

# A Fine-Grain Time-Sharing Time Warp System

ALESSANDRO PELLEGRINI and FRANCESCO QUAGLIA, Sapienza Università di Roma

Several techniques have been proposed to improve the performance of Parallel Discrete Event Simulation platforms relying on the Time Warp (optimistic) synchronization protocol. Among them we can mention optimized approaches for state restore, as well as techniques for load balancing or (dynamically) controlling the speculation degree, the latter being specifically targeted at reducing the incidence of causality errors leading to waste of computation. However, in state-of-the-art Time Warp systems, events' processing is not preemptable, which may prevent the possibility to promptly react to the injection of higher priority (say, lower timestamp) events. Delaying the processing of these events may, in turn, give rise to higher incidence of incorrect speculation. In this article, we present the design and realization of a fine-grain time-sharing Time Warp system, to be run on multi-core Linux machines, which makes systematic use of event preemption in order to dynamically reassign the CPU to higher priority events/tasks. Our proposal is based on a truly dual mode execution, application versus platform, which includes a timer-interrupt-based support for bringing control back to platform mode for possible CPU reassignment according to very fine grain periods. The latter facility is offered by an ad-hoc timer-interrupt management module for Linux, which we release, together with the overall time-sharing support, within the open source ROOT-Sim platform. An experimental assessment based on the classical PHOLD benchmark and two real-world models is presented, which shows how our proposal effectively leads to the reduction of the incidence of causality errors, especially when running with higher degrees of parallelism.

CCS Concepts: • **Computing methodologies** → **Discrete-event simulation**; *Massively parallel and high-performance simulations*; • **Software and its engineering** → **Scheduling**;

Additional Key Words and Phrases: Optimistic synchronization

## ACM Reference Format:

Alessandro Pellegrini and Francesco Quaglia. 2017. A fine-grain time-sharing time warp system. *ACM Trans. Model. Comput. Simul.* 27, 2, Article 10 (May 2017), 25 pages.

DOI: <http://dx.doi.org/10.1145/3013528>

## 1. INTRODUCTION

Parallel Discrete Event Simulation (PDES) is a universally recognized methodology for speeding up the execution of (very) large/complex discrete event models via the exploitation of hardware parallelism [Fujimoto 1990a]. It is based on partitioning the model into multiple simulation objects, historically referred to as Logical Processes (LPs), whose events are concurrently dispatched for execution.

One core problem in PDES is how to ensure causally consistent (namely, timestamp ordered) execution of the events at all the simulation objects, which is not trivial due to the dependencies that arise when different objects schedule events for each other. This is also known as the *synchronization* problem, for which different approaches and

---

Authors' addresses: A. Pellegrini, Dipartimento di Ingegneria Informatica, Automatica e Gestionale "Antonio Ruberti", Sapienza Università di Roma, Via Ariosto 25, 00185 Roma, Italy; email: [pellegrini@dis.uniroma1.it](mailto:pellegrini@dis.uniroma1.it); F. Quaglia, Dipartimento di Ingegneria Informatica, Automatica e Gestionale "Antonio Ruberti", Sapienza Università di Roma, Via Ariosto 25, 00185 Roma, Italy; email: [quaglia@dis.uniroma1.it](mailto:quaglia@dis.uniroma1.it).

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept., ACM, Inc., 2 Penn Plaza, Suite 701, New York, NY 10121-0701 USA, fax +1 (212) 869-0481, or [permissions@acm.org](mailto:permissions@acm.org).

© 2017 ACM 1049-3301/2017/05-ART10 \$15.00

DOI: <http://dx.doi.org/10.1145/3013528>

protocols have been presented in literature. The various proposals differ from each other by whether they entail (or not) the possibility to speculate along the simulated trajectory. If speculation is allowed, events are dispatched for processing at any simulation object as soon as they are available (at the underlying PDES environment) with no preliminary assessment of their safety. If causal inconsistencies arise, their effects are undone via rollback schemes. This is the Time Warp approach introduced in the seminal article by Jefferson [1985].

The relevance of speculative synchronization for PDES lies in that it allows for extremely high scalability. Recent results [Barnes-Jr. et al. 2013] have indeed shown how Time Warp systems exhibit the potential for scaling up to millions of processing units. On the down side, building Time Warp-based PDES platforms is far from being a trivial task because of two main reasons. One is related to the need for including the support for reversibility of the simulation model execution trajectory, an objective that should be pursued as transparently as possible to the overlying simulation application. Second, the actual performance delivered by Time Warp-based PDES platforms can be strongly affected by the rules according to which its worker threads CPU-dispatch the events to be processed.

