

**A Novel Approach to Multiprogrammed  
Multiprocessor Synchronization  
for Real-Time Kernels**

Hiroaki Takada and Ken Sakamura

January 29, 1997



DEPARTMENT OF INFORMATION SCIENCE  
FACULTY OF SCIENCE, UNIVERSITY OF TOKYO

7-3-1 HONGO, BUNKYO-KU, TOKYO, 113 JAPAN

<b>TITLE</b> A Novel Approach to Multiprogrammed Multiprocessor Synchronization for Real-Time Kernels	
<b>AUTHORS</b> Hiroaki Takada and Ken Sakamura	
<b>KEY WORDS AND PHRASES</b> multiprogrammed multiprocessor, synchronization, real-time system, SPEPP synchronization, wait-free synchronization, preemption-safe locking	
<b>ABSTRACT</b> In order to solve the problem of inopportune preemption in multiprogrammed multiprocessor synchronization, two strategies which are applicable to real-time systems have been investigated: preemption-safe locking and wait-free synchronization. Either of them, however, has a problem for use in the implementation of a real-time kernel. Preemption-safe locking has the drawback that the preemption overhead becomes quite large; while wait-free operations on complex data structures are generally very inefficient and are not practical. In this paper, we propose a novel approach to multiprogrammed multiprocessor synchronization, called the SPEPP (Spinning Processor Executes for Preempted Processors) synchronization, with which the preemption overhead can be reduced to zero, while operations on complex data structures can be realized with reasonable overhead. This paper presents the two algorithms of the SPEPP synchronization, and demonstrates its effectiveness through the performance measurements of a real-time kernel implemented with the SPEPP synchronization.	
<b>REPORT DATE</b> January 29, 1997	<b>WRITTEN LANGUAGE</b> English
<b>TOTAL NO. OF PAGES</b> 13	<b>NO. OF REFERENCES</b> 11
<b>ANY OTHER IDENTIFYING INFORMATION OF THIS REPORT</b>	
<b>DISTRIBUTION STATEMENT</b> First issue 45 copies. This technical report is available via anonymous FTP from ftp.is.s.u-tokyo.ac.jp (directory /pub/tech-reports).	
<b>SUPPLEMENTARY NOTES</b>	

# A Novel Approach to Multiprogrammed Multiprocessor Synchronization for Real-Time Kernels

Hiroaki Takada

Department of Information Science,  
School of Science, University of Tokyo

Ken Sakamura

The University Museum,  
University of Tokyo

Technical Report 97-01

January 29, 1997

## Abstract

In order to solve the problem of inopportune preemption in multiprogrammed multiprocessor synchronization, two strategies which are applicable to real-time systems have been investigated: preemption-safe locking and wait-free synchronization. Either of them, however, has a problem for use in the implementation of a real-time kernel. Preemption-safe locking has the drawback that the preemption overhead becomes quite large; while wait-free operations on complex data structures are generally very inefficient and are not practical.

In this paper, we propose a novel approach to multiprogrammed multiprocessor synchronization, called the SPEPP (Spinning Processor Executes for Preempted Processors) synchronization, with which the preemption overhead can be reduced to zero, while operations on complex data structures can be realized with reasonable overhead. This paper presents the two algorithms of the SPEPP synchronization, and demonstrates its effectiveness through the performance measurements of a real-time kernel implemented with the SPEPP synchronization.

## 1 Introduction

When a real-time system is constructed on a shared-memory multiprocessor, it is a common practice to use a real-time kernel<sup>1</sup> to execute multiple tasks preemptively on each processor. A real-time kernel designed for multiprocessor systems should also support some communication mechanisms among tasks on different processors. Various data structures within the real-time kernel are usually shared among processors and guarded with busy-wait mutual exclusion locks.

In realizing mutual exclusion locks used in a real-time kernel, the following difficulty exists. In order to bound the maximum execution time of a task that requires a lock, the turn that the task acquires the lock should be reserved when it begins waiting for the lock. Typically, this is realized by handing round the lock in an FCFS order. If the lock is handed to the task while the task is preempted, it cannot begin to execute the critical section immediately and holds the lock for an unnecessarily long duration. As the result, the other tasks waiting for the lock cannot acquire the lock in time. This difficulty has been discussed as one of the problems of multiprogrammed multiprocessor synchronization [1].

Two strategies have been investigated to solve this difficulty: preemption-safe locking and wait-free (or block-free) synchronization [2]. Wait-free synchronization guarantees that any task can complete an operation on a shared data structure in a finite number of steps, regardless of the execution speed of the other tasks [3], while block-free synchronization guarantees that one

---

<sup>1</sup>In this paper, the word “real-time kernel” indicates a basic run-time system for real-time systems supporting task management, priority-based preemptive scheduling, inter-task synchronization and communication, and some other functions, which is also called as a real-time monitor or a real-time executive.

of the contending tasks can. Block-free synchronization is not appropriate for real-time systems, because it cannot bound the worst-case execution time that a task executes an operation. Herlihy presented a general construction of wait-free synchronization operations, assuming a universal atomic primitive, such as compare\_and\_swap (CAS) or the pair of load\_linked and store\_conditional [3]. However, wait-free (or block-free) operations on complex data structures are very inefficient and are not practical [2].

