brought to you by T CORE

DAA/ LANGLEY

IN-62 64478-CR

P.417

Parallel Discrete Event Simulation: A Shared Memory Approach

> Daniel A. Reed Department of Computer Science University of Illinois Urbana, Illinois 61801

Allen D. Malony Center for Supercomputing Research and Development University of Illinois Urbana, Illinois 61801

Bradley D. McCredie Department of Electrical and Computer Engineering University of Illinois Urbana, Illinois 61801

ABSTRACT

With traditional event list techniques, evaluating a detailed discrete event simulation model can often require hours or even days of computation time. Parallel simulation mimics the interacting servers and queues of a real system by assigning each simulated entity to a processor. By eliminating the event list and maintaining only sufficient synchronization to insure causality, parallel simulation can potentially provide speedups that are linear in the number of processors. We present a set of shared memory experiments using the Chandy-Misra distributed simulation algorithm to simulate networks of queues. Parameters of the study include queueing network topology and routing probabilities, number of processors, and assignment of network nodes to processors. These experiments show that Chandy-Misra distributed simulation is a questionable alternative to sequential simulation of most queueing network models.

This work was supported in part by NSF Grant Number DCR 84-17948 and NASA Contract Number NAG-1-613.

(NASA-CR-180616) PARALLEL DISCHETE EVENT N87-26576 SIMULATION: A SHAFED MEMOBY AFFFCACH (Illinois Univ.) 47 p Avail: NTIS HC A03/MF A01 CSCL 09B Unclas G3/62 0064478

1864743.

Introduction

Historically, there have been two major techniques for modeling systems: queueing theory and discrete event simulation. When effective, queueing theoretic techniques can quickly provide mathematical insight into the behavior of systems over a broad range of parameter values. Their major limitation is the number of restrictive assumptions that must be satisfied to insure accuracy. Conversely, simulation models can mimic -a real-world system as closely as understanding permits and needs require. However, highly detailed simulation models can be computationally taxing. Computer systems simulations are particularly vexing because simulated events occur on a millisecond or microsecond time scale, often for many simulated minutes.

For example, simulating the behavior of processor executing a user or system program may involve millions or even tens of millions of events. In one architecture performance study, we recently examined the performance of allocation strategies for register windows in reduced instruction set computers (RISCs) [PaSe82, Patt85] as a function of multiprogramming level [KoRW86]. This analysis required instruction-level simulation for many different program mixes and consumed many hours of processor time.

Simulation of complex (VLSI) digital circuits for logic verification and fault analysis is another example of the computational constraints imposed on simulation of computing components. Although such simulations can consume *months* of machine time [Pfis82, FrWW84], designers have little choice; an untested design is unacceptable. Moreover, simulation complexity continues to increase dramatically; technology advances are doubling the number of circuits per chip every 1-2 years.

At a much higher level than logic design, we recently encountered difficulties while studying multicomputer networks Reed83, ReFu86, designed, ironically, to solve computationally

intensive problems. Briefly, a multicomputer network is a large number of interconnected computing nodes that asynchronously cooperate via message passing to execute the tasks of parallel programs.¹ Many design issues must be resolved before constructing a multicomputer network (e.g., the relative speeds of computation processors and internode communication links, topology of connecting communication links, buffer requirements for messages, and memory sizes). Although some of these issues can be attacked analytically, most are analytically intractable and can only be resolved via simulation. A parametric simulation study, whose individual simulation runs cover several minutes of simulated time, typically requires several hundred hours of processor time.

Although each of the three preceding examples, processor simulation, logic simulation, and network simulation, is very different, they share a common need for faster simulation techniques. Processor simulation reflects the fetch/decode/execute cycle of instruction execution and is, by its nature, sequential; parallelizing this application is the subject of architectural research. Circuit simulation, although clearly amenable to parallel processing [Pfis82], typically involves synchronous activation of many entities. In contrast, network simulation is typically asynchronous.

Prior Work

It might initially appear that evaluating models of many complex systems is both analytically and computationally intractable. However, recent developments have suggested that the computation time for *some* simulations can be reduced via either vector processing [ChBr83] or distributed simulation [PeWM79, ChMi81, JeSo85].

¹Hypercubes [Seit85] are a special case of multicomputer networks.

Vector Simulation

Chandak and Browne [ChBr83] recently proved an item of computing folklore — discrete event simulation models cannot always be vectorized. Specifically, they showed that any network of queues model containing feedback is not vectorizable. This result is quite negative: most interesting simulation models contain some type of feedback.

Given this result, we recently investigated the level of vectorization practically achievable [Reed85] by instrumenting a discrete event simulation of queueing network models on a Cray X-MP. Although we simulated a variety of workloads and queueing network models, the observed vectorization level never exceeded 5 percent. Even this fraction was primarily attributable to initialization code. Thus, the efficacy of vector simulation is in doubt.

Distributed Simulation

The inherently sequential nature of event list manipulation limits the potential parallelism of standard simulation models. The head of the event list must be removed, the simulation clock advanced, and the event performed (possibly causing new events to be added to the event list). Although techniques for performing event list manipulation and event simulation in parallel have been suggested [Comf82, Comf83], large scale performance increases seem unlikely. Only by eliminating the event list, in its traditional form, can additional parallelism be obtained; this is the goal of distributed simulation.

If one views a simulation model as a network of interacting servers and queues, distributed simulation maps each server/queue pair onto a processor of a multicomputer network. Each processor 'operates with its own simulation clock, and there is no global event list. Event occurrence times are transmitted across communication links to appropriate recipients (e.g., a message departing one server for another would carry with it its time of departure). Several distributed simulation techniques have been proposed, notably the Chandy-Misra algorithm [ChHM79, ChMi79, ChMi81] and the Time Warp algorithm [JeSo85]. The Chandy-Misra algorithm and Time Warp differ in their approach to time management. The former is pessimistic, advancing the processor simulation clocks only when conditions permit. In contrast, Time Warp assumes the simulation clocks can be advanced until conflicting information appears; the clocks are then *rolled back* to a consistent state, a so-called "time warp."

