

The Performance Implications of Thread Management Alternatives for Shared-Memory Multiprocessors

THOMAS E. ANDERSON, EDWARD D. LAZOWSKA, AND HENRY M. LEVY

Abstract—Threads (“lightweight” processes) have become a common element of new languages and operating systems. This paper examines the performance implications of several data structure and algorithm alternatives for thread management in shared-memory multiprocessors. Both experimental measurements and analytical model projections are presented.

For applications with fine-grained parallelism, small differences in thread management are shown to have significant performance impact, often posing a tradeoff between throughput and latency. Per-processor data structures can be used to improve throughput, and in some circumstances to avoid locking, improving latency as well.

The method used by processors to queue for locks is also shown to affect performance significantly. Normal methods of critical resource waiting can substantially degrade performance with moderate numbers of waiting processors. We present an Ethernet-style backoff algorithm that largely eliminates this effect.

Index Terms—Locking, multiprocessor, parallel computing, parallel software, performance, thread.

I. INTRODUCTION

THE purpose of this paper is to study the performance implications of thread management alternatives for shared-memory multiprocessors.

In traditional operating systems, a process, consisting of a single address space and a single thread of control within that address space, is used to execute a program. Within the process, program execution entails initializing and maintaining a great deal of state information. For instance, page tables, swap images, file descriptors, outstanding I/O requests, and saved register values are all kept on a per-program, and thus per-process, basis. The sheer volume of this information makes processes expensive to create and maintain.

Manuscript received March 1, 1989; revised July 15, 1989. This material is based on work supported by the National Science Foundation (Grants CCR-8619663, CCR-8703049, and CCR-8700106), the Naval Ocean Systems Center, U S WEST Advanced Technologies, the Washington Technology Center, and Digital Equipment Corporation (the Systems Research Center and the External Research Program). A preliminary version of this paper appeared in the Proceedings of the 1989 ACM SIGMETRICS and Performance '89 International Conference on Measurement and Modeling of Computer Systems, *Performance Evaluation Review*, vol. 17, no. 1 May 1989.

The authors are with the Department of Computer Science and Engineering, University of Washington, Seattle, WA 98195.
IEEE Log Number 8931176.

Threads, or “lightweight” processes, separate the notion of execution from the rest of the definition of a process. A single thread executes a portion of a program, cooperating with other threads concurrently executing within the same address space. Like processes, every thread must have a separate program counter and stack of activation records, describing the state of its execution. However, much of what is normally kept on a per-process basis can be maintained in common for all threads executing in a single program, with dramatic reductions in overhead.

Thread packages have become a common element of new languages and operating systems for both uniprocessor and multiprocessor architectures. Mach [1], Topaz [23], Psyche [21], DYNIX [22], and several extensions to UNIX [5], [11] are examples of operating systems that provide explicit support for concurrent or parallel execution of programs. Ada [19], CSP [12], Presto [7], Mesa [16], Concurrent Euclid [13], and Emerald [14] evidence equal interest within the language community.

On uniprocessors, threads are used as a program structuring aid or to overlap I/O with processing. The metric of goodness for these thread management implementations is simply processing cost per thread creation or context switch. No locking is needed inside thread routines, since only one routine can be executing at any one time.

Programs on multiprocessors use threads to exploit parallelism. The speedup achievable by any given application depends on the availability of thread management routines that provide low-cost facilities that are not a serial bottleneck. In Sequent's DYNIX operating system, for example, applications must use normal UNIX-like processes for parallelism [22]. Since process creation in DYNIX takes over 25 ms, only very coarse-grained parallelism can be exploited. As another example, the Topaz kernel provides relatively inexpensive thread creation and synchronization, but the routines are protected by a single lock [23]. While this may be appropriate for architectures with small numbers of processors, as the number of processors increases, the single lock could limit speedups for applications with fine-grained parallelism.

Our initial experience in the area of high-performance thread packages was with Presto, an application-level run time library that relies on the kernel only for processor allocation and memory management [7]. This work showed that there is

an order of magnitude performance advantage to using threads instead of DYNIX processes for exploiting parallelism. Drawing on this experience, we implemented a thread package that is, in turn, another order of magnitude faster than Presto. This basic package was then modified to implement each alternative we wanted to explore.

One consequence of the speed of our basic thread package is that small changes in the organization of data structures and locks have a significant impact on performance. Often, the choice involves a tradeoff between latency and throughput. Per-processor data structures can sometimes be used to avoid locking, however, improving latency and throughput at the same time.

Another consequence of the speed of our thread package is that its performance depends noticeably on the algorithm used to queue for locks. Earlier, we studied the relative performance of spinning and blocking locks [26]. In general, a thread that tries to acquire a lock that is already held can either spin ("busy-wait") until the lock is released, or relinquish the processor. However, within the thread management routines themselves, spinning is the only option. Thus, blocking at the user level may require spin-waiting in thread management routines. Spin-waiting has a cost not only to the processor waiting for the lock, but also to processors doing useful work. The degradation of other processors becomes substantial for moderate numbers of waiting processors, especially for small critical sections. We present an Ethernet-style backoff algorithm that largely eliminates this effect.

The following sections describe these issues in more detail. In Section II, we present an abstraction of a thread package: its objects, resources, and operations. Section III outlines the strategies for thread management that we examined and presents measurements of their relative performance. Section IV compares methods of queuing for locks. Section V combines these results in an analytical model. Section VI summarizes our experiences.

II. AN ABSTRACT THREAD PACKAGE

As noted in Section I, threads gain efficiency by separating the notion of execution from the rest of the definition of a process. The data structures needed by each thread are a program counter, a stack, and a control block. (The control block contains state information needed for thread management. Through the control block, the thread can be put onto lists and other threads can synchronize with it.) Another important data structure is the ready queue, which lists threads that are ready to run. Lampson and Redell [16] provide a good description of the functionality of a uniprocessor thread package.

Thread operations are shown in Table I. Creating a thread can be viewed as calling a procedure, except that the callee can execute in parallel with the caller. In both cases, the caller specifies a place to begin executing and some number of arguments. In fact, thread creation and startup is semantically equivalent to an asynchronous procedure call.

As Table I shows, a program can create a thread even if there is no idle processor available to run it. Because the parallelism cannot be immediately exploited in this case, it

TABLE I
THREAD OPERATIONS

Thread Creation	Allocate and initialize a control block, saving the initial PC. Allocate a stack and copy in the thread's arguments. Place the new thread on the ready queue.
Thread Startup	Remove the thread from the ready queue and begin to execute it.
Thread Block (wait on a blocking lock, monitor condition variable, or message)	Save register values and PC on the thread's stack. Place the thread on the condition queue for the event. Look for a thread in the ready queue, and start or resume it.
Signal a Blocked Thread	Remove the thread from the condition queue. Place the thread on the ready queue.
Thread Resume	Remove the thread from the ready queue. Restore registers. Continue executing it from the saved PC.
Thread Finish	Deallocate the stack and control block. Look for a thread in the ready queue, and start or resume it.

might seem that the overhead of thread creation should be avoided. The program may run faster by creating the thread, however, if at some future time there will be an idle processor that can be used to execute the thread. This idea of creating parallelism for future use is very powerful. Unfortunately, in the above framework, its space cost is prohibitive. Each thread must be initially allocated a large amount of space for its stack, since it is expensive to dynamically expand the space if the thread later runs out of it. In Table I, the thread is allocated space for a stack when it is created, but the space is largely wasted until the thread is actually started. Using virtual memory could remove the need to allocate physical memory to back the stack space until the thread begins to run; however, allocating extra virtual memory is itself expensive.

An important optimization to Table I, therefore, is to copy a thread's arguments into its control block when the thread is created. This way, the stack need not be allocated until thread startup; the arguments can be copied from the control block to the stack at that time. WorkCrews [24] and Presto [7] both take this approach.

Another important optimization is to store deallocated control blocks and stacks in free lists [7]. If these data structures were individually located out of the heap, thread overhead would include the cost of finding a free block of the correct size as well as possibly coalescing the block when it is returned to the heap. By using free lists, both allocation and deallocation can normally be simple list operations.

