r(1/T_1 - 1/T_r)$
Non-preemptive	$n(2^{1/n} - 1) - \max_{1 < r < n} C_r(1/T_1 - 1/T_r)$	$1 - \max_{1 < r < n} C_r(1/T_1 - 1/T_r)$
Mixed	$n(2^{1/n} - 1)$	

Figure 5: Theoretical results at a glance

threshold preemption is that tasks that have almost equal periods are assigned the same priority. It can be seen that threshold preemption does not always provide the locking benefits that delayed preemption does, unless synchronization is needed only among tasks that satisfy the threshold condition.

Table 5 gives an overview of the scheduling results that we derive in this paper. It shows the two preemption policies we have considered. For both these policies, it shows the modified schedulability tests for RM and EDF that account for the blocking caused by these preemption policies. The table also lists the extensions to non-preemptive scheduling (Section 7) and mixed (Section 8) scheduling.

## 5 Idealized Scheduler Simulation

We have developed idealized scheduler simulation (ISS) as a general method to analyze the effect of policy blocking on the schedulability test of a scheduling algorithm. The basic idea behind ISS is to simulate the essential operation of a blocking scheduler (the actual case) with an ideal scheduler (idealized case) that never blocks. By essential operation, we mean that task invocations in the idealized case complete at the same time as in the actual case. Thus if the tasks in the idealized case meet their deadlines, then they do so in the actual case as well. The schedulability test for the idealized case can therefore be adopted for the actual case. To go from the actual case to the idealized case in ISS, the computation requirements of the tasks are modified while keeping their periods the same. For each blocking event that occurs when a task delays its preemption (the blocking task), the ISS method does the following steps.

1. It determines which tasks need to be modified in the idealized case so as to simulate the essential operation of the actual scheduler. Lemma 1 shows that we only need to modify tasks that have priority higher than that of the blocking task.
2. For the higher priority tasks it determines by how much the original computation requirements must be changed and this is shown in Lemma 2.
3. Due to these changes, the requested utilization of the modified task set becomes artificially higher than that of the original task set. Since the upper bound on utilization is fixed by the scheduling policy (RM or EDF), some utilization becomes unavailable. Lemma 3 shows how to calculate this upper bound on the unavailable utilization.
4. Finally the schedulability test for the blocking case is obtained by taking the maximum of the unavailable utilization due to all blocking tasks, and subtracting this from the utilization bound of the idealized case.

In this section we state Lemmas 1, 2, 3, and discuss them informally. The proofs are only of technical importance and are in Appendix A.

We begin with some notation. A task  $J_i$  in the actual case has a corresponding task  $J'_i$  in the idealized case. The tasks  $J'_i$  have modified computation requirements as per the ISS method. We can imagine two schedulers, one that runs the tasks  $J_i$ , and the other that runs the tasks  $J'_i$ . We parameterize the time duration for which a running task  $J_r$  can block other tasks (the preemption interval) as  $\Delta_r$ . For example, in the case of threshold preemption we have  $0 \leq \Delta_r \leq C_r$ , and in the case of delayed preemption we have  $0 \leq \Delta_r \leq c_r$ . We now compare a running task  $J_r$ , and its counterpart  $J'_r$  in the following lemma.

**Lemma 1** *Suppose that corresponding invocations of  $J_r$  and  $J'_r$  are running (in their respective schedulers) at a given time, and a task with higher priority arrives. Suppose  $J_r$  continues to run for a duration  $\Delta_r$  thereby*

gaining a usage of  $\Delta_r$  compared to  $J'_r$ . Then at the instant that  $J_r$  resumes execution,  $J'_r$  would have resumed and regained exactly the amount  $\Delta_r$  in usage.

*Proof:* In Appendix A.

An important consequence of this lemma is that in both the actual and the idealized cases, the time at which the CPU becomes available to tasks with priority lower than  $J_r$  is the same. Hence only tasks with priority higher than  $J_r$  are affected by delayed preemption.

We next identify the set of higher priority tasks that actually get affected by the blocking due to  $J_r$ . Identifying this set is necessary to isolate and study the effect of the blocking caused by  $J_r$ . This set  $W_r$ , consists of tasks that arrive during the preemption interval of  $J_r$  and are blocked by  $J_r$ , but which run in the idealized case that is simulating the essential behavior of the actual case. While analyzing a particular preemption policy we will have to determine the set  $W_r$ . Examples are shown in Sections 6.2 and 6.1 where we apply ISS to analyze RM and EDF scheduling.

In general, we define the set  $W_r = \{J_{r_1}, \dots, J_{r_i}, \dots, J_{r_m}\}$ , and the set  $W'_r$  as the set that contains for each invocation  $J_{r_i}$  in  $W_r$ , its counterpart  $J'_{r_i}$ . The tasks in  $W_r$  are listed in order of decreasing priority. For each invocation  $J'_{r_i}$ , let the amount of usage that it obtains within the delay interval be  $d_{r_i}$ . It can be seen that the sum of these usages  $\sum d_{r_i} = \Delta_r$ .

**Lemma 2** Consider corresponding invocations in  $W_r$  and  $W'_r$ . At the end of the preemption interval, each invocation  $J'_{r_i}$  has got an extra  $d_{r_i}$  amount of usage compared to  $J_{r_i}$ . If the usage requirement of each  $J'_{r_i}$  is increased by an amount  $d_{r_i}$ , then both  $J'_{r_i}$  and  $J_{r_i}$  complete execution at the same time. In addition, if a task  $J'_{r_i}$  meets its deadline with this increased requirement, then so does the task  $J_{r_i}$ .