Preemption-safe locking is a class of spin lock algorithms with which a lock is not handed to preempted tasks. The MCS lock [4], which is an FCFS-ordered queueing spin lock algorithm and thus is suitable for real-time systems, has been extended to be preemption-safe [5, 6]. With the preemption-safe version of the MCS lock, when a task that is waiting for a lock is preempted, it modifies a shared variable and informs others that it is preempted. The task releasing the lock checks the state of the succeeding task in the waiting queue. If the succeeding one is preempted, its reservation to acquire the lock is canceled and the lock is handed to the next one in line. When a task finds that its reservation has been canceled while preempted, it re-executes the lock acquisition routine from the beginning. Obviously, the maximum time of this re-execution, called the preemption overhead below, is quite long depending on the number of contending tasks. In order to bound the preemption overhead with a constant time length, we have proposed an improved algorithm with which the reservation is not canceled but is postponed [7]. When the preempted task is scheduled again, it resumes waiting for the lock in its original position.

This paper proposes a novel approach to multiprogrammed multiprocessor synchronization, with which the preemption overhead can be reduced to zero, for use in a scalable real-time kernel for function-distributed multiprocessors [8]. In case of function-distributed multiprocessors, interrupt requests are raised on each processor, and the tasks on the processor can be preempted with them. The maximum time that a task waits for a lock is extended by the preemption overhead, whenever the task is preempted by an interrupt handler. Consequently, when the schedulability of the system is analyzed, the preemption overhead should be added to the interrupt service time, and thus reducing the preemption overhead is significant for improving the schedulability of the system. Though the preemption overhead becomes zero with the general construction of wait-free operations presented by Herlihy, shared data structures in a real-time kernel is too complex to be implemented efficiently with wait-free operations.<sup>2</sup>

Our approach is to incorporate the realization concept of the general construction of wait-free operations into the locking-based synchronization. In detail, when a task tries to acquire a lock, the operation that the task will execute in the critical section is posted to the waiting queue for the lock. If the turn that the task acquires the lock comes while it is preempted, the posted operation is executed by another processor that is spinning on the lock. With this method, named the SPEPP synchronization (*SPEPP* stands for “*Spinning Processor Executes for Preempted Processors*”), operations on a shared data structure are executed in an FCFS order regardless of the preemptions. If one of the tasks that are waiting for the lock is not preempted, the execution of the operations continues to make progress. If all the tasks are preempted, the execution ceases to progress. It is obvious, however, that the time duration that the execution is suspended is shorter than the duration that any one of the tasks is preempted. Therefore, the extension of the maximum time until the operation of a task is finished is shorter than or equal to the duration that the task is preempted. As the result, the preemption overhead is reduced to zero.

In Section 2, two algorithms of the SPEPP synchronization are presented. We first present the basic version assuming that more than one tasks on a same processor do not acquire a same lock. We then extend the algorithm to remove the assumption. In Section 3, the effectiveness of the SPEPP synchronization is investigated through performance measurements. We have implemented three versions of real-time kernels, two with preemption-safe lockings and one with the SPEPP synchronization, and compared their performance.

---

<sup>2</sup>To the best of our knowledge, no wait-free implementation of operating systems has been reported, while some block-free implementations have [9, 10]

