

Impact of Workload and System Parameters on Next Generation Cluster Scheduling Mechanisms

Yanyong Zhang, *Member, IEEE*, Anand Sivasubramaniam, *Member, IEEE*,
José Moreira, *Member, IEEE*, and Hubertus Franke, *Member, IEEE*

Abstract—Scheduling of processes onto processors of a parallel machine has always been an important and challenging area of research. The issue becomes even more crucial and difficult as we gradually progress to the use of off-the-shelf workstations, operating systems, and high bandwidth networks to build cost-effective clusters for demanding applications. Clusters are gaining acceptance not just in scientific applications that need supercomputing power, but also in domains such as databases, web service, and multimedia which place diverse Quality-of-Service (QoS) demands on the underlying system. Further, these applications have diverse characteristics in terms of their computation, communication, and I/O requirements, making conventional parallel scheduling solutions, such as space sharing or gang scheduling, unattractive. At the same time, leaving it to the native operating system of each node to make decisions independently can lead to ineffective use of system resources whenever there is communication. Instead, an emerging class of dynamic coscheduling mechanisms that attempt to take remedial actions to guide the system toward coscheduled execution without requiring explicit synchronization offers a lot of promise for cluster scheduling. Using a detailed simulator, this paper evaluates the pros and cons of different dynamic coscheduling alternatives while comparing their advantages over traditional gang scheduling (and not performing any coordinated scheduling at all). The impact of dynamic job arrivals, job characteristics, and different system parameters on these alternatives is evaluated in terms of several performance criteria. In addition, heuristics to enhance one of the alternatives even further are identified, classified, and evaluated. It is shown that these heuristics can significantly outperform the other alternatives over a spectrum of workload and system parameters and is thus a much better option for clusters than conventional gang scheduling.

Index Terms—Parallel scheduling, gang scheduling, dynamic coscheduling, clusters, simulation.

1 INTRODUCTION

SCHEDULING of processes onto processors of a parallel machine has always been an important and challenging area of research. Its importance stems from the impact of the scheduling discipline on the throughput and response times of the system. The research is challenging because of the numerous factors involved in implementing a scheduler. Some of these influencing factors are the parallel workload, presence of any sequential and/or interactive jobs, native operating system, node hardware, network interface, network, and communication software. The recent shift toward the adoption of off-the-shelf clusters/networks of workstations (called COWs/NOWs) for cost-effective parallel computing makes the design of an efficient scheduler even more crucial and challenging. Clusters are gaining acceptance not just in scientific applications that need supercomputing power, but also in domains such as

databases, web service, and multimedia, which place diverse Quality-of-Service (QoS) demands on the underlying system (not just higher throughput and/or lower response times). Further, these applications have diverse characteristics in terms of the computation, communication, and I/O operations which can raise complications when multiprogramming the system. Traditional solutions that have been used in conventional parallel systems are not adequately tuned to handle the diverse workloads and performance criteria required by cluster environments. This paper investigates the design space of scheduling strategies for clusters by extensively evaluating nine different alternatives to understand their pros and cons and compares them with a conventional solution (gang scheduling).

With the growing popularity of clusters, there have been several academic and commercial endeavors to develop networks, network interfaces, and messaging layers to provide efficient communication support for these systems. However, there is the more important and complex problem of managing and coordinating the resources across the nodes of a cluster. Optimizing communication alone may not necessarily translate to improved performance since a poor scheduling strategy could nullify any savings. For instance, a currently scheduled process on one node would experience a long wait for a message from a process not currently scheduled on another node, regardless of the

- Y. Zhang and A. Sivasubramaniam are with the Department of Computer Science and Engineering, Pennsylvania State University, University Park, PA 16802. E-mail: {yyzhang, anand}@cse.psu.edu.
- J. Moreira and H. Franke are with the IBM T.J. Watson Research Center, PO Box 218, Yorktown Heights, NY 10598. E-mail: {jmoreira, frankeh}@us.ibm.com.

Manuscript received 28 June 2000; revised 28 Feb. 2001; accepted 1 Mar. 2001.

For information on obtaining reprints of this article, please send e-mail to: tpds@computer.org, and reference IEEECS Log Number 112355.

low network latency. Scheduling and communication issues are thus closely intertwined and should be studied together. An efficient scheduling solution should understand the communication and I/O characteristics of the processes and coordinate the different activities across the nodes of a cluster without adding significant overheads.

Scheduling is usually done in two steps. The first step, spatial scheduling, consists of assigning processes to nodes. (A node can have one or more processors and runs a single operating system image.) The second step, temporal scheduling, consists of time multiplexing the various processes assigned to a node for execution by the processors of that node. There is a considerable body of literature regarding spatial scheduling [8] and we do not delve into this problem in this paper nor do we examine the issue of migrating processes during execution for better load balance. Without loss of generality, for this paper, we assume one processor (CPU) per node and an incoming job specifies how many CPUs it needs and executes one task (process) on each of the allocated CPUs.

The second scheduling step, temporal scheduling, is perhaps more important for a cluster. Just as common off-the-shelf (COTS) hardware has driven the popularity of clusters, it is rather tempting to leave it to the native (COTS) operating system scheduler to take care of managing the processes assigned to its node. However, the lack of global knowledge at each node can result in lower CPU utilization and higher communication or context switching overheads. As a result, there have traditionally been two approaches to address this problem. The first is space sharing [17], [36], which is a straightforward extension of batching to parallel systems. Multiple jobs can execute concurrently at different nodes, but each node is exclusively allocated to one process of a job which then runs to completion. Space sharing is simple to implement and is efficient from the viewpoint of reducing context switching costs as well as not having to hold the working sets of multiple applications in memory. However, space sharing in isolation can result in poor utilization (nodes can be free even though there are jobs waiting). Also, the lack of multiprogramming can hurt when a process performs a lot of I/O. The second approach is a hybrid scheme, called (exact) coscheduling or gang scheduling [22], [10], [9], [35], [11], [39], that combines time sharing with space sharing to remedy some of these problems. The processes of a job are scheduled on their respective nodes at the same time for a given time quantum, the expiration of which results in a synchronization between the nodes (using logical or physical clocks) to decide on the next job to run. This scheme usually requires long time quanta to offset high context switching and synchronization costs. Longer time quanta make the system less responsive for interactive and I/O intensive jobs (database services, graphics and visualization applications, etc.). In addition, strict gang scheduling also keeps the CPU idle while a process is performing I/O or waiting for a message within its allotted time quantum [11].

Recently, there has been interest in developing strategies that approximate coscheduling behavior without requiring explicit synchronization between the nodes (that still combine space and time sharing). We refer to this broad

class of strategies that approximate coscheduled execution as *dynamic coscheduling mechanisms*. The enabling technologies that have made this approach possible are the ability of the network interface card (NIC) and messaging layers to provide protected multi-user access to the network in conjunction with user-level messaging. Specifically, several user-processes at a node could concurrently send a message and an incoming message is directly transferred by the NIC to the corresponding destination process at that node (even if that process is not currently scheduled). It is thus not necessary to perform network context switching when the processes are switched out on their respective CPUs, as was necessary [35], [13] until recently. Dynamic coscheduling strategies try to hypothesize what is scheduled at remote nodes using local events (messaging actions/events in particular) to guide the native operating system scheduler toward coscheduled execution whenever needed. These strategies offer the promise of coscheduling, without the related overheads and scalability/reliability problems.

Prior to our work [20], [21], [32], there were only two suggestions [1], [31], [5] on how local messaging actions can be used to implement dynamic coscheduling. Both these mechanisms incur interrupts which can hurt performance under some situations. We have proposed two alternates, called Periodic Boost and Spin Yield, and have experimentally shown Periodic Boost to outperform the rest using an experimental cluster of eight Sun Ultra Enterprise servers running MPI applications [20], [21]. While our earlier study is a preliminary foray into this area, a comprehensive exercise exploring the pros and cons of these different alternatives is needed to answer several open and crucial questions:

- How do the different dynamic coscheduling alternatives compare when one considers dynamic job arrivals with different job sizes (number of CPUs) and execution times? Our earlier exercise considered only a few (constant) predetermined number of jobs, each demanding a constant number of CPUs and taking the same execution time. How does the arrival rate (load) of the jobs affect the average response times and throughput of the system?
- What is the impact of job characteristics on the performance of the system for the different scheduling alternatives? Specifically, how do the schemes compare as one varies the relative fraction of the computation (requiring CPU and memory resources only), communication, and I/O performed by a job? What is the impact of a multiprogrammed workload consisting of different job mixes? As mentioned earlier, clusters are intended to take on the demands of diverse applications, each with its own computation, communication, and I/O characteristics (for instance, a database application may be I/O intensive while a scientific application may be CPU or communication intensive), and should still meet the QoS requirements of each application. In addition, how does the work imbalance and skewness between the tasks (executing on different CPUs) affect the performance of each alternative?

- How do the system parameters, such as the multiprogramming level at each node, number of CPUs, and the operating system costs for context switching and interrupt processing, affect the relative performance of the schemes?
- Within the purview of the flexibility offered by the Periodic Boost mechanism that we have proposed, what coscheduling heuristics can we employ and how do these heuristics fare?

Investigation of these issues requires an extensive framework to model and evaluate the different scheduling alternatives. We have developed this framework with a detailed simulator that uses a multilevel feedback queue model (based on the Sun Solaris scheduler) as a core and builds the different scheduling alternatives around it. The choice of a simulation-based study stems from the need to vary different system parameters (some of which are difficult or even impossible to tweak on an actual system) in a controlled/nonintrusive manner. The simulator is driven by an abstraction of a real workload [11] that has been drawn from an actual supercomputing environment (Lawrence Livermore National Lab). Eight different dynamic coscheduling strategies are evaluated using this infrastructure and compared with exact coscheduling (as well as with not performing any coordinated scheduling at all) to draw revealing insights. To our knowledge, there has not been any prior work exploring this extensive design space using a spectrum of performance metrics (response time, wait time, utilization, and fairness) and dynamic workloads. In addition, this paper identifies several new heuristics that can be used in conjunction with the Periodic Boost mechanism. Some of these heuristics can significantly outperform the other alternatives in terms of response time, utilization, and fairness and are better equipped to handle the responsiveness demands of emerging applications.

The rest of this paper is organized as follows: The next section presents a short summary of User-level Networking, following which the different scheduling alternatives are discussed in Section 3. Section 4 gives details on the simulator itself, together with the simulation parameters and performance metrics under consideration. The performance results are presented in Section 5. Finally, Section 6 summarizes the observations and outlines directions for future research.

2 USER-LEVEL NETWORKING (ULN)

In the following discussion, we give a quick overview of user-level messaging since it is important for understanding the dynamic coscheduling strategies.

Traditional communication mechanisms have necessitated going via the operating system kernel to ensure protection. Recent network interface cards (NIC), such as Myrinet, provide sufficient capabilities/intelligence whereby they are able to monitor regions of memory for messages to become available and directly stream them out onto the network without being explicitly told to do so by the operating system. Similarly, an incoming message is examined by the NIC and directly transferred to the corresponding application receive buffers in memory (even

if that process is not currently scheduled on the host CPU). From an application's point of view, sending translates to appending a message to a queue in memory and receiving translates to (waiting and) dequeuing a message from memory. To avoid interrupt processing costs, the waiting is usually implemented as polling (busy-wait). Experimental implementations of variations of this mechanism on different hardware platforms have demonstrated end-to-end (application-to-application) latencies of 10-20 microseconds for short messages while most traditional kernel-based mechanisms are an order of magnitude more expensive. User-level messaging is achieved without compromising protection since each process can only access its own send/receive buffers (referred to as an *endpoint*). Thus, virtual memory automatically provides protected access to the network. Several ULNs [37], [38], [23], [2] based on variations of this paradigm have been developed.

User-level messaging, though preferable for lowering the communication overhead, actually complicates the issue from the scheduling viewpoint. A kernel-based blocking receive call would be treated as an I/O operation with the operating system putting the process to sleep. This may avoid idle cycles (which could be given to some other process at that node) spent polling for message arrival in a user-based mechanism. Efficient scheduling support in the context of user-level messaging thus presents interesting challenges.

