

## DSIM: SCALING TIME WARP TO 1,033 PROCESSORS

Gilbert Chen

The MathWorks, Inc.  
3 Apple Hill,  
Natick, MA 01760, U.S.A.

Boleslaw K. Szymanski

Department of Computer Science  
Rensselaer Polytechnic Institute  
110 8<sup>th</sup> Street  
Troy, NY 12180, U.S.A.

### ABSTRACT

This paper presents the design, implementation and performance of a Time Warp simulator, called DSIM, which targets clusters comprised of thousands of processors. DSIM employs a novel technique for GVT computation, called the Time Quantum GVT algorithm that requires no message acknowledgement, relies on constant-length messages and is efficient on clusters with large numbers of processors. Its implementation uses a technique called Local Fossil Collection to alleviate the overhead of memory reclamation and to support efficient event management. DSIM is also equipped with a simple programming interface to ease programming and debugging of simulations. Experimental results obtained on the PHOLD benchmark demonstrated that DSIM can process as many as 228 million events per second on 1033 processors.

### 1 INTRODUCTION

Clusters have steadily and gradually become the main stream of high performance computing platforms – the number of clusters on the top 500 supercomputer list (<http://www.top500.org/>) jumped from 208 in November 2003 to 291 in June 2004, an increase of 40% in merely 7 months. It is likely that this trend will continue in the future, as technological progresses on single processors and fast interconnection networks has made building clusters from commodity components easier and more cost-effective. The popularity of clusters brings a new challenge to PDES (Parallel Discrete Event Simulation) researchers designing efficient Time Warp simulators. Historically, at different stages of development, Time Warp simulators have always targeted the most widely available platforms at the time. For instance, the first generation Time Warp simulators (Jefferson et al.; Jefferson 1985; Baezner, Lomow and Unger 1994) were usually designed for parallel computers, especially MPP (Massively Parallel Processing) machines. When shared-memory computers become more prevailing in the nineties, the second generation Time Warp simulators,

such as GTW (Das et al. 1994) and ROSS (Carothers, Bauer and Pearce 2000), were developed to support this type of parallel computers. Today, as clusters are starting to dominate high-end computing, new Time Warp simulators capable of running efficiently on networks of distributed processors are particularly desired.

There are some Time Warp simulators such as WARPED (Martin et al. 2003) and ParaSol (Mascarenhas, Knop and Rego 1995) that were initially designed for distributed memory parallel computers. One may argue that these simulators can be ported to clusters with a little effort, as modern clusters may exhibit behavior that is close or comparable to distributed memory MPP computers, in aspects such as bandwidth, latency, and scale. However, little data have been available to prove these simulators' effectiveness on hundreds of processors. In fact, the largest Time Warp simulations that have ever been attempted, in terms of the number of processors used, to the best of our knowledge, were executed on 64-104 processors (Wieland et al. 1989; Fujimoto 1990; Kim and Jean 1996). Hence, the question remains whether Time Warp simulators running on modern large-scale clusters can attain satisfactorily good performance.

DSIM has been developed to support efficient Time Warp simulation on distributed clusters with thousands of processors. At the heart of DSIM is a new GVT (Global Virtual Time) algorithm, referred to as the Time Quantum GVT (TQ-GVT) algorithm that does not require message acknowledgement, and relies on short messages with constant length. This paper demonstrates that TQ-GVT is able to deliver a continuous stream of up-to-date GVT values during simulation. It achieves this by devoting one or more but less than 1% of the total number of processors to running the core of the GVT algorithm. In view of the large number of processors involved in the simulation on a large cluster, such a solution is more efficient than delaying all processors by a small fraction of time in a distributed implementation of the GVT computation.

In addition to the new GVT algorithm, DSIM uses a modified fossil collection mechanism called Local Fossil

Collection, as well as a careful implementation of event management mechanism. Another goal of DSIM design has been to demonstrate that the improvement of programmability of a parallel discrete event simulator can be accomplished without sacrificing its efficiency.

The primary goal of this paper is to introduce the basic features of DSIM to enable replication of our performance results. To evaluate the DSIM performance, the PHOLD model was chosen as an example of a non-trivial modeling problem. It needs to be demonstrated if the reported performance is repeatable for other types of models. Likewise, it is yet to be established which new features of DSIM are indispensable for high performance of Time Warp on large-scale clusters. Still, it is safe to predict that the new GVT algorithm has played a role in enabling DSIM superior performance.

The paper is organized as follows. Section 2 introduces the basic ideas behind the Time Quantum GVT algorithm. Section 3 presents DSIM implementation details, including Local Fossil Collection and the memory efficient event management mechanism. Section 4 describes the programming interface of DSIM in detail. Section 5 presents the experimental results of DSIM running on two different yet popular clusters. Section 6 concludes the paper by providing a summary of the basic features of DSIM and a vision of its future development.

## 2 TIME QUANTUM GVT ALGORITHM

The GVT is defined as the minimum of local simulated times on all processors and timestamps of all messages in transit (Fujimoto, 1989). A good GVT algorithm is critical to the overall performance of a Time Warp simulator. An event processed earlier than the GVT will not be subject to rollback under any circumstances, and therefore its associated memory can be released permanently. With more accurate GVT estimates, more obsolete memory can be reclaimed, decreasing the chance for performance losses stemming from the memory system bottleneck.