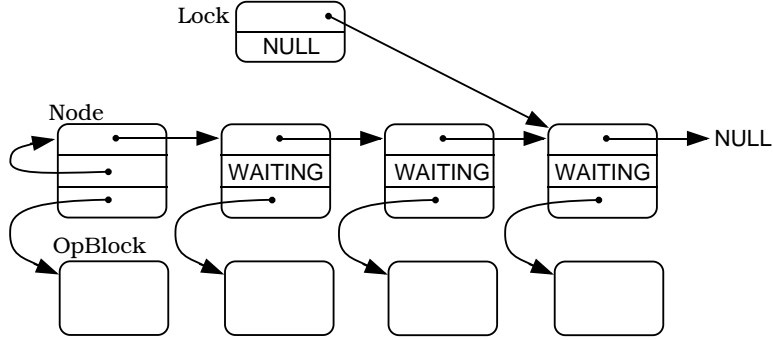


Figure 1: The Data Structures of the SPEPP Synchronization Algorithm

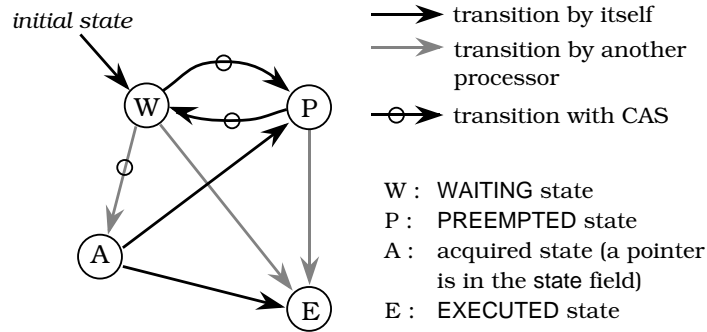


Figure 2: State Transition of a Queue Node

## 2 The SPEPP Synchronization

### 2.1 The Basic Algorithm

In this section, we present the basic version of the SPEPP synchronization algorithm with the restriction that each lock is acquired by at most one task on each processor. In the following, we assume that task switches are triggered by interrupt requests.

Pseudo-code for the algorithm appears in Figure 10, 11, and 12. In these figures, the lines beginning with “//” are comments. The keyword `shared` indicates that only one instance of the variable is allocated and shared in the system. Other variables are allocated for each task and located on the local memory of its host processor.<sup>3</sup> `SWAP` reads the memory addressed by the first parameter, returns the contents of the memory as its value, and atomically writes the second parameter to the memory. `CAS`, the abbreviation of `compare_and_swap`, first reads the memory addressed by the first parameter and compares its contents with the second parameter. If they are equal, the operation writes the third parameter to the memory atomically and returns true. Otherwise, it returns false.

`Lock` and `Node` are the data structures for a mutual exclusion lock and a queue node linked in the waiting queue, respectively. The last field of `Lock` points to the last node in the waiting queue, and the next field of `Node` points to the succeeding node in the queue. When the lock is idle and the waiting queue is empty, the last field is assigned `NULL`. The state field of `Node` includes a pointer to a queue node or one of the three special values (`WAITING`, `PREEMPTED`, and `EXECUTED`), and indicates the current state of the queue node. The `op` field points to the data block of the `OpBlock` type, in which the kind of operation and the parameters passed to it are stored. Also, the return values of the operation are written in this data block (Figure 1).

At first, the task trying to acquire a lock enqueues its node at the end of the waiting queue with a `SWAP` operation (Line 63). When the queue has been empty, the task succeeds in acquiring

<sup>3</sup>We assume that a processor can directly access its local memory without using the shared bus or interconnection network.

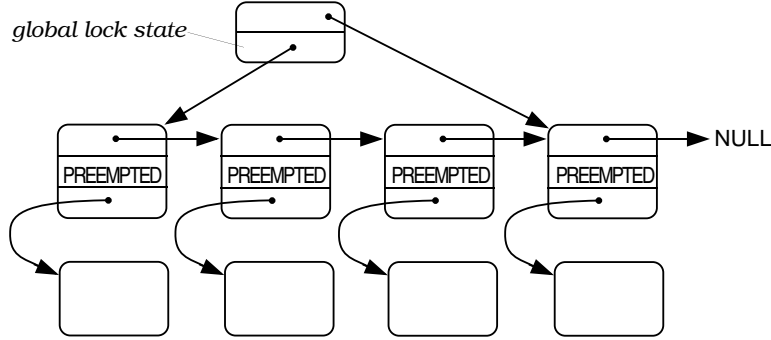


Figure 3: The Global Lock State

the lock. Otherwise, it initializes its state to `WAITING` (Line 69), makes the predecessor point to its queue node (Line 71), and begins to spin. When the task detects an interrupt request while spinning, it changes its state to `PREEMPTED` and services the request (Line 90). After the interrupt service, the task simply recovers its state from `PREEMPTED` to `WAITING` with a CAS operation (Line 82) and resumes to spin. When the lock is handed to the task, the state of the task is changed to the pointer to the head node of the waiting queue. This state is called the acquired state. If the task is at the head of the queue (in other words, if the state field points to its own queue node (Line 127)), the task executes its own operation (Line 142) and then releases the lock (Line 143). Otherwise, the task executes the operation posted to the head node (Line 129), and then proceeds to the succeeding node in the queue (Line 136), until the task's own node becomes the head of the queue (Line 127). When the operation of the task is executed by another processor, its state is changed to `EXECUTED`. In this case, the task simply exits the routine (Line 123). Figure 2 illustrates the state transition of the queue node of a task.

A difficult situation occurs when all the tasks waiting for a lock are preempted. In this case, there are no task to which to hand the lock, but the waiting queue should not be made empty. To cope with this difficulty, an additional field `glock`, called the global lock state, is introduced in the Lock data structure. When all the tasks in the waiting queue are preempted, the global lock state is assigned the pointer to the head node of the queue (Figure 3, Line 49). Each task which has resumed to spin must check the global lock state (Line 97). When a pointer is assigned to it, the task tries to acquire the lock by changing the global lock state to `NULL` atomically from the pointer (Line 108). Here, checking the global lock state only once is not sufficient. This is because at the moment of the checking, the task releasing the lock may be still trying to hand the lock to the succeeding ones. In order to inform that the global lock state must be checked again, a special value `UPDATING` is assigned to it (Line 38). In order to reduce the shared-bus traffic, we also adopt an exponential backoff scheme when checking the global lock state repeatedly (Line 103). Two constant parameters  $\alpha$  and  $\beta$  should be determined for each target hardware. It is also necessary to check the global lock state, just after a task enqueues its queue node at the end of the waiting queue.