We begin our study by assuming these optimizations. For simplicity, we will focus on the effect of thread management alternatives on the performance of only a few thread operations: creation, startup, and finish. These operations manipulate each of the three shared data structures: the ready queue, the stack free list, and the control block free list. Most of the discussion applies as well to the performance of block and resume operations.

III. THREAD MANAGEMENT ALTERNATIVES

In a parallel environment, access to shared data structures must be serialized to ensure consistency and correctness. Our thread package uses spin-locks for this purpose: when a processor tries to modify a data structure, it must first lock it to obtain exclusive access; if some other processor already holds the lock, the processor loops until the lock is released.

Locking implies dual concerns of latency and throughput [15]. Latency is the cost of thread management under the best case assumption of no contention for locks. Throughput, on the other hand, is the rate at which threads can be created, started, and finished when there is contention. If part of thread management must be done serially, then no matter how many processors work on a problem, there will be some maximum rate at which thread operations can be performed.

There are several ways of defining latency, with different implications for different types of applications. If an application keeps all of its processors continually busy, for instance by creating threads before they are needed, then any time spent in creating, starting, or finishing a thread is time that could have been spent doing other useful work. When a thread finishes, however, if there is no other work for the processor to do, the time spent deallocating the thread's data structures is unimportant. Instead, the relevant issues include how much a creating processor is delayed, since it has a thread to run, and how much time it takes for the created thread to begin running on a processor.

In the following subsections, we define five alternative thread management strategies, and describe some of the potential advantages and disadvantages of each approach. We then provide measurement and analytical comparisons of these alternatives.

A. Single Lock: Central Data Structures Protected by a Single Lock

The most obvious approach to thread management is to protect all data structures under a single lock. Once the lock is acquired by a processor, the processor is assured that it can modify any stored state. To perform a thread operation, a processor must first acquire the lock, then do what is needed to the shared data, and finally release the lock when done. In this way, only a single lock is needed per thread operation, but, since most of the thread management path is serialized, throughput is limited. In the typical scheme, idle processors loop checking the ready queue for work to do, causing useless contention for the ready queue lock; however, this can be avoided if idle processors check that the ready queue is not empty before acquiring the lock. (Ni and Wu [20] present a different approach.)

B. Multiple Locks: Central Data Structures Each Protected by a Separate Lock

A somewhat more modular approach to locking is to separately protect each data structure with its own lock [16]. Each operation on the data structure can then be surrounded by a lock acquisition and release. For thread management, this involves separately locking each enqueue and dequeue operation

on the ready queue, stack free list, and control block free list, the three shared data structures.

There is a basic tradeoff between latency and throughput in the choice between using a single lock or multiple locks in protecting shared data structures [15]. Since less of the total thread activity is in a critical section, and since it is split among several locks, the maximum rate of thread operations is higher with multiple locks than with a single lock. There is a cost to this increased throughput, however, more lock accesses are needed, increasing latency.

C. Local Free List: Per-Processor Free Lists without Locks; A Central Locked Ready Queue

One way of avoiding locking is to maintain as much state as possible locally, with each processor. If each processor maintains its own free lists of control blocks and stacks, these structures need not be locked, since only one processor will ever access them. As before, there is a single shared ready queue whose accesses are locked.

The tradeoff between latency and throughput can be largely avoided by using local free lists. Since fewer lock acquisitions are needed per thread, latency is lower than with multiple locks, yet since only accesses to the ready queue are serialized, throughput is better.

Local free lists need to be balanced. Control blocks and stacks can migrate between free lists if the thread is created or started on one processor and finished on another. This can happen, for instance, if a thread blocks and is resumed on a different processor. Thus, one free list can be empty, requiring the processor to obtain more space from the heap, while another free list has many entries. In the worst case, some processors only create and start threads (allocate structures), while other processors only finish them (deallocate structures). Without balancing, the deallocated structures are never reused; a separate stack and control block are needed for every thread. In contrast, with a centralized free list, only as many are needed as there are active (created or started, but not finished) threads.

It is inexpensive, however, to balance free lists by using global pool and a threshold T on the maximum size of each list. When the size of a free list reaches the threshold, half the list can be returned to the global pool; when a free list empties, $T/2$ entries can be claimed from the pool. The global pool must be locked, of course. For efficiency, it can be organized as a list of lists. The processing cost to balancing is thus one locked pool access amortized across at least $T/2$ free list accesses. Let P be the number of processors. An application using balanced local free lists will use no more than $O(P \times T)$ more space than one using a central list, even if one processor only allocates structures that other processors deallocate. The worst case occurs when the allocating processor's free list and the global pool are empty even though the other local free lists are almost full.

Thus, local free lists trade space for time. This tradeoff is practical for control blocks. Utilization of the pool lock is at most $O(P \times (R/T))$, where R is the rate of thread creation on a single processor. To ensure that the pool lock is not a source of contention (which would inflate the overhead per free list

access), we can set the threshold T to be equal to P . Control blocks are relatively small objects (in our implementation, roughly 100 bytes); provided P is not excessively large, using $100P$ bytes per processor is not onerous. If P is large, then a tree of pools could be used to limit the cost to balancing to $O(\log P)$ bytes per processor.

The tradeoff is not practical for stacks, however. Stacks are at least two orders of magnitude larger than control blocks. Even if sufficient memory were available, using that memory entails processing costs for initializing page tables and increased cache miss rates that could easily overwhelm the advantage gained from decreased locking. Instead, we let the local stack free lists contain at most one element. In this way, stacks need be allocated from the global pool only when a processor blocks a thread and then starts up a different thread, and deallocated only when a processor finishes a thread and then resumes another thread.

D. Idle Queue: A Central Queue of Idle Processors; Per-Processor Free Lists

None of the algorithms described so far exploit parallelism in creating threads. The creating processor allocates and initializes the control block; only when it is done is the starting processor allowed to allocate and initialize the stack. The cost of thread creation could be reduced if some of the work was done by idle processors in parallel with the creating processor.

In addition to a central queue of threads, we can maintain a central queue of idle processors. When there is a backlog of ready threads, there is no point to attempting parallel thread creation since all processors are already doing useful work. When a processor becomes idle and there is no backlog, it preallocates a control block and stack, puts itself on the idle queue, and spins on a local flag waiting for work. Thread creation then dequeues the idle processor, initializes the pre-allocated control block and stack, and sets that processor's flag, indicating that it now has a thread that is ready to run. Instead of processors searching for work, work searches for processors.

In fact, this approach does not alter the essentially sequential nature of thread creation. The idle processor must first queue itself before the creating processor can dequeue it, which in turn must set the flag before the idle processor can start running the thread. The critical path between the beginning of thread creation and when the thread starts running is reduced by doing some of the work (allocating structures, acquiring a lock, enqueueing) before the critical path begins. Since this adds complexity, and there is no benefit in the absence of idle processors, the effect is to trade off reduced latency when there are idle processors for increased latency when all processors are busy. Maximum throughput should be unchanged since two locked queue operations are still needed per thread life cycle. Wagner *et al.* [25] describe a different way of using idle processors to avoid work during blocking and resuming.

E. Local Ready Queue: Per-Processor Ready Queues; Per-Processor Free Lists

Once free lists are made local, the ready or idle queue lock can become a serial bottleneck as the number of processors or

the rate each processor schedules work increases [9]. One way of increasing throughput is to divide the load on a single lock among several locks. An application of this idea is to keep a ready queue per processor. In this way, enqueueing and dequeueing threads can occur in parallel, with each processor using a different queue. There is again a tradeoff between latency and throughput in the choice between using one or more ready queues.

Unlike the case of control block free lists, unlocked local ready queues are inefficient even if balanced through a global pool. Runnable threads are a scarce resource. An idle processor might have an empty queue, yet a ready thread that the processor could run is in some other processor's queue, while the global pool is empty. Performance can be arbitrarily bad in any scheme where a processor can be idle indefinitely while there is even one ready thread in some other queue. For instance, suppose P identical threads are created, but due to an imbalance, only $P - 1$ are started while one processor idles. The run time would then be twice as long as with any of the centralized queueing strategies.