Numerous GVT algorithms (Samadi 1985; Preiss 1989; Bellenot 1990; Baldwin, Chung and Chung 1991; Bauer and Sporrer 1992; Mattern 1993; Srinivasan and Reynolds 1993; Tomlinson and Garg 1993; D'Souza, Fan and Wilsey 1994; Das and Sarkar 1995; Steinman et al. 1995; Choe and Tropper 1998; Perumalla and Fujimoto 2001) have been proposed for distributed execution of Time Warp. Some of them depend upon message acknowledgement (Samadi 1985; Bellenot 1990; D'Souza, Fan and Wilsey 1994) to track those messages that are still in transit when the GVT computation is being conducted. Others (Bauer and Sporrer 1992; Mattern 1993; Tomlinson and Garg 1993) utilize Vector Clock or similar structures whose size is proportional to the number of processors. Still others (Preiss 1989; Baldwin, Chung and Chung 1991; Das and Sarkar 1995; Choe and Tropper 1998) take a dis-

tributed and passive approach, which requires the processor in need of memory to send out a special token message that will traverse all other processors. Synchronization-based GVT algorithms (Srinivasan and Reynolds 1993; Steinman et al. 1995; Perumalla and Fujimoto 2001) rely on global reductions to determine whether previously sent messages have all been received.

Instead of implementing one of these GVT algorithms, yet another GVT algorithm, referred to as Time Quantum GVT (TQ-GVT), is contributed to this already rich repository. Its uniqueness relies on allocating one processor, called the GVT master, exclusively to the GVT computation. The algorithm's name comes from the fact that the simulated time is divided into a sequence of non-overlapping time intervals (quanta) with equal width, somehow similarly as in (Perumalla and Fujimoto 2001). A GVT master monitors the numbers of sent and received messages within each time quantum for all simulating processors reporting to it. If the two numbers associated with a given time quantum are equal, then all messages sent during this time quantum must have already been received and they will not be counted when computing the new GVT value. This well-known method of accounting for transient messages can be traced back to Mattern (Mattern 1993).

The use of an exclusive processor for GVT computation may appear an obstacle to scalability, since it is basically a centralized approach. Nevertheless, one or more levels of intermediate GVT masters can be introduced, each of which keeps track of the number of transient messages in each time quantum for a subset of processors. These numbers must then be reported to the root GVT master, which in turn determines whether or not there are still messages in transit from each time quantum and calculates the GVT accordingly. Empirically, one GVT master can drive as many as 128 simulating processors, so an extra level of intermediate GVT masters could easily extend the number of simulating processors to 16,384 and if this is not sufficient, more levels of GVT masters can be added. A very small percentage of all processors, e.g., 129 out of 16,513, or merely 0.78 percent, will not be directly participating in the simulation. Hence, such a solution is suitable for clusters in which there are a large number of processors involved in a computation.

An essential feature that makes TQ-GVT scale so well is that it overlaps GVT computation with other tasks. Simulating processors are only engaged in processing events and transmitting and receiving messages. With regards to GVT computation, they just need to send report messages periodically, and receive GVT messages as they come. If GVT messages do not come on time, simulating processors are never delayed or blocked, as long as the amount of memory is adequate. As a result, only the masters are involved in a significant amount of computation of the GVT. The cost of GVT computation for simulating processors is non-blocking

send of a report message at the end of each quantum and non-blocking receive of the new GVT value if there is a GVT message. Thus, this solution eliminates two problems: (i) the large latency of the interconnection network often found in clusters, and (ii) the asynchrony that causes different processors to reach the synchronization point at vastly different wall-clock times when the number of involved processor exceeds, say, 1,000.

On the first glance, the TQ-GVT looks similar to the LBTS algorithm that chops the simulation time into bands or epochs (Perumalla and Fujimoto 2001). However, this superficial similarity disappears as soon as one considers the fact that the performance study in (Perumalla and Fujimoto 2001) was limited to 16 processors. Three essential differences distinguish TQ-GVT from the LBTS algorithm. First, multiple reductions may be needed by the LBTS algorithm within each band for all transient messages sent during the previous band to be received, while in TQ-GVT, each simulating processor is guaranteed to send only one report message and receive at most one GVT message per time quantum. Second, in the LBTS algorithm, during each band every processor must perform both global reduction and event processing. This, unfortunately, increases the latency of handling global reduction messages, because to maximize the event rates, each processor needs to execute a certain number of events consecutively before checking the message receive buffer. In contrast, in TQ-GVT, simulating processors are only engaged in event processing and message exchanges. It is the GVT master that carries out the global reduction, and consequently it can respond to GVT-related messages more quickly than simulating processors could, since this is its only job. Third, and perhaps the most important one, TQ-GVT never delays or blocks simulating processors at the end of the time quantum, and hence no temporary disturbances on communication can impact the progress of time quanta. In the LBTS algorithm, however, a new band may start only after all transient messages from the previous band have been accounted for.

Two methods can be used for advancing time quanta. The first method depends upon a hardware clock to advance the time quantum at fixed intervals. Clocks on different processors, however, need not be precisely synchronized, since a processor that is too far ahead or behind only causes some messages to be counted in the next time quantum. The other method simply requires that the GVT master broadcast a special message periodically, upon the receipt of which each simulating processor will advance the time quantum. The latter method requires more bandwidth.

The algorithm is briefly described below:

- At the beginning of the simulation, all processors, including the GVT master, perform barrier synchronization.