*Proof:* In Appendix A.

The above lemma determines the amount by which the usage requirements of tasks should be increased so that they complete (in the idealized case) at the same time as the corresponding tasks blocked by  $J_r$  in the actual case. We also note that if the usage requirement of the task  $J'_r$  is reduced by the amount  $\Delta_r$  then it will complete at the same time as  $J_r$ . With these changes, the essential operation of the simulated and actual cases would be identical. The unavailable utilization due to these modifications is given next.

**Lemma 3** The amount of unavailable utilization  $N_r$  due to blocking by  $J_r$  is bounded by  $\Delta_r(1/T_{r_1} - 1/T_r)$ .

*Proof:* In Appendix A.

The above result can be explained simply as follows. The ISS method increases the computation requirement of the highest priority task that can be blocked by  $J_r$  by an amount  $\Delta_r$  and reduces the computation requirement of  $J_r$  by  $\Delta_r$ . The unavailable utilization is the (artificial) increase in utilization of the blocked task minus the (artificial) decrease in utilization of  $J_r$ . Knowing the minimum period  $T_{r_1}$  (from  $W_r$ ), and the maximum value of  $\Delta_r$  for the task  $J_r$ , the value of  $N_r$  can be computed using the above formula. To determine the sufficient condition for the task set to be schedulable, the worst case value of  $N_r$  must be determined for all tasks in the task set and subtracted from the utilization bound of the scheduling algorithm used.

## 6 Analysis of Preemption Schemes

We now apply the ISS method to analyze various preemption policies introduced earlier. The general technique is to determine the set  $W_r$  for the particular scheduling policy (RM or EDF) and the preemption policy (TP or DP) being used. These can then be used to compute  $N_r$  in Lemma 3.

### 6.1 Delayed Preemption

In the case of delayed preemption, a task  $J_r$  can block any higher priority task. Thus  $W_r$  can potentially include all tasks with priority higher than  $J_r$ . In the rate-monotonic case this implies all tasks with periods smaller than that of  $J_r$  are part of  $W_r$ . The situation is the same in the earliest-deadline-first case as well. This is because tasks with periods bigger than that of  $J_r$  will have deadlines later than that of  $J_r$ . Thus the value of the smallest period in  $W_r$  is  $T_1$  in both the cases. Since the maximum time for which  $J_r$  can block others is  $c_r$ ,  $N_r$  is bounded by

$$N_r \leq c_r(1/T_1 - 1/T_r)$$

The schedulability test for RM is given by 1 and the test for EDF is given by 2.

$$U \leq n(2^{1/n} - 1) - \max_{1 \leq r \leq n} c_r(1/T_1 - 1/T_r) \quad (1)$$

$$U \leq 1 - \max_{1 \leq r \leq n} c_r(1/T_1 - 1/T_r) \quad (2)$$

The above results give an increased utilization bound compared to [13] where the reduction in utilization due to blocking would be  $c_{max}/T_1$ .

## 6.2 Threshold Preemption

In this scheme a task  $J_r$  can block another task for a maximum time  $C_r$ . Thus  $\Delta_r \leq C_r$ . To determine  $W_r$  we have to consider the RM and EDF policies as shown next.

**Rate Monotonic Policy** In RM an arriving task  $J_a$  has higher priority if  $T_a < T_r$ . However  $J_a$  does not preempt  $J_r$  if  $T_a \geq K_r T_r$ . Thus in this case the set  $W_r$  includes all tasks  $T_a$  with periods less than  $T_r$  such that  $T_a \geq K_r T_r$ . The period of the task with the smallest period in  $W_r$  is given by  $T_{r_1} \geq K_r T_r$ . Thus  $N_r$  is bounded by

$$N_r \leq C_r(1/K_r T_r - 1/T_r) = u_r(1/K_r - 1)$$

It must be noted here that if  $W_r$  is empty (as it is when  $r = 1$ ), then  $N_r = 0$  regardless of the value of  $U_r$ . The schedulability test can therefore be written as

$$U = \sum_{i=1}^n u_i \leq n(2^{1/n} - 1) - \max_{1 \leq r \leq n} u_r(1/K_r - 1) \quad (3)$$

It can be easily seen that for  $K_r = 1$  the test reduces to the standard result in [11]. It must be noted that the above result holds for any combination of task periods. For a particular task set, the value of  $T_{r_1}$  may be larger than  $K_r T_r$  and so we may get a smaller value of  $N_r$ , and consequently a better utilization bound.

**Earliest Deadline First Policy** We determine the set  $W_r$  as follows. To decide if a task  $J_a$  that arrives during  $J_r$ 's preemption interval should be included in  $W_r$  or not, we consider two cases—

$T_a \geq K_r T_r$ : We first consider the case when  $J_a$  has the highest priority among all runnable tasks (including  $J_r$ ). When it arrives,  $J_a$  is blocked by  $J_r$  because of the threshold condition. However in the idealized case,  $J_a$  preempts  $J_r$  and therefore it is included in  $W_r$ . On the other hand if  $J_a$  at the time of its arrival is not the highest priority task among the runnable ones, then it does not cause a preemption in the actual case. In the idealized case, the counterpart of the highest priority task in  $W_r$  should be running when  $J_a$  arrives. This is because the ISS method would have increased its usage requirement in such a way that it can complete only after the delay interval. Thus there is no preemption in the idealized case either, and so  $J_a$  should not be a member of  $W_r$ .

