

Operated by Universities Space Research Association



Institute for Computer Applications in Science and Engineering
NASA Langley Research Center
Hampton, VA 23681-0001

Contract NAS1-19480
December 1994

(NASA-CR-195011) ON EXTENDING
PARALLELISM TO SERIAL SIMULATORS
Final Report (ICASE) 18 p

David Nicol
Philip Heidelberger

ON EXTENDING PARALLELISM TO
SERIAL SIMULATORS

 ICASE

NASA Contractor Report 195011
ICASE Report No. 94-95

N95-18931

Unclass

03/61 0035843

4N-61
E5843
18P

On Extending Parallelism to Serial Simulators

David Nicol*

Department of Computer Science

The College of William and Mary

Williamsburg, VA 23185

Philip Heidelberger

IBM T.J. Watson Research Center

P.O. Box 704

Yorktown Heights, NY 10598

Abstract

This paper describes an approach to discrete event simulation modeling that appears to be effective for developing portable and efficient parallel execution of models of large distributed systems and communication networks. In this approach, the modeler develops sub-models using an existing sequential simulation modeling tool, using the full expressive power of the tool. A set of modeling language extensions permit automatically synchronized communication between sub-models; however, the automation requires that any such communication must take a non-zero amount of simulation time. Within this modeling paradigm, a variety of conservative synchronization protocols can transparently support conservative execution of sub-models on potentially different processors. A specific implementation of this approach, U.P.S. (Utilitarian Parallel Simulator), is described, along with performance results on the Intel Paragon.

*This work is supported in part by NSF grant CCR-9201195. It is also supported in part by NASA contract number NAS1-19480 while the author was a consultant at the Institute for Computer Applications in Science and Engineering (ICASE), NASA Langley Research Center, Hampton, VA, 23681.

1 Introduction

Few in the parallel discrete event simulation (PDES) community would argue with the assertion that PDES has not yet made a significant impact on discrete-event simulation practitioners. Indeed, this was the topic of a panel discussion at the 1994 PADS (Parallel and Distributed Simulation) Conference. During this panel discussion, a number of issues were raised and discussed:

1. This lack of impact is partially due to a lack of adequate PDES modeling tools, see also Fujimoto's discussion in [6] and Bagrodia's comment "The tools, stupid" [2]. The fact remains that efficient PDES requires the use of sophisticated techniques that are often application specific, platform specific, or both.
2. Fujimoto's "Holy Grail" does not currently exist. The "Holy Grail" provides arbitrary modeling capability and automatically parallelizes to yield significant speedups. Although one person out of a total attendance of approximately 80 disagreed, Fujimoto asserted that such a "Holy Grail" would not be available "in this century".
3. The best way for PDES to have impact in the near future is to focus on specific applications. One way to have maximum impact is to provide application specific PDES libraries or tools.

There is little to suggest that the "Holy Grail" will exist any time soon. PDES is an unusually tricky branch of parallel processing, and the overwhelming experience to date in parallel software development is that high performance requires hand-crafted attention to application and platform specifics, especially with respect to load-balancing, communication, and synchronization costs. Only recently have tools been developed to automate the generation of efficient communication and synchronization (e.g., High Performance Fortran), and even there the domain of applicability is limited to regular problems, and the user must map the workload to processors.

The existing parallel simulation literature is justifiably viewed from the outside as having little relevance to industrial simulation. The ideas found in that literature frequently lack a clear link to the computational models and the type of tool to which a simulation practitioner is accustomed. While one can understand that a researcher's goals and resources are different from an industrial simulation user, it comes as no surprise that parallel simulation theory has apparently not yet had an impact on industrial simulation practice.