3 SCHEDULING STRATEGIES

We present the gang coscheduling model, followed by the native operating system scheduler at each node that is modeled as the core around which the dynamic coscheduling mechanisms are structured. Finally, the details of the different coscheduling heuristics are presented. All the models have been designed and developed based on our implementation [20], [21] of these mechanisms on an actual Sun Solaris cluster connected by Myrinet.

3.1 Coscheduling or Gang Scheduling (GS)

Gang Scheduling ensures that the processes/tasks of a job are scheduled on their respective nodes at the same time. This usually requires some means of explicit or implicit synchronization to make a coordinated scheduling decision at the end of each time quantum. The CM-5 [35], Meiko CS-2 (gang mode) [3], Intel Paragon (gang mode) [14], AP/Linux [34], and the GangLL [11] on the IBM SP-2 at LLNL are examples of systems that use this approach.

The simulation model is based on the implementation of the GangLL scheduler [19], [11] on the Blue Pacific machine at Lawrence Livermore National Lab. The model uses an Ousterhout [22] matrix with the columns representing the CPUs and rows representing the time quanta (as many rows as the multiprogramming level). In an actual system, the multiprogramming level (MPL) will be set based on the available resources (such as memory, swap space, etc.) that can handle a certain number of jobs concurrently without significantly degrading performance. A job is allocated the required number of cells in a single row if available. Else, it is made to wait in an *arrival queue* (served in FCFS order) until there are enough free cells in a row. During each time

quantum, a CPU executes the assigned job for that row in the matrix and does not move to the next row until the next quantum (regardless of whether the process is waiting for a message, or performing I/O, or even finishes before the quantum ends). At the end of the quantum, a context switch cost is incurred. This not only includes the traditional costs, but also the cost for synchronizing between the nodes before it schedules the job for the next quantum (GangLL [11] actually uses physical clocks with large time quanta instead of explicit synchronization). Message receives are implemented as busy-waits (spinning), though some of this time could get hidden if the process is context switched out (quantum expires).

3.2 Local Scheduling

We refer to the system which does not make any coordinated scheduling decisions across the nodes as *local* scheduling. The native operating system is left to schedule the processes at each node. As in gang scheduling, each node can again handle a maximum of MPL processes at any time with the difference that an arriving job does not have to wait until free slots are found in a single row. Rather, a job can be scheduled to the corresponding CPUs that are not already operating at their full MPL capacity (can be a different row position for each column if one is to look at this problem as filling the Ousterhout matrix). If the job cannot find that many CPUs, it waits in an *arrival queue* (served in FCFS order) until it does. A brief description of the native scheduler (multilevel feedback queue) at each node, which closely resembles the Solaris scheduler, follows.

There are 60 priority levels (0 to 59 with a higher number denoting a higher priority) with a queue of runnable processes at each level. The process at the head of the highest priority queue is executed first. Higher priority levels get smaller time slices than lower priority levels that range from 20 ms for level 59 to 200 ms for level 0. At the end of the quantum, the currently executing process is degraded to the end of the queue of the next lower priority level. Process priority is boosted (to the head of the level 59 queue) when they return to the runnable state from the blocked state (completion of I/O, signal on a semaphore etc.) This design strives to strike a balance between compute and I/O bound jobs with I/O bound jobs typically executing at higher priority levels to initiate the I/O operation as early as possible. The scheduler, which runs every millisecond, ensures that lower priority processes are preempted if a higher priority process becomes runnable (the preemption may thus not take place immediately after priority changes). For fairness, the priorities of all processes are raised to level 59 every second.

The ULN messaging actions explained above are used as is in *local*; send is simply an append to a queue in memory and receive is busy waiting (spinning) in user-space for message arrival (consuming CPU cycles). This scheme has been considered as a baseline to show the need for a better scheduling strategy.

TABLE 1
Design Space of Dynamic Coscheduling Strategies

What do you do on message arrival?	How do you wait for a message?		
	Busy Wait	Spin Block	Spin Yield
No Explicit Reschedule	Local	SB	SY
Interrupt & Reschedule	DCS	DCS-SB	DCS-SY
Periodically Reschedule	PB	PB-SB	PB-SY

3.3 Dynamic Coscheduling Strategies/Heuristics

As mentioned earlier, these strategies rely on messaging actions to guide the system toward coscheduled execution and there is no coordinated effort explicitly taken to achieve this goal. Logically, there are two components in the interaction between a scheduler and the communication mechanism. The first is related to how the process waits for a message. This can involve: 1) just spinning (busy wait), 2) blocking after spinning for a while, or 3) yielding to some other process after spinning for a while. The second component is related to what happens when a message arrives and is transferred to application-level buffers. Here again, there are three possibilities: 1) Do no explicit rescheduling, 2) interrupt the host and take remedial steps to explicitly schedule the receiver process, and 3) periodically examine message queues and take steps as in 2. These two components can be combined to give a 3×3 design space of dynamic coscheduling strategies, as shown in Table 1 (a description of these strategies follows).

In the following discussion, we limit our explanations in order to familiarize the reader with these strategies and to explain how they are simulated, rather than give a detailed discussion of their implementation on an actual operating system. For a detailed description of the implementation of these different strategies on a Sun Solaris cluster connected by Myrinet, the reader is referred to our earlier work [20]. The simulation models and parameters are based on our earlier experimental exercises. All these strategies use the same scheme described above in *local* to assign the processes (tasks) of an arriving job to the different CPUs.

3.3.1 Spin Block (SB)

Versions of this mechanism have been considered by others in the context of implicit coscheduling [7], [1], [27] and demand-based coscheduling [31]. In this scheme, a process spins on a message receive for a fixed amount of time before blocking itself. The fixed time for which it spins, henceforth referred to as *spin time*, is carefully chosen to optimize performance. The rationale here is that if the message arrives in a reasonable amount of time (spin time), the sender process is also currently scheduled and the receiver should hold on to the CPU to increase the likelihood of executing in the near future when the sender process is also executing. Otherwise, it should block so that CPU cycles are not wasted.

The simulation model sets the spin time for a message to be slightly higher than the expected end-to-end latency (in the absence of any contention for network or node resources) of the message it is waiting for. If the corresponding message arrives within this period, the mechanism works the same way as the earlier scheduling

schemes (busy-wait). Else, the process makes a system call to block using a semaphore operation. On subsequent message arrival, the NIC firmware (having been told that the process has blocked) raises an interrupt, which is serviced by the kernel to unblock the process. As mentioned earlier, the process gets a priority boost (to the head of the queue of the highest priority level) on wakeup. Costs for blocking/unblocking a process, context switches resulting from these operations, and interrupt processing are modeled in the simulation.

3.3.2 Spin Yield (SY)

In SB, the process blocks after spinning. This has two consequences. First, an interrupt is required to wake the process on message arrival (which is an overhead). Second, the block action only relinquishes the CPU and there is no hint given to the underlying scheduler as to what should be scheduled next. In our earlier work [20], we have proposed Spin Yield (SY) as an alternative to address these problems. In this strategy, after spinning for the required spin time, the process does not block. Instead, it lowers its priority, boosts the priority of another process (based on the pending messages of the other processes at that workstation), and continues spinning. This avoids an interrupt (since the process keeps spinning, albeit at a lower priority) and gives hints to the scheduler as to what should be scheduled next.

In the simulation, the time spent spinning before yielding is again set similarly to the SB mechanism. On yielding, the process priority is dropped to a level that is one below the lowest priority of a process at that node and the priority of another process with a pending (unconsumed) message is boosted to the head of level 59 (the details of the algorithm that is used to select the candidate for boosting are described later in the context of the PB mechanism). The application resumes spinning upon returning from the yield mechanism (system call) and the scheduler is likely to preempt this process at the next millisecond boundary. System call costs, together with the overheads for manipulating the priority queues, are accounted for in the yield call.

3.3.3 Demand-Based Coscheduling (DCS)

Demand-based coscheduling [31] uses an incoming message to schedule the process for which it is intended and preempts the current process if the intended receiver is not currently scheduled. The underlying rationale is that the receipt of a message denotes the higher likelihood of the sender process of that job being scheduled at the remote workstation at that time.

Our DCS model is similar to the one discussed in [31]. Every one millisecond, the NIC finds out which process is currently being scheduled on its host CPU. The NIC uses this information to raise an interrupt (if the receiver process is not currently scheduled) on message arrival after transferring it to the corresponding endpoint. The interrupt service routine simply raises the priority of this receiver process to the head of the queue for level 59 so that it can possibly get scheduled at the next scheduler invocation (millisecond boundary). The application sends and receives (implemented as busy-waits) remain the same as in Local. The model again ensures that the costs for interrupts and queue manipulations are included based on experimental results.

3.3.4 Periodic Boost (PB)

We have proposed this as another interruptless alternative to addressing the inefficiencies arising from scheduling skews between processes. Instead of immediately interrupting the host CPU on message arrival as in DCS, the NIC functionality remains the same as in the baseline ULN receive mechanism. A kernel activity (thread) becomes active every one millisecond (the resolution of the scheduler activation), checks message queues, and boosts the priority of a process based on some heuristic. Whenever the scheduler becomes active (at the next millisecond boundary), it would preempt the current process and schedule the boosted process.

There are several heuristics that one could use within the PB mechanism to decide on who or when to boost and, in this paper, *we propose, classify, and evaluate 10 new heuristics*. For most of the experiments, the heuristic that is used (called Heuristic A) can be briefly explained as follows (further details are in Section 5.9): The PB mechanism goes about examining message queues in a round-robin fashion, starting with the current process, and stops at the first process performing a receive with the message that it is receiving present in the endpoint buffers (message has arrived but has not yet been consumed by the receive call). This process is then boosted to the head of the level 59 queue. If there is no such process, then, again going about it in a round-robin fashion, the mechanism tries to find a process which is not within a receive call (this can be incorporated easily into the existing ULN mechanism by simply setting a flag in the endpoint when the application enters a receive call and resetting it when it exits from the call). It then boosts this process if there is one. Else, the PB mechanism does nothing. This algorithm is used in finding a candidate for boosting in the SY mechanism as well.

The rationale behind this heuristic and a comparison with seven other variations is given later in Section 5.9. The simulation models the details of the PB mechanism, incorporating the costs associated with polling the endpoints to find pending message information and the subsequent costs of manipulating priority queues.

3.3.5 Combinations (DCS-SB, PB-SB, DCS-SY, PB-SY)

One could combine the choices for the two messaging actions, as was shown in Table 1, to derive integrated approaches that get the best (or worst) of both choices. As a result, there is nothing preventing us from considering the four alternatives—Demand-based coscheduling with Spin-Block (DCS-SB), Periodic Boost with Spin-Block (PB-SB), Demand-based coscheduling with Spin-Yield (DCS-SY), and Periodic Boost with Spin Yield (PB-SY)—as well.

4 EXPERIMENTAL PLATFORM

Before we present performance results, we give details on the simulation platform, the workload used to drive the simulator, the parameters that are varied in this exercise, and the performance metrics of interest.