With this algorithm, the maximum duration that interrupt services are inhibited can be bounded with  $T$  plus some lock handling overhead, where  $T$  is the maximum execution time of an operation. In order to realize this feature, interrupt requests are probed while spinning on the lock (Line 88). In addition, interrupt requests must be checked after the operation requested by another task is executed (Line 130). If an interrupt request is detected at this moment, the processor hands the lock to another one (Line 132) and services the interrupt request.

On the other hand, the maximum duration that a task spins on the lock can be bounded with  $(n-1) \cdot T$  plus lock handling overhead, where  $n$  is the number of contending processors for the lock, because of the restriction that at most one task on each processor acquires the lock. Consequently, the maximum execution time until a task finishes an operation<sup>4</sup> can be bounded with  $n \cdot T$  plus lock handling overhead regardless of the occurrence of preemptions.

<sup>4</sup>The period during which its host processor executes other tasks is not included in this execution time.

## 2.2 The Extended Algorithm

The basic algorithm still works, even if the restriction on the algorithm is removed. However, the maximum duration that a task spins on the lock becomes as long as  $(N - 1) \cdot T$  plus lock handling overhead, where  $N$  is the number of contending tasks for the lock. Because  $N$  is much larger than  $n$  in general, the maximum execution time until a task finishes an operation becomes very long. In this section, we present an extended algorithm of the SPEPP synchronization, with which the maximum execution time until a task finishes an operation can be regarded as being bounded with  $O(n)$  without the restriction.

Figure 13 presents the differences of the extended algorithm from the basic one. At first, the queue node should be prepared for each lock for each processor (instead of prepared for each task). The task trying to acquire a lock uses the queue node for the lock, even if a preempted task is using the node. When a task is scheduled again after preempted within the algorithm, it should check whether its operation has been executed during the preemption (Line 79.2). We assume that whether an operation is executed or not is recorded in the OpBlock structure and can be checked with the finished function. If the operation has not been executed, then the task checks the state of the queue node (Line 79.6). If it is EXECUTED, the queue node has been stolen by another task, and the former task should retry the operation. Otherwise, the task resumes to spin.

In order to estimate the maximum execution time until a task finishes an operation, we must consider the case when the queue node of a preempted task is stolen by another task. In this case, because the preempted task must re-execute the synchronization algorithm from the beginning, the maximum execution time to finish an operation is prolonged. Instead, the execution time that the preempting task is reduced. In more precise, when task  $\tau_1$  is preempted after spinning on a lock for the duration of  $T_s$ , the task  $\tau_2$  which steals the queue node of  $\tau_1$  can shorten its maximum waiting time for the lock by  $T_s$ . When the schedulability of the system is analyzed, we can safely count this stolen execution time of  $\tau_1$  as a part of the maximum execution time of  $\tau_2$ .

When a task detects that the state of its queue node is changed to EXECUTED, it must check with the finished function whether its operation has been executed (Line 122.1). This is because what has been executed may be the operation requested by another task on the same processor. In other words, a task must possibly wait for an execution of the operation requested by another task on the same processor. As the result, with this extended algorithm, the maximum execution time until a task finishes an operation can be considered to be bounded with  $(n + 1) \cdot T$  plus lock handling overhead. This is a significant improvement over the simple algorithm, when  $n$  is much smaller than  $N$ .

## 3 Performance Measurements

We have implemented a real-time kernel in which shared data structures within the kernel are operated with the extended version of the SPEPP synchronization algorithm presented in the previous section. The application program interface of the kernel is an extension of the  $\mu$ ITRON specification [11].

In this section, the performance of this real-time kernel is measured, and compared with the performance of two other versions of the real-time kernel implemented with preemption-safe lockings; one of them is implemented with the simple preemption-safe version of the MCS lock [5, 6] (called MCS/SPS locking in this section), and the other with the improved preemption-safe MCS lock [7] (called MCS/IPS locking).

### 3.1 Evaluation Environment and Performance Metric

A shared-bus multiprocessor system is used for the measurements. The shared bus is based on the VMEbus specification, and nine processor boards and a shared memory board are connected to the shared bus. Each processor board consists of a GMICRO/200 microprocessor, which is rated approximately at 10 MIPS, 1 MB of local memory, and some I/O interfaces. The local memory can be accessed from other processors through the shared bus. No coherent cache is equipped (Figure 4).