The common literature trend is to build Time Warp systems as user-space platforms that are seen by the application-level code as run-time environments offering a specific API (e.g., for cross-simulation-object scheduling of events) and, in the most advanced cases (see, e.g., Pellegrini et al. [2015]), providing application transparent support for reversibility of the actions performed by both the native application code and the invoked third-party (standard) libraries. Invocations to the latter, or side effects on the simulation state natively produced by the application code, are in fact transparently intercepted by the platform-level code (via wrapping and/or instrumentation [Pellegrini et al. 2015; Rönngren et al. 1996; West and Panesar 1996]), which runs reversible versions of the corresponding tasks.

Nonetheless, another common way of implementing Time Warp systems is the one where each CPU-dispatched simulation event is executed in *non-preemptable manner*. Consequently, the platform-level software is not allowed to re-evaluate CPU assignment until the completion of the last-dispatched event. This approach is not able to promptly react to the (system wide) dynamic generation and injection of events with higher priority, say, lower timestamps, compared to the one currently being processed by some CPU-core. Consequently, it is not fully optimized given that the generation of rollbacks, and the associated waste of computation, tends to increase when events are CPU-dispatched and processed according to a rule that does not fully fit the priorities associated with the dynamic generation of timestamped events [Quaglia and Cortellessa 2002]. We note that the reduction of rollback incidence cannot be fully tackled by solely relying on load balancing/sharing strategies (see, e.g., Carothers and Fujimoto [2000], Choe and Tropper [1999], Glazer and Tropper [1993], and Vitali et al. [2012]), since they operate as long term planners for fruitful CPU usage, thus being not suited for “prompt” response to punctual variations of the event priorities along time.

Clearly, the ideal approach to preempt the execution of a CPU-dispatched event would be to interrupt the thread execution flow right upon the delivery of some higher priority event (or anti-event), destined to one of the simulation objects managed by the same thread. This would require a mechanism to reflect the arrival of a new event (say, of a new message) into a change of the state of the CPU-core (e.g., the instruction pointer) so the running thread can change its execution path, thus enabling the higher priority event to be actually dispatched.<sup>1</sup>

---

<sup>1</sup>For the case of an incoming higher priority anti-event, the thread would dispatch the corresponding management operations, including the rollback of the target simulation object.

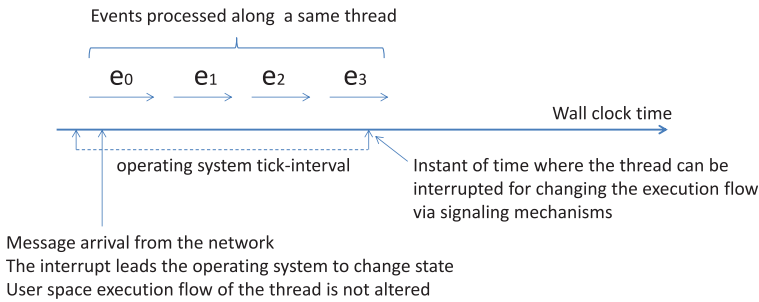


Fig. 1. Thread reaction to the injection of higher priority events on a conventional system.

For the case of network-based message passing in distributed memory systems, the arrival of a message from a network interface is reflected into a change of the operating system state, which makes the message accessible on I/O channels (e.g., via polling). But the thread itself does not change its execution flow, except if asynchronous signaling mechanisms are adopted. However, these would operate according to the time-granularity of conventional timer-interrupt mechanisms, say, 1 to 4ms on conventional operating system configurations. This is a granularity level that does not allow prompt preemption of events with common PDES workloads, where CPU requirements for processing simulation events are well below the order of milliseconds. An example scenario illustrating this problem is shown in Figure 1.

A similar problem still appears for the case of Time Warp systems running on top of shared-memory platforms. In more detail, if user space shared-memory support is used for exchanging messages across the threads, a sending thread will only post the new message on a shared-buffer, which will be checked by the destination thread according to a polling mechanism operating before (or after) the processing of an event. In fact, pure shared-memory based communication provides no effective mechanism for interrupting the event execution at some destination thread right upon the post of the new message. Even if a signaling mechanism was being used by the source thread, such as Posix user-defined signals, the time-granularity for the signal delivery to the destination thread would be still bound to the conventional operating system timer-interrupt interval configuration, thus resulting not adequate. Also, this approach would require the whole chain of signal management mechanisms to be passed through at the level of the operating system kernel, with consequent non-minimal overhead.