$T_a < K_r T_r$ :  $J_a$  cannot belong to  $W_r$  in this case because of the following reasons. If  $J_a$ , when it arrives, has the highest priority among all runnable tasks (including  $J_r$ ) then it preempts the running task both in the idealized, and in the actual case. In fact the delay interval  $\Delta_r$  ends at the instant this preemption occurs. So  $J_a$  does not belong to  $W_r$  in this case. On the other hand if  $J_a$  at the time of its arrival is not the highest priority task among the runnable ones, then it does not preempt  $J_r$  in the actual case by definition of the preemption policy. As in the previous case the counterpart of the highest priority task in  $W_r$  should be running in the idealized case and so  $J_a$  does not cause a preemption in the idealized case either. Thus  $J_a$  cannot be a member of  $W_r$ .

Hence the value of the smallest period in  $W_r$  is such that  $T_{r_1} \geq K_r T_r$ . Thus  $N_r$ , is bounded by

$$N_r \leq C_r(1/K_r T_r - 1/T_r) = u_r(1/K_r - 1)$$

As before if  $W_r$  is empty (as it is for  $r = 1$ ), then  $N_r = 0$  regardless of the value of  $U_r$ . The schedulability test can therefore be written as

$$U = \sum_{i=1}^n u_i \leq 1 - \max_{1 \leq r \leq n} u_r(1/K_r - 1) \quad (4)$$

### 6.2.1 Improvement over Existing Results

For both RM and EDF scheduling, the improvement in utilization compared to the result given in [13] is equal to  $u_r$ . As  $K_r$  approaches 1.0, the unavailable utilization given by [13] approaches  $u_r$  whereas in our case it approaches 0. Intuitively our result shows that in the limit when  $K_r$  approaches 1.0, there is no blocking and consequently there should be no reduction in utilization.

## 7 Analysis of Non-preemptive EDF Scheduling

In this section we show that non-preemptive scheduling can be treated as a special case of both the threshold preemption and delayed preemption schemes. Non-preemptive scheduling schemes have many practical advantages from the implementation standpoint. As stated in [5], they are easier to implement, and their overheads are easier to characterize. A necessary and sufficient condition is also derived in [5] to determine if a task set is schedulable under a non-preemptive EDF scheduling scheme. These conditions are stated as–

1.  $\sum_{i=1}^n C_i/T_i \leq 1$
2.  $\forall i, 1 < i \leq n; \forall L, T_1 < L < T_i : L \geq C_i + \sum_{j=1}^{i-1} \left\lfloor \frac{L-1}{T_j} \right\rfloor C_j$

The above tests are derived by computing upper bounds for the processor demand over an interval of time, and checking to see if this demand can be met by the EDF scheduling scheme within this interval.

The problem with the above test is that the time complexity of evaluating the second condition is polynomial in the maximum task period [5]. This is because the test has to consider all task arrivals for a duration equal to the maximum task period. We have obtained a much simpler test for non-preemptive EDF by considering it as a special case of the threshold and delayed preemption schemes. We illustrate this below:

- *Threshold preemption:* For each task  $J_i$ , we set its threshold  $K_i$  such that it cannot be preempted even by the task with the smallest period. Since  $T_1$  is the smallest period, if we set  $K_i$  to the value  $T_1/T_i$ , then no task can preempt  $J_i$ .
- *Delayed preemption:* For each task  $J_i$ , we set the quantum size  $c_i$  equal to  $C_i$ . Therefore  $J_i$  can never be preempted by any higher priority task.

Substituting these values into Equations 4, and 2 we obtain a simple sufficiency condition for non-preemptive EDF scheduling given by Equation 5 below.

$$U \leq 1 - \max_{1 \leq r \leq n} C_r(1/T_1 - 1/T_r) \quad (5)$$

Equation 5 has a time complexity of  $O(n)$  which is a vast improvement over the test in [5]. However our test only gives a sufficient condition for schedulability whereas the test in [5] provides a necessary and sufficient condition. In the next section we simulate task sets of different characteristics to see how close our sufficient condition is to the optimal condition given in [5].

### 7.1 Comparing the Sufficient and Optimal Conditions

We compare the simple utilization based sufficiency condition (Equation 5) to the optimal, but more complex necessary and sufficient conditions derived in [5]. We make the comparison by generating task sets based on some criteria, and determining the breakdown utilization [7] predicted by the two tests. Breakdown utilization is defined as follows. The execution time for each task is multiplied by a constant scaling factor  $\alpha$  while the periods remain fixed. The utilization point at which any further increase in  $\alpha$  causes the schedulability test to be violated is the breakdown utilization. The difference between the breakdown utilizations for the two tests is a measure of the amount by which the simple sufficiency condition underestimates the scheduling capacity. For evaluating the conditions in [5] we used a simplified method suggested in [6].

We vary two parameters while generating a task set—the number of tasks per task set, and the range of periods of tasks in a given task set. In the first simulation we fix the number of tasks per task set to be 5, and try out different ranges for task periods in each set. In the second simulation, the task periods range from 100 to 1000, but the number of tasks per set is varied.