Recognizing these problems, a number of tools for parallel simulation have been developed, with the goal of making the parallelism and synchronization more transparent to the user. The Army funded an effort at the Jet Propulsion Lab to develop the Time Warp Operating System (TWOS), to support the development of parallelized combat simulations. TWOS uses the event-oriented paradigm, and while the TWOS group no longer exists, TWOS is still available (without support) from NASA's COSMIC software clearing house. Jade Simulation developed sim++, a process-view system which was designed from the bottom to support Time Warp simulation. The Army also supported a port of the commercial simulation language ModSim (developed by CACI) onto TWOS, and then supported a port of ModSim to Jade's system. While we shall leave to the historians the task of analyzing the story of these early efforts, it is safe to say that the efforts were technically ambitious, but did not meet with the hoped-for level of success (e.g., see [12]). Even the most enthusiastic supporter of these efforts would agree with the assertion that issues related to synchronization (e.g., state-saving) ended up migrating to the modeler level by necessity, to escape unduly large overheads when purely automated means were employed. None of these products were widely used.

An optimistic simulation tool, SPEEDES[14], may be licensed from the Jet Propulsion Lab; SPEEDES is based in C++, and has the event-oriented world view. Little is known of its performance however, especially on large-scale parallel computers. Many of the details of synchronization are the responsibility of the modeler.

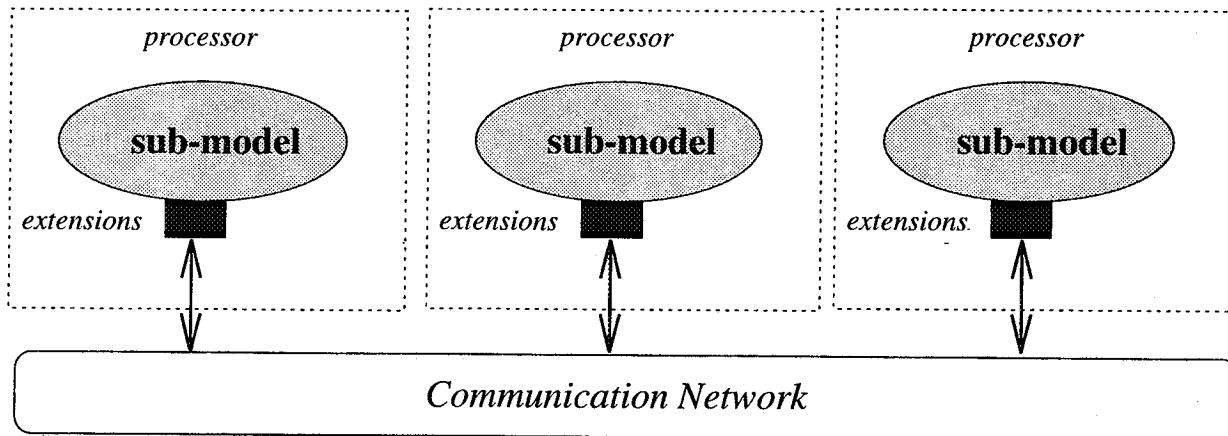


Figure 1: View of submodels and extensions

no knowledge of the simulation sub-model, except that which might be explicitly provided by the sub-model.

These considerations constrain the inter-processor synchronization to be conservative. This in turn constrains the models that may be parallelized, and constrains their partitioning. Furthermore, the burden of partitioning and mapping the simulation is the modelers. This latter requirement is still usual for parallel processing in general; as for the former requirement, there is one essential model characteristic which allows us to synchronize automatically—lookahead. In this context, lookahead is the ability to predict (or lower bound) when in the future one sub-model may affect another. A concrete example is a job entering service at a non-preemptive queue on sub-model i at time s and then being routed to another queue on sub-model j when its service is over at time t . Or, it may be that the routing is state-dependent, in which case i knows at time s that it may affect at time t any one of the queues to which it routes jobs, it simply doesn't know which one.

Lookahead is necessarily application dependent. However, in the context of computer and communication system simulations one generally uses standard abstractions such as queues, petri nets, or finite-state machines. The parallelizing extensions to such tools may provide special cases of these, designed in such a way that the lookahead calculation and dissemination may be automated.