In this article, we cope with preemptive events' processing in Time Warp systems to be run on shared-memory multi-core machines. Also, we target C-based software and Linux/x86-64 computing systems. Our solution overcomes the above depicted limitations by enabling the platform-level software to take back control and to re-evaluate CPU assignment with very fine grain period (on the order of tens of microseconds). To achieve this target we designed and developed an ad-hoc timer management Linux module, which allows for (periodical) control flow variations along any running thread with no intervention by the chain of kernel-level mechanisms used for supporting Posix signals, hence leading to minimal run-time overhead. This is achieved by dividing each operating system tick assigned to the thread into sub-ticks, each one leading to an execution flow variation that brings control to a fast priority check routine implemented at the Time Warp platform level. The latter accesses compact data posted on shared-memory to determine whether the currently processed event is still the highest priority one (across those bound to the simulation objects managed by that same thread). In the negative case, the thread preempts the execution of the current event and switches to the execution of some higher priority task, according to a *fine-grain time-sharing*

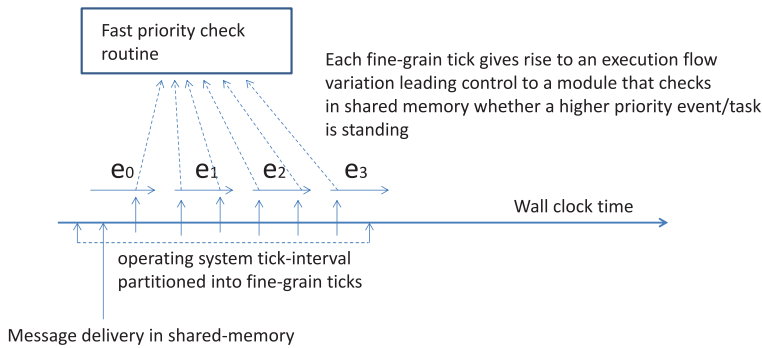


Fig. 2. Thread reaction to the injection of higher priority events in our fine-grain time-sharing system.

approach. The time-line of the execution of a thread with our approach is schematized in Figure 2. Clearly, the higher the frequency of fine-grain ticks' delivery, the higher the likelihood of prompt switch to some higher priority task, if any. But the overhead associated with the management of fine-grain ticks should be anyhow kept to a minimum value. To cope with this issue, we also provide a benchmark program that can be used to configure the frequency of fine-grain ticks in the target computing platform.

In our proposal, CPU assignment is also re-evaluated right before returning control back to the application code after the execution of an event has trapped into a platform level service, either explicitly or because of interception (aimed at making some action by the application modules reversible). Overall, the return to application code from platform level execution and the fine-grain ticks are exploited in a synergistic manner to maximize the opportunities to preempt events if higher priority ones have been delivered.

Our solution does not create any bias in terms of CPU assignment across threads running in the Linux system. In fact, the fine-grain tick mechanism we adopt does not alter the original operating system planning in terms of overall CPU time to be assigned both to the worker threads running within the Time Warp platform and to any other thread. This prevents impairing fairness when running our fine-grain time-sharing Time Warp system on top of a multi-user conventional platform.

Besides the ability to optimize CPU assignment depending on the (dynamic) priority of the tasks to be performed, our proposal has also the capability to address some specific liveness problems related to the speculative nature of Time Warp, such as application-level infinite loops that may arise when reaching a non-admissible state because of out of order events' execution [Nicol and Liu 1997]. These loops can be (timely) broken thanks to our fine-grain time-sharing approach, which can be exploited for supporting preemptive rollback operations leading to the squash of the non-admissible state trajectory.

The fine-grain time-sharing Time Warp architecture we have developed has been integrated within the open source ROOT-Sim package<sup>2</sup> (Pellegrini et al. [2011] and HPDCS Research Group [2012]) and operates in a fully transparent way to the overlying application code. Hence, the benefits from it come with no intervention by the application programmer. We also report experimental data for an assessment of our proposal, which have been collected by running three different test-bed applications—the PHOLD benchmark, a data store model and a personal communication system model—on top of an off-the-shelf 32-core machine.

<sup>2</sup>Available at <http://github.com/HPDCS/ROOT-Sim>.