One simple way of avoiding indefinite idling is to lock each ready queue; each idle processor can then scan the ready queues for work, beginning with its own [9]. If there is a ready thread, an idle processor will eventually find it. Processors can queue created threads locally, since balancing is achieved by idle processors. However, contention can still occur if a single processor enqueues (during create or signal) every ready thread; that processor's queue would operate as if it was a central ready queue, except that idle processors would have to waste time scanning for it. A simple way of avoiding this situation is to randomly choose a queue for each newly ready thread.

If each queue is equally likely to get a new ready thread, latency is bad when the number of runnable threads is near to the number of processors. There are two cases. Consider the cost of scheduling a thread onto a newly idle processor. If there are no ready threads, there is effectively no cost until a new thread is created. If there are ready but not running threads, any time spent finding a thread to run could have been spent running that thread. This time is small when there are many ready threads, because the idle processor will find the thread after scanning only a few queues; however, when there is only a single ready but not yet running thread, the processor will have to examine on average half of the queues in order to find it. The cost of scheduling a newly created thread onto an idle processor is similar: the thread will be found quickly if there are many idle processors and more slowly if there are only a few.

One reason to have a one-to-one correspondence between processors and ready queues is to maintain locality. There is an application-specific cost to migrating a newly resumed thread from the processor it last ran on, due to increased cache misses. The ability to avoid migration depends on the backlog of ready threads [10].

If maintaining locality is not important, then there is a tradeoff between latency and throughput in choosing the number of queues [20]. Up to some maximum, throughput is higher with more queues, but the number of queues that must be scanned

to find work, and thus the latency, is also higher. We set the number of queues equal to the number of processors for all measurements.

F. Measurement Results

To validate our intuitions about the relative merits of the alternative approaches, we implemented each on a Sequent Symmetry Model A shared-memory multiprocessor. All code was written in C and compiled with Sequent's standard compiler, with the exception of the locking and context switching code, which was programmed in assembler. Our Symmetry has 20 Intel 80386 processors, a shared bus, and a write-through cache coherency protocol [17]. The Symmetry has a timer with microsecond resolution that was used for all measurements. Table II contains times for sample Symmetry operations.

For all measurements, free lists were "warm started:" sufficient control blocks and stacks were preallocated for use by the benchmark. Our purpose was to measure the relative merits of each alternative, rather than the efficiency of the underlying memory management system. The cache was not warm-started, but we ran each benchmark long enough for this effect to become insignificant.

Fig. 1 is the principal performance comparison: it shows the elapsed time in seconds for each thread management alternative to create, start, and finish one million "null" threads, for varying numbers of processors. The one processor case shows the latency for a single thread in microseconds when there is no contention for locks.

Table III lists the code for this test. Initially, P threads are created; each recursively creates a thread then finishes, allowing that processor to start up one of the waiting threads. The test terminates when each processor has executed 1 million/ P threads; in practice, we found that the processors all completed at roughly the same time. For the multiple ready queue alternative, each newly created thread was added to a random queue to avoid biasing the results with the effect of locality. This test is not intended to be representative of a real parallel program, but it does expose the tradeoffs between the five alternatives. ("numberOfThreads" is a location private to each processor, initially set to 0.)

Fig. 2 shows the inverse graph for the same test: speedup as a function of the number of processors. We define speedup to be the ratio of the time for the fastest alternative on one processor (single lock) to the time each alternative takes on P processors. Speedup is proportional to throughput; this graph shows the maximum number of thread creates, starts, and finishes that can be performed in parallel for each alternative.

Before examining the relative performance of the five alternatives, we note that each of them has quite good performance. Threads are only an order of magnitude more expensive than a procedure call, and 500 times less expensive than normal DYNIX processes. Threads in Presto cost 600 μ s on the same Symmetry hardware, an order of magnitude worse than our threads although an order of magnitude better than DYNIX processes.

While Presto's speedup relative to DYNIX is due to using threads instead of processes, our speedup relative to Presto

TABLE II
RUN TIMES FOR SYMMETRY OPERATIONS (MEASURED)

Operation	Runtime (μ sec.)
Acquire and release a lock	5.6
Procedure call with no arguments	3.6
Each 4-byte argument	1
Iteration of null loop	2.5

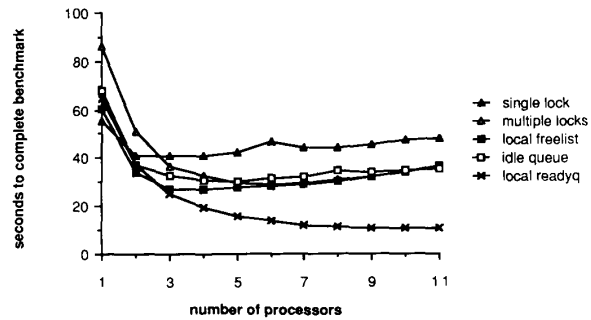


Fig. 1. Principal results for thread management—elapsed time to create, start, and finish 1 000 000 null threads (measured).

is due to attention to implementation details. We implemented Presto in C++; while this enhanced its ability to be modified [8], its C++ was first preprocessed into C, then compiled. This resulted in much less efficient code than could be achieved by direct coding in C. Another factor is that we stripped thread control blocks of all nonessential state, reducing the cost of initialization dramatically. We did not remove functionality: our thread package could be given Presto's user interface without sacrificing its performance.

Because our threads are inexpensive, the choice of alternatives has a large relative impact on both latency and throughput for applications with fine-grained parallelism. Specifically:

- Adding even a single lock acquisition into the thread management path can increase latency significantly. Locking each of the data structures separately results in a much higher latency than locking all data structures under the same lock. Using per-processor data structures to avoid locking is thus crucial to decreasing latency without sacrificing throughput.

- Additional complexity results in a noticeable increase in latency. There are on the order of 100 instructions in the thread management path; adding even a few extra instructions impacts performance. For example, the idle queue strategy checks for idle processors on thread creation. If the idle queue is always empty, as in the measurements in Fig. 1, it defaults to a normal ready queue. Even this simple check markedly increases in cost of threads. This implies that thread management routines must be kept simple; enhancements that would otherwise seem plausible but add complexity are unlikely to work, since there is little computation to save, and it is easy to swamp the savings with increased overhead.

- A large portion of the thread management path is locked, since little work is required beyond manipulation of shared data. When all data are kept under a single lock, throughput is limited by contention for this lock. However, even with local free lists, the lock on the ready queue limits throughput to

TABLE III
BENCHMARK: CREATE, START, AND FINISH 1 MILLION NULL THREADS

```

ThreadCycle () {
    numberOfThreads = numberOfThreads + 1;
    if (numberOfThreads == 1000000 / numberOfProcessors)
        Synchronize(); /* wait for all processors to reach here */
    else
        ThreadCreate (ThreadCycle);
}

BeginTimer ();
for (i = 1; i < numberOfProcessors; i++)
    ThreadCreate (ThreadCycle);
ThreadCycle ();
StopTimer ();

```

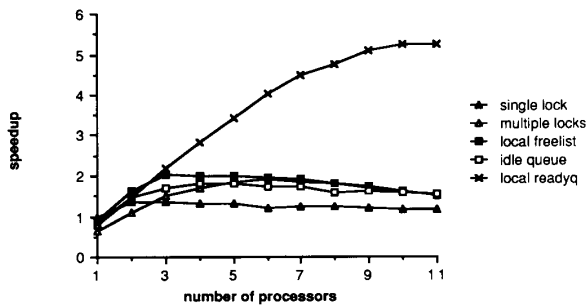


Fig. 2. Speedup to create, start, and finish 1 000 000 null threads (measured).

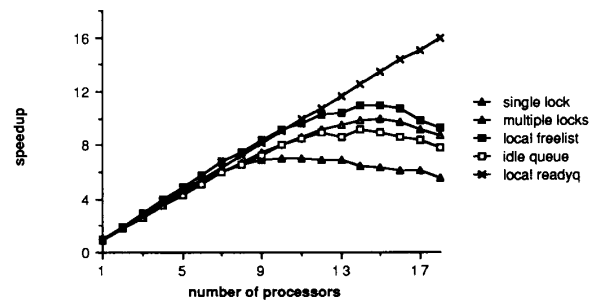


Fig. 3. Speedup, user work = 300 μ s (measured).

a few concurrent thread operations. Only local ready queues can support higher throughput.