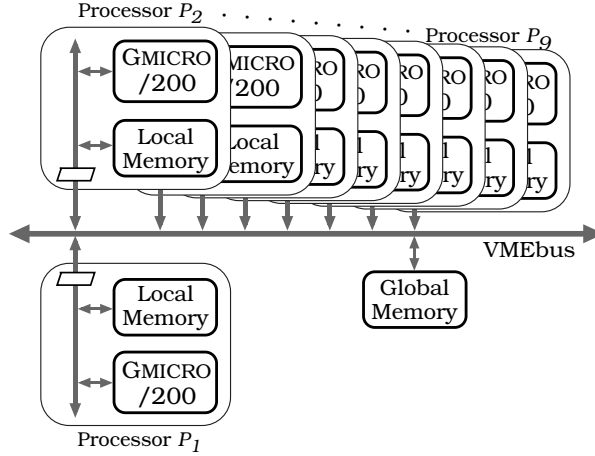


Figure 4: Evaluation Environment

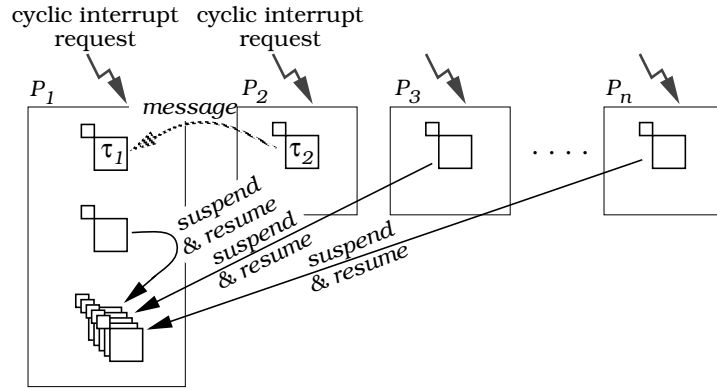


Figure 5: The Workload for the Measurements

The GMICRO/200 microprocessor supports the CAS instruction but does not support the SWAP operation. In the measurements, the SWAP operation is emulated using the CAS instruction and a retry loop. Because the VMEbus has only four pairs of bus request/grant lines, processors are classified into four classes by the bus request line they use. The round-robin arbitration scheme is adopted among classes and the static priority scheme is applied among processors belonging to a same class.

Because the worst-case behavior is the primary concern in real-time systems, the effectiveness of our proposals should be evaluated with its worst-case execution (or response) times. However, worst-case times cannot be obtained through experiments because of unavoidable non-determinism in multiprocessor systems. Moreover, our evaluation environment cannot bound the execution times of various kernel services, because the access time of the shared bus has no upper bound with our environment. Therefore, in place of a worst-case time, we have adopted a  $p$ -reliable time, the time within which a processor finishes to execute (or responds) with probability  $p$ , as the performance metric. In the following section, we use the 99.9%-reliable (i.e.  $p = 0.999$ ) execution (or response) times instead of the worst-case times.

### 3.2 Measurement Method

We have measured the execution times of a system call of the real-time kernels and the interrupt response times using a synthetic workload described below. A task  $\tau_2$  on processor  $P_2$  repeatedly invokes a system call that sends a message to a task  $\tau_1$  on processor  $P_1$ , and the execution times of the system call are measured. The execution times when an interrupt request is serviced during



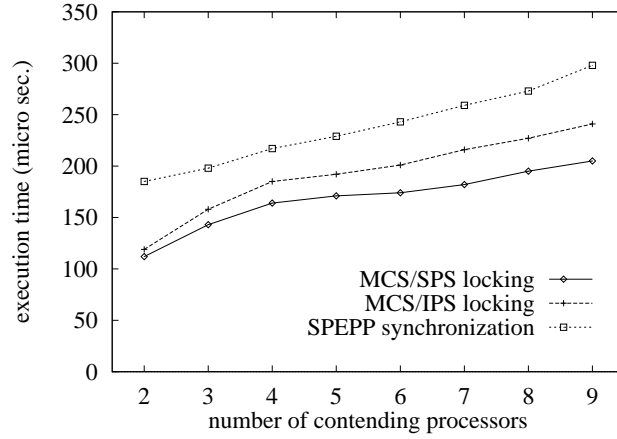


Figure 6: 99.9%-Reliable Execution Times of an Operation (when no interrupt request is serviced)

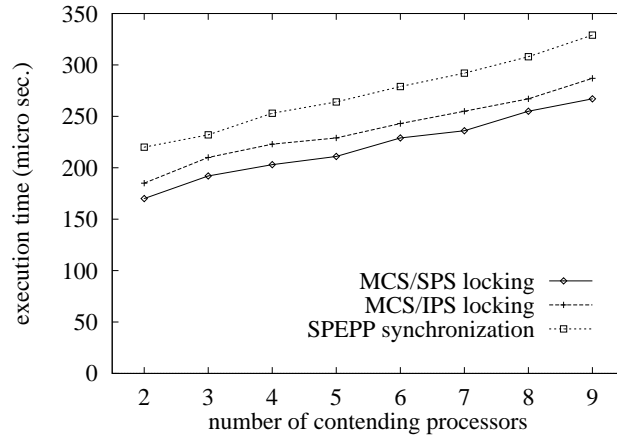


Figure 7: 99.9%-Reliable Execution Times of an Operation (when an interrupt request is serviced)

the execution are recorded separately.