- When a message is sent, it is always marked with the current time quantum of the sender.
- During the  $i$ th time quantum, the  $j$ th simulating processor keeps track of  $S_{ij}$ , the number of messages sent by itself,  $OT_{ij}$ , the smallest timestamp among all messages sent in this quantum, and  $\{R_{k,j}\}$ , an array of integers indicating the number of messages received marked with a time quantum  $k$ .
- At the end of the time quantum  $i$  (i.e., when it is time to move to the next time quantum), the  $j$ th processor reports  $S_{ij}$ ,  $OT_{ij}$ ,  $\{R_{k,j}\}$ , as well as  $i$ , the current time quantum index, and  $T_{ij}$ , the local virtual time, to the GVT master. As an example, suppose that a processor, in the time quantum 10, has sent 6 messages, with a smallest timestamp being 35.3, and received 9 messages. Among these 9 received messages, 3 were sent at time quantum 9, 2 at time quantum 10, 4 at time quantum 8. Besides, the local virtual time at the end of the time quantum 10 is 40.8 (the local virtual time must take into account the messages received in this time quantum). The processor will then report 10, 40.8, 35.3, 6, and  $\{3, 2, 4\}$  to the GVT master.
- The GVT master maintains three arrays. The first one is  $\{TM_i\}$ , for the number of transient messages from each time quantum. The second one,  $\{MVT_i\}$ , records the smallest timestamp among all messages sent from each quantum. The last one,  $\{LVT_j\}$ , stores the local simulated time of each processor.
- When receiving a report message for processor  $j$ , the GVT master first obtains the time quantum index  $i$ , and then updates its three arrays according to the following rules:
  - $TM_i += S_{i,j}$
  - For each  $k$  in  $\{R_{k,j}\}$ ,  $TM_k -= R_{k,j}$
  - If  $(OT_{i,j} < MVT_i)$   $MVT_i = OT_{i,j}$
  - $LVT_j = T_{i,j}$
- The new GVT value is the minimum among  $\{LVT_j\}$ , and  $\{MVT_i\}$  for such  $i$  that  $TM_i$  is a non-zero, or formally:
  - $GVT = \min(\min_i(MVT_i | TM_i > 0), \min_j(LVT_j))$

In the above version of the algorithm, the report message contains an array of integers each of which denotes the number of received messages marked with the corresponding time quantum. In the DSIM implementation, a maximum length is imposed on this array, by limiting the reporting of received messages to the oldest time quantum active at this processor (the smallest  $k$  such that  $R_{k,j} > 0$ ), and possibly several others that follow immediately. This change preserves the correctness of the algorithm; the only effect may be the more conservative GVT value computed.

This happens when at most a limited number of time quanta corresponding to the maximum length of the received time quantum array allowed in the report message can be removed from consideration at the end of a time quantum. In practice, an array of size 2 to 4 gives the optimal performance. Hence, only 2 or 4 such numbers need to be sent to the GVT master, no matter how many processors are being used.

The correctness of the algorithm can be proved by showing that any event message  $m_l$  received after a GVT value  $gvt$  is received must have a timestamp  $ts(m_l) \geq gvt$ .

**Proof.** We prove this property by contradiction using induction.

Let's assume that  $ts(m_l) < gvt$  and  $i_l$  is the processor that sent this message at time quantum  $s_l$ . Let  $j_l$  be the time quantum from which the last report message sent by processor  $i_l$  that was received by the GVT master before  $gvt$  was obtained. It cannot be that  $ts(m_l) \geq T_{i_l, j_l}$ , since  $T_{i_l, j_l}$  is taken into account in computing  $gvt$ , so  $T_{i_l, j_l} \geq gvt$ . Neither it could be that  $s_l \leq j_l$ , since then  $ts(m_l)$  would be reflected in the  $MVT$  value corresponding to  $s_l$ , contradicting our assumption that  $ts(m_l) < gvt$ . Hence,  $s_l > j_l$  and  $ts(m_l) < T_{i_l, j_l}$ , so there must be another message  $m_2$ , sent by a different processor  $i_2$ , which caused a rollback on processor  $i_l$ , and satisfies  $gvt < ts(m_l) < ts(m_2)$ , as a rollback never affects the events with the same or earlier timestamps than the timestamp of the rollback message itself. From that it follows that this message also satisfies  $s_2 > j_2$ , where  $s_2, j_2$  are analogs of  $s_l, j_l$ .

By induction, let's assume that there is a message  $m_k$ , sent by processor  $i_k$ , such that  $gvt > ts(m_k)$  and  $s_k \leq j_k$ , where  $j_k$  denotes the latest time quantum on processor  $i_k$  that is included in computing  $gvt$  and  $s_k$  is the time quantum at which message  $m_k$  was sent. It cannot be that  $ts(m_k) \geq T_{i_k, j_k}$ , since  $T_{i_k, j_k}$  is taken into account in computing  $gvt$ , so  $T_{i_k, j_k} \geq gvt$ . Hence,  $ts(m_k) < T_{i_k, j_k}$ , so there must be another message  $m_{k+1}$ , sent by a different processor  $i_{k+1}$ , which caused a rollback on processor  $i_k$ , and satisfies  $gvt < ts(m_k) < ts(m_{k+1})$ . We define  $s_{k+1}, j_{k+1}$  as analogs of  $s_k, j_k$ . For message  $m_{k+1}$ , it cannot be that  $s_{k+1} \leq j_{k+1}$ , since then  $ts(m_{k+1})$  would be reflected in  $MVT$  value corresponding to  $s_{k+1}$ , contradicting our conclusion that  $ts(m_{k+1}) < gvt$ .

Hence, by induction we conclude that our assumption that  $ts(m_l) < gvt$  implies that there is an infinite sequence of messages with timestamps smaller than  $gvt$ , which contradicts the basic premise that each simulation can generate only a finite number of messages in the finite simulation time  $gvt$ .

### 3 IMPLEMENTATION DETAILS

In this section, several techniques related to memory usage are described to provide the readers with an understanding of the inner workings of DSIM. Although some similar techniques may have been proposed in other Time Warp

systems, a combination of them used in DSIM defines the exact design under which the experiments presented in this paper have been run, thereby allowing replication of experimental results presented in this paper.

#### 3.1 Local Fossil Collection

In Time Warp, a processed event becomes a fossil when its timestamp is smaller than the GVT, and the operation of releasing memory allocated for such events is called *fossil collection*. Events cannot be immediately released after having been processed, because of the possibility of rollbacks. They must be kept in memory until fossil collection is performed. Usually, there are two ways of maintaining these processed events. The first approach uses a single processed event list for all LPs (Logical Processes) on one processor. The main drawback is that, when rollbacks occur on some LPs, it becomes difficult to keep the entire list sorted. If the list becomes unsorted, then the entire list has to be scanned in order to know which events are subject to fossil collection. GTW used an on-the-fly technique (Fujimoto and Hybinette 1997) to partially solve this problem, by checking events with local minimum timestamps and ignoring other events. However, there may exist some processed events with timestamps earlier than GVT that cannot be reclaimed by this technique.