- When lock contention is not a problem, the bandwidth of the bus limits thread management throughput. The curve in Fig. 2 levels out for the local ready queue alternative, even though there is no significant contention for locks. While the high bus demand per thread may be specific to the write-through cache protocol on the Symmetry, bus contention is likely to be a problem on any bus-structured shared-memory system.

In Figs. 1 and 2, threads do no work except to create other threads. It is natural to ask whether the performance implications of the thread management alternatives would still be significant in the presence of user-mode computing. We modified the test in Table III so that each thread performs an average of 300 μ s of user work, taken from a uniform distribution. This simulates the behavior of an application with fine-grained parallelism.

Fig. 3 graphs speedup for this modified benchmark. Differences appear as the number of processors increases.

Fig. 4 graphs thread cost in microseconds as a function of the number of runnable threads (parallelism). Thread cost was directly measured by taking timestamps before and after each thread was created and whenever a thread was started or finished. Multiple creations and completions were measured and averaged to improve accuracy; they were synchronized to avoid measuring lock contention.

There are a few things to note in Fig. 4:

- The curve for the central ready queue alternative jumps when the number of runnable threads reaches the number of

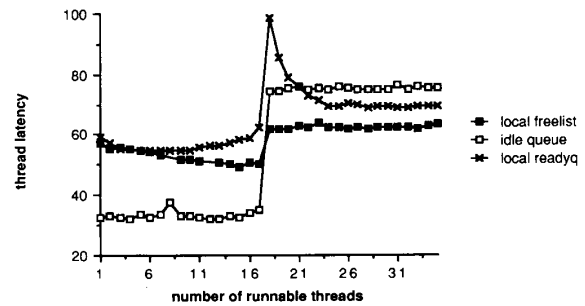


Fig. 4. Thread latency (μ s) versus number of runnable threads, 18 processors (measured).

processors because of a change in the definition of thread latency. When there are fewer threads than processors, thread cost is taken to be the time to create and start running a new thread. The time to finish a thread is unimportant if the idling processor has no work to do. When there are as many or more runnable threads as processors, the cost is the sum of the time to create a thread plus the time to finish it and start a new thread. The thread latency reported in Fig. 1 with one processor corresponds closely to the latency reported in Fig. 4 when there are more runnable threads than processors.

- Using an idle queue is faster when there are idle processors, but slower when there are more runnable threads than processors. Thread creation is faster if an idle processor can be used to do work before the thread is created, but checking the idle queue incurs overhead even if it is not used. Whether a particular application will run faster with an idle queue depends on how much time it spends in each case.

- The spike in the curve for per-processor ready queues shows that finding a ready thread among many queues is expensive when the parallelism of the application is near to the number of processors, but the expense fades when more ready threads or more idle processors are available. Since the height of the spike grows linearly with the number of queues, latency increases in proportion to the maximum throughput. However, the amount of user computing per thread needed to avoid contention for a central ready queue also grows linearly with the number of processors. Thus, the cost of searching for work among per-processor ready queues as a fraction of the time to do that work is not large unless the multiple ready queues are needed for additional throughput.

One area of further research is to examine hybrid thread management strategies to combine the advantages of some of the alternatives we have presented. For example, a per-processor version of the central idle queue could exploit locality and parallelize thread creation without compromising throughput. As another example, both central and per-processor ready queues could be used, by placing threads in a local queue if the lock on the central queue is busy. The drawback to any such approach is that complexity adds cost which may outweigh any benefits.

G. Analytical Explanation of Fig. 4

We now derive a formula that explains in detail the spike for the per-processor ready queue alternative in Fig. 4. When there are idle processors, we need to know the time between the queuing of a ready thread and the dequeuing of that thread by an idle processor; when there is a backlog of ready threads, we need to know how long it takes a newly idle processor to find one of the threads.

Let $E(r, q)$ be the expected number of queues examined by a newly idle processor to find one of r ready threads, which are randomly distributed among q queues. Without loss of generality, let the queues be numbered from 1 to q , let threads be numbered from 1 to r , let i_j be the queue containing the j th thread, and let the idle processor begin searching with queue 1. The idle processor must examine the number of queues equal to the lowest numbered nonempty queue. The number of ways of putting r threads into q queues is q^r .

$$E(r, q) = \frac{1}{q^r} \sum_{i_1=1}^q \sum_{i_2=1}^q \cdots \sum_{i_r=1}^q \text{minimum of } (i_1, i_2, \dots, i_r).$$

We can separately sum when each i_j is the minimum. When more than one thread is at the minimum, we count the value once in the sum for the least numbered thread. Thus, the value of i_j is counted whenever for all $k < j$, $i_k > i_j$, and for all $k > j$, $i_k \geq i_j$. There are $(q - i_j)^{j-1}$ values for the i_k , $k < j$ that satisfy the first condition and $(q - i_j + 1)^{r-j}$ values for the i_k , $k > j$ that satisfy the second.

$$E(r, q) = \frac{1}{q^r} \left(q + \sum_{j=1}^r \sum_{i=1}^{q-1} i(q-i)^{j-1}(q-i+1)^{r-j} \right) \quad (3.1)$$

By symmetry, (3.1) also holds when there are more processors than runnable threads. Let r be the number of idle

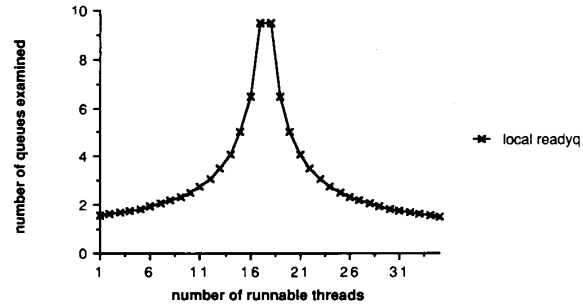


Fig. 5. Queues examined versus number of runnable threads, 18 processors (3.1).

processors, let i_j be the queue currently scanned by the j th idle processor, and let the newly created thread be put into queue 1. Then the processor that actually dequeues the thread will have to look through $E(r, q)$ queues, after the thread is queued, in order to find it.

Fig. 5 graphs (3.1) for 18 processors. In order to correspond to Fig. 4, the x -axis is the number of runnable threads, rather than the number of ready but not running threads or the number of idle processors. Since part of the spike in Fig. 4 is due to the difference in the measurements when there are idle processors or not, the curves in Figs. 4 and 5 correspond well.

The above analysis assumes that events occur one at a time. Since finding a ready thread among a number of queues can take a nontrivial amount of time, it is reasonable to consider what happens when another thread is created or another processor becomes idle during the interim. Suppose another thread is created before a newly idle processor finds one of the r ready threads. Let C be the cost (in number of queues examined) of finding a thread in this situation. C can be no better than if the new thread had been there all along and no worse than if the new thread is ignored. In other words, $E(r+1, q) \leq C \leq E(r, q)$. Similarly, if another processor becomes idle in the interim, provided $r \geq 2$, the combined cost for both processors to find threads is $E(r, q) + E(r-1, q)$, assuming the processors do not contend for the same queue, independent of which processor finds a ready thread first.

IV. SPIN-LOCK MANAGEMENT ALTERNATIVES

If a processor finds a thread management lock busy, it spin-waits for the lock to be released. An alternative would be to block, relinquishing the processor to do other useful work while the lock is busy. However, since thread management locks are held for only a short time, the overhead of performing this context switch would be prohibitive; in addition, any other work that the processor might do is controlled by a (or possibly the same) thread management lock. When an application lock is busy, a thread does have a choice between spin-waiting or blocking, but blocking at the user level may result in spin-waiting in a thread routine.

Spin-waiting has a hidden cost. Processors doing useful work may be slowed by processors that are merely waiting for a lock, due to bus contention. As a result, adding to the number of processors executing an application may in fact

slow it down by increasing the average number of spinning processors. Worse, the more spinning processors, the more the processor holding the lock is slowed, increasing the effective size of the critical section, resulting in even more waiting processors.