The discussion above identifies restrictions and limitations on the synchronization method used by the extensions. The base simulation tool must also have characteristics that permit one to graft on extensions. Foremost is the need for the simulation to provide a “trapdoor” to a base computer language, such as C. This is absolutely required, to provide access to the synchronization and communication primitives available on the tool's platform. Second, the tool must provide a means for the extension to cause an extension-defined synchronization function or routine to be invoked at a specific future point in simulation time. That function must be able to block all further simulation processing until (through the extension's synchronization activity) it is logically safe to do so. A third requirement is that a tool permit one to reset the random number generator seed. This is required to allow the extension to create separate independent random number streams on each sub-model. We know these requirements are met by several simulation tools.

A last requirement is that the simulation be partitioned so that a message from one sub-model carries all of the information needed by the recipient to correctly continue processing of the passed information. As a counter-example, RESQ has mechanisms to “split” a job into sub-jobs, and later “join” those sub-jobs. RESQ's mechanism for identifying the siblings of a split job is to have them contain a memory pointer

queue is FCFS, a job's departure time can be predicted exactly upon its arrival. Under other queueing disciplines (e.g. processor-sharing, round-robin) lower bounds on future departure times can be computed. Correct lower bounds can always be computed under any discipline such that a job's departure time cannot become earlier by virtue of an additional arrival to the queue. The answer to question (4) depends on the model. Sometimes the message reports a job transfer, with the identity of the job known at the time it enters service. Other times the message reflects the state of the sending sub-model at the time the job departs the server. The answer to question (5) is frequently "yes" in stochastic simulations where routings are chosen at random, independently of the state of the system. It is "no" if the routing is state-dependent, e.g., join the shortest queue among all destination possibilities. Question (6) asks whether special synchronization techniques can be applied. In particular, if the queue's output can be uniformized (which is true for a queue with a finite number of servers and Coxian phase-type service distributions), then techniques developed in [7] can be applied.

The degree to which a queue can predict its future behavior depends on the various considerations just listed. In our approach, when a sub-model specifies the carriers it will use, it also provides parameters that specify answers to the questions above. These parameters govern how lookahead is computed at the carrier. Only the author of the extensions need be concerned about the lookahead computations; the lookahead and synchronization based upon it are transparent to the simulation modeler.

There are at least three well-studied conservative synchronization protocols suitable for the extension. The selection of a protocol is constrained by the model characteristics. The YAWNS window-based protocol [8, 10] is appropriate when the message associated with a job completion can always be generated and sent a minimum of some $X > 0$ (which may be randomly sampled) simulation time units before the job completes service. Barring any further refined lookahead, the YAWNS protocol will cause all processors to synchronize globally every X_{\min} units of simulation, where X_{\min} is the minimum lookahead among all carriers in the simulation. YAWNS should be used if the connectivity between sub-models is high and the routing decisions cannot be predicted in advance. To attain good performance, this protocol does require that substantial sub-model simulation activity occur (on average) every $E[X_{\min}]$ units of simulation time, to amortize the cost of global synchronization.

The Appointments protocol [9] requires each carrier to maintain, for each sub-model to which it may route jobs, a lower bound (the appointment) on the next time at which it might next send a package there. The sub-model that creates the appointment is responsible for updating them. A sub-model that receives an appointment at time t will not advance past time t until the sub-model that established that appointment releases it to do so. Appointments should be used when a carrier's connectivity to other sub-models is low; they are especially effective if the service times and routing destinations can be pre-sampled.

The PUCS protocol [7] may be used when all carriers have Markovian service distributions that may be pre-sampled, and have routing distributions that may also be pre-sampled. The details are described elsewhere; essentially PUCS is the appointments protocol, with lookahead derived from the mathematical structure of the carriers.