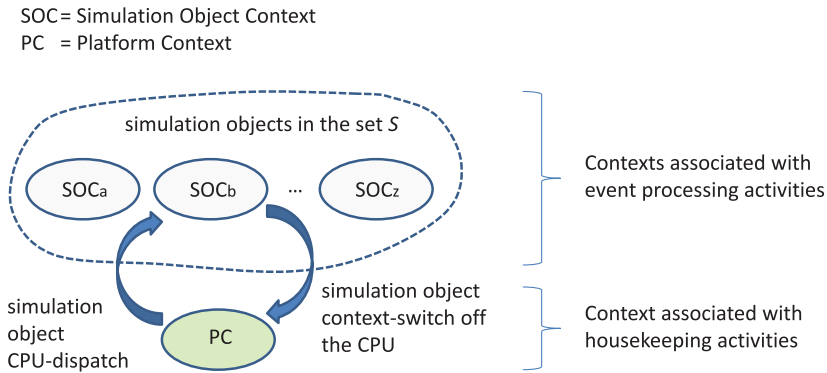


Fig. 3. Time-sharing Time Warp basics: execution contexts for an individual worker thread.

While presenting our proposal, we assume the reader is already familiar with Time Warp concepts, and we refer the less familiar readers to, for example, Jefferson [1985] and Jafer et al. [2013] for background information.

The remainder of this article is structured as follows. The fine-grain time-sharing Time Warp architecture is presented in Section 2. Experimental data are provided in Section 3. Related work is discussed in Section 4.

## 2. THE TIME-SHARING ARCHITECTURE

### 2.1. Basics

We assume the organization of the Time Warp PDES platform to adhere to the multi-thread paradigm, which has been recently shown to offer benefits (compared to counterpart implementations based on separate processes) for several aspects, such as the avoidance of simulation object migrations for well-balanced usage of resources [Vitali et al. 2012] and optimized data-exchange [Wang et al. 2014]. With this type of organization, any subset  $S$  of the simulation objects is (temporarily) bound to a specific worker thread, which is in charge of managing the corresponding event queues—each one associated with an object—and of CPU-dispatching its bound objects for event processing. Further, all the worker threads share platform level data structures, which plays a central role in how our fine-grain time-sharing architecture handles the detection of event priority variations at run-time, as we shall discuss.

The basic organization of our fine-grain time-sharing Time Warp system is schematized in Figure 3. Each simulation object belonging to the set  $S$  managed by a given worker thread has its own execution context (e.g., its own stack). Additionally, a platform context is included, thus each worker thread operates, at any time instant, either in the context of some simulation object or in platform context. In the real implementation these contexts, including the switch between them, are managed by relying on context management functions inspired to classical `setjump` and `longjump` functions provided by the Posix API, whose detailed description is provided in the appendix.

When running in platform context, the worker thread carries out housekeeping tasks, such as the check for incoming events (anti-events) destined to the simulation objects it is currently managing, and the actual CPU-dispatch of the simulation objects. The latter operation takes place according to the Lowest-Timestamp-First (LTF) policy [Lin and Lazowska 1991], which leads the worker thread to CPU-dispatch the simulation object (belonging to its bound set  $S$ ) whose next event to be processed is the one with the minimum timestamp among the ones already delivered (those already known to exist).

Third-party library functions accessible by the application code (e.g., `malloc`, `free`, and `printf`) are transparently intercepted via wrapping schemes, which enable running the corresponding reversible instances supported by the run-time environment,<sup>3</sup> conforming to what suggested by a few literature works (see, e.g., Antonacci et al. [2013] and Rönngren et al. [1996]). The same is true for the case of application transparent code injection (say, instrumentation) aimed at intercepting memory updates by the application code to make them reversible [Pellegrini et al. 2015; West and Panesar 1996]. The injected software brings control to platform level in a manner that is logically equivalent to the interception of external libraries' invocations by the application code. As for this aspect, in our fine-grain time-sharing Time Warp system, reversibility of the updates occurring within the (dynamically allocated) memory chunks forming the live state image of the simulation object is achieved by relying on the checkpoint support offered by the DyMeLoR library [Toccaceli and Quaglia 2008]. In what follows we focus our discussion on the interception of external libraries, with the implicitly assumption that the discussion also covers scenarios based on application level instrumentation. In our time-sharing architecture, each time an external library function is invoked, we say that the execution switches to platform mode, and then switches back to application mode as soon as the function returns. Clearly, when the worker thread operates in platform context, it also operates in platform mode. On the other hand, when it runs within the context of some simulation object, it can switch from application to platform mode multiple times, depending on the interactions between the application code and the intercepted external libraries. The wrapper that in our proposal encapsulates any intercepted function has the following structure:

```
return_type _function_name_wrapper(.. params ..){
    return_type ret;
    _enter_platform_mode();
    ret = function_name_reversible(.. params ..);
    _try_leave_platform_mode();
    return ret;
}
```