4.1 Simulator

We use a discrete-event simulator that has been built using the process-based CSIM package. The simulator has the

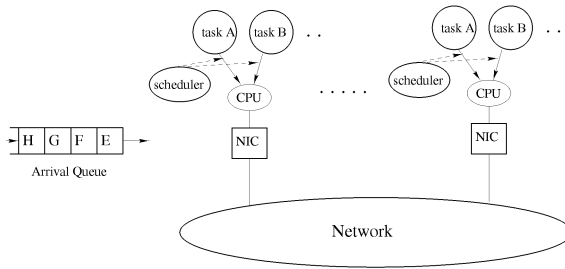


Fig. 1. Structure of the simulator.

following modules for each node in the system: network interface, operating system scheduler, and the application process. In addition, there is a network module that connects the different nodes. Since the focus of this paper is more on the scheduling mechanisms, we use a simple linear model for the network that is parameterized by the message size and do not consider network contention (though the contention at the network interface is modeled). The network interface module examines incoming messages from the network and deposits them into the corresponding endpoint. Similarly, it waits for outgoing messages and delivers them into the network module. Costs for these operations have been drawn from microbenchmarks run on our experimental platform discussed in our earlier work [20], [21], [33]. The core scheduler at each node uses a multilevel feedback queue that has been discussed in Section 3.2. The scheduler becomes active every one millisecond at each node (similar to the Solaris scheduler). At this time, it checks if the quantum has expired for the currently scheduled process and, if so, it preempts and reschedules another. Even if the quantum has not expired, the scheduler consults the feedback queues to check if there is a ready process with a priority higher than the currently scheduled one for possible preemption. There are two other components to the scheduler that correspond to interrupts and the periodic boost mechanism, respectively. The interrupt mechanism is used in SB, DCS, PB-SB, DCS-SB, and DCS-SY and becomes active immediately after the network interface module raises an interrupt. After accounting for interrupt processing costs, the scheduling queues may need to be manipulated in this mechanism. The periodic boost mechanism is used in PB, PB-SB, and PB-SY and becomes active every one millisecond to check the endpoints for messages and manipulate the scheduling queues as described in the previous section. Costs for the queue manipulations have again been drawn from our experimental studies. (See Fig. 1.)

Our simulator hides all the details of the scheduling models from the application process. The application interface allows the flexibility of specifying the computation time, communication events (sends or receives with message sizes and destinations), and I/O overheads. The development of the simulator has itself been a significant effort, but we do not delve further into the implementation details.

4.2 Workloads

We are interested in using realistic workloads to drive the performance evaluation. Toward this, we are interested in

capturing the dynamic behavior of the environment (i.e., dynamic job arrivals), different job execution times and job sizes (number of CPUs), and the characteristics of each job (computation, communication, and I/O fractions) as well.

To capture the dynamic behavior of the environment, our experiments use a workload that is drawn from a characterization of a real supercomputing environment at Lawrence Livermore National Lab. Job arrival, execution time, and size (number of CPUs—henceforth referred to as tasks) information of this environment have been traced and characterized to fit a mathematical model (Hyper-Erlang distribution of common order). The reader is referred to [11] for details on this work and the use of the model in different evaluation exercises [40]. However, due to the immense simulation details involved in this exercise (unlike in any of the previous studies), requiring the modeling of scheduling queues at granularities smaller than even a millisecond, it is not feasible to simulate very large systems. As a compromise, we have limited ourselves to clusters of up to 32 nodes and select jobs from the characterized model that fall within this limit.

As can be expected, the characteristics of each job in the system can further have an impact on the performance results. In particular, the time spent in the computation (only the CPU is required), communication, I/O activities, and the frequency of these operations can interact with the scheduling strategies in different ways. It is easy to draw false conclusions if one does not consider all these different artifacts in the performance evaluation. In reality, a job can be intensive in any one of the three components—computation (CPU), communication, or I/O—or can have different proportions of these components. To consider these different situations, we identify six different job types with different proportions of these components. Our evaluations use eight different workloads, termed *WL1* to *WL8*, with the first six using jobs of the corresponding class in isolation, as shown in Table 2. In the seventh workload (*WL7*), a job has an equal probability of falling in any of the six job types so that we consider the effect of a mixed load on the system. The last workload (*WL8*) considers an equal mix of the three job classes that are each intensive in one of the three components (CPU, I/O, and communication).

We consider four different communication patterns (that are typical of several parallel applications) between the tasks of a job that are shown in Fig. 2b. The considered patterns include *Nearest Neighbor (NN)* (process i communicates with $i-1$ and $i+1$, which is exhibited by applications such as Successive Over-Relaxation), *all-to-all (AA)* (process i sends a message to all other processes and waits for a message from everyone else, which is typical in many scientific computations, such as matrix transposes, multidimensional FFTs, etc.), *tree* (going up and coming down, similar to the tree-based barrier discussed in [18], which is a typical structure used for global reductions), and *linear* (process i sends a message to $i+1$ and waits for a message from $i-1$, which is a common structure in pipelined/systolic executions). All these patterns use a constant message size of 4,096 bytes.

Once the relative proportions and corresponding times in the three components (computation, communication, and

TABLE 2
Workloads

Job Type	Comp. (%)	I/O (%)	Comm. (%)
J1	35	15	50
J2	35	50	15
J3	35	35	30
J4	90	5	5
J5	35	5	60
J6	65	5	30

Workload	Job Types in Workload
WL1	J1
WL2	J2
WL3	J3
WL4	J4
WL5	J5
WL6	J6
WL7	equal mix of J1 thru' J6
WL8	equal mix of J2, J4 & J5

I/O) are derived for a given job, its tasks iteratively (as shown in Fig. 2a) go through a sequence of compute, I/O, sends/receives as per the communication pattern (Fig. 2b). By fixing the raw 1-way latency of a 4,096 byte message (from an experimental platform), the cost of communication per iteration in the ideal case (when everything is balanced) is known. Based on this and the relative proportion of the other two components (compute and I/O), which is determined by the job type, the computation, and I/O times per iteration, can be calculated. Together with the total job execution time (given by the characterized model), these individual times determine the number of iterations that a job goes through.

It is also important to note that the work imbalance/skewness between the tasks (and the resulting mismatch of the sends and receives) can have a significant impact on results. As a result, the skewness (s), which is expressed as a percentage of the computation and I/O fractions of the tasks, is another parameter that is varied in the experiments. Formally, the CPU and I/O components of each iteration of a task that were calculated earlier are each multiplied by a factor $(1 + \text{unif}(-s/2, s/2))$, where $\text{unif}(x, y)$ generates a random number between x and y using a uniform distribution. If s is set to 0, then all tasks spend the same computation and I/O times in each iteration and, thus, arrive at the communication events at the same time. In this

case, the execution time for this job will match the one picked for it from the characterized model [11] on a dedicated (nonmultiprogrammed) system. A larger skewness implies that tasks will arrive at the communication events at different times and the overall execution time per iteration will depend on who comes last to the send/receive calls. This also implies that the execution time is likely to get larger compared to that derived from the characterized model (each iteration can get elongated) with a larger s . The effectiveness of the scheduling mechanisms can be evaluated by how well they are able to hide the increase in execution times.

4.3 Parameters

Several parameters and costs can be varied in the simulator, and some of these (and the values that are used) are given in Table 3. Many of the times given in this table have been obtained from microbenchmarks on actual operating systems and Myrinet clusters with a ULN connecting the machines [20], [33].

4.4 Metrics

This exercise considers several metrics that are important from both the system and user's perspective:

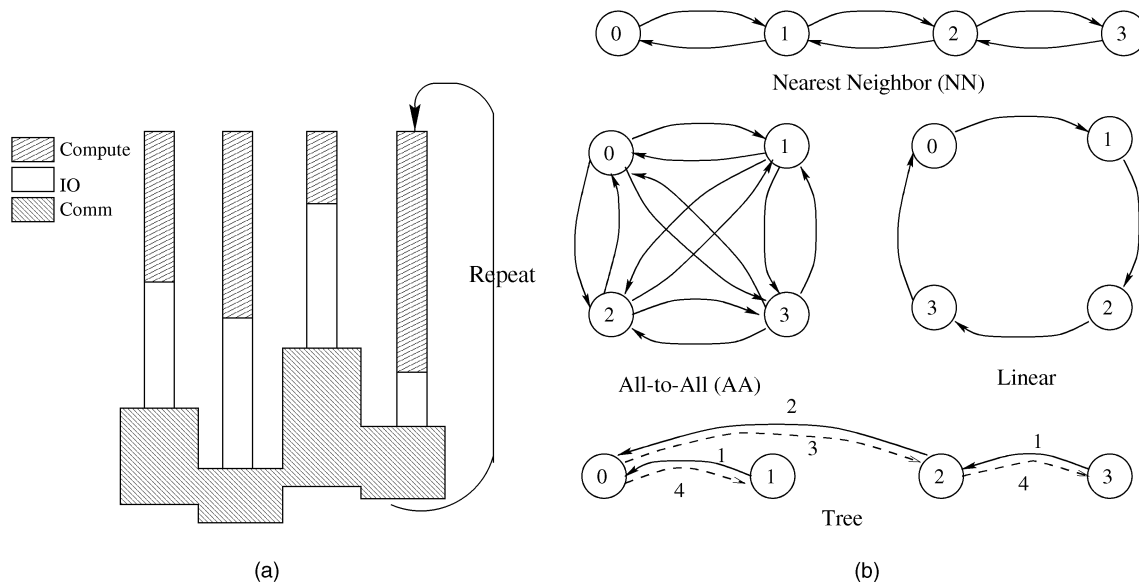


Fig. 2. Job structure. (a) Application model. (b) Communication pattern.

TABLE 3
Simulation Parameters and Values in Experiments

Parameter	Value(s)
p (# of nodes or CPUs)	32, 16
MPL (Multiprogramming Level)	2, 4, <u>5</u> , 16
s (Skewness)	150%, <u>20%</u>
CS (Context Switch Cost)	
Dynamic Coscheduling	<u>200</u> us, 100 us
GS	<u>2</u> ms, 1 ms
I (Interrupt Cost)	50 us, 25 us
Q (GS Time Quantum)	200 ms, 100 ms
Move between queues	<u>3</u> us
Check an endpoint	<u>2</u> us
Communication Pattern	NN, AA, Tree, Linear
Message Size	4096 bytes
1-way Message Latency	<u>185.48</u> us

Unless explicitly mentioned, the default (underlined) values are used.

- **Response Time:** This is the time difference between when a job completes and when it arrives in the system averaged over all jobs.
- **Wait Time:** This is the average time spent by a job waiting in the arrival queue before it is scheduled.
- **Execution Time:** This is the difference between *Response* and *Wait* times.
- **Slowdown:** This is the ratio of the *execution time* to the time taken on a system dedicated solely to this job. It is an indication of the slowdown that a job experiences when it executes in multiprogrammed fashion compared to running in isolation after it is allocated to the CPUs.
- **Utilization:** This is the percentage of time that the system actually spends in useful work.
- **Fairness:** The fairness to different job types (CPU, communication, or I/O intensive) is evaluated by comparing the slowdown of the individual job classes and its coefficient of variation (standard deviation divided by the mean) in a mixed workload. A smaller variation indicates a more fair scheme.
- **Performance Profile:** This is expressed as a graph showing the percentage of time that a CPU spends on the average in different components (compute, spinning/busy-wait, context switches, idling, and other overheads such as interrupt costs). Such a graph is useful to better understand the performance results and also to identify bottlenecks in the execution.

5 PERFORMANCE RESULTS

5.1 Impact of Load

As the load increases, the system is likely to be more heavily utilized. Consequently, jobs are likely to experience longer wait and response times. The workload characterization effort presented in prior research [11] provides a way of cranking the induced load by modulating the job arrival rates and execution times. The resulting effect on the response time of the system is plotted against the system utilization in Fig. 3 for the mixed workload (WL7). It should

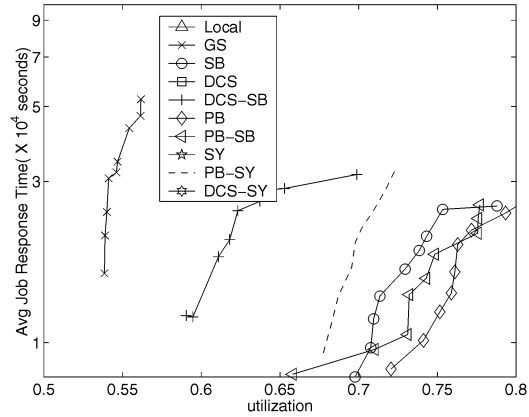


Fig. 3. Impact of load on response time (NN, WL7, $p = 32$, $MPL = 5$, $s = 20\%$).

be noted that utilization is calculated only using the useful work (computation time), as was mentioned earlier.