Both the Chandy-Misra algorithm and Time Warp have been simulated [Seet78, JeSo85], but, to our knowledge, no experimental results have yet been reported. In the remainder of this paper, we present the Chandy-Misra algorithm [ChMi81] and the results of an extensive study of its performance on a shared memory parallel processor when simulating queueing network models. Parameters of the study include queueing network topology and routing probabilities, number of processors, and assignment of queueing network servers to processors. We conclude with a summary of lessons learned and directions for future research.

The Chandy-Misra Distributed Simulation Algorithm

Consider some *physical* system composed of independent, interacting entities. A natural, distributed simulation of the physical system creates a topologically equivalent system of *logical* nodes. Interactions between two physical nodes are modeled by exchange of timestamped messages. The timestamp is the simulated message arrival time at the receiving node.

Each logical node is subject to some constraints. First, node interaction is only via message exchange; there are no shared variables. Second, each node must maintain a clock, representing the local simulated time. Finally, the timestamps of the messages generated by each node must be non-decreasing.

Intuitively, the distributed simulation has no single "correct" simulation time; each node operates independently subject only to those restrictions necessary to insure that events happen in the correct simulated order (i.e., *causality* is maintained). Independent events can be simulated in parallel even if they occur at different simulated times.

Message timestamps and node clocks are a manifestation of the need for causality: the behavior of a node P at its simulated time T cannot be influenced by any information transmitted to it after time T. This constraint has rather dramatic ramifications. Consider a node P that receives messages from two other nodes A and B. When a message arrives from node A, one would expect node P to interpret the message, perhaps producing a message as a consequence. However, if the arrival time of the message from A is greater than the arrival time of the *last* message from B, the message from A cannot be processed. Why? A message might later arrive from B with a smaller timestamp. Thus, a node with multiple inputs must wait until it receives messages from all inputs before selecting a message to interpret.

Although appealing, distributed simulation poses several pragmatic problems:

- Optimal assignment of nodes to processors is expensive.
- Only a subset of all discrete event simulation models are amenable to distributed simulation. As noted above, shared variables are not permitted. Hence, no events depending on the global system state are possible.
- Deadlocks can occur in most simulation models. Recall that a node must insure that no information received later can affect its output; this may require waiting for additional inputs. A cycle of waiting nodes results in deadlock.

The assignment problem deserves additional comment. The natural hardware realization of the network of nodes is a multicomputer network. Pragmatics dictate, however, that the multicomputer network have a fixed interconnection topology. Thus, a node network must be

mapped onto the multicomputer network. Unfortunately, finding an optimal mapping is known to be NP-complete [Chu80]. In practice, scheduling heuristics must be used; sub-optimal mappings are produced at considerably less computational expense. Even if an optimal mapping were found, the respective topologies of the multicomputer and node networks may be ill suited, resulting in either large communication delays or processor load imbalances.²

Like node scheduling, deadlock resolution, although difficult, is solvable. Chandy and Misra have described two distributed deadlock resolution techniques, avoidance and recovery [ChMi81]. For specificity's sake, we describe these techniques in the context of our RESQ implementation [SaMS80] for simulating queueing networks.

In the RESQ scheme, there are five node types: service, fork, merge, source, and -sink. Service nodes correspond to the interacting entities of a physical system (e.g., servers in a queueing network). In contrast, fork and merge nodes exist only to provide routing. Finally, source and sink nodes respectively create and destroy network messages. Thus, the central server model [Buze73] of Figure 1a would be represented, using the RESQ scheme, as shown in Figure 1b.

The RESQ notation for describing queueing network models has been widely used as an input language for sequential simulations. Using RESQ for parallel simulation entails modifying the semantics of some node node types. Specifically, distributed simulation with deadlock avoidance [ChMi79] requires fork, merge, and server nodes³ to send *null* messages under certain conditions. These null messages are time stamped and tell the receiving node that no *real* message will be forthcoming before the specified time. This enables the receiver to process outstanding messages with the assurance that its actions will not be revoked at a later time.

²Scheduling difficulties can be ameliorated by a shared memory implementation of message passing. This approach is discussed later.

³By definition, source and sink nodes can never be members of a deadlock set.

A fork node accepts a single stream of message inputs and distributes this stream across N outputs. Upon receiving a real or null input message, a fork node routes the message to the selected output and creates N-1 null messages, each with the same timestamp as the message received. One null message is routed to each destination not selected.

A merge node accepts N streams of message inputs and routes them in timestamp order to a single output. As noted earlier, the timestamp ordering forces the node to wait for messages, perhaps null, on all inputs before producing an output.

Finally, a *server* node accepts a single input stream and produces a single output stream. When the time of last message arrival is greater than the time of last message departure, and the server has no real messages to process, it produces a null message with a timestamp equal to the minimum time of next real message departure.

Although the null message technique provably avoids deadlocks, it does so at the price of potentially high overhead. In networks containing many fork/merge cycles, simulations have shown that the ratio of null to real messages can be very high [Seet78, Reed85]. The alternative to deadlock avoidance is deadlock detection and recovery. In this approach, the distributed simulation alternates between computation and recovery phases. As proposed by Chandy and Misra [ChMi81], the simulation runs until a distributed deadlock detection algorithm verifies deadlock. The simulation then enters a deadlock recovery phase and finally returns to active computation.

Although deadlock detection and recovery avoids null messages, it does so by *diverting* computation resources to detection and recovery. The performance advantage of this approach versus deadlock avoidance depends on the relative costs of synchronization and message passing.

In light of the many potentially performance-limiting problems with distributed simulation, it seems important to analyze the performance of distributed simulation in a realistic environment. Many such performance studies of traditional simulation algorithms have been conducted, and, based on these studies, new event list algorithms have been proposed [FrMa77, Wyma75]. Only limited *simulation* studies of distributed simulation have been reported [Seet78, JeSo85, Reed85]; little or no empirical data are available. In the remaining sections we discuss our experimental environment, implementation, and experimental results.

Experimental Environment

All simulation experiments were conducted on a Sequent Balance 21000 containing 20 processors and 16 Mbytes of memory. As shown in Figure 2, each Balance 21000 processor is a 10 MHz National Semiconductor NS32032 microprocessor, and all processors are connected to shared memory by a shared bus with a 80 Mbyte/s (maximum) transfer rate. Each processor has a 8K byte, write-through cache and an 8K byte local memory; the latter contains a copy of selected read-only operating system data structures and code.