In order to interfere the message sending operation, a task on  $P_1$  and tasks on the other processors alternately suspend and resume the execution of lower priority tasks on  $P_1$  at random intervals. The average interval is about  $1\text{ ms}$ . During the measurement, periodic interrupt requests are raised on each processor, and the interrupt response times are measured within the interrupt handler. The interrupt service time including the time for invoking and returning from the handler is about  $33\ \mu\text{s}$ . The interrupt period is around  $2\text{ ms}$  and is varied in 0–5% for each processor, in order that the timing of interrupt requests for each processor should not be synchronized. The relation of the tasks is illustrated in Figure 5.

### 3.3 Measurement Results

Figure 6 and 7 present the 99.9%-reliable execution times of a message sending operation when no interrupt request is serviced and those when an interrupt request is serviced during the execution, respectively. The number of contending processors (including  $P_1$ ) is changed from two (only one interfering task on  $P_1$ ) to nine (eight interfering tasks). The SPEPP synchronization exhibits worse performance than the preemption-safe lockings, mainly due to the overhead of storing the parameters passed to the operation into the OpBlock structure and reading the return values from

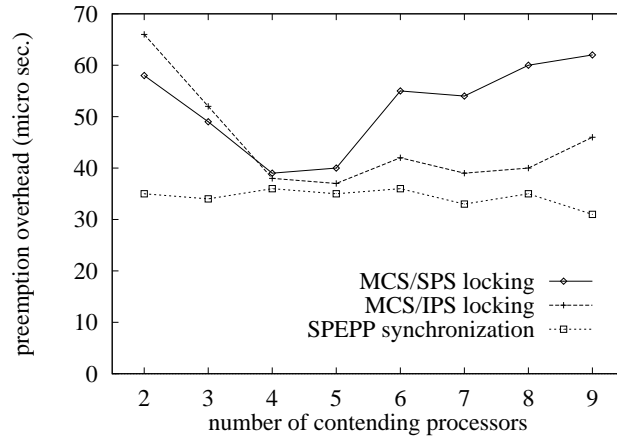


Figure 8: 99.9%-Reliable Preemption Overheads (an interrupt service time is included)

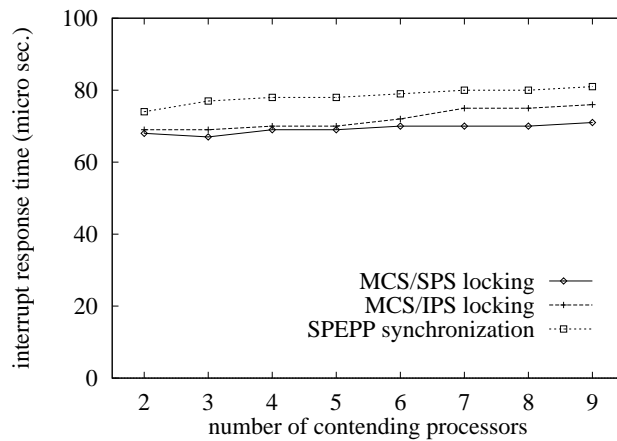


Figure 9: 99.9%-Reliable Interrupt Response Times

it.

Figure 8 presents the differences of these two graphs, i.e. the differences of the 99.9%-reliable execution times of an operation when an interrupt request is serviced and those when no interrupt request is serviced, which represent the measured preemption overheads including an interrupt service time. Considering that the interrupt service time is about  $33 \mu s$ , the preemption overheads with the SPEPP synchronization are almost zero. This demonstrates that the SPEPP synchronization has the expected property.

Finally, Figure 9 presents the 99.9%-reliable interrupt response times on processor  $P_1$ . With either method, the interrupt response times are independent of the number of contending processors. The response times are a little worse with the SPEPP synchronization than the other methods.

## 4 Conclusion and Future Work

In this paper, we have proposed a novel approach to multiprogrammed multiprocessor synchronization, called the SPEPP synchronization, for use in the implementation of a real-time kernel. The SPEPP synchronization has the merits of both of the two conventional approaches to multiprogrammed multiprocessor synchronization: preemption-safe locking and wait-free synchronization. In the concrete, the preemption overhead can be reduced to zero with the SPEPP synchroniza-

tion. It is also possible to operate on complex data structures with reasonable overhead. The two versions of the SPEPP synchronization algorithms have been described, and their timing behavior has been discussed. Finally, it is demonstrated through performance measurements that the SPEPP synchronization has the expected advantages, although its overhead is a little larger than preemption-safe lockings.