Local, SY, DCS, and DCS-SY, which all use spin-based receives, saturate even before the utilization reaches 50 percent and are thus not seen in Fig. 3. We find that GS can go only as high as 57 percent before saturation. The remaining schemes—DCS-SB, PB-SY, PB-SB, SB, and PB—perform significantly better than GS. Of these, PB, PB-SB, and SB perform the best, with the utilization reaching as high as 77 percent before saturation. In the rest of this section, we set the interarrival time of jobs to be 674.14 seconds with a mean execution time (including computation, I/O, and communication times) of 2.2250×10^3 seconds.

5.2 Impact of the Nature of Workload

Next, we examine the impact of the computation, communication, and I/O components of the workload on the performance of the different schemes. WL4, WL5, WL2, and WL3 are CPU intensive, communication intensive, I/O intensive, and evenly balanced (among the three components) workloads, respectively. The response times of the 10 scheduling alternatives for these four workloads are presented in Fig. 4. The response time is further broken down into the wait time in the arrival queue and the execution time. It should be noted that some of the bars which hit the upper boundary of the y-axis have been truncated (and the execution portion of these bars is not visible).

For the CPU intensive workload (WL4), Local and SY are simply not acceptable (this is true for the other workloads as well). Both DCS and DCS-SY have higher response times than the rest, mainly because of the longer wait in the arrival queue. GS has the next highest response time for this workload. For a high computation proportion, even an $s = 20\%$ skewness can make a difference in causing a load imbalance between the tasks, leading to inefficiencies in GS. There is not a significant difference between the response times of the other five schemes for this workload. When the computation is high, the inefficiencies due to inexact coscheduling are not really felt (since receives are infrequent) and any skewness is well hidden by the scheduling mechanism.

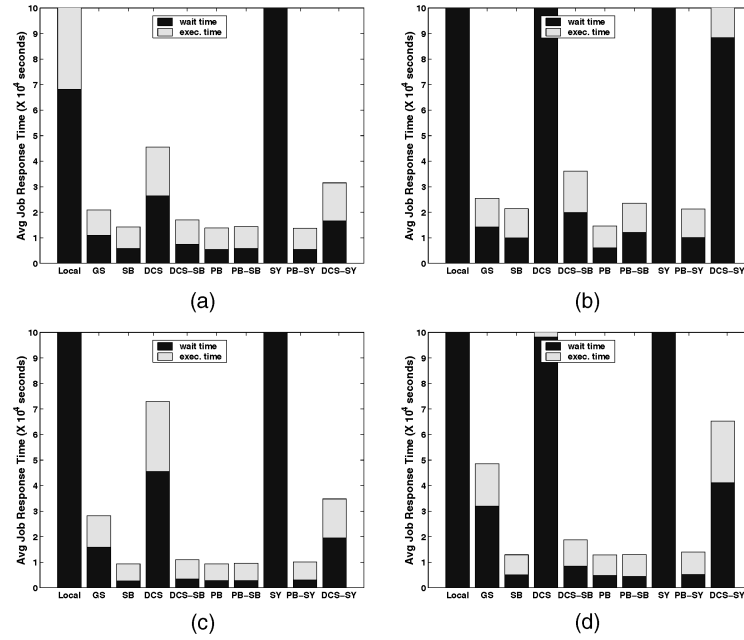


Fig. 4. Impact of the nature of workload on response time (NN, $p = 32$, $MPL = 5$, $s = 20\%$). (a) WL4. (b) WL5. (c) WL2. (d) WL3.

On the other hand, when one examines the communication intensive workload (WL5), the differences are more pronounced. DCS and DCS-SY become much worse here. The performance of GS, SB, PB-SB, and PB-SY is comparable, with PB giving the best performance. When the communication component becomes high, there is a stronger need for coscheduling. Further, the skewness between the executing tasks of an application gets lower with a lower compute fraction, making the inefficiencies of GS less important.

In the I/O intensive (WL2) or mixed workload (WL3), the occurrence of I/O activities within a time quantum keeps the CPU idle for the remainder of the quantum in GS. When the I/O portion is more intensive (WL2), the impact of noncoscheduled execution is not felt. As a result, there is not a significant difference between SB, DCS-SB, PB, PB-SB, and PB-SY. Finally, the evenly balanced workload (WL3) reiterates the observations made for the communication intensive workload, though to a lesser degree, while showing that GS is not as good as SB/PB/PB-SB/PB-SY because of the presence of I/O activities.

As in the previous section, we find that the PB, SB, PB-SY, and PB-SB schemes are uniformly good for all the workloads. Of these, PB is clearly the best for the communication intensive workload. DCS-SB falls in the next category. GS is reasonable as long as the execution does not lead to work imbalances. It should be noted that, in all four workloads, several of the dynamic coscheduling mechanisms have a shorter wait time component compared to GS. This is because allocating CPUs to a job in GS requires those CPUs to be available during the same time quantum (i.e., the same row of the Ousterhout matrix) while the others do not pose any such restrictions.

5.3 Impact of Multiprogramming Level (MPL)

As one increases the multiprogramming level, a larger number of jobs can be simultaneously accommodated in the

system, which works in favor of lowering wait times. A larger MPL also allows the system to find alternate work (that is useful) to do when processes block (or yield). However, larger MPLs also imply a larger number of context switches during execution, which stretches execution times. In addition, a larger MPL also decreases the likelihood of the tasks being coscheduled in the dynamic coscheduling mechanisms. Thus, it is interesting to study the impact of MPL to understand these factors and their interplay. Table 4 shows the change in response times normalized (with respect to $MPL = 2$) for WL2 with $MPL = 5$ and WL4 and WL5 with $MPL = 5$ and 16. Fig. 5 shows the performance profile graphs for the corresponding experiments, which can be used to observe the change in percentage of time spent in different activities when one increases the MPL level.

In Local, the overheads with a larger MPL dominate over other factors resulting in larger response times when we go

TABLE 4
Impact of MPL: Normalized Change in Response Time
with Respect to $MPL = 2$ (NN, $p = 32$, $s = 20\%$)

Scheme	WL2		WL4		WL5	
	MPL=5	MPL=16	MPL=5	MPL=16	MPL=5	MPL=16
Local	0.30	0.37	-	15.44	-	-
GS	-0.25	-0.20	-0.34	-0.30	-0.41	-
SB	-0.69	-0.36	-0.48	-0.44	-0.54	-
DCS	-0.27	-0.13	-0.34	0.18	-	-
DCS-SB	-0.68	-0.35	-0.44	-0.21	-0.38	-
PB	-0.65	-0.25	-0.38	-0.41	-0.52	-
PB-SB	-0.71	-0.34	-0.45	-0.43	-0.54	-
SY	-0.29	2.82	-	4.59	-	-
PB-SY	-0.64	-0.24	0.36	-0.28	2.78	-
DCS-SY	-0.58	-0.44	-0.81	-0.42	-0.94	-

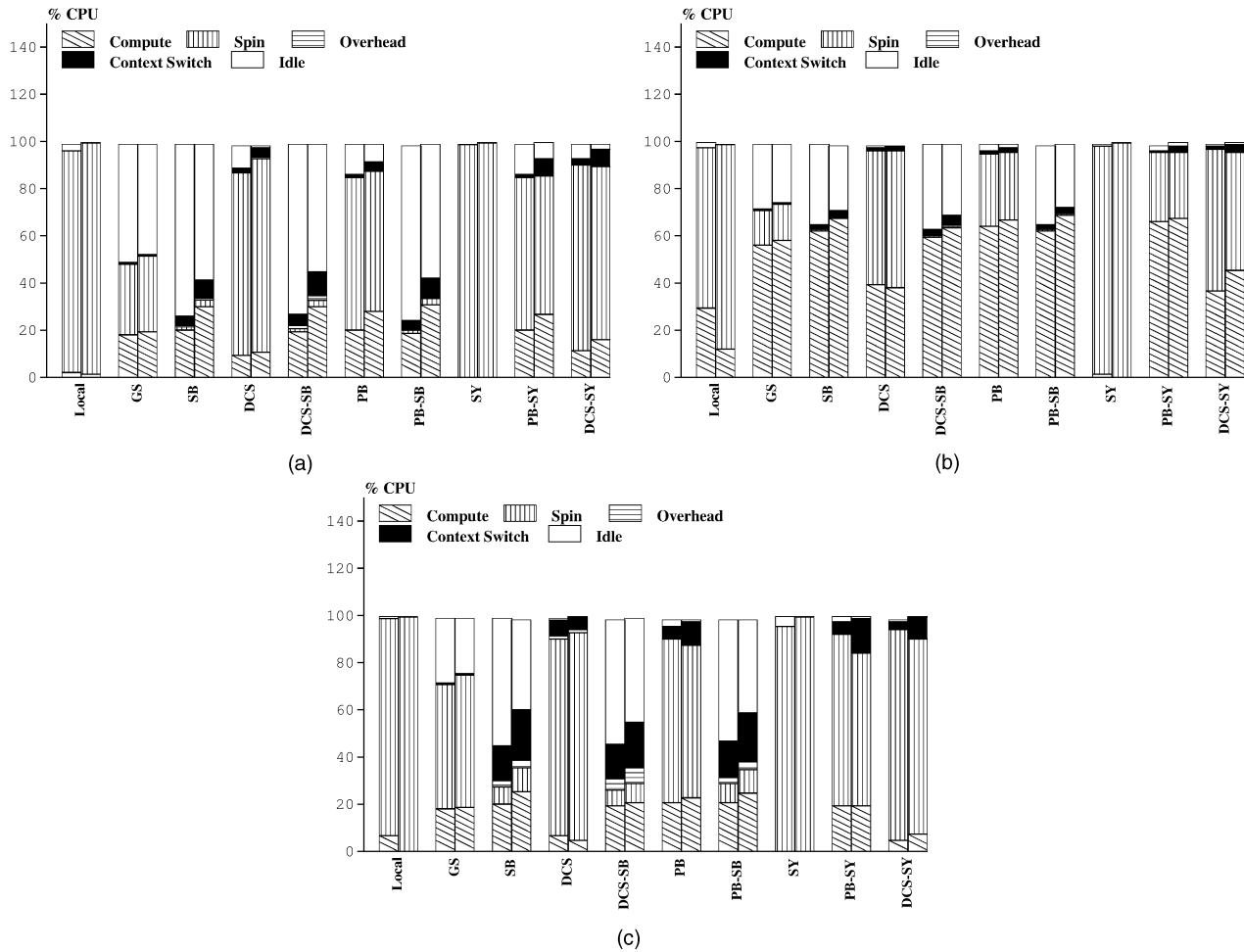


Fig. 5. Impact of MPL: performance profile (NN, $p = 32$, $MPL = 2$ (left bar) and 5 (right bar), $s = 20\%$). (a) WL2. (b) WL4. (c) WL5.

from MPL level 2 to 5 for all three workloads. We can glean this behavior from the performance profile graphs as well. Moving from MPL level 2 to 5 results in lowering the percentage of time spent in useful computation (increasing the spin component in the spinning alternatives). In these experiments again, we find that schemes incorporating PB and SB in some form or the other are able to provide a more scalable improvement in response times with increasing MPLs compared to the rest. At lower MPLs, the SB mechanisms are not able to keep the CPU fully occupied (idle time in the performance profile). This is even more significant for the I/O intensive and communication intensive workloads, which block more frequently. When the MPL is increased, while the context switch times do go up, it is seen that the reduction in the idling is more than adequate to compensate for the considered context switch overheads. With the PB mechanism, we find that it is better than SB at lower MPL levels (because it blocks less frequently), particularly for the communication intensive workload (WL5). SB really needs much higher MPL levels before its performance becomes comparable to PB. From the workload viewpoint, I/O and communication intensive workloads have larger changes (and will better benefit) with MPL compared to the CPU intensive workload.

It is to be noted that, after a point, one can expect response times to eventually go up (even though this is shown for only one of the schemes at $MPL = 16$ in Table 4) due to the overheads dominating any potential benefits.