In this section, we evaluate three different approaches to spin-waiting.

A. Hardware Description

On the Symmetry Model A, each processor has its own cache; provided all of its memory references can be satisfied out of that cache, a processor's progress is independent of the activity of other processors. Whenever a processor reads data that are not in its cache, it must wait for the data to come from memory via the bus; with a write-through protocol, a processor may also have to wait for writes to be sent to memory. In both cases, the wait can be longer and the processor's progress slowed because of bus contention.

The Symmetry has a basic test-and-set instruction, *xchgb* (exchange byte), that atomically reads a memory location and writes in a new value. The atomicity of the *xchgb* operation is enforced by the bus: a copy of the memory location is brought into the processor's cache, modified there, and then written back to memory. As the value is written to memory, all copies of the old memory value in other caches are invalidated. No comparison of the old and new values is performed; memory is written and other copies invalidated even if the value is unchanged. Any requests for that memory location in the interim are delayed until the processor is done modifying it [17].

The Sequent locking protocol is as follows. The lock is held if the value is a 1 and free if 0. To lock, a processor exchanges in a 1. If the old value was a 0, it got the lock; if the value was a 1, the lock was already held by someone else, and the processor must try again. In either case, the value is 1 afterwards. The lock is released by exchanging in a 0; this allows some other processor to get a 0 back in exchange for a 1. There are several potential protocols for spin-waiting, which are described below.

B. Spin on Xchgb

Perhaps the simplest way to implement spin-waiting is for each processor to loop on the *xchgb* instruction until it succeeds. The drawback to this approach is that every *xchgb* instruction consumes bus resources, whether or not it succeeds. As additional processors spin on the lock, the holder of the lock is slowed both because the bus is busier and because to free the lock it must contend for permission to update the lock value with the *xchgb*'s of processors uselessly trying to acquire the lock.

C. Spin on Read

Coherent caches seem to allow processors to spin without using bus cycles. A processor can try to acquire the lock once; if this fails, the processor can spin reading the lock memory location. As long as the value is 1, the lock is still held. This spinning is done in the cache, avoiding bus traffic. When the lock is released, the cache copy will be invalidated;

the spinning processor will see the value change to 0 and can then try to acquire the lock using an *xchgb* operation. Sequent's run-time library uses this implementation [22].

A problem arises when there are a number of processors waiting for a small critical section. When the lock is freed, every spinning processor's cache copy is invalidated, causing each processor to fetch the new value in turn. The first to try to acquire the lock succeeds; however, each processor that sees the value as 0 before this occurs will also, in turn, try to acquire the lock, fail, and go back to looping on a read. Unfortunately, each processor that does an unsuccessful *xchgb* operation invalidates all cache copies, forcing all processors that were looping to miss again. Thus, after each such operation, almost every spinning processor contends for the bus, some still waiting to do an *xchgb* and the rest to fetch the lock value. Eventually, each processor sees that the lock has been acquired and quiesces, looping in its cache.

For a given number of spinning processors, the performance of this algorithm is better for longer critical sections. After the lock is released and before quiescence, each spinning processor spends most of its time with a pending bus request; any normal bus request during this time will be correspondingly delayed. After quiescence, the spinning processors place no load on the bus, allowing the processor holding the lock to progress unhindered. With longer critical sections, the initial degradation is less significant. By contrast, spinning on the *xchgb* instruction degrades bus performance evenly throughout the critical section.

D. Ethernet-Style Backoff

The source of the difficulty is that there is a cost to attempting to acquire the lock. A generic solution to problems of this sort is to have each processor estimate its likelihood of success, and only try the lock when the probability is high. The estimate can be made from experience. The more times a processor has tried and failed, the more likely it is that many processors are spinning for the lock. When the lock is released, then, instead of every processor rushing to try to get it, each waits a period of time dependent on the number of past failures. If the lock is still free after this period, then the probability of success is high enough to try the lock. We used this algorithm for our measurements in Section III.

The analogy with Ethernet is revealing. In the Ethernet protocol, a processor can start a network transmission in any time slot that the network is free [18]. If two try to start transmitting in the same slot, both fail and must be retried later. To avoid further collisions, the length of time before retrying depends on the number of collisions encountered so far. In our case, when a number of processors simultaneously try to acquire a lock, one will succeed, but its progress will be slower than if there were no collisions.

The downside to Ethernet-style protocols is that they are unfair. A processor that has just arrived is more likely to acquire the lock (or network) than one who has been waiting, and failing, for some time. Spinning on a test-and-set instruction and spinning on a read of the lock location are both probabilistically fair; each spinning processor has an equal likelihood of getting the lock, even though the possibility of indefinite

starvation exists. Lock fairness is sometimes important to an application.

Another drawback of the backoff algorithm is that it takes longer for a spinning processor to acquire a newly free lock. The processor must check the lock value, delay, and check it again before trying the lock. Once the lock is acquired, however, the processor will proceed faster, relatively unimpaired by other spinning processors.

Even using this algorithm, there will be processor degradation when there are large numbers of spinning processors. When the lock is released, every spinning processor encounters a cache miss. After this initial miss, most processors delay locally until some other processor has acquired the lock, and then miss again to see that the lock has been acquired. With enough spinning processors, the bus can be saturated with these misses, slowing down the processor executing in the critical section.

These cache misses can be avoided. A processor can delay whenever it reads the lock value as busy. If the lock is not busy, the processor can immediately try to acquire it. Thus, spinning processors miss their cache every time the delay period expires, rather than every time the lock is released. This is analogous to the Ethernet notion of persistence [18]. A result of this variation is an even greater delay between when a lock is released and when a spinning processor will acquire the lock.

Processors can spin-wait (degrading other processors) for things other than locks. Agarwal and Cherian [2] apply backoff to spin-waiting for data to become available. Spin-waiting can also be a problem with idle processors polling a central or distributed ready queue. When a ready thread is queued, if every idle processor rushes to acquire the lock, bus saturation will result. Even if each idle processor delays after observing that a thread is queued, then makes sure that it is still queued, there is still a cache miss per idle processor, hurting performance for large numbers of idle processors.

If idle processors are kept on a queue, this problem does not occur. Each idle processor spins on a separate flag. When a thread becomes ready to run, only one processor's flag is modified; every other processor continues spinning without even a cache miss. The performance advantage of having work look for processors instead of processors looking for work will therefore be more important in systems with large numbers of processors. This effect can be seen in Fig. 4; the cost of the central ready queue is higher when there are only a few runnable threads, since there are more idle processors spin-waiting for work to appear in the ready queue.

E. Measurement Results

Fig. 6 shows the elapsed time for various numbers of processors to cooperatively increment and test a shared counter in a critical section 1 million times, for each method of spin-waiting. Each processor executes a loop: wait for the lock, increment the counter, and release the lock. We do not claim that this test is representative of the normal use of critical sections, but similar curves have been measured with more significant computation between lock accesses. Since there is

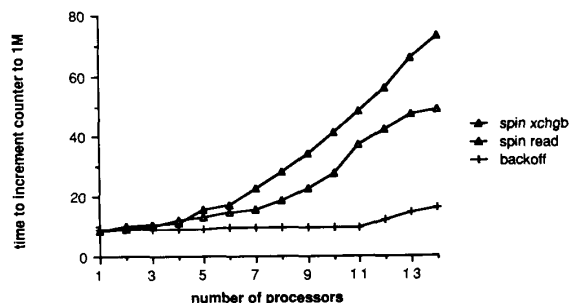


Fig. 6. Principal results for spin-waiting: elapsed time in seconds to increment a shared counter to 1 000 000 (measured).

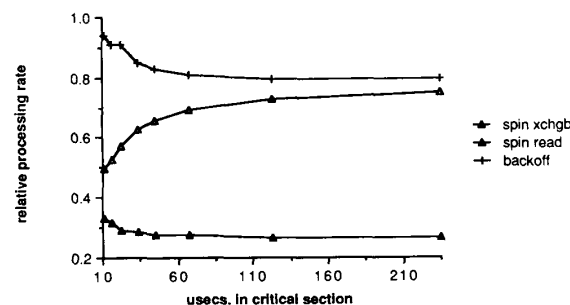


Fig. 7. Relative processor speed (8 processors to 1 processor) versus critical section size (measured).

little parallelism, if spinning processors did not slow the processor holding the lock, the curve would be flat.

