

CONTROLLING OVER-OPTIMISM IN TIME-WARP VIA CPU-BASED FLOW CONTROL

Vinay Sachdev, Maria Hybinette, Eileen Kraemer

Computer Science Department
415 Boyd Graduate Studies Research Center
The University of Georgia
Athens, GA 30602, U.S.A.

ABSTRACT

In standard optimistic parallel event simulation, no restriction exists on the maximum lag in simulation time between the fastest and slowest logical processes (LPs). Over-optimistic applications exhibit a large lag, which encourages rollback and may degrade performance. We investigate an approach for controlling over-optimism that classifies LPs as FAST, MEDIUM, or SLOW and migrates FAST and/or SLOW processes. FAST LPs are aggregated, forcing them to compete for CPU cycles. SLOW LPs are dispersed, to limit their competition for CPU cycles. The approach was implemented on distributed Georgia Tech Time Warp (GTW) (Das et al., 1994) and experiments performed using the synthetic application *P-Hold* (Fujimoto, 1990). For over-optimistic test cases, our approach was found to perform 1.25 to 2.75 times better than the standard approach in terms of useful work and to exhibit execution times shorter than or equal to the standard computation.

1 INTRODUCTION

Parallel simulation based on the optimistic Time Warp protocol (Jefferson and Sowizral, 1985) is widely used in large-scale simulations such as those of air-traffic control and the World Wide Web. A primary concern of such simulations is good performance. However, these simulations are subject to the problem of *over-optimism*, in which some logical processes (LPs) progress far beyond others. To facilitate discussion of the nature of inefficiencies that result from over-optimism, we first present some essential background on optimistic discrete event simulation using Time Warp.

A discrete event simulation consists of a collection of logical processes (LPs), which may execute on different

processors (PEs). The simulation is driven by the exchange of timestamped message by the LPs. Consistency in the processing of messages requires that all events be processed in timestamp order. Two main synchronization protocols exist: *conservative* and *optimistic*. The conservative protocol enforces consistency by avoiding the possibility of ever receiving an event in the past. In contrast, the optimistic protocol permits the receipt of an event from the past but responds by “rolling back” events that were optimistically processed too early.

Each LP in an optimistic simulation maintains a current *logical clock* (*local virtual time, or LVT*); whenever an LP receives a message with a timestamp earlier than LVT, it rolls back its execution to the time-stamp before that of the arrived message. Such an out-of-order message is called a *straggler message*. When rollback is necessary an LP reverts to the appropriate previous state and “un-schedules” any messages sent prior to the rollback. To support this, LPs maintain a history of state information and keep a record (an *anti-message*) for each message sent. In the case of rollback, the LP sends anti-messages, which annihilate the original messages sent. Anti-messages may cause additional rollbacks, called secondary rollbacks.

The earliest timestamp of any unprocessed or partially processed message in the system defines *Global Virtual Time (GVT)*. A message is guaranteed not to rollback if its timestamp is earlier than GVT. Thus, memory for events with timestamps *before* GVT, along with their corresponding state and anti-messages, can be reclaimed in a process known as *fossil collection*.

Over-optimism results in poor memory utilization because it creates a wide gap between GVT and the most recent timestamp in the system. State information, event histories, and anti-messages must be stored for all computation with a simulation time later than GVT. If over-

optimism is not controlled, the memory requirement may grow to the point that memory becomes exhausted.

Another inefficiency caused by over-optimistic behavior is that rollbacks may be quite long when slower LPs send a message to faster (over-optimistic) LPs. Previously processed events are rolled back and anti-messages are sent to cancel the events scheduled as a result of the events currently being rolled back. This nullifies all the useful computation done by both the over-optimistic LPs and the LPs to which the anti-messages were sent. Further, useful CPU cycles are wasted by undoing the work already done.

Finally, long and frequent rollbacks result in a large number of anti-messages being sent. These anti-messages utilize bandwidth that could otherwise have been used for useful communication. Thus, over-optimistic behavior of LPs can cause inefficiencies including **poor memory utilization**, **excessive rollbacks** and **communication overheads**.

Simulations that are susceptible to over-optimism include those operating in **heterogeneous environments**, those subject to **external workloads**, or those for which **application-specific characteristics** of the simulation promote over-optimistic behavior.

Networks of workstations (NOWs), an important platform for large scale simulations, are typically heterogeneous, with computing and memory resources varying among machines. Logical processes running on fast processors may progress faster in simulated time than logical processes running on slower processors.

A large scale simulation may run on a system that is shared among many users. Here, logical processes may compete with other applications for shared resources, causing some logical processes to run on more heavily loaded processors, while others run on less loaded processors. Logical processes running on heavily loaded processors make less progress in simulation time compared to logical processes on less loaded processors.