5.4 Impact of Skewness

Skewness (work imbalance) between the tasks of a job can determine the amount of time that a receiver spins or blocks for a message. Hence, it would be interesting to study how the mismatch of the sends and receives affects the performance of the scheduling mechanisms. Fig. 6 shows the response times and corresponding performance profile for the different schemes with two different skewness values (20 percent and 150 percent) for the CPU intensive workload.

Even if a job is running in isolation on a dedicated system, increasing the skewness would increase its execution time (as per the explanation for skewness given in Section 4.2). This effect is clearly seen for GS on the two workloads where the execution time goes up from $s = 20\%$ to 150% since GS does not have much scope for hiding the impact of this skewness. On the other hand, the dynamic coscheduling mechanisms can potentially better utilize the CPUs than GS for larger skewness values. We can see that most of these schemes are better able to hide the effect of the skewness compared to GS; the spin time increase with a

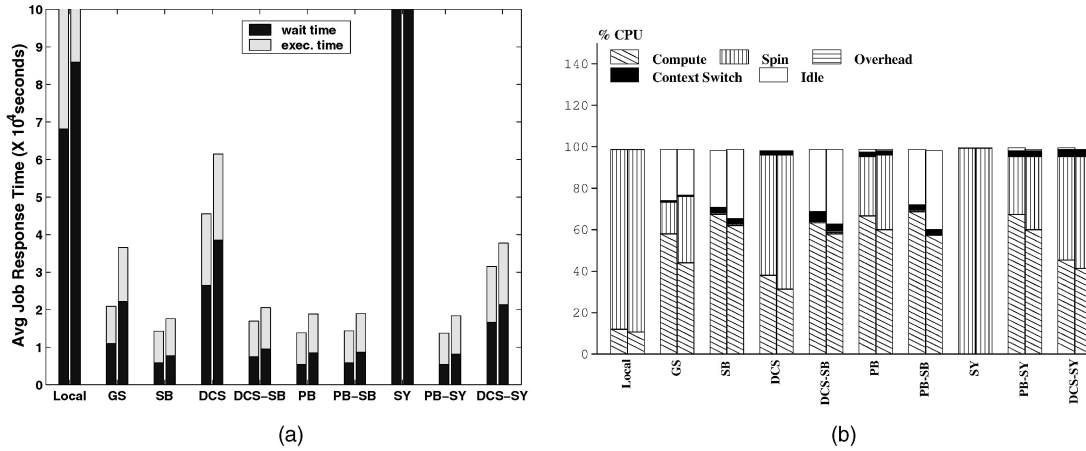


Fig. 6. Impact of skewness (NN, WL4, $p = 32$, $MPL = 5$, $s = 20\%$ (left bar) and 150% (right bar)). (a) Response time. (b) Performance profile.

larger skewness for the spinning mechanisms and the idle time increase for the blocking mechanisms are not as high as the increase in spin times for GS (as seen in performance profiles). Between the spinning and blocking mechanisms, we find that the latter are better able to hide the effect of skewness as can be expected. We have also tried a similar experiment for the communication intensive workload and we have found that the effect of skewness is less pronounced for this workload. This is because skewness is varied as a percentage rather than as an absolute value. As a result, the skew is higher in the compute intensive workload than in the communication intensive workload. The latter inherently synchronizes the tasks more often.

5.5 Impact of System Overheads

In the above discussions, the context switch and interrupt processing costs were kept constant at 200 and 50 microseconds, respectively, for all schemes except GS (where the context switch was set at two milliseconds to account for the synchronization between the CPUs). These parameters can also have an effect on the relative performance of the schemes and Table 5 shows these effects for SB, DCS, DCS-SB, PB, PB-SB, and GS. The response times in Table 5 have been normalized with respect to the first column. It should be noted that there are no interrupts in the PB and GS mechanisms.

TABLE 5
Impact of System Overheads on Response Time
(NN, WL5, $p = 32$, $MPL = 5$, $s = 20\%$)

Scheme	$Q = 200ms$, $CS = 2ms$	$Q = 100ms$, $CS = 1ms$	$Q = 100ms$, $CS = 2ms$
GS	1.000	0.896	0.977
	$CS = 200us$, $I = 50us$	$CS = 100us$, $I = 25us$	$CS = 100us$, $I = 50us$
SB	1.000	0.647	0.659
DCS	1.000	0.902	0.904
DCS-SB	1.000	0.710	0.804
PB-SB	1.000	0.670	0.687
PB	1.000	0.89	

The context switch times in GS are determined not only by the costs of swapping in/out processes at each node, but also due to the explicit synchronization between the nodes. As a result, context switch costs in GS are usually much higher than for the schemes which do not require any explicit synchronization. If we keep the ratio of the time quantum to the context switch overhead the same (100:1), we find that a smaller quantum helps this scheme. This is because, with a smaller quantum, the amount of time in the quantum that is wasted due to blocking (for I/O) or spinning (for a message) becomes smaller. By the time the process gets rescheduled, the operation may be complete, thus allowing better overlap with useful computation. As is to be expected, increasing the overhead percentage due to context switches (i.e., column 3 compared to column 2), extends the execution time (and thereby the response time) of a job.

With the overheads dropping (column 2 compared to column 1), the response times reduce for the different coscheduling heuristics. We find the benefits more significant for the blocking mechanisms compared to those that employ spinning. Blocking executions typically involve many more context switches (each block involves a switch) and, thus, can benefit from the lower associated costs. As interrupt costs go up with no changes to the context switch costs (column 2 to 3), the response times for the mechanisms which incur interrupts (DCS, DCS-SB, SB, and PB-SB) increase. PB makes rescheduling decisions at a much coarser granularity and, thus, is less likely to be affected by system overheads.

5.6 Impact of Communication Pattern

Fig. 7a, Fig. 7b, Fig. 7c, and Fig. 7d show the response times with the different schemes for the four communication patterns that have been considered. Examining the tree pattern, we find that the blocking schemes perform better than the other alternatives while this is not necessarily the case in the other patterns. With this pattern, sends and the corresponding receives at the other end are more likely to be mismatched (i.e., receive may happen much before the send). This is similar to what we observed with a higher skewness factor, making the blocking schemes perform rather well. With the all-to-all pattern, coscheduling

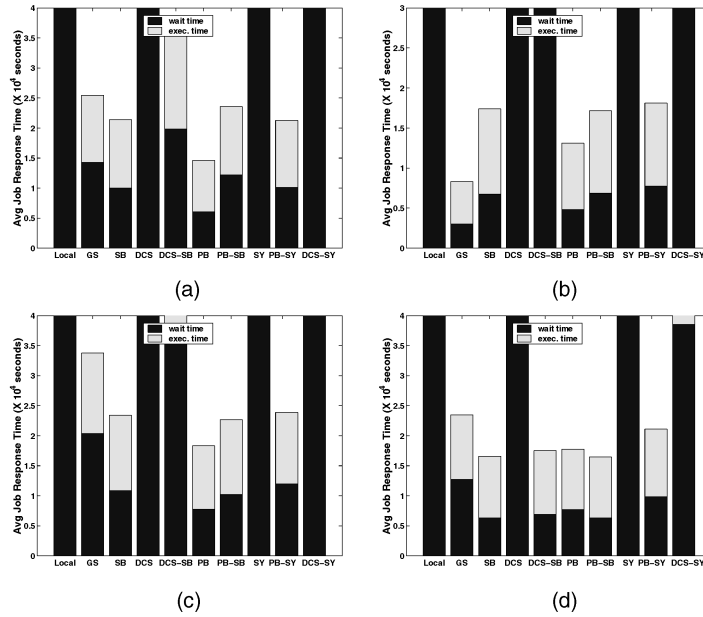


Fig. 7. Impact of the nature of the communication pattern on response time ($p = 32$, $MPL = 5$, $WL5$, $s = 20\%$). (a) NN. (b) AA. (c) Linear. (d) Tree.

becomes extremely important and one cannot tolerate any mismatch of sends and receives. As a result, we find the GS scheme performing the best. Linear and nearest-neighbor show comparable behavior for the different schemes and these patterns fall between the other two (tree and all-to-all) in terms of how important it is to coschedule the processes. Regardless of the communication pattern, we find the PB mechanism uniformly giving good performance.

5.7 Impact of System Size

It is interesting to find out what the improvement in response time would be as we increase the number of CPUs in the system (scalability issue). It is clear that the wait time in the arrival queue would go down. However, increasing system size (and keeping the same number of processes at each node) lessens the likelihood of coscheduled execution

due to the higher possibility of conflicting decisions being made across the nodes (stretching the execution time). The relative changes in these two components will determine the scalability of the system. To evaluate these issues, we conduct experiments with two system sizes (16 and 32 CPUs) and use the same load (with each job demanding at most 16 CPUs) on the two configurations. The resulting response times are shown in Fig. 8, together with its performance profile.

We observe that doubling the system size significantly (more than halving) lowers the wait time (which is the dominating factor of the response time) in the arrival queue. This more than adequately compensates for the minor increases in execution time (due to lower probability of coscheduled execution). Though the results are not explicitly shown here, it should be noted that these

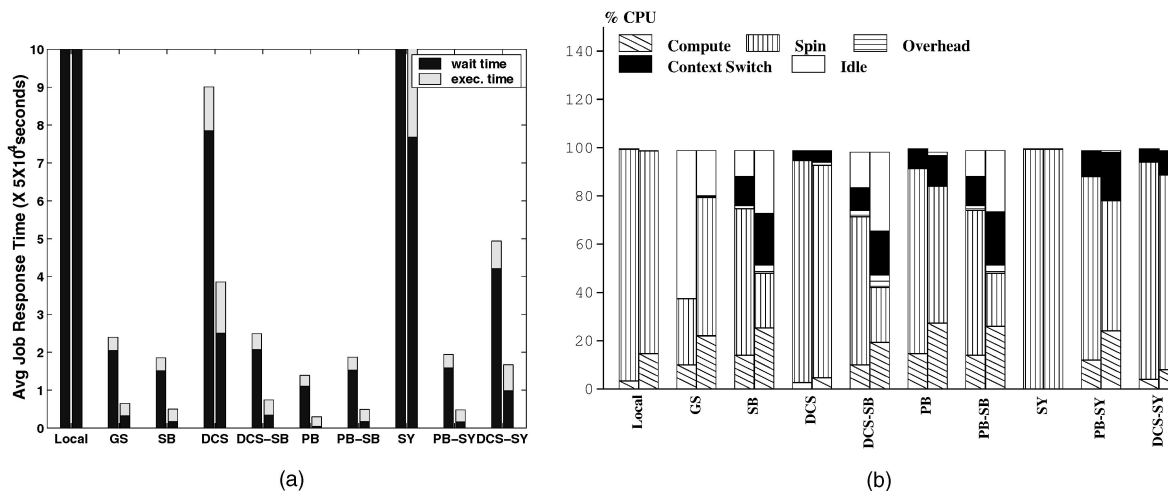


Fig. 8. Impact of the number of CPUs (NN, $WL5$, $p = 16$ (left bar) and 32 (right bar), $MPL = 5$, $s = 20\%$). (a) Response time. (b) Performance profile.

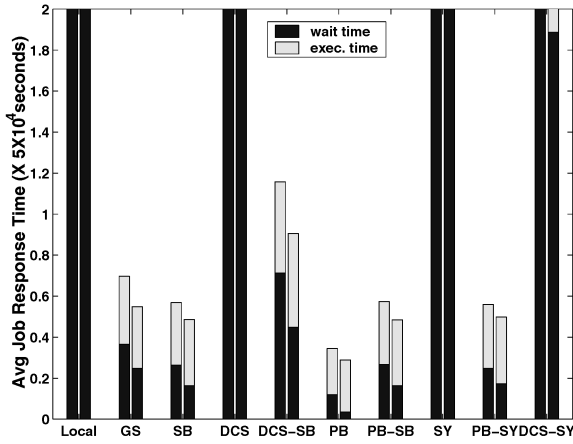


Fig. 9. Impact of partitioning on response time (NN, WL5, $p = 16$ (left bar) and 32 (right bar), $MPL = 5$, $s = 20\%$).

factors are accentuated with higher loads, where the wait times are likely to be even higher.