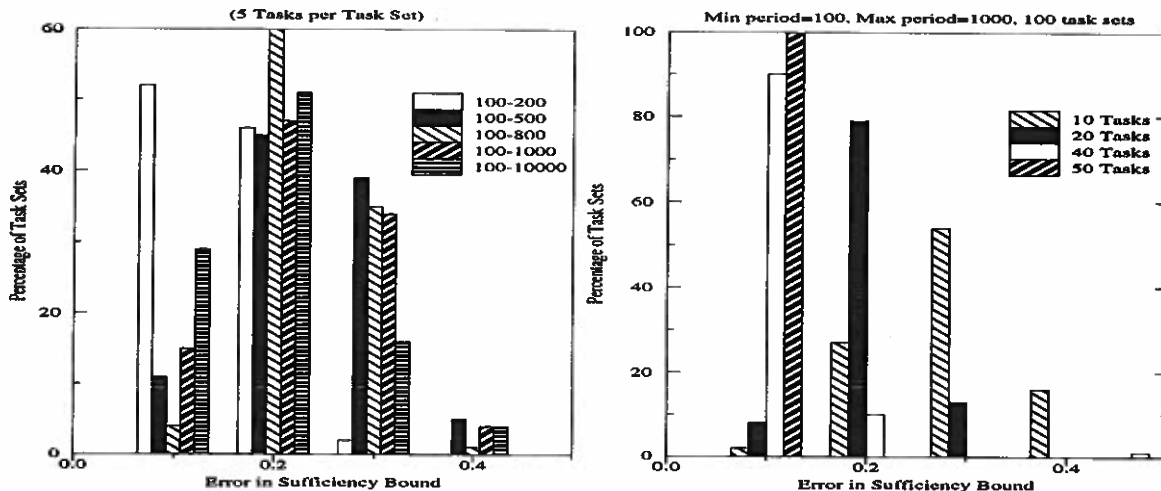


Figure 6: Accuracy of Sufficiency Condition

- *Simulation 1:* The ratio of the maximum task period to the minimum task period was varied over different runs in the range 2 to 10. For each run, 100 task sets were generated with task periods chosen according to a uniform distribution between the minimum and maximum periods. The computation requirement for each task was uniformly distributed between 1 and the minimum period. The breakdown utilization for the two schedulability tests was computed, and the difference between the two was obtained. This is shown in the first part of Figure 6. The x-axis has different values of this difference (or error). For a given error  $\epsilon$ , the y-axis shows for what percentage of task sets belonging to a particular period range, the error in our simple test was less than  $\epsilon$ . As can be seen, when the range of periods in the task set is small, the majority of the task sets have errors in the lower ranges. As the range of periods increase, higher error ranges are more prevalent. The absolute error however can be as high as 45%, which indicates that the sufficiency bound alone cannot give good utilization in all cases. A reasonable strategy would be to apply the simple test first while admitting a new task. This is a constant time operation. If the sufficiency condition fails, the more sophisticated optimal condition can be checked.
- *Simulation 2:* Here we investigate the effect of increasing the number of tasks in a task set on the accuracy of the sufficiency bound. The tasks were generated as before and the difference of breakdown utilizations was obtained. It can be seen (second part of Figure 6) that as the number of tasks are increased, the accuracy of the test increases steadily. When the number of tasks per set is 50 and above, the error is within 10% for all task sets. We therefore conclude that the sufficiency bound is very useful for larger task sets. It must be noted that the complexity of evaluating the necessary condition in [5] increases with the number of tasks, which makes our simpler test even more preferable.

## 8 Mixed Scheduling—A Modification to RM

In this section we use our earlier results to analyze a simple variation of the rate monotonic scheme. In this variation, tasks are assigned priorities according to the RM scheme. However, a task arrival that has higher priority than the running task preempts it only if its deadline is earlier than that of the running task. In other words, tasks are chosen to run according to the rate monotonic policy, but the preemption decision is based on an earliest deadline first policy. We refer to this as the *mized* policy. We show that any set of tasks that are schedulable by the RM policy, are schedulable by the mixed policy.

An interesting feature of the mixed policy is that it has fewer context switches compared to RM. In fact it is shown in [12] that EDF has fewer context switches than the RM for periodic tasks. One can see that when a preemption occurs in the EDF policy, the period of the arriving task must be smaller than the period of the running task, which implies that a preemption should occur in the RM policy as well. However, every preemption in RM does not imply a preemption in EDF. The same argument can be made for the mixed policy, and it turns out that it has almost as few context switches as EDF.

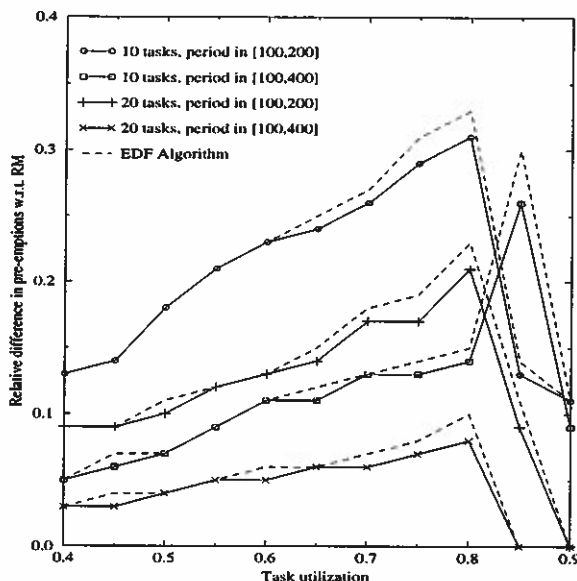


Figure 7: Reduction in Context Switching

### 8.1 Analysis of mixed scheduling

We show that any task set that is schedulable by RM is also schedulable by mixed scheduling. We use Lemma 1 which proves that if a running task  $J_r$  delays yielding the processor, then its usage at the time instant it resumes<sup>1</sup>, is the same as that of its counterpart in the idealized simulation. We consider two cases—