The magnitude of this effect is striking. Both spinning on the xchgb instruction and spinning on a read degrade processor performance badly for even a moderate number of spinning processors. For small critical sections, in either alternative, every spinning processor spends all of its time doing cache read misses or atomic xchgb operations, consuming bus cycles as fast as possible. By contrast, the backoff algorithm results in only slight degradation unless the number of spin-waiting processors exceeds ten.

Fig. 7 shows the effect of increasing the size of the critical section on each algorithm's performance. In addition to incrementing a counter, the critical section contains varying amounts of other work. We normalized the time for eight processors by the time for one processor. This measures relative processor speed. Again, if spin-waiting did not slow the processor holding the lock, one processor would not be faster than eight, and the relative processor speed would always be equal to 1. As expected, spinning on a read degrades performance less as the size of the critical section grows, while spinning on the xchgb instruction degrades performance evenly throughout the critical section.

To test the tradeoff between processor degradation and the delay in acquiring a newly released lock, we measured the elapsed time for a number of processors to each increment a shared counter within a critical section. Once a processor acquired the lock and bumped the counter once, it was set to loop until all processors were done. This test is indicative of the cost of using a lock for barrier synchronization. Fig. 8 shows the elapsed time divided by the number of processors.

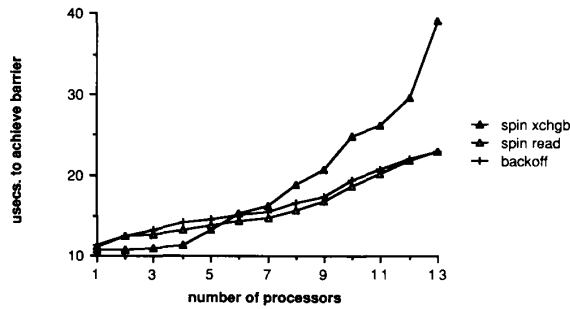


Fig. 8. Normalized time (microseconds per processor) to achieve barrier (measured).

If there is no processor degradation or delay in acquiring the lock, the elapsed time to achieve the barrier should increase linearly with each additional processor; the normalized curve in Fig. 8 should be flat.

Fig. 8 shows that for small numbers of processors, spinning on the xchgb instruction is fastest, since a processor immediately acquires the lock when it is released. As more processors are added, however, even though the lock is acquired faster, this is outweighed by the degradation of the processor holding the lock. The backoff algorithm shows a similar curve to spinning on a read, but for a different reason. Initially, many processors are queued for the lock; this leads spinning processors to guess large delay times. As more processors acquire the lock, there are fewer queued processors, and the delays become inappropriate.

Processors doing work are slowed proportional to the number of times they access the bus. Thus, the results of these tests depend somewhat on the content of the critical section. However, since the purpose of a critical section is to serialize modifications to shared data, its code is likely to be bus intensive. Our measurements indicate that almost half of the bus service demand of thread management is due to the critical section. Furthermore, thread management critical sections also tend to be small. For example, enqueueing or dequeueing a ready thread in a critical section both take less than 10 μ s, roughly the same as for Fig. 6.

F. Implications for Other Systems

In this section, we show that the performance of spin-waiting is of concern on architectures other than the Symmetry Model A.

Some multiprocessors do not provide hardware cache coherency; the BBN Butterfly [6] is an example. For these systems, every test of a lock value by a spinning processor requires a memory access. By inserting a delay between each test, the effect of spinning on busy processors can be reduced; backoff can be used to adapt the frequency of reads to the number of waiting processors.

The Symmetry Model A has a write-through protocol: when a processor modifies a location, the value is written to memory and all old copies of the location in other caches are invalidated. There is a cost to spin-waiting, even in architectures with a write-back cache coherency protocol. In a write-back protocol, the value is stored in the cache and later written to

memory when the cache block is replaced. There are two major approaches to keeping other caches consistent with the new value: all old copies in other caches can either be invalidated or updated with the new value (distributed-write) [4].

In the case of an invalidation-based write-back protocol, the spin-waiting alternatives have much the same effect as with write-through. If processors spin on the atomic test-and-set operation, the valid copy of the lock bounces from cache to cache, consuming bus resources. Provided test-and-sets invalidate all cache copies whether or not the lock value changes, spinning on a read does not help; there is still a cascade of repeated invalidations when the lock is released.

One possibility, then, is to add hardware to compare the old and new value of the lock on a test-and-set and to invalidate other copies only if they differ. While this would improve the performance of spinning on a read, it does not eliminate the problem. With P spinning processors, there are $O(P)$ bus requests per lock acquisition. Each processor must cache miss when the lock is released; it must also acquire the bus to ensure the atomicity of its subsequent test-and-set.

Systems with distributed-write coherency have similar performance. When a processor performs an atomic operation, every cache with an old copy is updated with the new value; thus no cache misses occur. If processors spin on a read, however, there will still be a rush of processors to try the lock when it is first released. Since the backoff algorithm reduces the number of unsuccessful lock attempts, it would reduce the bus load due to spinning even further.

Explicitly queueing spinning processors can further improve performance. Each processor in the queue spins on a separate flag; when a processor finishes with the lock, it passes control of the lock by setting the flag of the next one in the queue, without invalidating the flags of the other waiting processors.

In a related paper, we devised an efficient queue-based spin-waiting algorithm that uses only $O(1)$ bus transactions per execution of the critical section [3]. Each arriving processor does an atomic read-and-increment to obtain a unique sequence number. When a processor finishes with the lock, it taps the processor with the next highest sequence number; that processor now owns the lock. Since processors are sequenced, no atomic read-then-write instruction is needed to pass control of the lock. Table IV lists the code for this approach ("myPlace" is a location private to each processor). Measurements of this algorithm appear in that paper.

V. ANALYTICAL RESULTS

We developed a simple queueing network model for our thread package to demonstrate that the combination of processor degradation due to bus contention and the effect of lock contention can account for our measurements. We then used the validated model to project the performance of our thread package under varying conditions.

Our model is hierarchical. The low-level model represents the effect of bus contention on processor speed. The high-level model represents the effect of lock contention on throughput and response time. Since processor speed affects the amount of lock contention and the number of spinning processors af-

TABLE IV
QUEUE-BASED ALGORITHM FOR SPIN-WAITING

Init	<code>flags[0] = HAS_LOCK;</code> <code>flags[1..P-1] = MUST_WAIT;</code> <code>queueLast = 0;</code>
Lock	<code>myPlace = ReadAndIncrement(queueLast);</code> <code>while (flags[myPlace mod P] == MUST_WAIT)</code> <code>;</code> <code>flags[myPlace mod P] = MUST_WAIT;</code>
Unlock	<code>flags[(myPlace + 1) mod P] = HAS_LOCK;</code>

fects bus contention and thus processor speed, we iterate between levels to convergence. We describe the two submodels in more detail below.

A. Modeling Bus Contention

In the low-level model, we represent each processor as a customer in a closed queueing network. The network has two service centers: a queueing center for the bus and a delay center for nonbus activity. Each processor spends some of its time referencing memory through the bus and thus contending with other processors also using the bus, and some of its time processing out of its cache, independent of the activity of other processors. Processor speed is degraded by the percentage of time spent queueing, but not in service, at the bus.

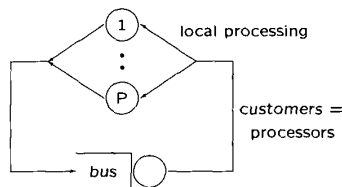


Diagram 5.1 Low-level model of bus contention.

This model is an approximation of the real bus mechanism, which is considerably more complex [17]. At moderate loads, our model will be pessimistic by predicting more contention than is actually experienced. Because of the regularity of the time each processor spends computing between accesses to the bus, if two processors collide at the bus, they are unlikely to collide at their next visit. Our model assumes that arrivals are more nearly independent.