The performance profile graph also sheds interesting insights on scalability of the system. For the blocking schemes (SB, DCS-SB, and PB-SB), we find that the idle time component of the execution goes up. This is because the same load is exercising a larger system, resulting in lower load (tasks assigned or effective MPL) per node. Lower effective MPLs imply a larger idle time for the blocking mechanisms (since they run out of work) and a lower spin component in the spinning mechanisms due to the higher likelihood of coscheduled execution. In GS, the idle component of the execution goes down with the system size, which may appear somewhat counterintuitive. This is because there is an optimization being performed wherein tasks are replicated across time quanta whenever there are idle slots as long as it does not affect the coscheduled execution of other jobs [40]. Thus, a job could occur in more than one row of the Ousterhout matrix provided it does not require any migration across nodes. Since the overall response times drop, this optimization contributes to lowering the idleness component of the system.

Another exercise that would be useful to a system administrator/designer toward configuring a given system for maximum rewards is: Does it make sense to

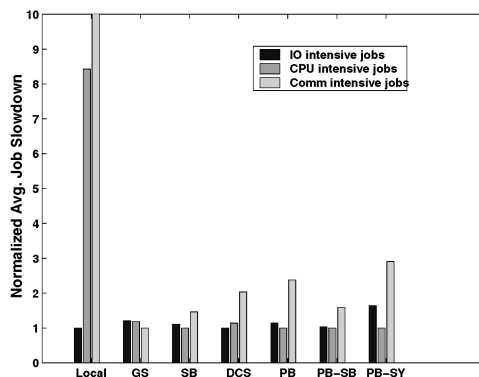
partition a given cluster into smaller fragments and specifically assign jobs to these fragments or would it be better to treat it as one big machine? To investigate this issue, we run two sets of experiments, first with a single cluster of 32 nodes and second with two clusters of 16 nodes each, and exercise them with the same workload. Although there has been prior work on better job assignment schemes using their execution time [6], in this study, we just assign each job to the partition with the least load since we do not have execution times beforehand.

Fig. 9 compares the response times with these two executions (right bar for 32-node cluster and left bar for the partitioned system). We can observe that, in nearly all cases, the savings in the wait time offset any stretching of the execution time, suggesting that it is not a good idea to partition the cluster. Similar results were observed for other system sizes/configurations as well.

5.8 Fairness

Until now, we have purely looked at performance issues for the different scheduling mechanisms. Fairness to different jobs is another important criteria from both the system administrator and user perspectives. There are several ways of analyzing the fairness of the system. In this paper, we look at the fairness issue from the viewpoint of how much the scheduler is biased/discriminatory based on the nature of the job. If a scheduler is biased towards/against any particular type of job, then it is likely to experience a lower/higher slowdown. So, comparing the differences in slowdown (or, rather, the coefficient of variation of slowdown) to the different job classes in a mixed workload can be used to argue the fairness of a scheduling strategy. To examine the fairness for different job types, we run workload WL8 (with $s = 0$ to ensure that all job types have the same average execution time on a dedicated system) on a 32-node system and classify the jobs within this workload into those that are CPU intensive, communication intensive, and I/O intensive. Fig. 10 shows the average job slowdown (normalized for each with respect to the job class with the least slowdown) and its coefficient of variation over these three job classes for Local, GS, SB, DCS, PB, PB-SB, and PB-SY.

As expected, GS is the most fair, with all three job classes experiencing equal slowdowns (low coefficient of



(a)

Scheme	Coeff. of Var.
Local	1.632
GS	0.099
SB	0.203
DCS	0.403
PB	0.504
PB-SB	0.274
PB-SY	0.524

(b)

Fig. 10. Fairness: slowdown of different job types (NN, WL8, $p = 32$, $MPL = 5$, $s = 0\%$). (a) Slowdown. (b) Coefficient of variation of slowdown.

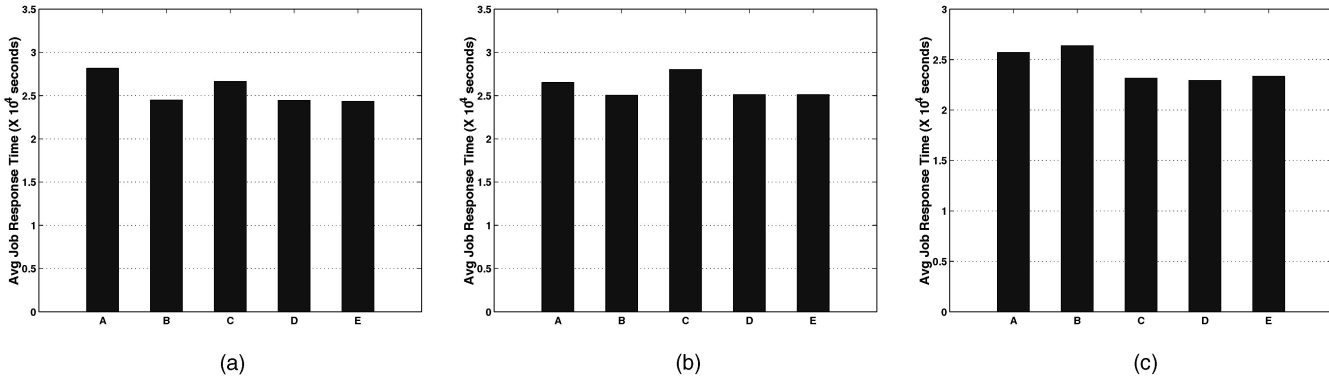


Fig. 11. PB heuristics (NN, $p = 32$, $MPL = 5$, $s = 20\%$). (a) WL4. (b) WL8. (c) WL5.

variation). At the other end of the spectrum is Local, where I/O intensive jobs get boosted much more often than the others (I/O completion is the only criterion for boost in priority in the Solaris scheduler at each node). Of the remaining schemes, the blocking schemes (SB and PB-SB) are more fair than their spinning counterparts (DCS, PB, and PB-SY). We will revisit the fairness issue for the PB mechanism again in the next section after we discuss the different heuristics that one can use for PB.

5.9 PB Heuristics

The PB mechanism is a powerful one, allowing us to employ several heuristics to decide on who to boost before the next scheduling decision. Until now, we have used only one heuristic for this mechanism and it will be interesting to explore alternatives for this choice (particularly if we can do even better). When the PB mechanism is invoked, a ready/executing process can be in one of four states:

- S1. It is in the compute/send phase and there are no pending (unconsumed) messages.
- S2. It is in the compute/send phase and there is at least one pending message.
- S3. It is in the receive phase and the message for which this receive is posted has arrived (and is not yet consumed).
- S4. It is in the receive phase and the message for which it is waiting has not yet arrived.

Identifying these states can help us decide whether to boost a process and to decide on a priority order if there are choices. Of these four states, it intuitively makes the most sense to give the highest priority to S3 and S2 since a pending message implies that the receiver should soon be scheduled as per the DCS heuristic. Another observation is that a process in S4 should never get boosted since it will waste CPU cycles spinning redundantly (so we can confine ourselves to states S1, S2, and S3). Based on the relative priority of these states, we identify the following heuristics:

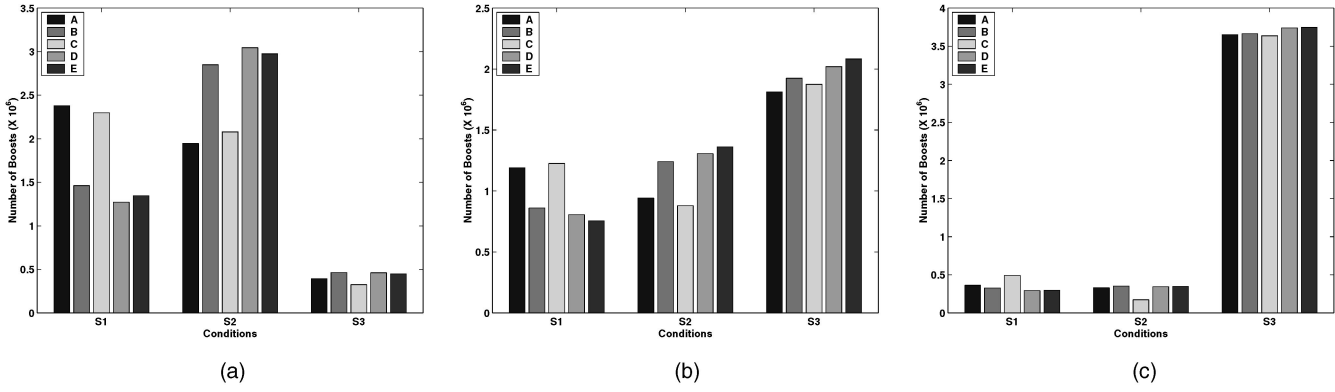
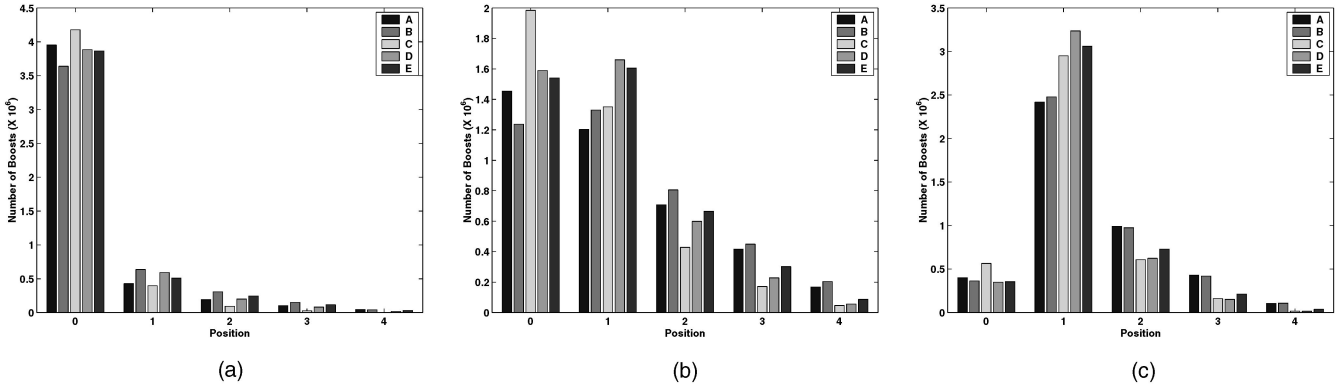
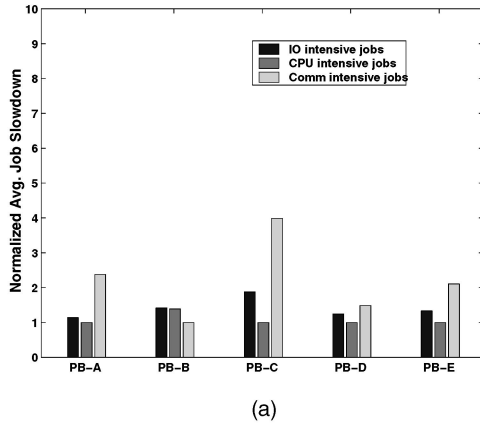
- A : $S3 \rightarrow \{S2, S1\}$,
- B : $S3 \rightarrow S2 \rightarrow S1$,
- C : $\{S3, S2, S1\}$,
- D : $\{S3, S2\} \rightarrow S1$,
- E : $S2 \rightarrow S3 \rightarrow S1$,

where \rightarrow denotes the priority ordering. For instance, scheme A denotes that all processes are first checked to see if any of them is in state S3 in round-robin order, stopping (and boosting) at the first process in this state. Only when none of them are in state S3 are the conditions S2 and S1 evaluated for all of them (again in round-robin order) to see if a process is in either of these two states (and, if so, stopping and boosting that process). The other schemes can be similarly described. In all the previous performance results, heuristic A has been used.