The most important future work for us is to extend our algorithm for nested locks, in order to implement a scalable full-fledged real-time kernel for function-distributed multiprocessors. The difficulty in its implementation is described in [8].

Another work to be done is to apply our approach to multiprogrammed multiprocessor synchronization of non-real-time systems. When the maximum execution times need not be bounded, various optimizations are possible. For example, in case of cache coherent multiprocessor systems, it is a promising approach that once a processor acquires a lock, it executes the operations of all the tasks that are waiting for the lock. In case of NUMA architecture, on the other hand, if the processor to which the target shared data structure is local is spinning on the lock, the lock should be handed to the processor and the processor executes for all of the waiting tasks.

## References

- [1] R. W. Wisniewski, L. Kontothanassis, and M. L. Scott, "High performance synchronization algorithms for multiprogrammed multiprocessors," in *Proc. 5th ACM Symposium on Principles and Practice of Parallel Programming*, pp. 199–206, July 1995.
- [2] M. M. Michael and M. L. Scott, "Relative performance of preemption-safe locking and non-blocking synchronization on multiprogrammed shared memory multiprocessors," in *Proc. 11th Int'l Parallel Processing Symposium*, Apr. 1997. (to appear).
- [3] M. Herlihy, "Wait-free synchronization," *ACM Trans. Programming Languages and Systems*, vol. 13, pp. 124–149, Jan. 1991.
- [4] J. M. Mellor-Crummey and M. L. Scott, "Algorithms for scalable synchronization on shared-memory multiprocessors," *ACM Trans. Computer Systems*, vol. 9, pp. 21–65, Feb. 1991.
- [5] H. Takada and K. Sakamura, "A bounded spin lock algorithm with preemption," Tech. Rep. 93-2, Department of Information Science, University of Tokyo, July 1993.
- [6] R. W. Wisniewski, L. Kontothanassis, and M. L. Scott, "Scalable spin locks for multiprogrammed systems," in *Proc. 8th Int'l Parallel Processing Symposium*, Apr. 1994.
- [7] H. Takada and K. Sakamura, "Predictable spin lock algorithms with preemption," in *Proc. Real-Time Operating Systems and Software*, pp. 2–6, May 1994.
- [8] H. Takada, C.-D. Wang, and K. Sakamura, "Issues for realizing a scalable real-time kernel for function-distributed multiprocessors," in *Proc. Work in Progress Session of 17th IEEE Real-Time Systems Symposium*, pp. 23–26, Dec. 1996.
- [9] H. Massalin and C. Pu, "A lock-free multiprocessor OS kernel," Tech. Rep. CUCS-005-91, Department of Computer Science, Columbia University, 1991.
- [10] M. Greenwald and D. Cheriton, "The synergy between non-blocking synchronization and operating system structure," in *Proc. 2nd Symposium on Operating Systems Design and Implementation*, pp. 126–136, Oct. 1996.
- [11] K. Sakamura, ed.,  *$\mu$ ITRON 3.0 Specification*. Tokyo: TRON Association, 1994. (can be obtained from "ftp://tron.um.u-tokyo.ac.jp/pub/TRON/ITRON/SPEC/mitron3.txt.Z").

```

1:      // data structure for a queue node.
2:      type Node = record
3:          next: pointer to Node;
4:          state: pointer to Node;
5:          op: pointer to OpBlock
6:      end;
7:
8:      // data structure for a lock.
9:      type Lock = record
10:         last: pointer to Node;
11:         glock: pointer to Node
12:     end;
13:
14:     // globally shared variable.
15:     shared var L: Lock;    // L.last and L.glock should be initialized to NULL.
16:
17:     // the subroutine to release the lock.
18:     // top is the head node of the queue before the release.
19:     procedure release_lock(lock: pointer to Lock, top: pointer to Node);
20:         var entry, next: pointer to Node;
21:     begin
22:         next := top→next;
23:         if next = NIL then
24:             // tries to make the waiting queue empty.
25:             if CAS(&(lock→last), top, 0) then
26:                 // the waiting queue is empty, now.
27:                 top→state := EXECUTED;
28:                 return
29:             end;
30:             repeat next := top→next until next ≠ NULL
31:         end;
32:         top→state := EXECUTED;
33:
34:         // waits if lock→glock is not NULL.
35:         repeat until lock→glock = NULL;
36:
37:         entry := next;
38:         lock→glock := UPDATING;
39:         repeat
40:             // tries to hand the lock to entry.
41:             if CAS(&(entry→state), WAITING, next) then
42:                 lock→glock := NULL;
43:                 return
44:             end;
45:             // proceeds to the next node in the queue.
46:             entry := entry→next
47:         until entry = NULL;
48:         // no task to which to hand the lock.
49:         lock→glock := top
50:     end;