There are three components to bus utilization. A processor can be executing user code, thread management code, or spin-waiting, each with different service demands on the bus. Given these service demands and the ratio of time each processor spends in each type of activity, we determine the aggregate service demands at the bus and at the delay center and use these aggregate demands to solve the model.

Since it is difficult to analytically determine the bus demand of a section of code, we determine a portion of it inductively from measurements. We provide each processor with its own copy of all data structures; we then run the code in parallel on each processor. Since there is no shared data, there can be no contention for software resources; any delay experienced by a processor relative to when it is running the code by itself must be due to contention for hardware resources, such as for memory or the bus. We then match a curve from our model of the bus to the measured curve and use the result as the service

demand for that section of code. The curves matched well in practice, deviating only at moderate loads, as expected.

Since bus contention may disproportionately impact the critical section execution time, affecting lock contention in the high-level model, we used this approach separately for the critical section and noncritical section code within thread management. The critical section code turns out to account for much of the bus demand of thread management.

Even though it could affect bus usage, we did not include in our model the effect of different numbers of processors on cache hit ratios. When a processor writes a location, the Symmetry updates both memory and that processor's cache. As a result, on a single processor, data that are both written and read will tend to stay in the cache, avoiding cache misses. When multiple processors read and write shared data, the cache copies of the data will be repeatedly invalidated as different processors update it, resulting in more cache misses than in the single processor case. Our model therefore underestimates bus demand, making it optimistic, especially as the bus nears saturation.

The bus demand of spinning processors was also determined inductively. P processors were set to run the critical section with separate copies of the data structures; by the experiment described above, we know the bus service demand of these processors. Q processors were set to run a shared copy of the critical section; one of these processors has the normal bus service demand, and $Q - 1$ spin-wait. By measuring the processor degradation of the P copies, we can determine the aggregate bus demand of the $Q - 1$ spinning processors. A two-class model was used, one class representing processors executing critical sections and one representing spinning processors. Only the response time of the processors executing the critical section is important.

The bus demand, at least for the backoff algorithm, is linear with the number of processors. While there is no *a priori* reason for this, it intuitively makes sense. The effect of adding a spinning processor with the backoff algorithm is to add two cache misses per execution of the critical section. The bus demand of other processors is relatively unaffected. While this invariance would also hold for the spin on xchgb algorithm, it is less true when processors spin on memory reads, because the cascade of cache misses is longer for every processor when more processors are spinning. Note that the graphs in Section IV could be used to infer the bus demand of spinning processors. We did not choose this approach because there is a correlation between when the processor holding the lock and when the processors spinning on the lock use the bus. The curve for the backoff algorithm in Fig. 6, for example, is similar to that of an optimistic asymptotic bound.

B. Modeling Lock Contention

In the high-level model, we represent each lock in the thread management path by a separate queueing center. Processing time spent not holding a lock is modeled as a delay center. Service demands were directly measured, then the part of each service demand due to bus accesses was inflated by the bus response time of the low-level model. As in the low-level model, each processor is represented as a single customer in

a closed class. By solving this model, we can determine the average amount of time each processor spends spin-waiting for a lock versus executing thread operations or user code. This ratio is then used as an input to the low-level model. (Note that it is a simple matter in this model to add queueing centers if the application-level code does further locking.)

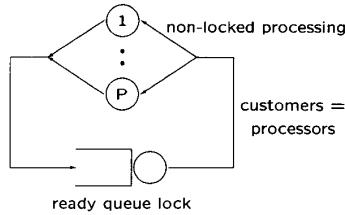


Diagram 5.2 High-level model of lock contention for the local free list alternative.

If the time between thread operations is deterministic, our model is pessimistic at moderate loads. As for the bus, if two processors collide at a lock, the effect of deterministic processing times is to reduce the likelihood that they will collide at the next visit. Fig. 2 shows this effect. The curves are similar in shape to asymptotic optimistic bounds, since the processing time to do each thread operation is deterministic. Fig. 3 does not show this effect, since the user computation for each thread was randomly chosen from a uniform distribution.

Our model does not explicitly represent an application's distribution of parallelism, although Fig. 4 shows that this affects performance. We chose not to include this in our model since the distribution, and more importantly the effect of lock queueing delay on that distribution, are almost always application-dependent.

Given the distribution, the model could be evaluated separately for each population of threads; these separate evaluations could then be averaged, weighted by the proportion of time for that population. The population of the high-level model should be the minimum between the number of processors and the number of threads, reflecting the number of active processors. The population of the low-level model should be set similarly, except that since idle processors consume bus resources, a second class should be added to represent them.

This method of separate evaluations ignores the fact that lock contention can only occur when the parallelism is being incremented or decremented; we believe that any distortion introduced by the adaptive nature of the mechanism will be outweighed by the effects of lock and bus contention. Ni and Wu [20] also discuss this issue.

C. Comparison to Measured Results, and Projections

Fig. 9 compares our model results with our measurement results previously reported in Fig. 3. We modeled two alternatives: per-processor ready queues (local readyq) and per-processor free lists with a central ready queue (local freelist). Our model agrees well with the measurements, within 5 percent except for the central ready queue with 18 processors. The model predicts the shape of the curve, but is somewhat optimistic; this appears to be due to underestimating the bus demand, which is important in determining the effective size

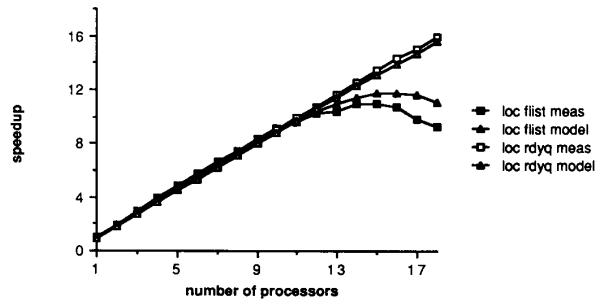


Fig. 9. Comparison of analytic and measured results from Fig. 3.

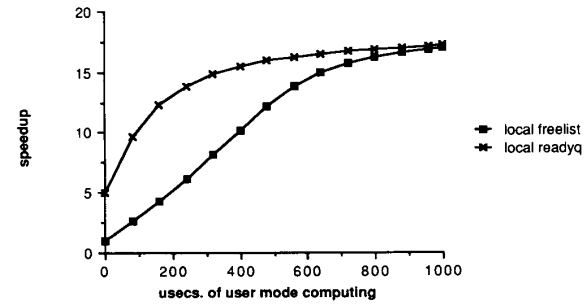


Fig. 10. Speedup versus microseconds of user computation per thread, 20 processors, bus load = 5 percent (analytic).

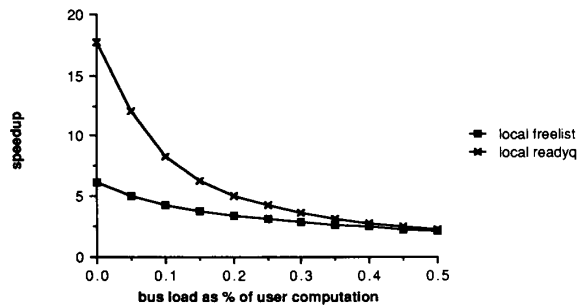


Fig. 11. Speedup versus bus load, user work = 200 μ s, 20 processors (analytic).

of the critical section. The model does capture the difference between the alternatives.

Having validated our model, we used it to investigate the effect of varying key parameters. Fig. 10 shows speedup of a hypothetical application with 20 processors as a function of the amount of user computation per thread. As we would expect, as an application uses finer grained parallelism (smaller amounts of computation per thread), the central lock on the ready queue becomes a bottleneck. For sufficiently coarse-grained parallelism, the performance of the thread package ceases to matter. In the limit, even DYNIX processes could be used.

Contention for the bus can also reduce the difference between the alternatives. Fig. 11 shows speedup as a function of the percentage usage of the bus by each thread. As the bus usage increases, the bus limits the speedup with local ready queues, but it also limits the speedup with the central ready queue, since bus contention inflates the critical section time.