1. *All higher priority tasks that arrive have deadlines later than that of  $J_r$ :* In this case  $J_r$  runs to completion before yielding the processor. When  $J_r$  resumes by definition all the blocked tasks should have completed. By Lemma 1  $J_r'$  also completes at this time. Since  $J_r'$  meets its deadline, so do the blocked tasks since they have deadlines later than  $J_r'$ .
2. *Some tasks have deadlines earlier than that of  $J_r$ :* An arrival that has an earlier deadline would also have a smaller period than that of  $J_r$  (since deadlines are equal to periods for periodic tasks). When  $J_r$  resumes, by definition all higher priority blocked tasks would have completed, and by Lemma 1  $J_r$  will have the gotten the same usage as  $J_r'$ . After this instant either both  $J_r$  and  $J_r'$  complete execution without further higher priority task arrivals, or there is a higher priority task arrival and the above two cases repeat again.

**Performance of Mixed Scheduling** We simulated randomly generated task sets to compare the number of context switches in the RM, EDF, and mixed scheduling policies. The task periods and computation times were uniformly distributed as before. The utilization of the task sets was increased in each run until the breakdown utilization was reached. A count of the number of context switches for each policy was recorded in each run. The count for each utilization point was averaged over several task sets. The averaged number of context switches saved in the EDF and mixed policies compared to RM, normalized with respect to the number of context switches in the RM policy is plotted at different utilization points in Figure 7. Each task set was simulated long enough so that these values stabilized.

It can be seen that the number of context switches saved in EDF is only slightly more than in the mixed policy. These two curves track each other closely in the utilization range shown. We also see that as the number of tasks in the task set increases, the relative savings in context switching diminish for both the EDF and mixed policies. As expected the savings in context switches increase at higher utilizations. The difference between the EDF and the mixed policies is more at higher utilizations but maximum difference did not exceed 2%. From this we conclude that the mixed policy is almost always preferable to the RM policy if reducing context switching is

<sup>1</sup>If it completes execution when it yields the processor, we assume that it resumes but runs for zero time



	Pentium	Sparc-1
Inter-process	35	373
Intra-process	14	114

Figure 8: Measured Context Switch Time (microseconds)

of concern. The only extra operation in the mixed scheme compared to RM is the deadline comparison which is a constant time operation. In EDF on the other hand the arriving task must be inserted in a queue sorted by deadline and this takes  $O(\log n)$  time. The mixed policy therefore combines the simplicity of RM with the context switching performance of EDF.

## 9 Accounting for Context Switching Costs

Our analysis of preemption policies shows that the blocking caused by them reduces utilization. However our experimental results with RMDP indicate that we can support task sets with higher utilization than the case when we use immediate preemption. It follows that the reduction in the number of involuntary context switches offsets the loss of utilization due to blocking. We have verified this fact experimentally on the Sparc and Pentium platforms. We first describe our experiment to measure the context switch time using the RTU mechanism. We then measure the rate at which involuntary context switches occur when the task set is run. This rate depends on the clock interrupt frequency. Finally we calculate the loss of utilization due to context switching and how much we gain by using RMDP.

### 9.1 Measuring Context Switch Time

The RTU facility helps us measure context switch time using user level programs. We ran two experiments, one to measure context switch time between RTUs in two different processes (inter-process), and the other to measure the time to switch between RTUs in the same process (intra-process). We created two RTUs, one with a period of 10 msec and the other with a period of 20 msec. In one case these RTUs were created in different processes and in the other case they were in the same process. The 20 msec handler goes into a loop polling the flag in its shared memory structure which is set by the kernel when the 10 msec handler becomes runnable. When the handler detects the flag to be set, it records the time and yields the CPU. The kernel saves the state of the yielding process and switches in the process that has the 10 msec RTU. The handler of the 10 msec RTU is written to return as soon as it is called. The 20 msec handler that had yielded earlier is still in the run queue at this point. The kernel switches in the process that contains the 20 msec RTU and invokes its handler. When the handler resumes, it records the resumption time. The difference between these times gives the time for two context switches, *i.e.*, from the 20 msec RTU to the 10 msec RTU, and back. Half of this gives the context switch time. The context switch times as measured on a 100 Mhz Pentium and on a 25 Mhz Sparc-1 are shown in Table 8. We can see that with a 1 msec clock rate, preemptions on the Sparc-1 will impose too much overhead. We expect that newer Sparc machines will not have this problem. For now, all experiments are reported for the Pentium machine.

### 9.2 Measuring Number of Context Switches

We now go back to our experiment in Figure 4 to explain why the 90 msec RTUs were missing their deadlines when we used immediate preemption (IP) but were able to meet their deadlines when we used non-preemptive (NP) scheduling. Consider the IP case shown in Figure 9. A 90 msec RTU  $J_a$  arrives at  $t_a$  but is unable to complete before  $t_a + 90$ . Let  $t_{prev}$  be the time at which the previous invocation of  $J_a$  completes.  $J_a$  can be affected due to preemption in the following two ways:

- If  $P_1$  preemptions occur between  $t_{prev}$  and  $t_a$ , the higher priority RTUs (dashed lines) that are in the run queue at  $t_a$  get delayed by the time for  $P_1$  preemptions. Since  $J_a$  can run only after the higher priority tasks complete,  $J_a$  also gets delayed by  $P_1$  preemptions.
- If  $P_2$  preemptions occur between  $t_a$  and  $t_a + 90$  then  $J_a$  gets further delayed by time for  $P_2$  preemptions.