The manner in which a particular application is implemented can also influence over-optimistic behavior. Applications that exhibit *self-instantiation* and *uneven granularity of load* per LP may demonstrate over-optimistic behavior. *Self-instantiation* means that an LP schedules events to itself rather than to a remote LP. *Degree of self-instantiation* refers to the number of messages an LP sends to itself before sending a remote message. Applications that consist mainly of LPs with a high degree of self-instantiation communicate with other LPs infrequently. Because of this infrequent communication, whenever an LP that is far behind sends a message to an LP that is far ahead, it causes long rollbacks due to out-of-order messages. For example, the implementation of a Personal Communication Systems (PCS), described in (Carothers and Fujimoto, 1994), that includes LPs having a high degree of self-instantiation, has been shown to exhibit over-optimistic behavior.

Another characteristic that may cause over-optimistic behavior is uneven granularity of load per LP in an application. This happens when some LPs incorporate more work and take more time to process their event set than the other LPs. Again, when the LPs that are far behind in the simulation communicate with LPs that are far ahead, long rollbacks and anti-messages result. Implementations of Asynchronous Transfer Networks (ATMs) are prone to such an uneven granularity of load per LP (Hao, et al., 1996).

A goal of our approach is to use the CPU as a flow control mechanism for over-optimistic execution. We present a process migration scheme that controls over-optimistic processes by isolating their impact on other processes, while promoting the progress of slower, less optimistic processes. We evaluate our approach using a synthetic benchmark application called *P-Hold*.

2 BACKGROUND

Prior work addressing the problem of over-optimism falls into three broad categories: protocols using limited optimism, memory management protocols and adaptive techniques.

2.1 Optimism Limiting Protocols

Blocking is a commonly used technique for reducing the amount of rollback (Reiher et al., 1989) that limits the progress of over-optimistic LPs through use of a time window of size W . LPs are prevented from progressing beyond $GVT + W$, and are blocked until LPs that are far behind catch up. Window size may be determined statically dynamically.

The aggressive no-risk protocol (Dickens and Reynolds, 1990) avoids sending a message until it is guaranteed that it won't cause rollbacks. Messages sent by an LP are stored in its PE's buffer and not sent until GVT advances beyond the send timestamp of the message, assuring that the messages will not be rolled back later.

The *look ahead information* approach (Lubachevsky et al., 1989), also may be used to decide whether it is safe to process a given message. A hybrid conservative and optimistic protocol is employed that begins with the conservative protocol to determine which events are safe to process and later adds optimistic synchronization features to "unsafe" events.

Another approach introduces *additional rollbacks* at stochastically selected intervals (Madiseti et al., 1983). These additional rollbacks prevent overly optimistic execution of LPs that could be rolled back to GVT if a rollback decision were determined for that LP. Probability vectors are used to determine if the LP should be rolled back.

The *breathing time protocol* (Steinman, 1983) limits the number of events a particular LP can process beyond

GVT and involves determination of the minimum time stamp among events that will be produced in the future.

2.2 Memory Management Protocols

Two protocols used to limit the memory utilization in an over-optimistic simulation are *artificial rollback* (Jefferson, 1990) and *cancel back* (Lin and Preiss, 1991). These protocols are used when the system runs out of memory and fossil collection attempts cannot reclaim the memory needed for the simulation to progress. These schemes roll back some of the logical processes and utilize the freed memory to continue. *Artificial rollback* works by identifying the most over-optimistic LPs (those furthest ahead in simulation time) and then rolling them back. *Cancel back*, on the other hand, achieves the same effect by sending back certain messages to the sender LP, rolling back the sender.

2.3 Adaptive Techniques

S. Das and R. Fujimoto proposed an adaptive technique combining memory management and limited optimism synchronization protocols (Das and Fujimoto, 1997). The amount of memory allocated to a Time Warp simulation automatically limits the amount of optimistic execution, i.e., the degree to which processes may advance ahead of other processes. This protocol seeks to provide sufficient memory for Time Warp to execute efficiently, but not so much memory that overly optimistic execution can occur. The protocol attempts to simultaneously address rollback thrashing and memory management issues. The approach is adaptive; it monitors the execution of the Time Warp program and automatically adjusts the memory provided to the parallel simulator. An adaptive protocol was necessary because the the synchronization and memory management protocol parameters depend on characteristics of the application such as symmetry and homogeneity among the simulation processes and memory required to execute the program using Time Warp.

A load distribution system for background execution of Time Warp (Carothers and Fujimoto, 2000) is designed to use free cycles of a collection of heterogeneous machines to run a Time Warp simulation. The load management policy involves both processor allocation and load balancing. The processor allocation policy dynamically determines the set of processors to be used for a Time Warp simulation. LPs are grouped into clusters by the application to reduce migration overhead. Clusters, rather than individual LPs are migrated, with the goal of goal of equalizing the progress of all the processors, taking into consideration the external and internal workload, processor speeds, etc. The metrics for classifying the processors and individual clusters are PAT (Processor Advance Time) and CAT (Cluster Allocation Time), respectively.