On the other hand, the central ready queue lock can

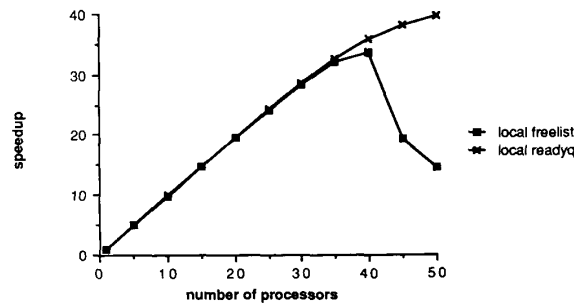


Fig. 12. Speedup versus number of processors, user work = 2 ms (analytic).

again limit speedup even for more coarsely-grained parallelism, given a sufficient number of processors. Fig. 12 shows speedup as a function of the number of processors when threads each compute for 2 ms. The sharp dropoff for the central ready queue alternative shows the inherent instability of a system where spinning processors consume resources.

VI. CONCLUSIONS

Threads have become a common element of new languages and operating systems. Efficient thread management is critical to achieving good performance from parallel applications. We have studied the performance implications of several thread management and locking alternatives. We showed the following.

- It is possible to implement a fast thread package. Simplicity is crucial for this.
- For fine-grained parallelism, small changes in data structures and locking have a significant effect on both latency and throughput.
- Per-processor data structures can be used to improve throughput; if a resource is not scarce, localizing data can avoid locking, improving latency as well.
- Spin-waiting can delay not only the processor waiting for a lock, but other processors doing work. This appears to be independent of the cache coherence protocol.
- The cost of spin-waiting can be reduced by using an Ethernet-style backoff or a queue-based algorithm.
- A simple queueing model can accurately predict the effect of a combination of factors on the performance of shared-memory multiprocessors.

An area of future research is to determine the extent to which our results, developed in the context of thread management systems, also apply to application programs that exploit fine-grained parallelism on shared-memory multiprocessors.

ACKNOWLEDGMENT

We would like to thank D. Wagner for suggesting that an Ethernet-style algorithm might solve the spin-waiting problem.

REFERENCES

- [1] M. Accetta, R. Baron, W. Bolosky, D. Golub, R. Rashid, A. Tevianian, and M. Young, "Mach: A new kernel foundation for UNIX development," in *Proc. Summer 1986 USENIX Tech. Conf. Exhibition*, pp. 93-112.
- [2] A. Agarwal and M. Cherian, "Adaptive backoff synchronization techniques," in *Proc. 16th Int. Symp. Comput. Architecture*, June 1989, pp. 396-406.
- [3] T. E. Anderson, "The performance implications of spin-waiting alternatives for shared-memory multiprocessors," in *Proc. 1989 Int. Conf. Parallel Processing*, Aug. 1989.
- [4] J. Archibald and J.-L. Baer, "Cache coherence protocols: Evaluation using a multiprocessor simulation model," *ACM Trans. Comput. Syst.*, vol. 4, no. 4, pp. 273-298, Nov. 1986.
- [5] M. J. Bach and S. J. Buroff, "Multiprocessor UNIX operating systems," *AT&T Bell Labs. Tech. J.*, vol. 63, no. 8, pp. 1733-1749, Oct. 1984.
- [6] BBN Laboratories, *Butterfly Parallel Processor Overview*, 1985.
- [7] B. Bershad, E. Lazowska, and H. Levy, "Presto: A system for object-oriented parallel programming," *Software: Practice and Experience*, vol. 18, no. 8, pp. 713-732, Aug. 1988.
- [8] B. Bershad, E. Lazowska, H. Levy, and D. Wagner, "An open environment for building parallel programming systems," in *Proc. ACM/SIGPLAN PPEALS 1988*, pp. 1-9.
- [9] K. W. Dritz and J. M. Boyle, "Beyond 'speedup': Performance analysis of parallel programs," Tech. Rep. ANL-87-7, Math. and Comput. Sci. Division, Argonne Nat. Lab., Feb. 1987.
- [10] D. Eager, E. Lazowska, and J. Zahorjan, "Adaptive load sharing in homogeneous distributed systems," *IEEE Trans. Software Eng.*, vol. 12, no. 5, pp. 662-675, May 1986.
- [11] J. Edler, J. Lipkis, and E. Schonberg, "Process management for highly parallel UNIX systems," *Ultracomputer Note 136*, Apr. 1988.
- [12] C. A. R. Hoare, "Communicating sequential processes," *Commun. ACM*, vol. 21, no. 8, pp. 666-677, Aug. 1978.
- [13] R. Holt, "A short introduction to concurrent Euclid," *SIGPLAN Notices*, vol. 17, pp. 60-79, May 1982.
- [14] E. Jul, H. Levy, N. Hutchinson, and A. Black, "Fine-grained mobility in the Emerald system," *ACM Trans. Comput. Syst.*, vol. 6, no. 1, pp. 109-133, Feb. 1988.
- [15] B. Kumar and T. Gonsalves, "Modelling and analysis of distributed software systems," in *Proc. 7th ACM Symp. Oper. Syst. Principles*, Dec. 1977, pp. 2-8.
- [16] B. W. Lampson and D. D. Redell, "Experiences with processes and monitors in Mesa," *Commun. ACM*, vol. 23, no. 2, pp. 104-117, Feb. 1980.
- [17] T. Lovett and S. Thakkar, "The Symmetry multiprocessor system," in *Proc. 1988 Int. Conf. Parallel Processing*, pp. 303-310.
- [18] R. Metcalfe and D. Boggs, "Ethernet: Distributed packet switching for local computer networks," *Commun. ACM*, vol. 19, no. 7, pp. 395-404, July 1976.
- [19] D. A. Mundie and D. A. Fisher, "Parallel processing in Ada," *IEEE Comput. Mag.*, pp. 20-25, Aug. 1985.
- [20] L. Ni and C.-F. Wu, "Design tradeoffs for process scheduling in tightly coupled multiprocessor systems," *IEEE Trans. Software Eng.*, vol. 15, no. 3, pp. 327-334, Mar. 1989.
- [21] M. Scott, T. LeBlanc, and B. Marsh, "Design rationale for Psyche, a general purpose multiprocessor operating system," in *Proc. 1988 Int. Conf. Parallel Processing*, Aug. 1988.
- [22] Sequent Computer Systems, Inc., *Symmetry Technical Summary*.
- [23] C. Thacker, L. Stewart, and E. Satterthwaite Jr., "Firefly: A multiprocessor workstation," *IEEE Trans. Comput.*, vol. 37, no. 8, pp. 909-920, Aug. 1988.
- [24] M. Vandevoorde and E. Roberts, "WorkCrews: An abstraction for controlling parallelism," *Int. J. Parallel Programming*, vol. 17, no. 4, pp. 347-366, Aug. 1988.

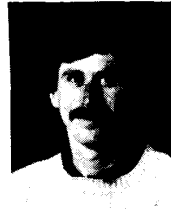
- [25] D. Wagner, E. Lazowska, and B. Bershad, "Techniques for efficient shared-memory parallel simulation," *Distribut. Simulation 1989*, Society for Computer Simulation, pp. 29-37.
- [26] J. Zahorjan, E. Lazowska, and D. Eager, "Spinning versus blocking in parallel systems with uncertainty," in *Proc. Int. Seminar Perform. Distribut. Parallel Syst.*, North Holland, Dec. 1988.



Thomas E. Anderson received the A.B. degree in 1983 from Harvard University, Cambridge, MA.

Since 1987, he has pursued a doctoral degree in the Department of Computer Science, University of Washington, Seattle. His research interests include multiprocessor operating systems and performance modeling.

Mr. Anderson won an IBM Graduate Fellowship in 1989.



Henry M. Levy is Research Associate Professor in the Department of Computer Science at the University of Washington, Seattle. His research involves computer architecture, distributed and parallel operating systems, object-oriented systems, and performance evaluation. Formerly, he was a consulting engineer with Digital Equipment Corporation. He is author of the books *Computer Programming and Architecture: The VAX* (second edition), and *Capability-Based Computer Systems*.