The Dynix $\[Box]$ operating system for the Balance 21000 is a variant of UC-Berkeley's 4.2BSD Unix $\[Box]$ with extensions for processor scheduling. Because Dynix schedules all processes from a common pool, a process may execute on different processors during successive time slices. However, as long as the number of active processes is less than the number of processors, each process will execute on a separate processor. In this case, process and processor are equivalent notions. To the time-sharing user, the Balance 21000 appears as a standard Unix system, albeit with better interactive response time.

Parallel programs consist of a group of Unix processes that interact using a library of primitives for shared memory allocation and process synchronization. Shared memory is implemented by mapping a region of physical memory into the virtual address space of each process. This mapping can be done only once during program execution, typically at the

beginning. Once mapped, shared memory can be allocated to specific variables as desired.

Access to the shared memory region is controlled by software spin *locks* and *barriers*. These locks, semantically equivalent to binary semaphores, provide mutual exclusion. Barriers are used to synchronize a group of processes; a process reaching a barrier is forced to wait until all processes in the specified group have reached the barrier.

In summary, the Balance 21000 is a "standard" Unix system with minimal extensions for parallel programming. Consequently, many parallel operations are dominated by operating system overhead. For comparison with later discussion, Table 1 shows the elapsed times for typical operations.

Shared Memory Implementation of Distributed Simulation

A shared memory multiprocessor, such as the Balance 21000, provides a flexible testbed for studying the performance of distributed simulation. The problems associated with mapping a node network onto a multicomputer network are removed; the shared memory processors are, effectively, completely connected. By implementing message passing using shared memory, communications costs are the same for all processors. However, a shared memory implementation of distributed simulation requires special consideration for synchronization of shared message queues, processor allocation, and deadlock management.

In a shared memory implementation of distributed simulation, all node state information, including input message queues, resides in shared memory. Message-based communication between nodes is implemented via shared access to the message queues of each node. Each message queue is protected by a synchronization lock to guarantee mutual exclusion. Synchronization is only necessary, however, if the communicating nodes execute on separate processors. Before a node can send a timestamped message to another node, it must first acquire a free message from a shared free message list. A lock is necessary to prevent simultaneous access to the free message list. After retrieving a free message, the node timestamps it and writes it to the destination node's message queue, using synchronization primitives to lock and unlock the queue, if necessary. A message is returned to the free message list once it has been processed by the destination node. Because only messages are used for internode communication, the requirement that no simulated events depend on the global system state is still satisfied.

Processor Allocation

There are two basic approaches to processor allocation in a shared memory implementation of distributed simulation. The first approach, *static node assignment*, fixes the assignment of nodes to processors for the duration of the simulation. When the number of network nodes equals the number of allocated processors, each node is assigned to a separate processor. Otherwise, nodes must be clustered, and these clusters are assigned to individual processors. Several clusterings are possible when the number of nodes exceeds the number of processors; each such clustering exhibits different performance. One advantage of static node assignment is that communication between nodes in a cluster can be done "locally" without the overhead for locking message queues. However, intercluster message transmissions require queue locking.

The second approach, dynamic node assignment, assigns nodes to processors during the simulation. Idle processors obtain work from a shared queue of unassigned network nodes. This shared node work queue must be locked before a processor can be allocated an unassigned network node. When a processor obtains a node, it satisfies any outstanding work for the node before returning the node to the tail of the node work queue. Because processors are assigned only one node at any time, all communication between nodes must be synchronized to guarantee exclusive access to shared message queues.

With dynamic node assignment, nodes must wait on the work queue until assigned to a processor. The length of this delay depends on size of the node work queue and can be quite large for large networks. However, not all nodes on the work queue have outstanding work (i.e., there are input messages that will generate output messages when processed). In deadlock avoidance mode, for example, those nodes awaiting input can only generate null messages if processed. A natural strategy for improving performance places only those nodes with outstanding work on the work queue. This reduces the size of the node work queue and the waiting delay.

Our implementation of the above node waiting strategy is conservative. When a processor identifies a node with no outstanding work, it sets a "waiting" flag in the node's state and does not place the node on the work queue. When a message is sent to a waiting node, the processor sending the message will reset the waiting flag of the waiting node and place it on the work queue. The implementation is conservative because the new message may not actually instigate any new work for the node.

To investigate the effects of this node waiting strategy, we also implemented a *no node waiting* scheme. In this approach, a node is immediately placed at the tail of the work queue after it has been processed, even if it has no outstanding work.

Although static node assignment is efficient for nodes within a cluster, the node assignment cannot be changed to balance network load. Conversely, dynamic node assignment naturally adjusts to network load but incurs synchronization overhead not only for all messages but also for access to the node work queue. Which implementation is best for a particular simulation model depends on the relative costs of synchronization and the beneficial effects of load balancing.

Deadlock Avoidance and Recovery

Our deadlock avoidance approach is a straightforward implementation of the algorithm described earlier [Seet78, ChMi81]. In contrast, deadlock recovery merits further discussion.

As described by Chandy and Misra, distributed simulation with deadlock detection and recovery alternates between simulation and distributed deadlock detection and recovery phases. The presence of shared memory obviates the need for most of the protocol for distributed deadlock detection [ChHM83]. Instead, each processor sets a flag in global memory when it believes it is deadlocked. A guardian processor monitors the global system state and forces the processors to rendezvous at a synchronization barrier when they all report potential deadlock. The deadlock recovery algorithm is then invoked.

Notice, however, that all processors reporting an inability to progress is a necessary but not sufficient condition for deadlock. Between the time a processor P reports potential deadlock and the time the guardian processor sees this report, processor P may have received messages enabling it to progress. Consequently, the processors may appear deadlocked when they are not. To reduce the probability of detecting such false deadlocks, the guardian uses a *backoff* algorithm that must re-verify a potential deadlock before invoking deadlock recovery. This algorithm, controlled by an input parameter, weighs the relative cost of forcing synchronization and deadlock recovery for a false deadlock against the lost time when detection of a real deadlock is delayed.