Details of how one constructs carriers and communicates messages and appointments (and/or synchronizes globally) will vary with simulation system, and execution platform. However, the essential ingredients needed to support these extensions are access to an underlying computer language (e.g., C or C++), and the ability to schedule future events at known simulation times. Access to an underlying language permits access to communication and synchronization routines. An ability to schedule future events at known time-stamps gives a sub-model the ability to block at an appointment or synchronization instant. The sub-model just schedules an event that (in the case of appointments) waits for a message signaling release by the scheduling sub-model, or (in the case of YAWNS) engages in a global reduction to synchronize and establish the next

name of the associated station, and the identity of the processor on which that station resides. Thus just as CSIM requires a **mailbox** call for each mailbox in the model, U.P.S. requires a **remote-mailbox** call for each remote-mailbox in the model. When a process gives a package to a carrier it may specify the remote-mailbox that describes the destination. U.P.S. then takes care of the actual package transmission and receipt.

It is also possible to associate a *routing distribution* with a carrier. In this case a process cedes to the carrier the selection of package destination, but provides a probability distribution the carrier uses. This mechanism supports multi-casts; each element of the random routing table is a list of remote-mailboxes. When selected, every station associated with a remote-mailbox in the list receives a copy of the package.

In order to send a package, a CSIM process builds a package and calls an U.P.S. routine to cause this package to be delivered. In the case of a carrier without an associated routing distribution, this call is simply **ship**(carrier,msg,len,rbox) where carrier is the identity of the carrier, msg is the starting address of the package, len is the package's length, and rbox is the name of the remote-mailbox to which the package is to be sent. In CSIM, messages are sent to ordinary mailboxes with a **send**(mailbox,msg), thus the U.P.S. call naturally extends the CSIM call. As regards time advancement, U.P.S. mimics normal CSIM usage by not returning control to the process until the delivery time of the package. Functionally then, passing a package to a carrier is no different from a CSIM process running through the sequence of acquiring a facility, holding it for the service duration, delivering a message in a mail-box at the service end, and releasing the facility. In order to receive a package sent by a process in a different sub-model, a CSIM process accesses a station just as though it is an ordinary CSIM mailbox, because it is an ordinary CSIM mailbox! The only purpose for distinguishing a station from a mailbox at declaration time is to allow U.P.S. to maintain a list of stations and their names, so that incoming messages can be properly delivered.

The U.P.S. handling of packages deserves comment. A sub-model can accept inter-processor messages only from within an U.P.S. extension routine. To avoid buffering problems we must ensure that an extension routine is called frequently to look for newly arrived messages, and move them into the simulator's address space. This is easily accomplished by creating a CSIM process that just loops through two activities (i) call **hold**(*t*) for some time *t*, (ii) check for, receive, and deal with any new messages. Currently, with YAWNS, this process is invoked only at end of windows at which time global synchronization is required; clearly this could (and should) be modified so as to pick up messages more frequently should there be heavy message traffic. Messages can also be sought from within any other U.P.S. routine. The routine that looks for messages just creates a CSIM process for that message, passing to the process the arrival time and other particulars. The newly arrived process does a **hold**() to suspend itself until the arrival instant, deposits the package into the appropriate station, and then departs. Appointments are handled slightly differently. A sub-model maintains a list of all exterior carriers that send it appointments. An "in" and an "out" appointment count is maintained for each. An appointment message always serves both to release the sub-model from the last appointment made by the carrier, and to establish a new appointment. Receiving an appointment message, a sub-model increments the "in" count associated with the sending carrier, holds until the appointment time, then enters a loop where it waits for the "in" count to increase beyond the "out" count. Inside of this loop the code looks for new messages, and calls **hold**(0), which gives any other process with the same activation time a chance to be processed. On exiting the loop, the carrier's "out" count is incremented, and the process terminates.