The other approach to maintaining processed events is to keep a separate processed event list on each LP. The problem of unsorted lists no longer exists because it is always the head (the latest one) of the processed event list that needs to be rolled back first. However, it introduces another problem. As explained in (Carothers, Bauer and Pearce 2000), it is costly to search through all those processed event lists during fossil collection, especially when the number of LPs is high. ROSS (Carothers, Bauer and Pearce 2000) solved this problem by grouping LPs into kernel processes, thereby helping to reduce the number of processed event lists. Another technique (Vee and Hsu 2002) is to sort these processed event lists further by their tails (the earliest ones) so that lists with a tail larger than the GVT can be completely skipped.

DSIM adopts the separate processed event list approach, in which fossil collection is not carried out when the GVT is updated. Instead, each LP checks if the GVT has been updated *before* it is about to process an event. If so, it will then compare the earliest processed event with the GVT. Only if the GVT has been recently updated and if the earliest processed event is earlier than the GVT will the LP invoke the fossil collection procedure. Otherwise fossil collection will be bypassed. Although this technique does not decrease the number of operations, it improves the locality of memory references, since the event memory released in the fossil collection procedure can be immediately reused in the processing of the new event (if there are new events to be scheduled).

Local Fossil Collection comes with its own drawbacks. First, if an LP suddenly becomes inactive after a highly active period and before the GVT is updated, it will have no new events to process and consequently will not be able to perform fossil collection. Second, delaying fossil collection until an LP is processing an event may increase memory usage, as more processed events will stay in the memory.

### 3.2 Event Management

Event allocation and deallocation are often the most frequently executed operations during discrete event simulation. Low-level system calls cannot be directly used, since such calls cannot be guaranteed to complete in a constant time. To minimize the cost of these operations, DSIM creates a customized event allocator for each event type. The rationale is that in C/C++, for objects of the same type, the memory footprint is always the same. The actual amount of memory used may vary from one object to another, but any extra memory must be explicitly obtained, and this is not a responsibility of the simulator.

Since event allocators only handle events of equal sizes, they can pre-allocate a number of memory buffers in a free buffer pool. To handle a request for a new event, the event allocator simply retrieves one buffer from the free buffer pool. If the pool is empty, it will acquire more buffers for it, using the low-level memory allocation function. When an event is to be released, its buffer is returned to the free buffer pool.

In DSIM, an unprocessed event becomes a processed one after it has been processed, and, conversely, a processed event becomes an unprocessed one after it has been rolled back. To save memory and to avoid unnecessary memory operations, however, both events are represented by the same data structure, with a flag denoting the status of the event. Therefore, within this single event data structure, there are two data blocks, one for the processed event and the other for the unprocessed one. DSIM is capable of overlying one over the other, to save memory further, as these two blocks are never needed at the same time during execution.

DSIM adopts the direct cancellation approach proposed by Fujimoto (Fujimoto 1989) for intra-processor events and extends it to inter-processor events. When a new inter-processor event is to be created and scheduled for an LP in another processor, the event is packed into a positive message. A stub event, containing only the timestamp of the new event, the message identifier, and the receiving LP identifier, is left behind to preserve the event dependency. The stub event, treated in the same way as other intra-processor events, will be inserted into the dependent list of the event that scheduled the new inter-processor event.

When receiving a positive message, an LP must unpack the message to restore the inter-processor event, and store the sender LP identifier and the message identifier into an input queue. To cancel an inter-processor event, an anti-message is created from the corresponding stub event. When the receiver LP sees the anti-message, it will check its input queue for presence of the inter-processor event to be canceled by comparing the message identifier and the sender LP identifier. In this way, the output queue can be completely avoided, while the input queue is still retained, but only for inter-processor events.

## 4 PROGRAMMING INTERFACE

DSIM comes with a simple programming interface that hides many implementation details while providing much freedom for programmers in making design decisions. To build a simulation, a DSIM programmer must first define the types of events, and then implement the LPs that comprise the simulation, and finally complete several auxiliary functions according to certain rules.

### 4.1 Event Declaration

The following syntax is used to create a new event type named *new\_event\_t*:

```
typedef tw_event_t < positive_data_t,
anti_data_t, type_id> new_event_t;
```

In the above statement, *positive\_data\_t*, the positive data type, is the type of the data stored in the unprocessed event, while *anti\_data\_t*, the anti-data type, is the type of the data stored in the processed event. The third parameter *type\_id* is the identifier of the event type. It must be different for each event type so that the LP can determine what event has happened by checking this field of the event.

### 4.2 State Saving versus Reverse Computation

The positive data type is decided by the simulation model semantics. The determination of the anti-data type is more difficult. Normally, it is determined by what data are needed to make each event reversible. DSIM supports two styles of undoing events, one is traditional state saving (infrequent state saving is not supported) and the other is based on *reverse computation* (Carothers, Perumalla and Fujimoto 1999). If the former is taken, the anti-data must basically store any change made by the processing of an event. For reverse computation, generally, much less information needs to be stored in the anti-data. For example, if an ordinary random number generator is used during the event processing, then the anti-data must contain the random seed before the event arrives, which can be written back to the random number generator if the event is to be

reversed. This is traditional state saving. In contrast, a reversible random number generator (Carothers, Perumalla and Fujimoto 1999) can be used, which can go back to the original state after the reverse function is called.

### 4.3 Implementing LPs

An LP class must be derived from the base class *tw\_lp*. There are five functions of *tw\_lp* that can be overridden: *Start*, *Stop*, *Process*, *Undo*, and *Commit*.