A load balancing technique for the Time Warp distributed system for object-oriented simulation (Burdorf and Marti, 1993) distributes objects across nodes and provides optimistic concurrency control. The scheme consists of static and dynamic load balancing monitors. The static monitor determines pre-assignment of objects to processors. The dynamic load balancing module monitors load imbalance and initiates migration of objects, using knowledge of simulation time (LVT) to reduce rollback. That is, it minimized the distance between the simulation time of the farthest ahead object and the furthest behind object.

Another interesting approach that applies load balancing and optimism limiting protocol (Jones and Das, 1998) combines the throttling of over-optimistic processes and scheduling (or load balancing) to control over-optimistic behavior. Throttling is implemented by the *moving time window* protocol described above. In the scheduling component LPs are remapped to processors so the N slowest LPs in simulation time are mapped to different processors, where N is the number of processors on which the simulation is run.

2.4 Comparison

S. Das's adaptive memory management technique that uses memory as a flow control for controlling over-optimism is similar to our approach, which uses CPU as a flow control mechanism. Unlike Carothers and Fujimoto's approach, which uses available CPU cycles to load balance logical processes, we use the CPU itself to control over-optimism. By aggregating over-optimistic LPs to one CPU, we force these LPs to compete with each other for available CPU cycles. This slows down their progress while isolating their impact on other LPs. In addition, we spread out the less optimistic LPs, limiting their competition for CPU cycles.

S. Das and K. Jones approach of using throttling with scheduling is similar to our approach in terms of making the less optimistic LPs progress to catch up with LPs that are ahead in simulated future. But for LPs that are far ahead in simulation time, Das and Jones use blocking, which wastes CPU cycles, unlike our approach which isolates these LPs on one PE to slow their progress. Too much throttling is harmful as too few events are admitted for processing. Another important difference is that their approach was implemented on a simulated distributed system, whereas our implementation is deployed on a real distributed system on real processors.

3 OUR APPROACH

The idea behind our approach is to aggregate fast LPs onto one processor so that they must compete for processor cycles, slowing their progress. Slow logical processes are dispersed to different processors to limit their competition

for CPU cycles. Migration costs are minimized by reducing the number of LPs moved. This is done by selecting a *fast repository* (CPU on which the fast LPs will be aggregated) that already has the most fast LPs mapped to it. The cost of moves are justified in our approach in that we move only LPs that are either too fast and waste work, or too slow and likely require more resources.

Our approach to controlling over-optimism involves the following steps:

(i) Ranking of LPs: LPs are ranked according to how far ahead in simulation time they are compared to the rest of the simulation. We use a metric called *GVTlag*, the difference between an LP's current simulation time and GVT ($GVTlag = LVT - GVT$). LPs are ranked according to their *GVTlag* value.

(ii) Classification of LPs: LPs are then classified as one of *FAST*, *MEDIUM* and *SLOW*. The $k * N$ LPs that have the highest value of *GVTlag* are selected to be labeled as *FAST*, where k is a scaling factor and N is the number of processors on which the GTW kernel is currently running. The factor k is an experimental parameter and depends on the migration cost in the system. A selected number of the slowest LPs are classified as *SLOW*; the remaining LPs are classified as *MEDIUM*. In our experiments, the number of *SLOW* LPs is fixed at $N-1$.

(iii) Identification of the FAST-REPOSITORY: The *FAST-REPOSITORY* is the CPU on which the over-optimistic (*FAST*) LPs will be aggregated. The CPU containing the greatest number of LPs labeled as *FAST* is selected as the *FAST-REPOSITORY*. This reduces the number of LPs that will be moved in subsequent steps.

(iv) Isolation of FAST LPs: All *FAST* LPs not currently mapped to the *FAST-REPOSITORY* are migrated to it. This is done to isolate the effects of the over-optimistic LPs from the rest of the LPs in the simulation and to force these optimistic LPs to compete for CPU cycles.

(v) Spreading of SLOW LPs: *SLOW* LPs are redistributed to the processors other than the *FAST-REPOSITORY*, one per processor. The goal of this step is to limit the contention of the least optimistic LPs for CPU cycles, in the hopes of allowing them to "catch up" with the rest of the simulation.

3.1 Implementation

Our load balancing algorithm is implemented on the distributed Georgia Tech Time Warp system (GTW) (Das et al., 1994), which is a parallel and distributed discrete event simulation executive based on Jefferson's Time Warp (Jefferson, 1985). GTW runs on both shared memory machine and distributed memory machines. Details of GTW can be found in (Fujimoto et al., 1997).

Distributed GTW employs an additional thread on each machine (Carothers and Fujimoto, 2000) that handles all the external communication with other machines. The

Parallel Virtual Machine (PVM) communication library is used for remote communication. We implemented our process migration algorithm in a separate thread called *MonitorOPT* on top of distributed GTW.