Presently U.P.S. supports the YAWNS and PUCS protocols. An U.P.S. statement declares which one is to be used. We are actively incorporating the appointments protocol, and considering how to allow mixtures within the same simulation. We are also implementing support for queries, and for "remote-storages" which will allow one sub-model to contend for and acquire CSIM storage variables on another sub-model. Finally,

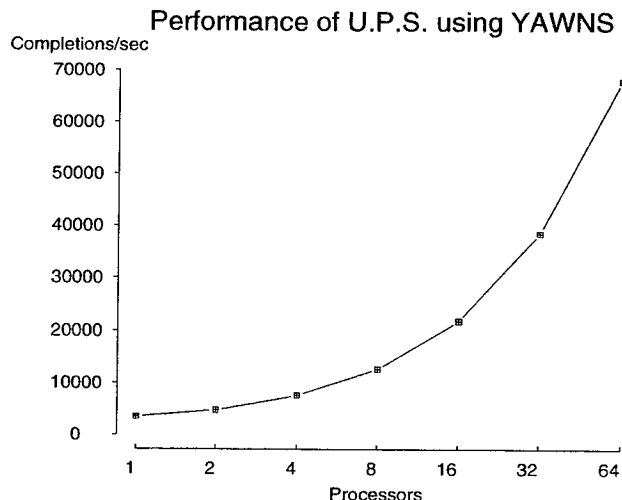


Figure 3: Performance of U.P.S. using the YAWNS protocol. I/O service distribution is 0.5 plus an exponential with mean 0.5.

server I/O queue, and mailboxes where a CPU queue looks for job arrivals. Next we modified this model for U.P.S. just by making the multi-server I/O queue a carrier, and by setting up remote mailboxes and stations for every central server subsystem. Then, to investigate how well the YAWNS protocol might work under advantageous conditions, we assigned each I/O device a service time of 0.5 plus an exponential with mean 0.5. Figure 3 then plots the aggregate job completion rate as we vary the number of processors through powers-of-two, keeping the load on each processor fixed at one basic network size, as described earlier. The total problem size varies through the same values as were studied for the serial CSIM case. The 1 processor U.P.S. rate provides insight into the overheads of executing YAWNS on this problem, overheads due primarily to increased computation for lookahead and window calculations, and to extra logic in U.P.S. that is executed at the point a job leaves service at an I/O device. The model studied has a very high U.P.S. component, and the overhead slows the U.P.S. execution rate to about 70% of native CSIM. As we simultaneously increase the problem size and number of processors, we observe that the aggregate job completion rate increases. Compared with the corresponding native CSIM runs, we obtain speedups of 1.02, 1.8, 3.22, 6.20, 13.3, and 42.5 on power-of-two numbers of processors between 2 and 64.

By changing the I/O service distribution to be exponential we may use the PUCS synchronization mechanism. The results of doing so are presented in Figure 4 for two cases. The “light load” case is the same as the problem studied under YAWNS, except for the service times. The overheads in this situation are quite large for PUCS. Among the underlying synchronization messages, the ratio of overhead messages to messages that carry jobs is 10 to 1. By increasing p_r to 0.999 and quadrupling the total number of jobs, we achieve rather better performance. This is due both to a better computation to communication ratio (as a result of increasing p_r) and a better ratio of overhead messages to job-carrying messages, (1.5). The wide variation in performance shows the sensitivity of the method to problem characteristics.

We also used YAWNS to simulate the light workload model, and found it to be very disadvantageous. No problem/processor configuration achieved an aggregate rate larger than than 2000 completions per second.

A second problem is a network of switches, such as are used to establish long-distance phone calls. The

Performance of U.P.S. on telephone switching network

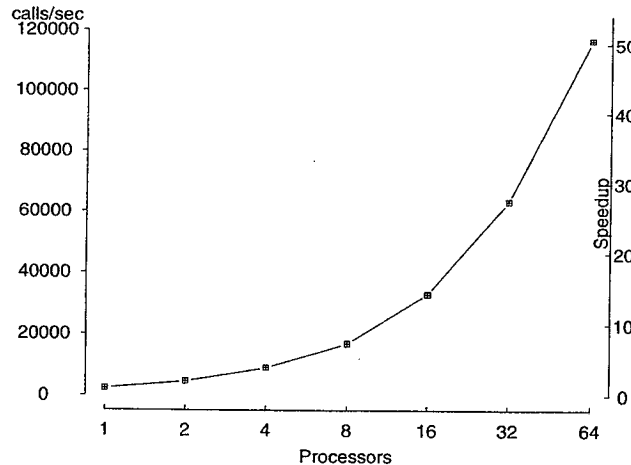


Figure 5: Performance of telephone switching network

output port. In addition, there is a propagation delay for transferring packets from one switch to another. This propagation delay is modeled by a simple `hold` statement in the case that the sending and receiving switches are co-resident on the same processor while carriers are used if the switches are not co-resident. These propagation delays are deterministic; depending only on the distance between the switches. Thus the carriers have an infinite number of servers with the constant c equal to the propagation delay (and the random part equal to zero). Thus the model uses the YAWNS protocol.