We need to measure  $P_1$  and  $P_2$  in the worst case. We measured  $t_{prev}$  to be  $t_a - 70$  msec. Therefore we need to measure the number of preemptions that occur in the worst case in an interval of length  $70 + 90 = 160$  msec.

We measured the number of preemptions in any interval  $I$  msec in the following manner. The scheduler maintains a moving window with  $I$  slots, one for each tick. Each slot records a “hit” if a switch occurred on that

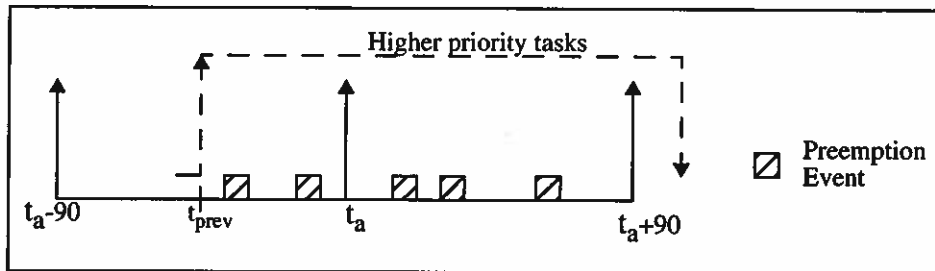


Figure 9: Delays due to Context Switching

clock tick. By moving this window across the “tickline” we counted the maximum number of hits in any  $I$  msec interval. For  $I = 160$  this number was as high as 110. Each preemption costs 2 context switches. Using the values shown in Table 8 for Pentium, each preemption costs 0.07 msec which results in a delay of  $110 * .07 = 7.7$  msec.

### 9.3 Interpreting Experimental Results

The above measurements show that preemptions can cause delays which could lead to missed deadlines. By reducing the number of preemptions using a scheduling scheme such as RMDP, we can offset the reduction in utilization that is predicted by our analysis. In the experiment above, the reduction due to blocking given by Equation 1 was  $0.54(1/40 - 1/90) = 0.75\%$ . On the other hand by avoiding the overhead due to preemption we can increase available utilization by upto  $7.7/90 = 8.5\%$ . Our goal is to be able to predict the gain due to reduction in preemptions not only for the non-preemptive case shown above, but also for other cases. We are currently looking into ways to do this.

**Estimating gains from saving system calls** The measurements presented earlier only consider gains due to reduction in preemptions. We can also model the gains due to saving on locking system calls in a simple manner. If we assume from [1] that each system calls is 770 cycles, then one pair of lock set/release operations should take (assuming cycle time is 10 nanoseconds)  $770 * 2 * 10 = 0.015$  msec. Even if one lock is acquired and released per period by each RTU the total utilization due to locking is  $.015 * 16 * (1/40 + \dots + 1/90) = 2.4\%$ . Typically locking is done per packet (or a batch of packets), and thus the cost of locking increases linearly with the number of batches processed per period. If we assume a 10Mbps link speed, at least 1 1KB packet must be processed per msec to fully utilize the link. In this case a lock per packet represents 1.5% of CPU time. This cost increases linearly with link speeds assuming packet sizes remain the same. Since speeds of 100 Mbps are not uncommon, this linear increase can reduce throughput seen by applications. Thus schemes such as RMDP that do not incur additional per packet costs will scale well to higher speeds.

## 10 Conclusions

With the rapid deployment of multimedia applications, the need for real-time scheduling in general purpose computers continues to increase. We have implemented the RMDP scheduling mechanism so that protocol processing can be done efficiently. RMDP lowers the costs of context switching, and eliminates the need for locking operations in protocol processing. To be able to use RMDP and other practical scheduling schemes that operate in the region between strict preemption and non-preemption, we have developed analytical tools based on the idealized scheduler simulation methodology. Using this we were able to derive simple schedulability tests for rate-monotonic and earliest-deadline-first scheduling that apply in cases when preemption is not immediate. The use of preemption policies such as delayed preemption has the potential to reduce context switching costs. To quantify these savings, we measured context switching costs on the Sparc-1 and Pentium platforms. From the measured values we conclude that the Sparc-1 platform will perform poorly specially when clock interrupts occur every 1 msec. We then estimate the gains due to savings in system calls during protocol processing. We show that these costs are substantial and increase with transmission speeds, and so eliminating them makes RMDP scale well to high speeds.