The software architecture of distributed GTW, including our thread, is shown in Fig.1. GTW provides the APIs for the simulation application to exchange information with GTW regarding the number of LPs, number of processors, event handlers for each LP and other information. Once GTW has the application-specific information it sets up various data structures to carry out the execution of the simulation. Distributed GTW consist of two libraries: the *kernel library* and the *kernel communication library*. The kernel library consists of the core functionalities of GTW, including the state saving mechanism, scheduler of events, mechanism for computing GVT and communication thread. The kernel communication library consists of various methods that invoke PVM calls. The kernel library calls the method in the communication library in order to execute PVM functionality.

We implemented an optimism controlling module consisting of a central monitoring process called "*MonitorOPT*" that is heart of our algorithm. The process runs on a dedicated machine. *MonitorOPT* executes periodically to collect statistics from other processing elements (PEs). The specific period is an experimental parameter and could be varied. *MonitorOPT* computes the moves of the LPs based on the collected statistics of the whole system and uses this for controlling over-optimism of the LPs. The *MonitorOPT* process responsible for migration decisions communicates with the kernel communication library in order to exchange messages with the communication thread of the kernel library.

Once the new mapping information of where each LP will be moved is computed a move-list is created having information about what LPs will be needed to move, the source PE from which it will be moved and to what destination PE will it be moved to. If *MonitorOPT* determines some moves to be made, a *HALT message* is sent to each processor. This is done to stop each PE's computation and roll back all LPs to GVT. This helps to synchronize all PEs and to perform load redistribution effectively, reducing the cost of LP migration as the history information of the LPs will not be moved from source processor to the destination.

After all the processors roll back their computations to GVT, they send an acknowledgment message, *HALT-ACK*, to the *MonitorOPT* process. The *MonitorOPT* process then sends a *MOVE message* for each move in the move-list to the source processor from which a particular LP will be moved. The source PE transfers all the information regarding the LP to the destination PE, and then sends a *MOVE-LP message* to the destination PE. The destination PE updates all the necessary information and then sends a *MOVE-LP-ACK message* to the source PE, after which the source PE sends a *MOVE-ACK* to the *MonitorOPT*.

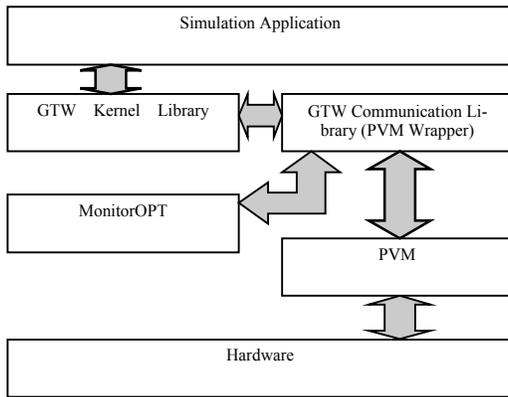


Fig 1. A Software architecture for Distributed GTW

Once the MOVE messages have been sent for each move in the move-list, the MonitorOPT process computes the new mapping information which is the updated mapping of all LPs to PEs to update the mapping structures used for migration of LPs in the future. It then sends a *MAPPING message* to all the processors to reflect this new mapping. Once all PEs update their local mapping structures for the mapping message they send a *MAPPING-ACK message* to MonitorOPT. Next, a *RESUME message* is sent so that the normal computation can resume on all processors. The load balancing interval starts with this message and at the interval expiry new computations and data collection can begin. This process continues until the end of the simulation.

4 EXPERIMENTS AND DISCUSSION

Experiments were performed to compare the performance of our approach with that of the standard Time Warp simulation without process migration. The experiments presented in this study were performed in a heterogeneous environment consisting of 8 machines: one SGI Origin 2000 with sixteen 195 MHz MIPS R 10,000 processors running IRIX version 6.5 and seven 167-MHz Sun Sparc Ultra -1 workstations running version 2.5 of the Solaris operating system. In all experiments the GTW kernel was executed on a total of 8 processors with one processor

from each of the machines involved in the experiment. MonitorOPT was executed on the SGI machine. The benchmark application used to perform these experiments is P-Hold.

4.1 The Benchmark application: P-Hold

P-Hold is a benchmark application using a synthetic workload model (Fujimoto, 1990). The benchmark uses a fixed message population. Processing of each message takes a finite amount of time, after which a new message is sent to another LP with a specified time stamp increment. The initial messages have a timestamp that is exponentially distributed between 0 and 1. As the messages are forwarded their new timestamp increments are fixed at 1. The total number of LPs for these experiments was fixed at 256. These LPs were initially evenly distributed on 8 processors, with 32 LPs per processor. The experiments used a fixed message population of 6,400.

In order to effect an imbalance, P-Hold was instrumented with two distinct synthetic workloads: null and one millisecond. In the null case, event processing is made as small as possible; in the one millisecond case event processing includes a 1 msec delay loop.

To make the application more imbalanced, P-Hold was also configured as self-instantiated, and the degree of self-instantiation of an event was varied. The experiments used two degrees of self-instantiation: 50 and 200. When an event is processed in which the source LP differs from the destination LP, the destination LP schedules the next d generations of the event to it. After d generations of the event have been produced, a new destination LP is randomly selected.