```

Figure 10: The Basic Algorithm of the SPEPP Synchronization (Part 1)

```

51:      // shared variables for each task.
52:      var I: Node;
53:      var opblock: OpBlock;
54:
55:      // unshared variables for each task.
56:      var pred, top, entry, next: pointer to Node;
57:      var interval, i: integer;
58:
59:      // opblock is assumed to be prepared beforehand.
60:      disable_interrupts;
61:      I.next := NULL;
62:      // inserts I at the end of the queue.
63:      pred := SWAP(&(L.last), &I);
64:      if pred = NULL then
65:          // when the queue has been empty.
66:          goto acquired
67:      end;
68:      // when the queue has not been empty.
69:      I.state := WAITING;
70:      I.op := &opblock;
71:      pred→next := &I;
72:      goto spin;
73:
74:  preemption:
75:      // branches here after assigning PREEMPTED to I.state
76:      //          when an interrupt request should be serviced.
77:      enable_interrupts;
78:      // interrupt requests are serviced here.
79:      disable_interrupts;
80:      // tries to recover I.state from PREEMPTED to WAITING.
81:      // does not recover I.state if it is already EXECUTED.
82:      CAS(&(I.state), PREEMPTED, WAITING);
83:  spin:
84:      i := 1; interval :=  $\alpha$ ;
85:      entry := I.state;
86:      // spins while I.state is WAITING.
87:      while entry = WAITING do
88:          if interrupt_requested then
89:              // tries to change I.state from WAITING to PREEMPTED.
90:              if CAS(I.state, WAITING, PREEMPTED) then
91:                  goto preemption
92:              end
93:          end

```

Figure 11: The Basic Algorithm of the SPEPP Synchronization (Part 2)

```

94:         i := i - 1;
95:         if i = 0 then
96:             // checks the global lock state.
97:             top := L.glock;
98:             if top = NULL then
99:                 i := ∞
100:            else if top = UPDATING then
101:                // must check the global lock state again
102:                // with exponential backoff.
103:                i := interval; interval := interval × β
104:            else begin
105:                entry := I.state;
106:                if entry ≠ WAITING then
107:                    break
108:                else if CAS(&(L.glock), top, NULL) then
109:                    // succeeds in acquiring the global lock.
110:                    entry := top;
111:                    goto execute
112:                else
113:                    i := ∞
114:                end
115:            end
116:        end;
117:        entry := I.state
118:    end;
119:
120:    // I.state (= entry) is not WAITING, now.
121:    if entry = EXECUTED then
122:        // my operation has been executed by another processor.
123:        goto finish
124:    end;
125:    execute:
126:        // succeeds in acquiring the lock.
127:        while entry ≠ me do
128:            // executes the operation posted to entry.
129:            execute(entry→op);
130:            if interrupt_requested then
131:                I.state := PREEMPTED;
132:                release_lock(lock, entry);
133:                goto preemption
134:            end;
135:            // proceeds to the next node in the queue.
136:            repeat next := entry→next until next ≠ NULL;
137:            entry→state := EXECUTED;
138:            entry := next
139:        end;
140:    acquired:
141:        // executes my operation, now.
142:        execute(&opblock);
143:        release_lock(lock, me)
144:    finish:
145:        enable_interrupts

```

Figure 12: The Basic Algorithm of the SPEPP Synchronization (Part 3)

```

51:      // shared variables for each lock for each processor.
52:      var I: Node;
52.1:    // shared variables for each task.
53:      var opblock: OpBlock;

59:      // opblock is assumed to be prepared beforehand.
60:      disable_interrupts;
60.1:  retry:
61:      I.next := NULL;
62:      // inserts I at the end of the queue.
63:      pred := SWAP(&(L.last), &I);

74:  preemption:
75:      // branches here after assigning PREEMPTED to I.state
76:      //      when an interrupt request should be serviced.
77:      enable_interrupts;
78:      // interrupt requests are serviced here.
79:      disable_interrupts;
79.1:  entry := I.state;
79.2:  if finished(&opblock) then
79.3:      // opblock has already been executed.
79.4:      goto finish
79.5:  end;
79.6:  if entry = EXECUTED then
79.7:      // the queue node has been stolen.
79.8:      goto retry
79.9:  end;
79.10: // steals the queue node (if another task is using it).
79.11: I.op := &opblock;
80:      // tries to recover I.state from PREEMPTED to WAITING.
81:      // does not recover I.state if it is already EXECUTED.
82:      CAS(&(I.state), PREEMPTED, WAITING);
83:  spin:

120:     // I.state (= entry) is not WAITING, now.
121:     if entry = EXECUTED then
122:         // the operation has been executed by another processor.
122.1:     if ¬finished(&opblock) then
122.2:         // executed operation has been the request of another task.
122.3:         goto retry
122.4:     end;
122.5:     // executed operation has been my request.
123:     goto finish
124:     end;
125:  execute:

```

Figure 13: The Extended Algorithm of the SPEPP Synchronization