## References

- [1] Chen, J.B., et. al., "The Measured Performance of Personal Computer Operating Systems," 15<sup>th</sup> ACM SOSP, Dec. 1995.
- [2] Gopalakrishnan R., Parulkar G.M., "Real-time Upcalls: A Mechanism to Provide Real-time Processing guarantees," Tech. Rep. WUCS-95-06, Washington University, St.Louis, 1995.
- [3] Gopalakrishnan R., Parulkar G.M., "A Generalized Preemption Model for Real-time Scheduling," Tech. Rep. WUCS-96-04, Washington University, St.Louis, 1996.
- [4] Gopalakrishnan R., Parulkar G.M., "A Real-time Upcall Facility for Protocol Processing with QoS Guarantees," (Poster) 15<sup>th</sup> ACM SOSP, Dec. 1995.
- [5] Jeffay, K., Stanat, D.F., Martel, C.U., "On Non-Preemptive Scheduling of Periodic and Sporadic Tasks," 12<sup>th</sup> IEEE Real-Time Systems Symposium, Dec 1991.
- [6] Jeffay, K., Stone, D.L., "Accounting for Interrupt Handling Costs in Dynamic Priority Task Systems," 14<sup>th</sup> IEEE Real-Time Systems Symposium, Dec 1993.
- [7] Katcher et. al., "Engineering and Analysis of Fixed Priority Schedulers," IEEE Transactions on Software Engineering, Sep 1993.
- [8] Katcher, D.I., "Engineering and Analysis of Real-time Operating Systems," PhD Thesis, Carnegie Mellon University, 1994.
- [9] Khanna, S., et. al., "Realtime Scheduling in SunOS5.0," USENIX, Winter 1992, pp.375-390.
- [10] Lehoczky, J.P., Sha, L., "Performance of Real-Time Bus Scheduling Algorithms," ACM Performance Evaluation Review, Vol.14, No.1, May 1986.
- [11] Liu, C.L. and Layland, J.W., "Scheduling Algorithms for Multiprogramming in a Hard-Real-Time Environment," JACM, Vol. 20, No. 1, January 1973.
- [12] Pingali, S., "Protocol and Real Time Scheduling Issues for Multimedia Applications," PhD Thesis, University of Massachusetts, Amherst, Sep 1994.
- [13] Sha, L. et. al., "Priority Inheritance Protocols: An Approach to Real-Time Synchronization," IEEE Transactions on Computers, Vol.39, No.9, Sep 1990.
- [14] Sha, L., Lehoczky, J.P., Rajkumar, R., "Solutions for Some Practical Problems in Prioritized Preemptive Scheduling," IEEE Real-Time Systems Symposium, Dec 1986.
- [15] Tokuda, H., Nakajima, T., Rao, P., "Real-Time Mach: Towards Predictable Real-time Systems," USENIX Mach Workshop, Oct 1990.

## A Appendix

This section contains the proofs of Lemmas 1, 2, 3 on which the ISS method is based. We restate the lemmas that appeared in Section 5.

**Lemma 1** *Suppose that corresponding invocations of  $J_r$  and  $J'_r$  are running (in their respective schedulers) at a given time, and a task with higher priority arrives. Suppose  $J_r$  continues to run for a duration  $\Delta_r$  thereby gaining a usage of  $\Delta_r$  compared to  $J'_r$ . Then by the time  $J_r$  resumes execution,  $J'_r$  would have resumed and regained exactly the amount  $\Delta_r$  in usage.*

*Proof:* We assume that both  $J'_r$  and  $J_r$  resume execution after the higher priority task completes execution. (If  $J_r$  were to complete its execution in the time  $\Delta_r$ , it would normally not resume. However for the sake of uniformity, we assume that  $J_r$  does resume, but it runs for 0 time before completing).

Figure 10 shows the times  $b$  and  $a$  at which  $J_r$  and  $J'_r$  are preempted, and the times  $d$  and  $c$  at which they first resume execution. We note that the usage requirement of higher priority arrivals is the same in both schedulers. There are two cases to consider.