For the first set of experiments, two classes of LP/event restrictions are generated:

Class [200: 1 msec]: This includes events with degree of self-instantiation 200 and LPs that take a 1 millisecond delay (event granularity of 1 msec) to execute every event.

Class [50: null]: This includes events with degree of self-instantiation 50 and LPs that execute events without any delay (null event granularity).

The percentage of LPs in Class [200:1msec] was varied from 0 to 100 percent in increments of 20 percent. The remaining LPs in the simulation operated under Class [50 : null]. By varying the percentage of LPs executing in either of the two classes, we created different levels of optimism and imbalance. LPs under Class [200: 1msec] progress slowly as their event granularity is larger than event granularity of Class [50: null]. Additionally, whenever the LPs in Class [200:1msec] send remote messages to LPs in Class [50:null], it may roll back the computation of the Class [50:null] LPs.

A second set of experiments was performed for comparison. While keeping all other operating parameters the same, we swapped the delay of event execution on the two

classes of LPs. LPs under these experiments were made to execute under the two classes of operating restrictions,

Class [50: 1msec]: Events with degree of self-instantiation 50 and LPs execute event with 1 millisecond delay (1 msec event granularity).

Class [200: null]: Events with degree of self-instantiation 200 and LPs execute events without any delay (null event granularity).

An empirically determined sample period of 50 seconds was used. Results represent an average of four trials.

4.2 Experiments under Set 1

In figure 2 below, we show the performance of three variations of our migration algorithm: (1) with fast repository and slow LP spreading, (2) without fast repository but with slow LP spreading, (3) with fast repository but without slow LP spreading. We measure performance as the ratio of the useful work with process migration to the useful work without process migration. Useful work is the ratio of events that were committed (not rolled back) to all events processed. The three variations are described below:

Both Fast Repository and Slow LP spreading: The number of FAST LPs is $N * k$, where $k = 5$ and $N = 8$. The numbers of LPs classified as SLOW is fixed to 7 ($N-1$, where N is the number of processors in the system).

Slow LP spreading only: The number of FAST LPs is 0; accomplished by setting $k=0$. We classify 8 ($N=8$) LPs as SLOW and the remaining LPs as MEDIUM.

Fast Repository only: 40 LPs ($k=5$) are classified as FAST. The remaining LPs are classified as MEDIUM.

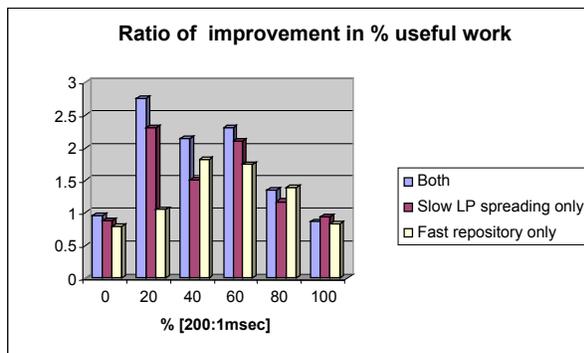


Fig. 2 Ratio of percentage improvement in useful work as compared to without process migration, (Set 1)

From figure 2 we see that in cases where over-optimism exists the simulation performs better with process migration (“Both”) than without process migration in terms of percentage of useful work, by a factor of 1.25 to 2.75 over useful work without process migration. The best

performance is at data point 20%, where “Both” (Fast repository and slow spreading) performs approximately 2.75 times better than the simulation without process migration. None of the variations of process migration perform better than the simulation without process migration at data points 0 and 100. At these points the simulation is balanced; all LPs execute under one class of event granularity, so it is likely that it is not over-optimistic. At data point 0, all the LPs execute under Class [50 : null], whereas at data point 100 all LPs operate under Class [200:1 msec]. The LPs all have the same event granularity and progress at the same pace, causing less rollback. The ratio of useful work of less than one in these cases in which over-optimism is absent is a result of the overhead of process migration.

When the application is unbalanced, process migration(Both) demonstrates maximum benefits performing by a factor of 1.25 to 2.75 times better than the results without process migration. When the application is overly optimistic, with no migration, the simulation spends much of its time rolling back its computation rather than progressing forward and hence results in less useful work. On the other hand, due to process migration, the over-optimistic LPs are isolated on a separate processor slowing their progress down and less optimistic LPs are allowed to make progress by redistributing them on different processors. This helps in reducing the differences among the progress of the LPs and the number of rollbacks.

The two variations of process migration also perform better than no process migration when the application is unbalanced (data points 20% – 80%). In some cases the fast repository (which limits fast LPs) produces the benefit, while in other cases the spreading of slow LPs (which promotes progress by the slow LPs) produces the benefit. It depends upon the application state which type of LPs dominates the simulation. To complement each of these two approaches, we combined the two variations in our implementation of the process migration algorithm. In this combined approach, we control both the less optimistic and over-optimistic LPs, consistently providing better performance in cases where over-optimism exists.