where the preamble `_enter_platform_mode` and the tail `_try_leave_platform_mode` are macros that set/unset a per-worker thread flag indicating the current running mode. Further, as we shall discuss, the `_try_leave_platform_mode` macro is also used for implementing part of the event preemption logic leading to switch the current execution context, if needed. This is the reason for the “try” prefix in the macro.

The switch between application and platform mode (and vice versa) occurs in our architecture not only because of synchronous invocations to intercepted external functions issued by the application level software (when running in the context of some simulation object). Rather, a timer-interrupt handler operating in user space is used to bounce control to platform mode periodically. We refer to this handler as **extra-tick-manager** given that, as hinted, the time-sharing architecture leads a single operating system tick interval assigned to a worker thread to be partitioned into multiple fine-grain tick intervals just leading to extra-ticks hitting that thread. Overall, the state diagram for any worker thread operating within the fine-grain time-sharing Time Warp system is the one depicted in Figure 4.

The execution of **extra-tick-manager** is triggered by the ad-hoc timer management logic we have embedded within our Linux module, which allows delivering fine-grain

<sup>3</sup>In our view, reversibility also means that the behavior of the intercepted libraries is guaranteed to be piece-wise-deterministic, to allow optimized state restore schemes based on infrequent checkpointing and coasting forward.

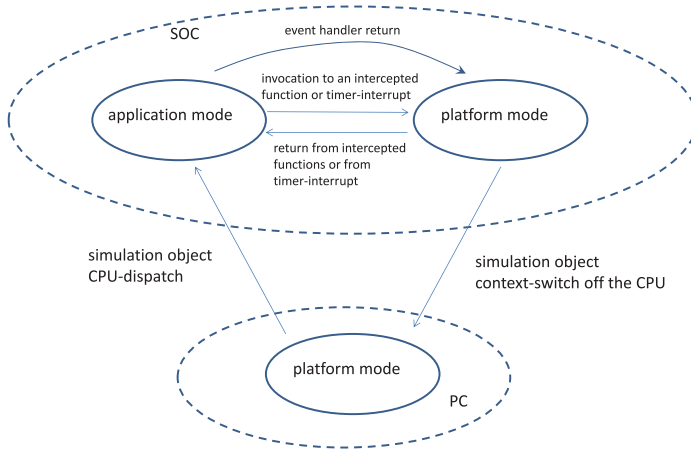


Fig. 4. Worker thread state diagram.

ticks to any worker thread running within the time-sharing Time Warp platform. The details of the implementation are provided in the appendix. For the abstraction level of the discussion in the main body of this article, the important point is that a thread that wishes to be interrupted according to fine-grain ticks needs to register itself via an `ioctl` command on the `dev_extra_tick` device file we support with our Linux module.

## 2.2. Run-Time Detection of Priority Variations and Event Preemption Logic

As well known, when CPU-dispatching of the simulation objects in  $S$  is carried out by the worker thread according to the LTF algorithm, priority variations of the currently executed events (say, a decrease of the priority of the last CPU-dispatched event caused by the delivery of some event—or anti-event—with a lower timestamp destined to some object belonging to  $S$ ) may only arise due to communication between simulation objects belonging to different subsets, say,  $S$  and  $S'$ , which are bound to different worker threads. Consequently, the architectural organization of the communication facility within the multi-thread Time Warp platform plays a relevant role in the run-time determination of the priority variation (if any) of the currently processed event.

As hinted, we focus on shared-memory communication, and we consider a scenario conforming to the indications in Vitali et al. [2012], where the exchange of messages/anti-messages across different worker threads does not take place by directly incorporating the corresponding information into the destination object event queue. Rather, messages are exchanged according to a top-half/bottom-half approach oriented to scalability. In particular, each worker thread manages a set of bottom-half queues (one for each simulation object belonging to the set  $S$  it is currently handling), such that any other worker thread in the system can notify the presence of new data to be ultimately incorporated into the destination object's event queue via the corresponding bottom-half queue. This is done via the execution of a top-half data record (tail) insertion into the bottom-half queue. Checking whether some new data is present into a bottom-half queue, and actual processing of the data with (timestamp-ordered) incorporation into the destination event queue, is carried out exclusively by the worker thread in charge of (currently) handling the destination object.