One may well ask why this deadlock detection technique was used, rather than a variation of graph reduction [Fink86] or the distributed deadlock detection proposed by Chandy *et al* [ChHM83]. Simply put, the number and frequency of deadlocks in a distributed simulation is potentially enormous. Hence, deadlock detection and recovery must be fast. To obtain a consistent state for graph reduction, the processors must either exchange messages or synchronize. The overhead of the first is near that for deadlock avoidance. The latter is as expensive as detecting false deadlocks. Thus, detecting some false deadlocks using a backoff mechanism seems a reasonable compromise.

Simulation Experiments

Experimental evaluation of distributed simulation requires not only an implementation but also a set of test cases. This is particularly important in light of earlier simulation studies [Seet78, Reed85], which showed that the performance of distributed simulation is extremely sensitive to the topology of the simulated network. Simple tests (e.g., tandem queues) have easily interpretable results, but do not reflect typical simulations. Conversely, simulation of complex queueing networks, although realistic, make it difficult to interpret the sources of performance degradation in distributed simulation.⁴

As a compromise, we selected several simple queueing networks and a few complex ones.

- tandem networks (1, 2, 4, 8, and 16 server nodes)
- general, feed-forward networks (6, 10, and 14 nodes),
- cyclic networks (2, 4, and 8 nodes)
- central server networks (5 nodes), and
- cluster networks (10 and 18 nodes).

The tandem and feed-forward networks are open networks and contain no cycles. With potentially linear speedup, they represent the best-case performance of distributed simulation. The cyclic networks show the performance degradation of tandem networks when they are closed. As an often used model of computer systems, central server networks have pragmatic importance [Buze73]. In addition, they have nested cycles, a more restrictive constraint than the simple cyclic networks. Finally, the cluster networks illustrate the effects of decomposability on simulation performance.

⁴We distinguish between the performance of the Chandy-Misra simulation and the performance measures for the simulated network. The former are the subject of this study.

Each of these networks was simulated for a variety of workloads, (e.g., routing probabilities, arrival rates, and service times) using six variations of a Chandy-Misra implementation: static node assignment with deadlock avoidance, static node assignment with deadlock recovery, dynamic node assignment with deadlock avoidance, dynamic node assignment with deadlock recovery, dynamic node assignment with waiting and deadlock avoidance, and dynamic node assignment with waiting and deadlock recovery. In all cases, we varied the number of processors from one to the number of nodes in the simulated network. Together, these simulations represent approximately two weeks of computation time on the Sequent Balance 21000. Figures 8-25 and Tables 3-5, discussed below, show the results of a portion of these experiments. All such figures and tables show 95 percent confidence intervals about mean values.

Speedup, defined as

$$S_p = \frac{T_1}{T_p},$$

where T_1 and T_p are the respective execution times using one and p processors, is the performance metric used to compare all experimental results. All speedups are shown relative to a one processor distributed simulation using static assignment with deadlock recovery. In this case, all simulated nodes execute on one processor. Consequently, no synchronization is needed during queue insertion and deletion. Deadlocks *can* occur with one processor. This increases the value of T_1 and, consequently, increases the apparent speedup. Although it might seem preferable to use an event list oriented simulation as the point of reference, this would color the results with the idiosyncrasies of *two* implementations. For comparison, we conducted equivalent event list simulations on the Balance 21000 using SMPL, a portable simulation package. These results show that a single processor distributed simulation always executes more slowly than the equivalent sequential simulations. Thus, the speedups presented can be viewed as upper bounds on the speedup achievable with distributed simulation.

In addition to speedup, we also use deadlock recovery and null message fractions as performance measures. These are defined as

$$F_D = \frac{number of deadlock recoveries}{number of message transmissions}$$

and

$$F_N = \frac{number \ of \ null \ message \ transmissions}{number \ of \ message \ transmissions}$$

respectively. The deadlock recovery and null message fractions measure the amount of *useful* computation performed by each simulation.

Tandem Networks

Tandem networks are a feasibility test of distributed simulation. If distributed simulation cannot achieve good pipelined speedups for tandem networks, there is little prospect for success for networks containing cycles.

Figure 8 shows the speedups for both deadlock avoidance and recovery when nodes are statically assigned to processors. Because there are no cycles, no deadlocks occur, and there is little distinction between deadlock recovery and avoidance. Recall that deadlock avoidance must continually verify that no null messages need be sent. Conversely, deadlock recovery does nothing until deadlock is detected. Thus deadlock avoidance incurs a small overhead even if no null messages need be sent. This difference is magnified as the number of nodes increases, leading to a small, but perceptible difference at 16 server nodes.

Figure 8 shows a linear speedup for a small number of nodes, and a decrease in the slope of the speedup curve for additional nodes. This sublinear speedup for a larger number of nodes . arises from memory and bus contention, as well as synchronization overhead. By comparison, Figure 9 shows speedups when deadlock recovery is used, and nodes are retrieved from a work queue.⁵ Dynamic node assignment yields greater speedup than static assignment, but it too suffers from memory contention. Using half as many processors as nodes results in near linear speedup, albeit a smaller speedup than that obtained with maximal parallelism. This is the final confirmation of the effects of memory contention.

Waiting (i.e., placing nodes on the work queue only when they can profitably be evaluated), is ineffective because, in a tandem network, all nodes are always active. The additional overhead simply reduces the speedup, as shown in when maximal parallelism is used.

The previous discussion, with one exception, assumed the number of processors equaled the number of nodes. When the number of nodes exceeds the number of processors, nodes must, with static assignment, be clustered onto processors. Table 2 shows the effects of this clustering for a tandem network containing 16 server nodes. For static assignment, speedup declines precipitously as the number of processors is reduced (e.g., reducing the number of processors from 18 to 12 reduces the speedup from approximately 9 to 4). In contrast, dynamic node assignment allocates processors to nodes based on their need for evaluation. When the number of nodes exceeds the number of processors, dynamic assignment is the method of choice.

Finally, we note that the sequential execution time, 113 seconds, compares favorably to the single processor distributed simulation. This suggests that the overhead for distributed simulation, other than that for deadlock avoidance or recovery, compares favorably to that for event-driven simulation.

General Feed-forward Networks

⁵In Figure 9, "half parallelism" means that the number of processors used is equal to one half the number of network nodes.