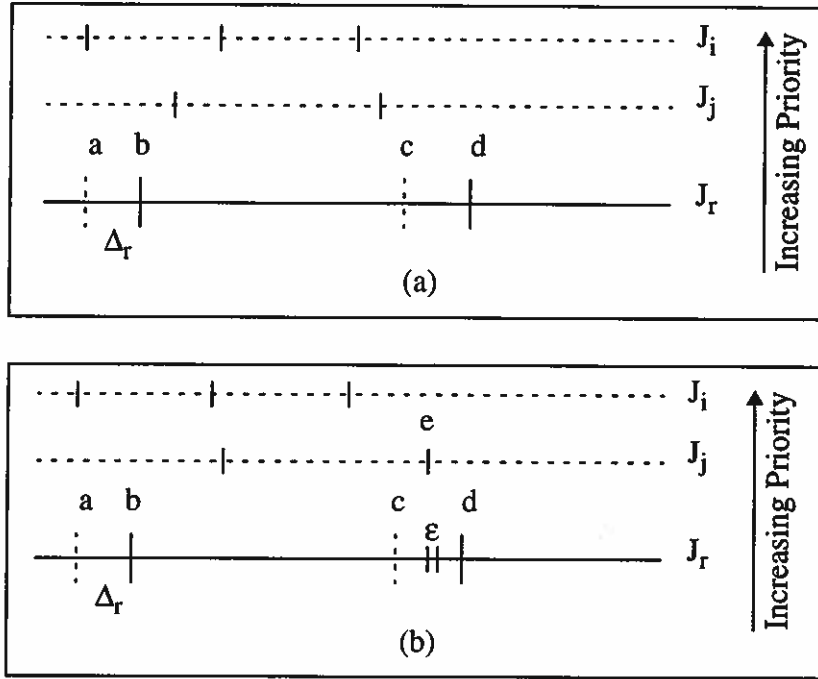


Figure 10: Effect of Delaying Preemption

- There is no higher priority arrival between  $c$  and  $d$ . In this case the amount of demand for CPU by higher priority tasks between  $a$  and  $c$  is the same as that between  $a$  and  $d$ . Since  $J'_r$  resumes at  $c$ , the time for higher priority tasks to complete their execution is  $c - a$ . Hence  $d - b$  must equal  $c - a$  and hence  $d - c = b - a = \Delta_r$ . Therefore at time  $d$ ,  $J'_r$  has got back exactly  $\Delta_r$  of CPU time after resumption. This also means that  $J_r$  and  $J'_r$  have obtained the same usage between  $a$  and  $d$ .
- There are higher priority task arrivals between  $c$  and  $d$ . This case is shown in Figure 10(b) where there is a higher priority task arrival at  $e$  with a service time requirement of  $\epsilon$ . After  $J'_r$  resumes at  $c$ , it runs until  $e$  and then resumes after an interval  $\epsilon$ . However,  $J_r$  must also wait for the arrival at  $e$  to complete its work. Hence  $d - b = c - a + \epsilon$  and hence  $d - c = \Delta_r + \epsilon$ . Since  $J'_r$  resumes at  $c$  and it takes  $\Delta_r + \epsilon$  time after this for it to accumulate  $\Delta_r$  amount of usage, the time at which this happens is  $d$  which is the time at which  $J_r$  resumes. Therefore in the interval  $[a, d]$  both  $J_r$  and  $J'_r$  have got the same amount of usage namely  $\Delta_r$ .

□

One implication of this lemma is that the relationship between the CPU usage of  $J_r$  and  $J'_r$  that exists at the time when a higher priority task arrives is restored once the higher priority arrival completes its execution requirement and both invocations resume execution. Another consequence of the above lemma is that regardless of whether a task blocks other tasks, the time at which the CPU becomes available to lower priority tasks is the same. Hence we only need to consider higher priority tasks that get blocked because of the blocking caused by  $J_r$ .

The next lemma uses the set  $W_r$ , and  $W'_r$  that was explained in Section 5. The tasks in  $W_r$  are listed in order of decreasing priority. For each invocation  $J'_{r_i}$ , let the amount of usage that it obtains within the delay interval be  $d_{r_i}$ . It can be seen that the sum of these usages  $\sum d_{r_i} = \Delta_r$ . We will refer to the time at which an invocation  $J'_{r_i}$  completes its execution to be the “normal” completion time. We relate the members of  $W_r$  and  $W'_r$  in the following lemma.

**Lemma 2** Consider corresponding invocations in  $W_r$  and  $W'_r$ . At the end of the delay interval, each invocation  $J'_{r_i}$  has got an extra  $d_{r_i}$  amount of usage compared to  $J_{r_i}$ . If the usage requirement of each  $J'_{r_i}$  is increased by

an amount  $d_{r_i}$ , then both  $J'_{r_i}$  and  $J_{r_i}$  complete execution at the same time. In addition, if a task  $J'_{r_i}$  meets its deadline with this increased requirement, then so does the task  $J_{r_i}$ .

*Proof:* We prove this by induction on the index  $i$  in  $r_i$ . We first consider  $i = 1$ . After the preemption interval, the task  $J_{r_1}$  begins execution, and at this point  $J'_{r_1}$  would have got a usage of  $d_{r_1}$  since its arrival. We first consider the case when there is no higher priority invocation arrival until  $J'_{r_1}$  completes. In this case,  $J_{r_1}$  will have to run for a time  $d_{r_1}$  after the normal completion time in order to satisfy its usage requirement. However if the usage requirement of  $J'_{r_1}$  were to be increased by an amount  $d_{r_1}$ , then both  $J'_{r_1}$  and  $J_{r_1}$  will have the same CPU requirement at the end of the delay interval and they will therefore complete at the same time. If the total utilization is such that this increase in usage requirement does not increase the total utilization beyond the bound given by the schedulability test, then the task  $J'_{r_1}$  meets its deadline and hence  $J_{r_1}$  also meets its deadline. In addition they both complete at the same time.

If we allow higher priority arrivals before  $J'_{r_1}$  completes, then  $J_{r_1}$  could delay yielding the CPU by an amount  $\Delta_{r_1}$ . However according to lemma 1, the relationship in usage between  $J_{r_1}$  and  $J'_{r_1}$  is restored after they both resume. Therefore the argument in the previous case applies also to this case. Thus  $J_{r_1}$  meets its deadline and it completes at the same time as  $J'_{r_1}$  even when we allow preemptions.

We now assume that the result holds for all invocations upto  $J_{r_{i-1}}$ . Since  $J_{r_{i-1}}$  and  $J'_{r_{i-1}}$  complete at the same time,  $J_{r_i}$  and  $J'_{r_i}$  get the CPU at the same time. At this time the difference in usage between them is exactly  $d_{r_i}$ . Therefore the argument for  $i = 1$  again applies, and as before holds good even if we consider preemptions. Therefore they complete at the same time and the invocation  $J_{r_i}$  meets its deadline provided  $J'_{r_i}$  meets its deadline.

□

**Lemma 3** *The amount of unavailable utilization  $N_r$  due to blocking by  $J_r$  is bounded by  $\Delta_r(1/T_{r_1} - 1/T_r)$ .*

*Proof:* The unavailable utilization due to  $J_r$  is the amount by which the utilization of the task set has been artificially increased by modifying the usage requirements. This is given by

$$\begin{aligned}
 N_r &= (d_{r_1}/T_{r_1} + \dots + d_{r_m}/T_{r_m}) - \Delta_r/T_r \\
 &\leq \frac{1}{T_{r_1}}(d_{r_1} + \dots + d_{r_m}) - \Delta_r/T_r \\
 &= \Delta_r(1/T_{r_1} - 1/T_r)
 \end{aligned} \tag{6}$$