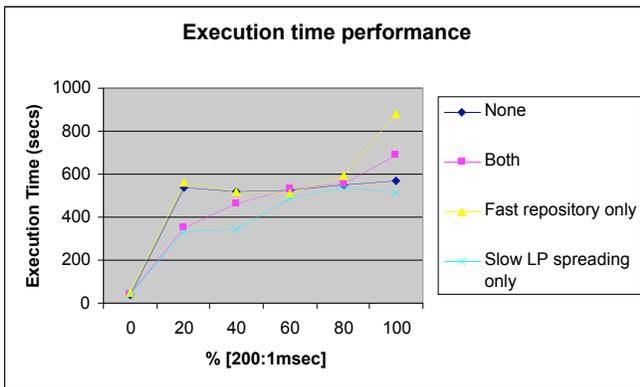


Fig. 3 Execution time performance (Set 1)

Figure 3 shows the execution time performance for the above experiment, comparing the simulation run with process migration and its various variations against the simulation without process migration. It was observed that process migration did not provide any benefit in terms of execution time when the optimism is balanced (data point 0 and 100 when all LPs have same event granularity). The most benefit in terms of execution time is obtained at 20% and 40% of Class [200:1msec]. At these two points there is also better percentage of useful work as shown in fig. 2. For data points 60 and 80 it is observed that process migration takes almost equal execution time compared to without process migration. The percentage of useful work performed by the process migration run for data point 60 is approximately 2 times better than without process migration and for data point 80 is 1.15 times better than without process migration. That is, although we do not benefit in terms of execution time at these two data points, we do benefit in terms of useful work. This is the result of migration overhead. If this overhead consists largely of communication costs, then it may be possible for external workloads to benefit from this reduced work, as the CPU cycles not used in performing wasted work would be available. However, we have yet to characterize the computation/communication trade-offs of process migration, and plan to do so in future work. For point 100, without proc-

ess migration (None) performs better for both the percentage of useful work in figure 2 and the execution time in figure 3. This is due to the application being balanced at this instant as all LPs belong to one class of self-instantiation and event granularity and thus progress at the same pace. The overheads of process migration exceed the benefits obtained when the application is the balanced and least optimistic.

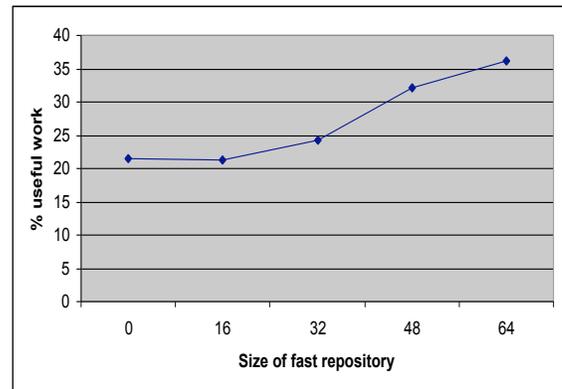


Fig. 4 FAST LPs scaling performance

Fig. 4 shows the results of varying the value of k , used to determine the number of LPs to be classified as FAST. This experiment is performed for the data point 40 of the previous experiment (60% [200:1], 40% [50:0]). It is observed that as the number of fast LPs is increased, the percentage of useful work also improves. However, if k is large migration cost will dominate.

5.2.2 Experiment under Set 2

Another set of experiments was performed, in which the self-instantiation degree of messages and the delays of LP execution were swapped in the two classes of operating restrictions to create imbalance in the application. With all the other experimental parameters being the same as set 1, but instead of varying the messages with self-instantiation of 200 with LPs having event granularity one millisecond as in previous experiments, we combined messages with self-instantiation 50 with LPs having event granularity of 1 millisecond. Again in this set of experiments we vary the percentage Class [50: 1msec] from 0 to 100 in increments of 20 percent. The remaining LPs in the system execute with null event granularity (Class [200: null]).

Fig. 5 shows the performance in terms of percentage of useful computation for simulation runs with and without process migration. Each of the two scenarios were plotted to show the performance of the simulations with various levels of optimism created by varying LPs with different event granularity. As shown in Fig 5, it was observed that process migration performs better than without it at all the levels except at data point 100. This is because the simula-

tion application is balanced at data point 100. At data point 100, all the LPs had one millisecond event granularity. This caused all LPs to progress at the same rate. In this scenario the overhead of process migration exceeded the benefits obtained by it and hence, the simulation without

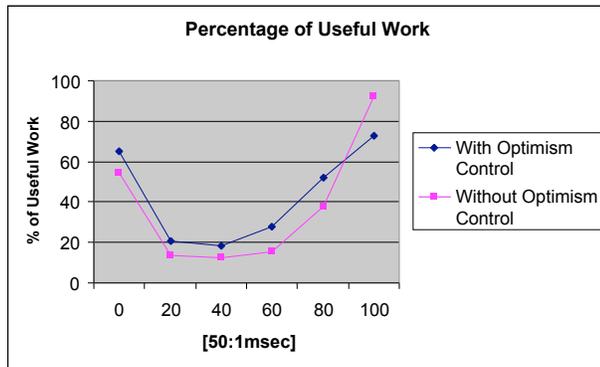


Fig. 5 Percentage of Useful Work (Set 2)

process migration performed better. In all the other instances, when the simulation is unbalanced, simulation with process migration performs 1.2 to 1.79 times better than without process migration.