Among the simplest generalizations of a tandem network are those containing forks and joins. Figures 4 and 5 show two such networks of differing complexity. Tables 3 and 4 show the corresponding speedups as a function of node clustering and deadlock technique.

Unlike the tandem networks, where deadlock avoidance and recovery are indistinguishable, feed-forward networks with forks necessarily distinguish between the two deadlock techniques. Because this network is open and contains no cycles, no deadlocks can occur, and deadlock detection detects none. In contrast, deadlock avoidance requires that null messages be sent at each fork node. This overhead is the reason for the difference in the performance of the two deadlock techniques.

Although speedups are *not* linear in the number processors, the fork and join nodes do not require as much processing time as the server nodes. Because of this, a linear speedup from a sequential simulation cannot be expected.

Cyclic Networks

The closed equivalent of a tandem network is the cyclic queue of Figure 6. Unlike the tandem network, where the interarrival time at the source node does not affect the execution time of the simulation, the cyclic network depends on the simulated population. Figure 10 shows the speedup obtained for a four node cyclic network; similar results also hold for larger cyclic networks. In contrast to the tandem networks, the cyclic network does not show linear speedup as a function of the number of processors.

Initially, one might suspect deadlock avoidance or recovery caused this decrease in performance. However, examination of the simulations shows that deadlock avoidance sent *no* null messages. Instead, messages circulate in large groups or *trains*; a node processes a train of messages and waits until they return on their next cycle. This suggests that fewer processors, dynamically assigned to the nodes, would achieve most of the potential speedup. Figure 11 confirms this supposition: two, dynamically assigned processors, achieve nearly 80 percent of the speedup obtained with four processors. When dynamic assignment is augmented with waiting, as in Figure 12, the difference grows even smaller.

The second important conclusion drawn from simulating cyclic queues is the extremely high cost for deadlock recovery. As noted above, deadlock avoidance sent no null messages. In contrast, deadlock recovery detected a small number of potential deadlocks. At population 40, the deadlock recovery fraction F_D was 0.0015. This corresponds to approximately 250 deadlock recoveries in 160,000 message transmissions. As discussed earlier, deadlock detection recovery forces all processors to synchronize at a barrier before invoking the deadlock recovery algorithm. Execution profiling showed that the deadlock recovery routine and the barrier primitive comprised a negligible fraction of the simulation time. Synchronizing is expensive, but only because there is a significant interval between the arrival of the first processor at the barrier and the last. The parallelism declines as each processor reaches the barrier. Were the transition from computation to deadlock recovery abrupt, deadlock recovery would be inexpensive. As Figure 11 shows, dynamic node assignment improves the performance of deadlock recovery, primarily because fewer processors are actively evaluating nodes, dynamic node assignment with waiting further reduces the rendezvous delay; see Figure 12.

Central Server Networks

Central server networks have long been used as models of computer systems [Buze73], and consequently have pragmatic importance. Because they contain nested cycles, central server networks are susceptible to deadlock in a distributed simulation. Hence, they are a more realistic test of distributed simulation. Figures 13 and 14 show the speedup obtained for a central server network containing three servers; see Figure 1 for the network topology. Even with five processors, the speedup barely exceeds unity. Moreover, this is using the single processor, static node assignment case as the basis for calculating speedup. As Table 5 shows, the parallel implementation rarely completes more quickly than the sequential implementation. Indeed, static node assignment with deadlock avoidance runs 16 times more slowly than the sequential implementation. Consequently, the speedups over an event-driven simulation are much lower than Figures 13 and 14 suggest.

Unlike the simple cyclic network, where both deadlock avoidance and recovery were rare, the central server network frequently forces the simulation to either send null messages or attempt deadlock recovery. Figures 15 and 16 respectively show the deadlock recovery and null message fractions for static node assignment. Although not shown, the fractions for dynamic node assignment are similar.

With only one circulating message, nearly fifty null messages are transmitted for each movement of the real message. Although the null fraction decreases as the number of circulating messages increases, it converges to approximately twenty null messages per real message transmission.⁶ In contrast, the deadlock recovery fraction converges to 0.35. Although these deadlock recoveries are expensive, as the analysis of cyclic networks showed, their number is so small compared to the number of null messages sent during deadlock avoidance that deadlock recovery is significantly faster.

Finally, we must emphasize that these results are *significantly* more negative than earlier simulated results [Reed85]. A sequential simulation of a network, by its nature, imposes some sequential ordering on the evaluation of network nodes. When those nodes are not being evaluated, they do not generate null messages, nor can they deadlock. In contrast, in a fully

⁶This value, twenty, seems independent of the network routing probabilities. Removing the nested cycle from Figure 1 neither increases the observed speedup nor decreases the null fraction. We hypothesize that the value is a function of the relative speeds of the processors and memory.

parallel implementation, all nodes are always active. Thus, they continue to receive and generate null messages while awaiting receipt of real messages. Thus, the overhead is *higher* than suggested by a sequential simulation of distributed simulation.

Cluster Networks

Cluster networks were the most complex simulated during our experimental study. As Figure 7 shows, a cluster network is composed of several tightly clustered subnetworks. This has two important ramifications. First, the network is nearly decomposed and should yield significant speedups with parallel simulation. Secondly, the clustering increases the expected execution time of the simulation. Why? The clocks of all nodes must reach the terminating value before the simulation completes. With only a few circulating messages, some nodes may be idle for long periods of time. Only when a message "escapes" from a subcluster will the clocks of other nodes advance. This asynchrony means that the clocks of some nodes may run far past the terminating simulation value.

Figure 17 shows the speedup of the cluster network with static node assignment. For small populations, deadlock recovery is significantly faster than deadlock avoidance. As with the central server network, this is attributable to the large number of null messages sent. As the population increases, the null fraction decreases precipitously, and deadlock avoidance becomes the method of choice. (Figure 18 shows the deadlock and null fractions.) Interestingly, the speedup obtained with deadlock recovery is relatively insensitive to the simulated population. As Figure 18 shows, the deadlock fraction is negligible but constant. With a large number of processors, the delay to synchronize at a barrier is prohibitive; this overhead is the reason for the poor performance of deadlock recovery.