Fig. 11 compares the response times of the five heuristics for the WL4 (CPU intensive), WL8 (mixed), and WL5 (communication intensive) workloads. For WL4, we find that heuristics B, D, and E outperform the other two. Even though the workload is CPU intensive, it is still important to boost processes with one or more pending messages (to increase the likelihood of coscheduling). Of the five heuristics, B, D, and E are likely to give more boosts to such processes than the others. The likelihood of S3 being the chosen condition for boosting a process is rather low in this workload (being CPU intensive). In A and C, S1 has a higher chance (than in B, D, and E) of being the condition used for boosting (because of the priority order). This intuition is also confirmed by the statistics from the experiments, shown in Fig. 12a, which gives the number of boosts given by the different heuristics based on the condition causing the boosts.

We find a similar argument (Fig. 11b) holding for the mixed workload (WL8) as well, where we find that heuristics B, D, and E give more boosts to a job with a pending message (states S2 or S3) than the other two (Fig. 12b).

Moving to WL5, we find that heuristics C, D and E are better than the first two which can be explained as follows. For this workload, S3 is most likely to be the chosen condition for boosting a process since it is very likely that the current process is in the S4 state when the periodic boost mechanism is invoked and the next immediate process in round-robin order is chosen (which is likely to be in the S3 state). This intuition is borne out by the statistics shown in Fig. 12c, where S3 gives the maximum number of boosts, and in Fig. 13c, which shows the number of boosts given to a process that is i positions away (on the x-axis) from the currently scheduled process in round-robin order. In this workload, it is more important to adhere to some order that is common across the different CPUs to increase the

Fig. 12. Conditions causing boosts (NN, $p = 32$, $MPL = 5$, $s = 20\%$). (a) WL4. (b) WL8. (c) WL5.Fig. 13. Position of process being boosted (NN, $p = 32$, $MPL = 5$, $s = 20\%$). (a) WL4. (b) WL8. (c) WL5.

Scheme	Coeff. of Var.
PB-A	0.504
PB-B	0.184
PB-C	0.671
PB-D	0.196
PB-E	0.384

(b)

Fig. 14. PB fairness: slowdown of different job types (NN, WL8, $p = 32$, $MPL = 5$, $s = 0\%$). (a) Slowdown. (b) Coefficient of slowdown.

likelihood of the processes belonging to the same job being scheduled at the same time (since they are all likely to be in the same state). Of the five heuristics, we find the latter three have a higher likelihood of selecting the next process in round-robin order and, thus, have a lower response time.

In essence, two factors are to be considered in the performance of these heuristics. First is the condition causing the boost and the second being the position (from the current process in round-robin order) of the process that is boosted. In WL4 and WL8, the first factor is more important since the processes are likely to be in different states while the latter factor is more important in WL5, where all (or most) the other processes at a node are likely to be in the same state (namely S3).

Overall, we find that heuristics **D** and **E** perform the best across all three workloads. We have also tried heuristics that are based on the number of pending messages and time-stamp of message arrivals, but we have still found **D** and **E** outperforming those as well.

We next examine the fairness of the different heuristics using the same experiment mentioned in Section 5.8. As before, Fig. 14 shows the normalized job slowdown and its corresponding coefficient of variation, respectively, for heuristics **A** through **E**. We find that **B** and **D** are much more fair than the others, with their coefficient of variation comparable to that of **SB** (Fig. 10). We noticed that **D** was giving low response times as well, suggesting that we use this heuristic for the PB mechanism.

TABLE 6
PB': Coefficient of Variation of Slowdown
(NN, WL8, $p = 32$, $MPL = 5$, $s = 0\%$)

Scheme	Coeff. of Var.
PB-A'	0.101
PB-B'	0.123
PB-C'	0.082
PB-D'	0.058
PB-E'	0.280

The fairness issue suggests an even further improvement to the five PB heuristics, where, instead of just checking the condition (S1 through S4) for boosting, we also ensure that the selected process does not get more than a fair share of the CPU (else an alternate is chosen) [4]. So, we can revise the five heuristics with this additional check and we refer to these new mechanisms as PB-A' through PB-E', respectively. For instance, PB-E' can be explained as follows. First, we find all processes that are in the S2 state. If there is at least one such process, then we use the process with the minimum ratio of the CPU time that it has been allocated to the time since it entered the system (i.e., percentage of the CPU that it has received since it started) as the target for boosting. However, if no process is in the S2 state, then we check for the S3 state and do likewise. Table 6 shows the coefficient of variation of these five variations and it can be seen that the first four of these are as fair as GS (compare to Fig. 10). Further, the response times of the enhanced mechanisms are plotted in Fig. 15 in relation with the plain A through E heuristics to show that we are not losing out on response times at the cost of fairness.

6 CONCLUDING REMARKS AND FUTURE WORK

Advances in user-level networking (ULN) allows us to explore a new domain of dynamic coscheduling mechanisms that offer potential improvement over conventional scheduling strategies (such as space sharing or gang scheduling) for parallel machines. These mechanisms become even more important for cluster environments where applications with diverse characteristics and QoS requirements coexist. Until now, the understanding and knowledge of the relative performance of dynamic coscheduling mechanisms is rather limited [32], [20], [7], [1], [31]. While our previous work [20] did a preliminary examination of these dynamic coscheduling mechanisms, the study was rather limited due to the workloads that were considered, the experiments that were conducted, the inflexibility of the underlying system in modulating parameters, and the performance metrics that were evaluated. For the very first time, this paper has presented a comprehensive evaluation study of the different dynamic coscheduling alternatives using realistic and dynamic workloads with varying job sizes, execution times, and characteristics, in studying the impact of different system and workload parameters on numerous performance metrics. This exercise has required the development of a comprehensive and flexible simulator, which has itself been a substantial effort. Using this

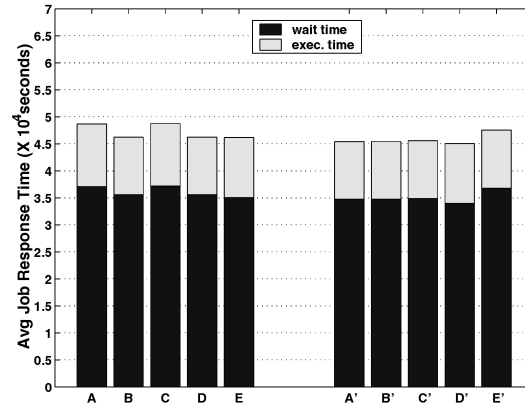


Fig. 15. PB': Response time (NN, WL8, $p = 32$, $MPL = 5$, $s = 20\%$).

simulator, we have been able to examine the suitability of different scheduling mechanisms under varying conditions.

There is clearly a great need for some coordinated scheduling effort between the nodes to accommodate the parallel jobs. This is evidenced by the poor performance of Local, which leaves it to the individual operating systems (on these nodes) to independently make their own decisions. Regardless of the operating condition (unless there is little or no communication at all), Local is not a viable option.

We have also observed that the dynamic coscheduling mechanisms are, in general, a better alternative to the conventional solution, Gang Scheduling (GS), employed today in many parallel machines [35], [3], [14], [34], [11]. In the first place, GS requires slots to be available on the required number of CPUs within the same time quantum (row of the Ousterhout matrix), which is not necessary in the dynamic coscheduling mechanisms. As a result, wait times in the arrival queue are typically higher in GS. Further, GS executions do not fare as well as the dynamic coscheduling alternatives when:

1. The skewness between the tasks of a job is high (GS is not able to hide the spin times as well as the dynamic coscheduling mechanisms),
2. The jobs are I/O intensive (GS wastes the rest of the time quantum),
3. The multiprogramming level is high (system overheads become more important), and
4. When the costs for explicitly synchronizing the nodes between time quanta becomes high (perhaps, for large systems).

The one advantage that GS has over its dynamic coscheduling counterparts is that it is fair across different job types. It should be noted that GS could be enhanced further, as others [12], [41] have suggested, which may tip the balance in favor of GS when the above factors are not overwhelming.

Comparing the Spin Block (SB), Spin Yield (SY), and Demand-based Coscheduling (DCS) approaches within the dynamic coscheduling domain, we find that blocking (SB) seems to be a better option most of the time. This is particularly true for: 1) higher multiprogramming levels (since the probability of finding alternate useful work on

blocking is higher), 2) higher skewness (since they are less sensitive to skews between tasks), and 3) communication intensive workloads (since it better approaches coscheduling). The reason for the last observation was also made in [20], where it was mentioned that it is important to schedule the receiver of a message as early as possible while reducing spin times. DCS only addresses the first goal and SY addresses only the second goal. SB, on the other hand, gives the blocked process a priority boost on message arrival to attempt achieving both goals.

Uniformly, we find PB outperforming SB (and the other schemes). It is able to get the advantages of spinning (to avoid interrupt processing costs) and is able to relinquish the CPU whenever needed. While the performances of PB and SB are comparable for CPU and I/O intensive workloads, PB is much better when workloads are communication intensive. In addition, PB is preferable over SB for: 1) smaller MPLs (SB becomes comparable to PB only when the MPL is larger) and 2) higher system overheads (as switching/interrupt costs increase, the advantage of PB over SB is amplified since SB incurs more system overheads). It is only when the skewness gets large that SB has a slight edge and, in such cases, one could resort to a combination (PB-SB) of these mechanisms.

In this paper, we have also opened a new topic for further research that offers a lot of potential for improving PB even further. We have been able to identify a set of heuristics that can be used to decide on when and whom to boost, developed a classification of these heuristics, and presented results showing the performance potential of this avenue of research. These heuristics take message arrival and process state into consideration in picking a candidate for a priority boost. Further, one could also augment these heuristics with a fair share (of the CPU) mechanism without compromising on the performance (response time) benefits. The resulting heuristic (PB-D') has significantly lower response times than the other coscheduling alternatives and is even more fair than GS.

While this paper is a comprehensive and novel exercise in the evaluation of numerous parallel coscheduling mechanisms using a unified framework, there are several interesting directions for future research. It would be interesting to theoretically model the scheduling mechanisms/heuristics so that we can understand how well they perform (and how close they are to the optimal). Recently, a related study [25] showed how buffering and latency tolerance (separating the posting of a receive from the time when the receiver should actually block) can be used to minimize the impact of noncoscheduled execution. However, their study assumes that all the communication latency can be hidden, in which case coscheduling is not important. We propose exploring the possibility of exploiting communication slackness (gap between posting of a receive and the time the data is actually needed) within the PB heuristics for even better performance. We have not considered memory constraints in this work and incorporating such constraints to determine the multiprogramming level and associated overheads on the dynamic coscheduling strategies would also be very interesting [24], [15], [26].

This study has taken an important step in improving the responsiveness of a high performance (parallel) computing platform which has traditionally been optimized for throughput (for scientific applications). This can make it more suitable for the emerging use of clusters for web and database services. We propose exploring this issue using workloads from the emerging applications. Further, guaranteed (soft or hard) service [29], [30], [28], [16] is also expected to become an important Quality-of-Service demand from such an environment as clusters take on the challenges of multimedia and real-time applications. Extending the PB mechanism using well-known rate-based and real-time scheduling techniques is part of our ongoing and future work.

ACKNOWLEDGMENTS

This research has been supported in part by US National Science Foundation grants CCR-9988164, CCR-9900701, DMI-0075572, Career Award MIP-9701475, and equipment grant EIA-9818327.