Packets arrive from the outside according to a batch Poisson process; the batch size is uniformly distributed between one and some maximum value. This permits some burstiness in the arrival process. When a batch of packets arrives, a destination switch is chosen uniformly among the other switches. This is something of a worst case as far as parallel simulation goes since it results in both a single packet crossing multiple processor boundaries and induces a load imbalance. (Switches in the center of the mesh are more heavily loaded than those at the edges.) More elaborate arrival processes and flow control (e.g., leaky buckets) could be added to the model; these would tend to increase the computation to communications ratio of the model.

Even at the speed of light, propagation delays down the fibers provide excellent lookahead. For example, suppose the output ports operate at 100 megabits/second and the switches are 50 miles apart (the parameters used in our model). Then it takes about 4 microseconds to put a 53 byte packet on the fiber, while it takes about 270 microseconds to transmit it from one switch to the next at the speed of light. Thus intra-switch events happen about 67 times faster than inter-switch events. This ratio provides excellent lookahead for the purpose of parallel simulation.

A purely sequential version of this model (one that does not use any U.P.S. constructs) took several days to develop. Converting this model to incorporate the U.P.S. constructs took approximately another day. This primarily involved extending the model's data structures to support the notion of sub-models, which are basically just rectangular arrays of switches with one carrier and one station. This version of the model was debugged on a single processor. When put on the Paragon, it took an additional 15 minutes of debugging before the model ran correctly; the error was a simple indexing mistake involving the mapping of sub-models to processors. In developing this model, most of the attention was paid to model correctness aspects. (CSIM

Performance of U.P.S. on ATM Model (16x16 network)
normalized procs x time

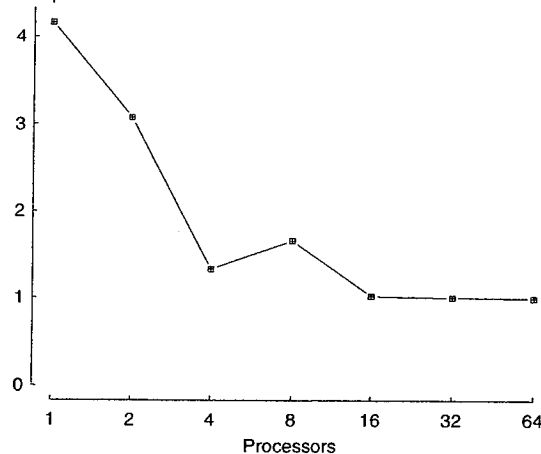


Figure 7: Normalized processor \times time product, illustrating super-linear accelerations on large numbers of processors.

partitions the model into sub-models, one per processor, and incorporates the extension constructs into those sub-models. The extensions are designed to provide essentially the same functionality as basic constructs in the base simulator, but to do so in a way that the inter-processor synchronization and communication is automated. The distributed model looks very much like a non-distributed model. In U.P.S. , a small set of additional calls provide the parallel simulation capability. Almost all of the model is ordinary C and CSIM.

The methodology very much involves the modeler, who must partition the model, and must consider the performance tradeoffs between different synchronization strategies when performing that partition. However, the simulation modeler does not have to implement the communication and synchronization. That is done once, by the extension's author. This provides industry with a lower risk path to parallel processing, and an ability to experiment with it with low software cost. In view of the general perception of parallel simulation as tricky business, we believe this approach proves an important first step towards making parallel simulation safe for the masses.