When nodes are dynamically assigned to processors, Figure 19, the performance of both deadlock avoidance and deadlock recovery increase significantly. As noted earlier, messages are

often "trapped" in network subclusters and many nodes are often idle. With deadlock avoidance and static node assignment, many nodes continually generate null messages. These messages simply cause additional overhead and memory contention. With dynamic assignment, a node must migrate from the tail to the head of the node work queue before being evaluated. This additional delay between evaluations reduces the number of null messages and is the source of the additional speedup with dynamic node assignment.

Figure 19 also shows that a smaller number of processors yields a marginally larger speedup than that obtained with maximal parallelism. Although the difference is not statistically significant in this case, Figure 20 shows that the difference is large when node waiting is introduced. The reason is, as before, the presence of many idle nodes. By suspending nodes that cannot productively contribute to the simulation, contention for the node work queue is reduced. and those nodes with work can proceed without interference.

In summary, the cluster network shows that distributed simulation can produce significant speedups *if* the network is decoupled, and subclusters interact infrequently.

Summary

Distributed simulation has been the subject of several *simulated* performance studies; little or no experimental data have heretofore been available. Obtaining such data was the primary goal of this work. Using queueing networks as the simulation application, we simulated a variety of such networks with varying workloads using several variations of the Chandy-Misra algorithm on a shared memory machine.

These experiments show that, with rare exception, the Chandy-Misra approach to distributed simulation is not a viable approach to parallel simulation of queueing network models. There are two primary reasons for this. First, a single processor implementation of the Chandy-Misra algorithm is sometimes slower than the equivalent sequential, event-driven simulation. Thus, multiple processors are needed just to recoup the loss due to the inefficiency. Second, networks with cycles require deadlock avoidance or recovery techniques. These techniques are extremely costly, and there is little prospect that they can be reduced to acceptable levels.

Because queueing network simulation requires little processing by server nodes, nodes interact frequently in real time. Because of this, queueing networks are a stress test for distributed simulation. In simulations that require extensive computation between node interactions, distributed simulation is analogous to a group of decoupled processes. In such cases, distributed simulation should prove more attractive.

Acknowledgments

Jack Dongarra and the Advanced Computing Research Facility of Argonne National Laboratory graciously provided both advice and access to the Sequent Balance 21000.

References

- [Buze73] J. P. Buzen, "Computational Algorithms for Closed Queueing Networks with Exponential Servers," Communications of the ACM, Vol. 16, No. 9, September 1973, pp. 527-531.
- [ChHM79] K. M. Chandy, V. Holmes, and J. Misra, "Distributed Simulation of Networks," Computer Networks, Vol. 3, No. 1, February 1979, pp. 105-113.
- [ChMi79] K. M. Chandy and J. Misra, "Distributed Simulation: A Case Study in Design and Verification of Distributed Programs," *IEEE Transactions on Software Engineering*, Vol. SE-5, No. 5, September 1979, pp. 440-452.
- [ChMi81] K. M. Chandy and J. Misra, "Asynchronous Distributed Simulation via a Sequence of Parallel Computations," Communications of the ACM, Vol. 24, No. 4, April 1981, pp. 198-206.
- [ChHM83] K. M. Chandy, L. M. Haas, and J. Misra, "Distributed Deadlock Detection," ACM Transactions on Computer Systems, Vol. 1, No. 2, May 1983, pp. 144-156.
- [ChBr83] A. Chandak and J. C. Browne, "Vectorization of Discrete Event Simulation," Proceedings of the 1983 International Conference on Parallel Processing, August 1983, pp. 359-361.
- [Chu80] W. W. Chu, et al, "Task Allocation in Distributed Data Processing," IEEE Computer, Vol. 13, No. 11, November 1980, pp. 57-69.
- [Comf82] J. C. Comfort, "The Design of a Multi-microprocessor Based Simulation Computer -I," Proceedings of the Fifteenth Annual Simulation Symposium, March 1982, pp. 45-53.
- [Comf83] J. C. Comfort, "The Simulation of a Master-Slave Event Set Processor," Simulation, Vol. 42, pp. 117-124.
- [Fink86] R. A. Finkel, An Operating Systems Vade Mecum, Prentice-Hall, 1986.
- [FrWW84] M. A. Franklin, D. F. Wann, and K. F. Wong, "Parallel Machines and Algorithms for Discrete-Event Simulation," 1984 International Conference on Parallel Processing, August 1984, pp. 449-458.
- [FrMa77] D. Franta and W. Maly, "An Efficient Data Structure for the Simulation Event Set," Communications of the ACM, Vol. 20, No. 8, August 1977, pp. 596-602.

- [Heid86] P. Heidelberger, "Statistical Analysis of Parallel Simulations," Proceedings of the 1986 Winter Simulation Conference, to appear.
- [JeSo85] D. Jefferson and H. Sowizral, "Fast Concurrent Simulation Using the Time Warp Mechanism, Distributed Simulation 1985, The 1985 Society for Computer Simulation Multiconference, San Diego, California.
- [KoRW86] M. B. Konsek, D. A. Reed, and W. Watcharawittayakul, "Context Switching with Multiple Register Windows: A RISC Performance Study," in preparation.
- [PaSe82] D. A. Patterson and C. H. Sequin, "A VLSI RISC," *IEEE Computer*, Vol. 15, No. 9, September 1982, pp. 8-21.
- [Patt85] D. A. Patterson, "Reduced Instruction Set Computers," Communications of the ACM, Vol. 28, No. 1, January 1985, pp. 8-21.
- [PeWM79] J. K. Peacock, J. W. Wong, and E. G. Manning, "Distributed Simulation Using a Network of Processors," Computer Networks, Vol. 3, No. 1, February 1979, pp. 44-56.
- [Pfis82] G. F. Pfister, "The Yorktown Simulation Engine: Introduction," ACM IEEE 19th Design Automation Conference, June 1982, pp. 51-54.
- [Reed83] D. A. Reed, "A Simulation Study of Multimicrocomputer Networks," Proceedings of the 1983 International Conference on Parallel Processing, August 1983, pp. 161–163.
- Reed85] D. A. Reed, "Parallel Discrete Event Simulation: A Case Study," Record of Proceedings: 18th Annual Simulation Symposium, March 1985, pp. 95-107, invited paper.
- [ReFu86] D. A. Reed and R. M. Fujimoto, Multicomputer Networks: Message Based Parallel Processing, submitted to MIT Press.
- [SaMS80] C. H. Sauer, E. A. MacNair, and S. Salza, "A Language for Extended Queueing Networks," *IBM Journal of Research and Development*, Vol. 24, No. 6, November 1980, pp. 747-755.
- [Seet78] M. Seethalakshmi, "Performance Analysis of Distributed Simulation," M.S. Report, Computer Science Department, University of Texas, Austin, Texas, 1978.
- [Seit85] C. L. Seitz, "The Cosmic Cube," Communications of the ACM, Vol. 28, No. 1, January 1985, pp. 22-33.
- [Wyma75] P. F. Wyman, "Improved Event Scanning Mechanisms for Discrete Event Simulation," Communications of the ACM, Vol. 18, No. 4, April 1975, pp. 221-230.