The *Start* and *Stop* functions are called when the simulation is started and stopped, respectively. The *Process* function processes an unprocessed event, while the *Undo* function revokes a processed event and brings the LP to the previous state.

The *Commit* function is called when a given processed event memory is to be reclaimed (that is when the event time is smaller or equal to the current GVT). This gives the LP an opportunity to perform irreversible operations, such as printing to the standard output or some others I/O actions. If the event contains some pointer to allocated memory, it is also in this function that the memory can be released.

The *Commit* function is useful for debugging as well. In a PDES program, most programming errors would cause the total number of processed events to be different when the number of processors varies. Locating the first event that is incorrect is difficult, since a processed event may be rolled back later. A customized *Commit* function would allow only those committed events (processed events earlier than the GVT) to be printed out. By comparing the outputs of a program running on different numbers of processors, incorrect events can be quickly identified.

### 4.4 Auxiliary Functions

Programmers must also implement several auxiliary functions in order for the simulation to work properly. These functions include *id\_to\_proc*, which returns the identifier of the processor where a given LP resides, and *id\_to\_lp*, which converts the identifier of the given LP to the pointer to this LP. In DSIM, each LP has a unique integral identifier that is used to address them. For instance, to schedule an event for an LP, the identifier of that LP must be provided. DSIM will then call the *id\_to\_proc* function to determine whether the destination LP is in the same processor. If so, it will call the *id\_to\_lp* function to obtain the address and forward the event. Otherwise, DSIM will simply pack the event into a message and then send it to the destination processor.

Another function that must be implemented by the programmer is *tw\_main*. In this function, a predefined DSIM simulation engine must first be initialized. This simulation engine is then used to create a set of LPs and to initialize them one by one. After setting the parameters of the simulation engine and calling its *Setup* function, the

*Run* function, which contains the main scheduler loop, must be called to run the simulation. The simulation terminates after the return from the *Run* function.

### 4.5 Shared vs. Distributed Memory Support

Although DSIM was designed for distributed memory clusters, it supports shared-memory machines as well. The reason for dual mode support is that shared-memory machines, while limited in the number of processors available, may be a more convenient platform to program, debug, and test PDES programs. The only difference between these two modes is in the creation of LPs. In the distributed mode, the program is only responsible for creating LPs that are assigned to the current processor, while in the shared-memory mode the program must create all LPs no matter which processor they are assigned to. The GVT algorithm proposed by Xiao et al. (Xiao et al. 1995) is adopted for the shared-memory version. As integral identifiers rather than addresses, are used to address LPs, programmers can completely ignore the differences between these two modes except when they have to create LPs, which can be accomplished by merely a few lines of code.

### 4.6 Artificial Rollback

Another feature of DSIM is that it provides a special debug mode to facilitate the debugging of PDES programs. If the program is compiled with a macro *DSIM\_AR*, artificial rollbacks will be introduced even though the simulation is executed on a single processor. This is achieved by using an error-prone priority queue that does not always find the earliest unprocessed event. Instead, with an adjustable probability, this queue may retrieve a randomly chosen, arbitrary event in the queue. Therefore, LPs will receive out-of-timestamp order events even when there is only one processor being used. The advantage of this technique is that errors are now reproducible, so the sequence of rollbacks remains the same every time the program is executed. In contrast, a PDES program with a regular priority queue running on multiple processors tends to produce different orders of execution, and therefore errors may appear or disappear during different runs, making the program extremely hard to debug.

## 5 EXPERIMENTAL RESULTS

Experiments have been performed on two different clusters, in order to evaluate the performance and scalability of DSIM with respect to the PHOLD model (Fujimoto 1990). The first set of experiments was carried out on a cluster of 40 nodes, where each node is an IBM Netfinity server with two 700-MHz Intel Pentium III processors. Half of these nodes are connected by fast Ethernet and the other by Gigabit Ethernet.

In the PHOLD model (Fujimoto 1990), each event stays at an LP for an exponential time and then departs to one of four nearest neighbors randomly chosen. In all experiments presented here the number of initial events per LP is always 16. LPs are organized into a two dimensional  $X$  by  $Y$  grid, where  $X$  is the number of columns and  $Y$  is the number of rows.  $N$  denotes the number of processors used.

Strip partitioning by columns is selected so that each processor has a sub grid of  $X/N$  by  $Y$  to simulate. One metric, the ratio of remote events (or inter-processor events), is of particular importance in asserting the performance of any Time Warp simulations. For the PHOLD model that is partitioned by columns, the ratios of remote events can be easily deduced from  $X$  and  $N$  as follows. Only two columns in each processor, the leftmost and the rightmost, may generate events that must be sent to a different processor. However, on average only 1 out of 4 such events produced by LPs on these two columns will actually depart. Therefore, the fraction of remote events is  $N/(2X)$ .

It must be noted that strip partitioning does not result in fewest remote events. It is the block partitioning which divides the entire grid into two-dimensional tiles, that leads to the lowest percentage of remote events when  $X=Y$ . However, one goal of the performance studies here is to demonstrate the scalability of DSIM with different ratios of remote events, so block partitioning was not implemented.

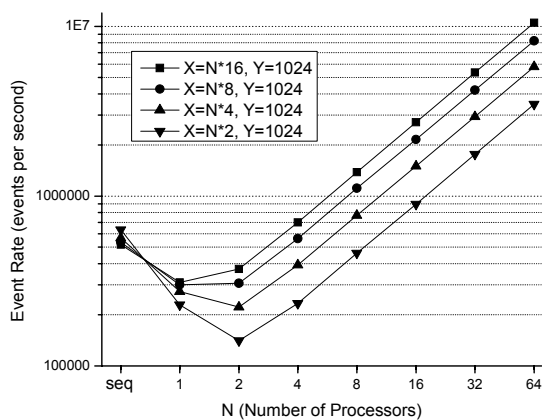


Figure 1. Performance of DSIM on a Netfinity Cluster with Fixed Percentages of Remote Events

It is also worth to note that although the PHOLD model may seem to be a toy example, it is indeed difficult to be parallelized, because there is no lookahead and the event granularity is extremely low. Some real-world models, such as PCS networks (Carothers, Fujimoto and Yi-Bing 1995), exhibit similar event behaviors as PHOLD, only with higher ratios of local events and coarser event granularity. Successful parallelization of the PHOLD model may provide a lower bound on the DSIM performance on this class of real-world applications.