In fig. 6, we compare the execution time performance for the above experiment with simulation runs with and without process migration. As mentioned above, due to the simulation application being balanced at data point 0 and 100 as all the messages and LPs belong to either of the two classes of self-instantiation and event granularity respectively, we do not observe any benefits in execution time. At data point 0, the simulation without process migration performs better than one with process migration both in terms of execution time and percentage of useful work. Whereas at data point 100, we obtain benefits in terms of percentage of useful work but not in terms of execution time. A similar scenario exists at point 80, where we obtain the benefits in terms of useful work but not in execution time. For the remaining proportions of Class [50:1msec], benefits in both execution time and computation performance were obtained when the simulation was run with process migration as compared to without process migration. When the application is unbalanced, the im-

provements in execution time by doing process migration were approximately 10 to 20 percent.

We observed that our algorithm performs better in terms of useful work in all the experiments except when the application is balanced. We observed that the migration cost of LPs do incur a penalty in terms of execution time in some cases. One of the reasons for this penalty is due to the PVM communication which is expensive. Due to resource limitations, the experiments were performed on older hardware. By updating the hardware we might reduce this migration cost to some extent. Also it was observed that a higher value of k improves the performance in terms of useful work, but may incur higher migration cost in terms of execution time. Hence, an appropriate value of k should be selected as a tradeoff between useful work and migration cost.

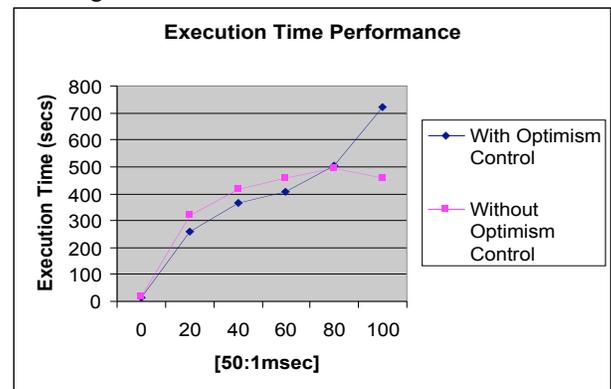


Fig. 6 Execution Time Performance (Set 2)

5 CONCLUSION AND FUTURE WORK

In a standard optimistic parallel event simulation, there are no restrictions on how far an LP can proceed ahead of other LPs. This imbalance may degrade performance as LPs may spend more time rolling back than progressing forward. In order to control such over-optimistic behavior, we implemented a process migration system on distributed Georgia Tech Time Warp (GTW). The process migration system classifies LPs with respect to optimism as FAST, MEDIUM and SLOW depending on how far each LP is making progress in simulation time compared to other LPs. The statistics collected to classify an LP is called *GVTlag*, which is the difference of an LPs current simulation time and the systems GVT. We aggregate the FAST LPs on one processor called the FAST-REPOSITORY to make them compete with each other for available CPU cycles to slow their progress and consequently isolate their impact on other logical processes. At the same time we spread the SLOW LPs, by redistributing each SLOW LP to a different processor. This provides SLOW LPs with sufficient CPU cycles to make progress and help them to catch up with the FAST LPs.

We performed various experiments with a synthetic benchmark application called *P-Hold*. We implemented P-Hold using two computation event granularities: null and 1msec. Additionally, P-Hold was configured to be self-instantiating. In order to create imbalance in the application LPs and events were classified in two classes having either self-instantiation degree of 50 or 200 and with and without 1 msec delay in event processing.

It was observed that whenever the application is unbalanced simulation with process migration performs 1.25 to 2.75 times better in terms of useful computation than without process migration. Also, the execution time under unbalanced application conditions is either better or equivalent to the time taken without process migration.

On the other hand, when the application is balanced and over-optimism is minimal, simulations without process migration perform better in terms of useful work and execution time. This is due to the overheads of process migration exceeding the benefits obtained when the application is balanced.

Due to high memory requirements and resource limitations, we could not experiment with more than 256 LPs in the simulation. We propose to experiment with more than 256 LPs in the future and observe how the process migration performs when the number of logical process increases in the simulation. Also, we propose to modify the process migration algorithm to make the process of selecting the number of fast LPs dynamic and dependent on how the simulation is performing. This could help in improving the amount of useful work done and execution time compared to the simulation without process migration.

A problem our approach may encounter in extreme cases is that logical processes might aggregate on the same processors leaving other processors sparsely populated. For example, we may encounter a situation in which a million LPs are on one processor and one logical process is on each of the remaining processors. Another problem could be the movement of logical processes that are already balanced. One of the possible solutions to alleviate such scenarios is to use a threshold to determine the number of processes classified as fast.