Table 1Typical Operation Times for the Sequent Balance 21000

Operation	Time (µsec)
Lock/unlock	60
Subroutine call/return	60
System call	400
Context switch	1000
Process creation	60000

STATIC		DYNAMIC				
Clustering Case	Recovery -	Avoidance	Recovery	Recovery w/ Waiting	Avoidance	Avoidance w/ Waiting
Α	$9.24 \pm 0.41\%$	$8.73 \pm 0.50\%$	$10.22 \pm 1.92\%$	9.97 ± 0.94%	$10.19 \pm 2.49\%$	$9.74 \pm 1.25\%$
В	$8.53 \pm 0.31\%$	$8.23 \pm 0.56\%$				
С	$4.58\pm8.03\%$	$\boldsymbol{4.27 \pm 11.29\%}$				
D	$\boldsymbol{3.87 \pm 7.61\%}$	$3.28 \pm 23.16\%$				
Е			$7.50\pm0.99\%$	$7.52 \pm 1.32\%$	$7.46 \pm 1.48\%$	$7.41 \pm 1.71\%$
F	$3.22\pm3.19\%$	$2.83 \pm 8.11\%$				
G	$3.60\pm0.23\%$	$3.51\pm0.86\%$				
Н	$2.41 \pm 2.24\%$	$2.38 \pm 1.69\%$				1
I	$1.87 \pm 0.21\%$	$1.83 \pm 0.62\%$				
J	$1.00\pm0.33\%$	$0.98 \pm 0.37\%$	$0.91 \pm 0.17\%$	$0.98 \pm 0.37\%$	$0.91 \pm 0.39\%$	$0.90\pm0.77\%$

Table 2Speedups for tandem network with 16 server nodes

Parameter	Value
Node Service Time	0.0625
Confidence Level	95%
Speedup Base	One processor static deadlock recovery
Mean Base Execution Time	117.88 seconds
Mean Sequential Execution Time	113.45 seconds
Cluster case A	(1)(18) ⁻
Cluster case B	(1 17) (2)(15) (16 18)
Cluster case C	(1 17) (2 3) (4)(13) (14 15) (16 18)
Cluster case D	$(1 \ 17) (2 \ 3) (4 \ 5) (6)(11) (12 \ 13) (14 \ 15) (16 \ 18)$
Cluster case E	(1 17 2) (3 4 5) (6 7 8) (9 10 11) (12 13) (14)(16) (18)
Cluster case F	(1 17 2) (3 4 5) (6 7 8) (9 10 11) (12 13 14) (15 16 18)
Cluster case G	(1 17 2 3 4) (5 6 7 8) (9 10 11 12) (13 14 15 16 18)
Cluster case H	(1 17 2 3 4 5) (6 7 8 9 10 11) (12 13 14 15 16 18)
Cluster case I	(1 17 2 3 4 5 6 7 8) (9 10 11 12 13 14 15 16 18)
Cluster case J	(1 17 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 18)

Node numbers refer to Figure 3. Parenthesized node groups execute on one processor.

~

	STATIC		DYNAMIC				
Clustering Case	Recovery -	Avoidance	Recovery	Recovery w/ Waiting	Avoidance	Avoidance w/Waiting	
A	$3.34 \pm 1.00\%$	$2.57\pm0.62\%$	$3.74 \pm 0.00\%$	$3.66 \pm 1.40\%$	$2.90 \pm 1.80\%$	$2.81 \pm 1.35\%$	
В	$2.16\pm0.64\%$	$1.81 \pm 1.67\%$					
с	$2.17\pm0.38\%$	1.36 ± 4.95					
D	$1.57 \pm 1.40\%$	$1.21\pm0.36\%$	$2.47 \pm 1.13\%$	$2.47 \pm 1.12\%$	$1.86 \pm 1.31\%$	$1.83 \pm 2.09\%$	
E	$1.55\pm0.46\%$	$1.25 \pm 1.46\%$					
F	$1.00\pm0.82\%$	$0.76 \pm 0.51\%$	$0.88 \pm 0.48\%$	$0.87\pm0.47\%$	$0.65 \pm 0.43\%$	$0.65 \pm 0.40\%$	

 Table 3

 Speedups for generalized feed forward network with 6 nodes

Parameter	Value
Node Service Time	1.0
Confidence Level	95%
Speedup Base	One processor static deadlock recovery
Cluster case A	(1)(6) ⁻
Cluster case B	(1) (2) (3 4) (5 6)
Cluster case C	(1 2) (3 4) (5 6)
Cluster case D	(1) (2) (3 4 5 6)
Cluster case E	(1 5 6) (2 3 4)
Cluster case F	(1 2 3 4 5 6)

Node numbers refer to Figure 4. Parenthesized node groups execute on one processor.