We are exploring this methodology by writing an extension library, U.P.S. , for the commercial tool CSIM (now distributed by Mesquite Software). It is a natural match, as CSIM modelers must be C programmers, and should easily adapt to the interacting sub-models view of the simulation. A CSIM modeler needs to incorporate only a few U.P.S. constructs into their models. To use U.P.S. , it helps for a CSIM modeler to understand the fundamentals of SPMD programming, e.g., partitioning and mapping of submodels onto a parallel processor. U.P.S. currently runs on the Intel Paragon. This paper provides preliminary performance results of its performance on three models: a network of central server computer systems, a telephone switching network, and an ATM network. We see that good performance is delivered on models where there is ample parallelism which is easily abstracted.

We are continuing to extend the capabilities of U.P.S. , and seek to expand its user base and to port it to other platforms. We are implementing an MPI (Message Passing Interface) version; MPI is the newly emerging standard for message-passing applications. When operational, U.P.S. will execute on any platform supporting both CSIM and MPI, which we expect to include most parallel systems as well as networks of

- [13] H. Schwetman. CSIM : A C-based, process oriented simulation language. In *Proceedings of the 1986 Winter Simulation Conference*, pages 387–396, 1986.
- [14] J.S. Steinman. Speedes: Synchronous parallel environment for emulation and discrete event simulation. In *Advances in Parallel and Distributed Simulation*, volume 23, pages 95–103. SCS Simulation Series, Jan. 1991.

REPORT DOCUMENTATION PAGE			Form Approved OMB No. 0704-0188	
Public reporting burden for this collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington, VA 22202-4302, and to the Office of Management and Budget, Paperwork Reduction Project (0704-0188), Washington, DC 20503.				
1. AGENCY USE ONLY (Leave blank)	2. REPORT DATE December 1994	3. REPORT TYPE AND DATES COVERED Contractor Report		
4. TITLE AND SUBTITLE ON EXTENDING PARALLELISM TO SERIAL SIMULATORS		5. FUNDING NUMBERS C NAS1-19480 WU 505-90-52-01		
6. AUTHOR(S) David Nicol Philip Heidelberger				
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) Institute for Computer Applications in Science and Engineering Mail Stop 132C, NASA Langley Research Center Hampton, VA 23681-0001		8. PERFORMING ORGANIZATION REPORT NUMBER ICASE Report No. 94-95		
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES) National Aeronautics and Space Administration Langley Research Center Hampton, VA 23681-0001		10. SPONSORING/MONITORING AGENCY REPORT NUMBER NASA CR-195011 ICASE Report No. 94-95		
11. SUPPLEMENTARY NOTES Langley Technical Monitor: Michael F. Card Final Report Submitted to Workshop on Parallel and Distributed Simulation				
12a. DISTRIBUTION/AVAILABILITY STATEMENT Unclassified-Unlimited Subject Category 60, 61		12b. DISTRIBUTION CODE		
13. ABSTRACT (Maximum 200 words) This paper describes an approach to discrete event simulation modeling that appears to be effective for developing portable and efficient parallel execution of models of large distributed systems and communication networks. In this approach, the modeler develops sub-models using an existing sequential simulation modeling tool, using the full expressive power of the tool. A set of modeling language extensions permit automatically synchronized communication between sub-models; however, the automation requires that any such communication must take a non-zero amount of simulation time. Within this modeling paradigm, a variety of conservative synchronization protocols can transparently support conservative execution of sub-models on potentially different processors. A specific implementation of this approach, U.P.S. (Utilitarian Parallel Simulator), is described, along with performance results on the Intel Paragon.				
14. SUBJECT TERMS simulation; parallel processing		15. NUMBER OF PAGES 17		
		16. PRICE CODE A03		
17. SECURITY CLASSIFICATION OF REPORT Unclassified	18. SECURITY CLASSIFICATION OF THIS PAGE Unclassified	19. SECURITY CLASSIFICATION OF ABSTRACT	20. LIMITATION OF ABSTRACT	