The first set of experiments was designed to demonstrate how DSIM scales with different percentages of remote events (Figure 1). The number of processors marked on the horizontal axis does not count the GVT master. Here, the problem size was increased linearly as more processors were used. The percentages of remote events ranged from 3.125%, to 6.25%, 12.5%, and 25%. The actual percentages of remote events were slightly higher, since there were anti-messages. It is evident from Figure 1 that the performance of DSIM drops as remote events increase, due to the extra time needed for sending and receiving messages.

In Figure 1, DSIM is also compared with a sequential discrete event simulator that uses the same simulation engine but with facilities to handle rollbacks and to exchange message being removed. The overhead of parallelization becomes manifest when comparing this simulator with the sequential execution of DSIM, which is at least twice slower. More interestingly, with decreasing numbers of columns, the performance of the sequential simulator improves while that of DSIM on one processor drops. The former happens because of shrinking memory footprints with fewer LPs. The latter results from the use of a particular parameter called *event batch* which controls the number of consecutive events that can be processed in a single batch. With more remote events, this parameter must be decreased, and therefore the performance declines.

The performance gap between 1 processor and 2 processors hints the overhead of remote events. The event-processing rate may grow or drop when 2 processors are used as opposed to 1, depending on the percentage of remote events. However, in all cases starting from 2 processors on, DSIM maintains almost linear increases in event rates, implying an excellent scalability of the simulator.

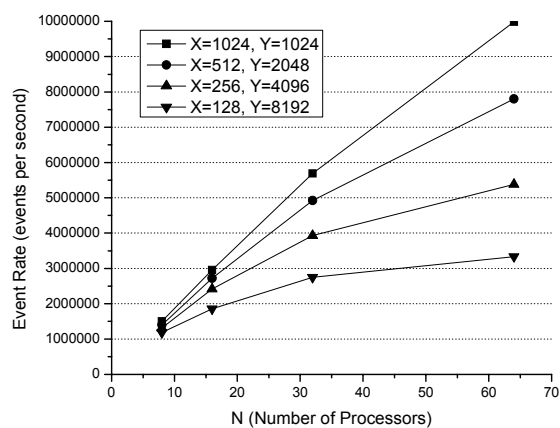


Figure 2. Performance of DSIM on a Netfinity Cluster with Fixed Problem Sizes

In the second set of experiments running on the same cluster, the number of LPs simulated was fixed but the

numbers of processors varied (Figure 2). Four configurations, namely, 1024 by 1024, 512 by 2048, 256 by 4096, and 128 by 8192, give different ratios of remote events. As these results indicate, DSIM is capable of simulating about 10 millions events per second using 64 processors, with 3.25% remote events (the 1024 by 1024 configuration). When the ratio of remote events increases to 25% (the 128 by 8192 configuration), DSIM is still able to process more than 3 million events per second. For this configuration, the improvement of using 64 processors over 32 processors is minor, but with 32 processors the ratio of remote events drops from 25% to 12.5% as in this set of experiments the problem size is fixed.

What is the speedup on 64 processors? Unfortunately, for the presented performance studies, the term *speedup* loses its precise meaning. All four configurations, each of which contained 1 million LPs and 16 million events shown in Figure 2 could not be executed on less than 8 processors due to the memory requirement (even if the efficient sequential simulation engine were to be used), so it was impossible to measure how fast a sequential simulation of the same problem size would be. If, on the other hand, the data presented in Figure 1 are to be used to derive speedups, the results would be biased against parallel execution, since models of different sizes are simulated, and larger sizes mean slower sequential speed, as evident in Figure 1. Even so, the speedup of parallel execution on 64 processors of a 1024 by 1024 grid versus efficient sequential execution of a much smaller size (16 by 1024) is still 20.3.

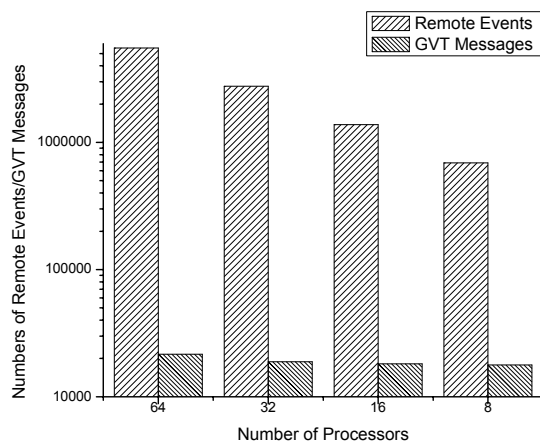


Figure 3. Numbers of Remote Events/GVT Messages in the 1024 by 1024 Configuration

Figure 3 shows the numbers of remote events and GVT messages in the 1024 by 1024 configuration. The numbers of GVT messages stay roughly the same for different numbers of processors as they are proportional to the product of the execution time and the number of processors, which in turn is proportional to the workload. This is because the frequency of GVT updates is controlled by the

width of time quanta. In all experiments presented in this paper, this width is 0.1-0.2 seconds, meaning that the GVT value can be updated 5-10 times per second. Notice the logarithmic scale used in the vertical axis – the ratios of GVT messages to remote events are 0.39%, 0.68%, 1.3%, and 2.6% respectively. These numbers illustrate the extremely low overhead of the TQ-GVT algorithm even with a high GVT update frequency.