STATIC		DYNAMIC				
Clustering Case	Recovery -	Avoidance	Recovery	Recovery w/ Waiting	Avoidance	Avoidance w/ Waiting
Α	$6.05 \pm 1.46\%$	$2.36 \pm 0.45\%$	$7.19 \pm 2.27\%$	$6.93 \pm 2.69\%$	$2.53 \pm 2.56\%$	$2.25 \pm 2.54\%$
В	$4.86\pm0.91\%$	$2.32\pm0.93\%$				
С	$1.82\pm0.42\%$	$\boldsymbol{0.97\pm0.19\%}$	$5.44 \pm 0.96\%$	$5.42\pm0.81\%$	$\boldsymbol{3.04\pm0.64\%}$	$2.98 \pm 0.45\%$
D	$1.83 \pm 0.65\%$	$0.96 \pm 0.31\%$				
Е	$1.79 \pm 4.69\%$	$0.99 \pm 0.52\%$		· · ·		
F	$1.00\pm0.91\%$	$0.57 \pm 0.76\%$	$0.90\pm0.56\%$	$0.89 \pm 0.13\%$	$0.48\pm0.46\%$	$0.48\pm0.47\%$

Table 4Speedups for generalized feed forward network with 14 nodes

Parameter	Value
Node Service Time	1.0
Confidence Level	95%
Speedup Base	One processor static deadlock recovery
Cluster case A	(1)(14) -
Cluster case B	(1) (2) (3 9) (4 14) (5 6) (7) (8) (10 13) (11) (12)
Cluster case C	(1) (2) (7) (8) (11) (12) (3 4 5 6 9 10 13 14)
Cluster case D	(1 2) (7 8) (11 12) (3 4 5 6 9 10 13 14)
Cluster case E	(1 2 7 8 11 12) (3 4 5 6 9 10 13 14)
Cluster case F	(1 2 7 8 11 12 3 4 5 6 9 10 13 14)

 $^{\rm L}$ Node numbers refer to Figure 5. Parenthesized node groups execute on one processor.

Table 5	
Sequential and parallel mean execution	time
for five node central server	
(time given in seconds)	

		STATIC PARALLEL		DYNAMIC PARALLEL			
Popu- lation	SEQUENTIAL	Recovery	Avoidance	Recovery	Recovery w/ Waiting	Avoidance	Avoidance w/ Waiting
1	26.32	28.97	491.85	33.62	35.47	569.00	662.30
2	42.80	44.67	510.71	50.19	56.70	619.56	655.15
3	51.44	52.73	490.48	59.89	61.95	599.29	632.47
4	56.96	56.92	477.92	63.64	67.02	580.51	623.12
10	67.22	67.20	471.20	76.66	82.69	601.95	602.09
20	74.42	70.58	450.91	83.47	88.70	628.46	620.91
40	87.78	74.74	1419.22	86.08	93.38	602.21	595.28

Parameter	Value
Routing Probability	(1) 0.10, (4) 0.45, (5) 0.45
Clustering case (5 PEs)	(1) (2) (3) (4) $(5)^{-1}$

Node numbers refer to Figure 1. Parenthesized node groups execute on one processor.



Figure 1a Central Server Queueing Model



Figure 1b RESQ Representation of Central Server Model



Figure 2 Sequent Balance 21000 Configuration













Figure 7 Cluster network





Parameter	Value
Node Service Time	1 / number of server nodes
Confidence Level	95%
Speedup Base	One processor static deadlock recovery

Figure 9 Speedup for tandem queue with deadlock recovery (dynamic node assignment)

- Speedup



Parameter	Value
Node Service Time	1 / number of delay nodes
Confidence Level	95%
Speedup Base	One processor static deadlock recovery

Figure 10 Speedup for four node cyclic queue (static node allocation)



Parameter	Value
Node Service Time	0.25
Confidence Level	95%
Speedup Base	One processor deadlock recovery
Cluster case A	$(1) (2) (3) (4)^{-1}$
Cluster case B	(1 2) (3 4)

 $^{-}$ Node numbers refer to Figure 6. Parenthesized node groups execute on one processor.





Parameter	Value
Node Service Time	0.25
Confidence Level	95%
Speedup Base	One processor deadlock recovery





Parameter	Value
Node Service Time	0.25
Confidence Level	95%
Speedup Base	One processor deadlock recovery





Parameter	Value
Routing Probability	(1) 0.10, (4) 0.45, (5) 0.45
Confidence Level	95%
Speedup Base	One processor static deadlock recovery
Cluster case A	$(1) (2) (3) (4) (5)^{-1}$
Cluster case B	$(1\ 2)\ (3)\ (4\ 5)$
Cluster case C	(1 2) (3 4 5)

Node numbers refer to Figure 1. Parenthesized node groups execute on one processor.



Speedup



Parameter	Value
Routing Probability	(1) 0.10, (4) 0.45, (5) 0.45
Confidence Level	95%
Speedup Base	One processor static deadlock recovery

Figure 15 Deadlock fractions for five node central server (static node assignment)



Parameter	Value
p0 Routing Probability	(1) 0.0, (4) 0.5, (5) 0.5
p1 Routing Probability	$(1) \ 0.10, \ (4) \ 0.45, \ (5) \ 0.45$
Confidence Level	95%
Speedup Base	One processor static deadlock recovery
Cluster case A	(1) (2) (3) (4) (5)
Cluster case B	(1 2) (3) (4 5)
Cluster case C	(1 2) (3 4 5)

Node numbers refer to Figure 1. Parenthesized node groups execute on one processor.



Fraction



Parameter	Value
p0 Routing Probability	(1) 0.0, (4) 0.5, (5) 0.5
p1 Routing Probability	(1) 0.10, (4) 0.45, (5) 0.45
Confidence Level	95%
Speedup Base	One processor static deadlock recovery
Cluster case A	(1) (2) (3) (4) (5)
Cluster case B	(1 2) (3) (4 5)
Cluster case C	(1 2) (3 4 5)

 $^{-}$ Node numbers refer to Figure 1. Parenthesized node groups execute on one processor.



Speedup



Parameter	Value
Node Service Time	1.0
Confidence Level	95%
Speedup Base	One processor deadlock recovery
Cluster case A	(1) (2) (18)
Cluster case B	(1 4) (2 3) (5 8) (6 7) (9 10)
Cluster case C	(1 2 3 4) (5 6 7 8) (9 10)

Node numbers refer to Figure 7. Parenthesized node groups execute on one processor.

Figure 18 Deadlock and null fractions for four block cluster (static node assignment)



Node numbers refer to Figure 7. Parenthesized node groups execute on one processor.



Speedup



Parameter	Value
Node Service Time	1.0
Confidence Level	95%
Speedup Base	One processor deadlock recovery



Speedup



Parameter	Value
Node Service Time	1.0
Confidence Level	95%
Speedup Base	One processor deadlock recovery