We also propose to cluster LPs similar to (Carothers and Fujimoto, 2000), in order to reduce the process migration cost. This would need a different metric for evaluating process migration. We propose to evaluate such a possibility and observe its performance compared to without process migration.

REFERENCES

- Jefferson, D. R., and H. Sowizral. 1985. Fast Concurrent Simulation Using the Time Warp Mechanism, Part – I : Local Control, *ACM Transactions Programming Languages and Systems*, 7(3): 404-425.
- Briner, J. 1997. Fast Parallel Simulation of Digital Systems, *Advances in Parallel and Distributed Simulation*, pp 20-28, June 1997.
- Wieland, F., E. Blair and T. Zukas, Parallel Discrete Event Simulation (PDES): A Case Study in Design, Development and Performance Using SPEEDES, *In Proceedings of the 9th Workshop on Parallel and Distributed Simulation (PADS '95)*, pages 103 -110, June 1995.
- C. Carothers and R. Fujimoto, Distributed Simulation of Large-Scale PCS Networks, *In Proceedings of the Second International Workshop on Modeling, Analysis, and Simulation of Computer and Telecommunication Systems(MASCOT'94)*, pages 2-6, 1994 .
- F. Hao, K. Wilson, R. Fujimoto and E. Zegura, *Logical Process Size in Parallel Simulations*, *In Proceeding of the Winter Simulation Conference* , pp 645-652, December 1996
- R. M. Fujimoto, *Parallel and Distributed Simulation Systems*, Wiley Series on Parallel and Distributed Computing, Reading Pg 137 – 138, 2000.
- P. L. Reiher, F. Wieland, and D. R. Jefferson, *Limitation of Optimism in the Time Warp Operating System*, *In Proceedings of the Winter Simulation Conference*, pages 765-770, 1989.
- P. M. Dickens and Jr. P. F. Reynolds, SRADS with local rollback, *In Proceeding of the SCS Multiconference on Distributed Simulation*, pp. 161-164, vol 22, January 1990.
- B. D. Lubachevsky A. Shwartz and A. Weiss, Rollback Sometimes Works...if Filtered, *In Proceedings of the Winter Simulation Conference Proceedings*, pp. 630 – 639, Dec 1989.
- Madisetti V. K, Hardaker, D.A. , and Fujimoto, R. M. 1983, The MIMIDIX Operating System for Parallel Simulation and Supercomputing, *J. Parallel and Distributed Computing*, vol 18, no. 4, pp. 473-483, August 1993.
- Steinman, J. S. 1983, Breathing Time Warp, *In Proceedings of the Seventh Workshop on Parallel and Distributed Simulation*, vol 23, pp. 109-118, May 1993.
- Jefferson, D. R., Virtual Time II: The Cancel Back Protocol for Storage Management in Distributed Simulation, 1990. *In Proceedings of the 9th Annual ACM Symposium on Principles of Distributed Computing*, pp. 75-90, Aug 1990.
- Lin, Y.-B., Preiss, B.R, Optimal Memory Management for Time Warp Parallel Simulation, *ACM Transactions on Modeling and Computer Simulation*, pp 283 – 307, October 1991.
- K. Jones and Samir R. Das, Combining Optimism Limiting Schemes in Time Warp based Parallel Simulations, *In Proceedings of the 1998 Winter Simulation Conference*, pages 499--505, Dec. 1998 .

- C. D. Carothers and R. M. Fujimoto, Efficient Execution of Time Warp Programs on Heterogeneous, NOW Platforms. *IEEE Transactions on Parallel and Distributed Systems*, Vol. 11, No. 3, March 2000.
- D. R. Jefferson, Virtual Time, *ACM Trans. Programming Languages and Systems*, vol. 7, no. 3, pp. 404-425, July 1985.
- S. Das, R. Fujimoto, K. Panesar, D. Allison and M. Hybinette, GTW: A Time Warp System for Shared Memory Multiprocessors, In *Proceedings of the Winter Simulation Conference 1994*, pp. 1,332-1,339, December 1994.
- R. M. Fujimoto, Performance of Time Warp under Synthetic Workloads, In *Proceedings of the SCS Multi-conference on Distributed Simulation*, vol. 22, pp. 23-28, Jan. 1990.
- C. Burdorf and J. Marti, Load Balancing Strategies for Time Warp on Multi-User Workstations, *The Computer Journal*, vol 36, no.2, pp. 168-176, 1993.
- R. M. Fujimoto, S. R. Das, K.S. Panesar, M. Hybinette and C. Carothers, Georgia Tech. Time Warp Programmer's Manual for Distributed Network of Workstations, Technical Report GIT-CC97-18, College of Computing, Georgia Inst. of Technology , July 1997.
- S. Das and R. M. Fujimoto, Adaptive Memory Management and Optimism Control in Time Warp, *ACM Transactions on Modeling and Computer Simulation (TOMACS)*, Volume 7, Issue 2, 1997.