Another set of experiments was run on an Alphaserer cluster consisting of 750 nodes, connected by a Quadrics interconnection network, located at Pittsburgh Supercomputing Center. This cluster was ranked 34<sup>th</sup> on the top 500 supercomputer list as of November 2004. Each node is a SMP with 4 1-GHz processors and 4 Gbytes of memory. Figure 4 depicts the event processing rates of DSIM on up to 1024 simulating processors. The numbers alongside each data point denote the numbers of extra processors that were used for GVT masters. It has been empirically determined that a GVT master in the TQ-GVT algorithm can drive as many as 128 processors without degrading the performance noticeably (whether more processors can be supported is yet to be established). Therefore, 2, 4, and 8 intermediate GVT masters were introduced for 256, 512, and 1024 processors respectively.

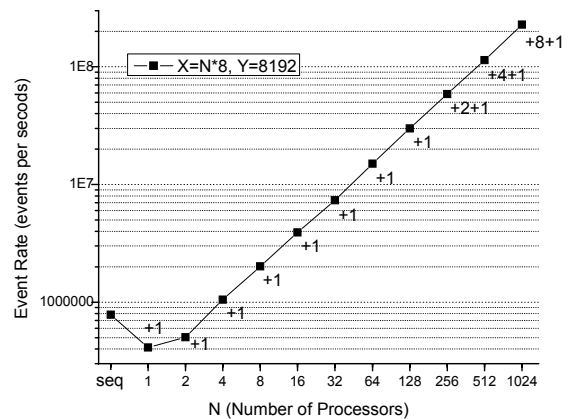


Figure 4. Performance of DSIM on the PSC Cluster with Fixed Percentages of Remote Events

The number of columns increases as more simulating processors are added, to maintain a constant percentage of remote events of 6.25%. A performance curve similar to those in Figure 1 is observed in Figure 4. In the 1024 processor case, which actually used 1033 processors, 67,108,862 LPs were simulated, yielding an event-processing rate of 228 million events per second, and a speedup of 296 (again, this speedup is somehow unfair to parallel execution since it is impossible for any single CPU to execute a simulation this large). Since each LP was assigned 16 events initially, at any moment during the simulation there were total of 1,073,741,824 unprocessed events.



## 6 CONCLUSION AND FUTURE WORK

As clusters become more prevailing, more and more PDES applications will be executed on this type of parallel computers. DSIM was developed under this premise. The major difficulty of porting PDES applications to clusters is designing an efficient and scalable GVT algorithm. The Time Quantum GVT algorithm adopted by DSIM meets these requirements. Furthermore, various improvements, such as Local Fossil Collection and efficient event management, enable DSIM to run with an unprecedented speed of 228 million events per second. As the same time, the programmability of DSIM has been an equally important design consideration to ensure that a programmer can quickly get familiar with the simulator.

DSIM is freely available with complete source code at <http://www.cs.rpi.edu/~cheng3/dsim>. More simulations will be implemented to verify its performance for various applications. The Time Quantum GVT algorithm will continue to be improved in order to enable Time Warp simulations on tens of thousands or even millions of processors. Moreover, as DSIM is an open source project, it is hoped that it will provide a standard simulation platform for researchers to implement and test various PDES algorithms.

## ACKNOWLEDGMENTS

The work presented in this paper has been done while the first author was a postdoctoral researcher associate at the Center for Pervasive Computing and Networking, RPI, Troy, NY. The largest simulations were carried out on the National Science Foundation Terascale Computing System at the Pittsburgh Supercomputing Center under the grant of time PSC IRI05001P. The authors wish also to express their gratitude to Professors Christopher Carothers, Carl Tropper and Adelinde Uhrmacher for their valuable comments on an early version of this paper.

## REFERENCES

- Baezner, D., G. Lomow and B. Unger. 1994. Parallel simulation environment based on time warp. *International Journal in Computer Simulation*, 4(2): 183.
- Baldwin, R., M. J. Chung and Y. Chung. 1991. Overlapping window algorithm for computing GVT in Time Warp. *11th International Conference on Distributed Computing Systems* (Cat. No.91CH2996-7), 20-24 May 1991, 534-41, Arlington, TX: IEEE Comput. Soc. Press.
- Bauer, H. and C. Sporrer. 1992. Distributed logic simulation and an approach to asynchronous GVT-calculation. *Proceedings of the 1992 SCS Western Simulation MultiConference on Parallel and Distributed Simulation*, 205-208, Newport Beach, CA: SCS.
- Bellenot, S. 1990. Global Virtual Time Algorithms. *Proceedings of the SCS Multiconference on Distributed Simulation*, 122-127, San Diego, CA: Soc. for Computer Simulation Int.
- Carothers, C. D., D. Bauer and S. Pearce. 2000. ROSS: a high-performance, low memory, modular Time Warp system, 53, Los Alamitos, CA: IEEE.
- Carothers, C. D., R. M. Fujimoto and L. Yi-Bing. 1995. *A case study in simulating PCS networks using time warp*, 87, Lake Placid, NY: IEEE Comput. Soc. Press.
- Carothers, C. D., K. S. Perumalla and R. M. Fujimoto. 1999. Efficient optimistic parallel simulations using reverse computation. *ACM Transactions on Modeling and Computer Simulation* 9(3): 224.
- Choe, M. and C. Tropper. 1998. *An Efficient GVT Computation Using Snapshots*. CSMA 98, 33-43.
- D'Souza, L. M., X. Fan and P. A. Wilsey. 1994. pGVT: an algorithm for accurate GVT estimation. *Proceedings of 8th Workshop on Parallel and Distributed Simulation*, 102-109, Edinburgh, UK: SCS.
- Das, S., R. Fujimoto, K. Panesar, D. Allison and M. Hybinette. 1994. GTW: a time warp system for shared memory multiprocessors. *Proceedings of Winter Simulation Conference*, 1332-1339, M. S. Manivannan and J. D. Tew, Piscataway, NY: Institute of Electrical and Electronics Engineers.
- Das, S. K. and F. Sarkar. 1995. A hypercube algorithm for GVT computation and its application in optimistic parallel simulation. *Proceedings of Simulation Symposium*, 51-60, Phoenix, AZ: IEEE Comput. Soc. Press.
- Fujimoto, R. M. 1989. Time warp on a shared memory multiprocessor. *Transactions of the Society for Computer Simulation* 6(3): 211-239.
- Fujimoto, R. M. 1990. Performance of time warp under synthetic workloads. *Distributed Simulation. Proceedings of the SCS Multiconference*, 23-28, San Diego, CA: SCS.
- Fujimoto, R. M. and M. Hybinette. 1997. Computing global virtual time in shared-memory multiprocessors. *ACM Transactions on Modeling and Computer Simulation*, 7(4): 425-46.
- Jefferson, D., B. Beckman, F. Wieland, L. Blume and M. Diloreto. Time warp operating system. *Proceedings of the eleventh ACM Symposium on Operating systems principles*, 77-93, ACM Press.
- Jefferson, D. R. 1985. Virtual time. *ACM Transactions on Programming Languages and Systems*, 7(3): 404-25.
- Kim, H. K. and J. Jean. 1996. *Concurrency preserving partitioning (CPP) for parallel logic simulation*, 98, Los Alamitos, CA: IEEE.
- Martin, D. E., P. A. Wilsey, R. J. Hoekstra, E. R. Keiter, S. A. Hutchinson, T. V. Russo and L. J. Waters. 2003. Redesigning the WARPED simulation kernel for analysis and application development. *Proceedings*