REFERENCES

- [1] A.C. Arpaci-Dusseau, D.E. Culler, and A.M. Mainwaring, "Scheduling with Implicit Information in Distributed Systems," *Proc. ACM SIGMETRICS 1998 Conf. Measurement and Modeling of Computer Systems*, 1998.
- [2] M. Banikazemi, V. Moorthy, D.K. Panda, L. Herger, and B. Abali, "Efficient Virtual Interface Architecture Support for the IBM SP Switch-Connected NT Clusters," *Proc. Int'l Symp. Parallel and Distributed Processing Systems*, pp. 33-42, May 2000.
- [3] E. Barton, J. Cownie, and M. McLaren, "Message Passing on the Meiko CS-2," *Parallel Computing*, vol. 20, no. 4, Apr. 1994.
- [4] P. Brinch-Hansen, "An Analysis of Response Ratio Scheduling," *Int'l Federation Information Processing Congress*, vol. TA-3, pp. 150-154, Aug. 1971.
- [5] M. Buchanan and A. Chien, "Coordinated Thread Scheduling for Workstation Clusters under Windows NT," *Proc. USENIX Windows NT Workshop*, Aug. 1997.
- [6] M. Crovella, M. Harchol-Balter, and C. Murta, "Task Assignment in a Distributed System: Improving Performance by Unbalancing Load," *Proc. ACM Sigmetrics Conf. Measurement and Modeling of Computer Systems*, pp. 268-269, June 1998.
- [7] A.C. Dusseau, R.H. Arpaci, and D.E. Culler, "Effective Distributed Scheduling of Parallel Workloads," *Proc. ACM SIGMETRICS 1996 Conf. Measurement and Modeling of Computer Systems*, pp. 25-36, 1996.
- [8] D.G. Feitelson, "A Survey of Scheduling in Multiprogrammed Parallel Systems," Technical Report Research Report RC 19790(87657), IBM T.J. Watson Research Center, Oct. 1994.
- [9] D.G. Feitelson and L. Rudolph, "Coscheduling Based on Run-Time Identification of Activity Working Sets," Technical Report Research Report RC 18416(80519), IBM T.J. Watson Research Center, Oct. 1992.
- [10] D.G. Feitelson and L. Rudolph, "Gang Scheduling Performance Benefits for Fine-Grained Synchronization," *J. Parallel and Distributed Computing*, vol. 16, no. 4, pp. 306-318, Dec. 1992.
- [11] H. Franke, J. Jann, J.E. Moreira, P. Pattnaik, and M.A. Jette, "Evaluation of Parallel Job Scheduling for ASCI Blue-Pacific," *Proc. Supercomputing*, Nov. 1999.
- [12] G. Ghare and S. Leutenegger, "The Effect of Correlating Quantum Allocation and Job Size for Gang Scheduling," *Proc. Fifth Ann. Workshop Job Scheduling Strategies for Parallel Processing*, Apr. 1999.
- [13] A. Hori, H. Tezuka, and Y. Ishikawa, "Global State Detection Using Network Preemption," *Proc. Int'l Parallel Processing Symp. Workshop Job Scheduling Strategies for Parallel Processing*, pp. 262-276, Apr. 1997.
- [14] Intel Corporation, *Paragon User's Guide*, 1993.

- [15] N. Islam, A.L. Prodromidis, M.S. Squillante, L.L. Fong, and A.S. Gopal, "Extensible Resource Management for Cluster Computing," *Proc. 17th Int'l Conf. Distributed Computing Systems*, pp. 561-568, 1997.
- [16] H. Kaneko, J.A. Stankovic, S. Sen, and K. Ramamritham, "Integrated Scheduling of Multimedia and Hard Real-Time Tasks," *Proc. 17th IEEE Real-Time Systems Symp.*, June 1996.
- [17] D. Lifka, "The ANL/IBM SP Scheduling System," *Proc. Int'l Parallel Processing Symp. Workshop Job Scheduling Strategies for Parallel Processing*, pp. 295-303, Apr. 1995.
- [18] J.M. Mellor-Crummey and M.L. Scott, "Algorithms for Scalable Synchronization on Shared-Memory Multiprocessors," *ACM Trans. Computer Systems*, vol. 9, no. 1, pp. 21-65, Feb. 1991.
- [19] J.E. Moreira, H. Franke, W. Chan, L.L. Fong, M.A. Jette, and A. Yoo, "Gang-Scheduling System for ASCI Blue-Pacific," *Proc. Seventh Int'l Conf. High-Performance Computing and Networking (HPCN '99)*, pp. 831-840, Apr. 1999.
- [20] S. Nagar, A. Banerjee, A. Sivasubramaniam, and C.R. Das, "A Closer Look at Coscheduling Approaches for a Network of Workstations," *Proc. 11th Ann. ACM Symp. Parallel Algorithms and Architectures*, pp. 96-105, June 1999.
- [21] S. Nagar, A. Banerjee, A. Sivasubramaniam, and C.R. Das, "Alternatives to Coscheduling a Network of Workstations," *J. Parallel and Distributed Computing*, vol. 59, no. 2, pp. 302-327, Nov. 1999.
- [22] J.K. Ousterhout, "Scheduling Techniques for Concurrent Systems" *Proc. Third Int'l Conf. Distributed Computing Systems*, pp. 22-30, May 1982.
- [23] S. Pakin, M. Lauria, and A. Chien, "High Performance Messaging on Workstations: Illinois Fast Messages (FM) for Myrinet," *Proc. Supercomputing '95*, Dec. 1995.
- [24] V.G.J. Peris, M.S. Squillante, and V.K. Naik, "Analysis of the Impact of Memory in Distributed Parallel Processing Systems," *Proc. ACM SIGMETRICS Conf. Measurement and Modeling of Computer Systems*, pp. 5-18, May 1994.
- [25] F. Petrini and W. Feng, "Buffered Coscheduling: A New Method for Multitasking Parallel Jobs on Distributed Systems," *Proc. Int'l Parallel and Distributed Processing Symp.*, pp. 439-444, May 2000.
- [26] F. Petrini and W. Feng, "Time-Sharing Parallel Jobs in the Presence of Multiple Resource Requirements," *Proc. Sixth Ann. Workshop Job Scheduling Strategies for Parallel Processing*, pp. 71-92, May 2000.
- [27] R. Poovendran, P. Keleher, and J.S. Baras, "A Decision-Process Analysis of Implicit Coscheduling," *Proc. Int'l Parallel and Distributed Processing Symp.*, pp. 1115-120, May 2000.
- [28] K. Ramamritham, C. Shen, O. Gonzalez, S. Sen, and S.B. Shirkurkar, "Using Windows NT for Real-Time Applications: Experimental Observations and Recommendations," *Proc. Fourth IEEE Real-Time Technology and Applications*, June 1998.
- [29] K. Ramamritham and J.A. Stankovic, "Scheduling Algorithms and Operating Systems Support for Real-Time Systems," *Proc. IEEE*, pp. 55-67, Jan. 1994.
- [30] C. Shen, O. Gonzalez, K. Ramamritham, and I. Mizunuma, "User Level Scheduling of Communication for Real-Time Tasks," *Proc. Fifth IEEE Real-Time Technology and Applications*, June 1999.
- [31] P.G. Sobalvarro, "Demand-Based Coscheduling of Parallel Jobs on Multiprogrammed Multiprocessors," PhD thesis, Dept. of Electrical Eng. and Computer Science, Massachusetts Inst. of Technology, Jan. 1997.
- [32] M.S. Squillante, Y. Zhang, A. Sivasubramaniam, N. Gautam, H. Franke, and J. Moreira, "Analytic Modeling and Analysis of Dynamic Coscheduling for a Wide Spectrum of Parallel and Distributed Environments," Technical Report CSE-01-004, Computer Science and Eng. Dept., Pennsylvania State Univ., Feb. 2001.
- [33] R. Subrahmaniam, "Implementing Coscheduling Heuristics for Windows NT Clusters," master's thesis, Dept. of Computer Science and Eng., Pennsylvania State Univ., October 1999.
- [34] K. Suzuki and D. Walsh, "Implementing the Combination of Time Sharing and Space Sharing on AP/Linux," *Proc. Int'l Parallel Processing Symp. Workshop Job Scheduling Strategies for Parallel Processing*, pp. 83-97, Mar. 1998.
- [35] Thinking Machines Corp., *The Connection Machine CM-5 Technical Summary*, Oct. 1991.
- [36] L.W. Tucker and G.G. Robertson, "Architecture and Applications of the Connection Machine," *Computer*, vol. 21, no. 8, pp. 26-38, Aug. 1988.
- [37] Specification for the Virtual Interface Architecture, <http://www.viarch.org>. 2001.
- [38] T. von Eicken, A. Basu, V. Buch, and W. Vogels, "U-Net: A User-Level Network Interface for Parallel and Distributed Computing," *Proc. 15th ACM Symp. Operating System Principles*, Dec. 1995.
- [39] F. Wang, M. Papaefthymiou, and M. Squillante, "Performance Evaluation of Gang Scheduling for Parallel and Distributed Multiprogramming," *Proc. Int'l Parallel Processing Symp. Workshop Job Scheduling Strategies for Parallel Processing*, pp. 277-298, Apr. 1997.
- [40] Y. Zhang, H. Franke, J. Moreira, and A. Sivasubramaniam, "Improving Parallel Job Scheduling by Combining Gang Scheduling and Backfilling Techniques," *Proc. Int'l Parallel and Distributed Proc. Symp.*, pp. 133-142, May 2000.
- [41] B. Zhou, R. Brent, C. Johnson, and D. Walsh, "Job Re-Packing for Enhancing the Fine and Coarse Grain Parallel Processes," *Proc. Fifth Ann. Workshop Job Scheduling Strategies for Parallel Processing*, Apr. 1999.



Yanyong Zhang received the BS degree in computer science from the University of Science and Technology of China in 1997. Since 1997, she has been pursuing the PhD degree in computer science and engineering at Pennsylvania State University. She spent the summers of 1999 and 2000 interning at the IBM T.J. Watson Research Center. Her research interests are in operating systems and performance evaluation with an emphasis on parallel scheduling and clusters. She is a recipient of a 2001 IBM Research Fellowship and is a member of the IEEE and the ACM.

Anand Sivasubramaniam received the BTech degree in computer science from the Indian Institute of Technology, Madras, in 1989, and the MS and PhD degrees in computer science from the Georgia Institute of Technology in 1991 and 1995, respectively. He has been on the faculty at The Pennsylvania State University since Fall 1995, where he is currently an associate professor. His research interests are in computer architecture, operating systems, performance evaluation, and applications for high performance computer systems. More recently, he has also been examining resource-constrained computing issues. His research has been funded by the US National Science Foundation through several grants (including the CAREER award) and support from industries, including IBM and Unisys Corp. He has published extensively in several leading journals and conferences, and is currently serving as an associate editor of the *IEEE Transactions on Computers*. He is a member of the IEEE, IEEE Computer Society, and ACM.



José Moreira received BS degrees in physics and electrical engineering in 1987 and the MS degree in electrical engineering in 1990, all from the University of Sao Paulo, Brazil. He received the PhD degree in electrical engineering from the University of Illinois at Urbana-Champaign in 1995. He is a research staff member and manager, Blue Gene System Software, at the IBM Thomas J. Watson Research Center. Since joining the IBM T.J. Watson Research Center in

1995, he has been involved in several high-performance computing projects, including the Teraflop-scale ASCI Blue-Pacific and ASCI White. He is the author of more than 30 publications on high-performance computing. He participated in the development of a system for dynamic reconfiguration of HPF and MPI applications on a parallel computer (DRMS), the development of a gang-scheduling system for large parallel supercomputers (GangLL), and the development of compilers and libraries for high-performance technical computing in Java. He is currently the spec lead for the Java Community Process proposal to add multidimensional arrays. In the Blue Gene project, he is co-leader of the efforts to develop system software for a Petaflop-scale machine. He is a member of the IEEE.



Hubertus Franke received a summa cum laude Diplom Informatik degree in 1987 from the Technical University of Karlsruhe, Germany. He received the MS and PhD degrees in electrical engineering from Vanderbilt University, in 1989 and 1992, respectively. In 1988, he was the recipient of the computer science achievement award of the Computer Science Research Center, Karlsruhe. Since 1993, he has been a research staff member at the IBM T.J. Watson Research Center, where he currently manages the Enterprise Linux Group. Since joining IBM Research he has been involved in various high-performance computing, operating systems, and computer architecture projects. He was the initial developer of the IBM SP2 MPI communications software and gang scheduling system. He participated in the K42 operating system, NUMA partitioning, memory compression and Linux scalability projects. He is the coauthor of more than 50 publications in these fields and is a member of the IEEE.

► For more information on this or any computing topic, please visit our Digital Library at <http://computer.org/publications/dlib>.