- 36th Annual Simulation Symposium (ANSS-36 2003), 216-23, Orlando, FL: IEEE Comput. Soc.
- Mascarenhas, E., F. Knop and V. Rego. 1995. *ParaSol: a multithreaded system for parallel simulation based on mobile threads*, 690, Piscataway, NJ: IEEE.
- Mattern, F. 1993. Efficient algorithms for distributed snapshots and global virtual time approximation. *Journal of Parallel and Distributed Computing* 18(4): 423-34.
- Perumalla, K. and R. Fujimoto. 2001. Virtual time synchronization over unreliable network transport. *Proceedings 15th Workshop on Parallel and Distributed Simulation*, 129, Lake Arrowhead, CA: IEEE Comput. Soc.
- Preiss, B. R. 1989. The Yaddes Distributed Discrete Event Simulation Specification Language and Execution Environment. *Proceedings of the SCS Multiconference on Distributed Simulation*, 139-144.
- Samadi, B. 1985. *Distributed Simulation, Algorithms and Performance Analysis*. Computer Science Department, University of California, Los Angeles.
- Srinivasan, S. and P. F. Reynolds, Jr. 1993. Non-interfering GVT computation via asynchronous global reductions. *Proceedings of 1993 Winter Simulation Conference*, G. W. Evans, M. Mollaghasemi, E. C. Russell and W. E. Biles, eds. 740-749, Piscataway, NJ: Institute of Electrical and Electronics Engineers.
- Steinman, J. S., C. A. Lee, L. F. Wilson and D. M. Nicol. 1995. Global virtual time and distributed synchronization. *Proceedings 9th Workshop on Parallel and Distributed Simulation (ACM/IEEE)*, 14-16 June 1995, 139-48, Lake Placid, NY, USA, IEEE Comput. Soc. Press.
- Tomlinson, A. I. and V. K. Garg. 1993. An algorithm for minimally latent global virtual time. *1993 Workshop on Parallel and Distributed Simulation*, 16-19 May 1993, 35-42, San Diego, CA, USA, SCS.
- Vee, V.-Y. and W.-J. Hsu. 2002. Pal: a new fossil collector for time warp. *Proceedings 16th Workshop on Parallel and Distributed Simulation*, 35-42, Washington, DC, USA, IEEE Comput. Soc.
- Wieland, F., L. Hawley, A. Feinberg, M. Di Loreto, L. Blume, P. Reiher, B. Beckman, P. Hontalas, S. Belenot and D. Jefferson. 1989. *Distributed combat simulation and time warp. The model and its performance*, 14, Tampa, FL, USA, Publ by Soc for Computer Simulation Int, San Diego, CA, USA.
- Xiao, Z., F. Gomes, B. Unger and J. Cleary. 1995. A fast asynchronous GVT algorithm for shared memory multiprocessor architectures. *Proceedings 9th Workshop on Parallel and Distributed Simulation (ACM/IEEE)*, 203-208, Lake Placid, NY: IEEE Comput. Soc. Press.

## AUTHOR BIOGRAPHIES

**GILBERT CHEN** is a simulation design engineer in The MathWorks, Inc. His research interests include parallel discrete event simulation, simulation architecture, and wireless sensor networks. Prior to joining The MathWorks, he had been a postdoctoral research associate in the Center for Pervasive Computing and Networking at Rensselaer Polytechnic Institute, where he obtained his PhD in Computer Science.

**BOLESŁAW K. SZYMANSKI** is the Director of the Center for Pervasive Computing and Networking and a Professor of Computer Science at Rensselaer Polytechnic Institute. He received his Ph.D. in Computer Science from the National Academy of Sciences in Warsaw, Poland, in 1976. Prior to joining RPI in 1985, he was a faculty member at the Department of Computer and Information Sciences at University of Pennsylvania. He is an author and co-author of more than two hundred fifty scientific publications and an editor of four books. Dr. Szymanski is also an Editor-in-Chief of *Scientific Programming* and an Area Editor of *Simulations*. He is also an IEEE fellow and a member of the IEEE Computer Society, and the Association for Computing Machinery. Dr. Szymanski's interests include distributed and parallel computing and system modeling and simulation. His recent work includes sensor networks, on-line network simulation, and network security.