In our time-sharing Time Warp system, the above scheme has been extended along the following lines, to support early detection of priority variations. First, each worker thread  $t$  has been associated with a  $BH_{min}_t$  record, which represents at any time

instant the minimum timestamp of a message/anti-message that has been recorded in any of the bottom-half queues associated with the simulation objects that  $t$  is currently managing, since the last flush operation of these queues. In other words,  $BH\_min_t$  represents the minimum value among the timestamps of data in transit (if any), destined to some simulation object belonging to the set  $S$  handled by  $t$ .

This record is initialized to a special macro `INFINITE` when the worker thread  $t$  accesses its bound bottom-half queues and flushes the data into the corresponding event queues. Whenever a different worker thread inserts a bottom-half record into any of the bottom-half queues associated with the simulation objects managed by  $t$ , the reduction  $BH\_min_t = \text{Min}(BH\_min_t, T)$  is performed, where  $T$  represents the timestamp of the message/anti-message that is being placed into the destination bottom-half queue. In our implementation, this reduction is performed via an atomic Compare-And-Swap (CAS) instruction. This allows manipulating  $BH\_min_t$  while not requiring worker threads that concurrently access two distinct bottom-half queues associated with  $t$  to execute a conflicting critical section.<sup>4</sup>

Another record, called *current\_time<sub>t</sub>*, is associated with each worker thread  $t$ . It is used to keep track of the timestamp of the current simulation event, if any, that has been CPU-dispatched along  $t$ —this is the lowest-timestamp event according to LTF. The value of *current\_time<sub>t</sub>* is set to the special value  $-1$  if thread  $t$  is not currently processing any event, which means that it is running housekeeping operations in platform context. The values of *current\_time<sub>t</sub>* and  $BH\_min_t$  are used in combination to determine whether some higher priority task (compared to the one currently processed along thread  $t$ ) needs to be CPU-dispatched. In particular, the platform level function that executes the check and determines whether some higher priority task needs to be executed along thread  $t$ , which needs therefore context switch between simulation object and platform contexts (thus enabling CPU reassignment via platform level actions), is structured as follows:

---

*void* **check-and-switch()**

1. **if** (*current\_time<sub>t</sub>*  $\leq$   $BH\_min_t$ )
  2.     **return**;
  3. **else**
  4.     \_enter\_platform\_mode();
  5.     switch\_to\_platform\_context();
- 

The above structure allows changing the current execution flow along thread  $t$ , by pushing it to platform-context (and also to platform mode, if not already operating with this mode) if:

- (1) The simulation object currently dispatched for event execution along  $t$  needs to rollback, since it is the recipient of a message or an anti-message in its past— $BH\_min_t$  corresponds to the timestamp of a message/anti-message destined to the currently running simulation object. In this case the rollback operation will take place according to a preemptive mode.
- (2) Any simulation object belonging to the set  $S$  managed by  $t$  dynamically gains a priority higher than that of the currently running one, since it becomes the recipient of some message or anti-message with a timestamp lower than the one of the last event that was CPU-dispatched according to LTF. The case of an incoming anti-message is again representative of a causal inconsistent execution at the

---

<sup>4</sup>In fact, each of them needs to temporarily lock a different bottom-half queue for data insertion, which helps not hampering concurrency [Vitali et al. 2012].



destination simulation object, given the adopted LTF rule for CPU-dispatching the events.

In either case, control must return to the Time Warp platform layer, so the higher priority task (either a rollback operation or not) can be promptly executed. On the other hand, if no higher priority task needs to be executed, the **check-and-switch** function simply returns control back to its caller. Clearly, if the simulation object that is context-switched off the CPU still runs on a consistent path, the preempted event will be resumed (with no loss of already performed work) when LTF will again find it as the highest priority one.

The last aspect to discuss is related to how the **check-and-switch** function is integrated with the `_try_leave_platform_mode` macro and with the **extra-tick-manager** module. The integration with `_try_leave_platform_mode` takes place as follows:

---

```

_try_leave_platform_mode()
1. check-and-switch();
2. _leave_platform_mode(); //reset of the flag indicating platform mode execution
3. return; //regular return from an intercepted function

```

---

By the above pseudo-code, the **check-and-switch** function—which possibly leads to context-switch off the CPU the currently running simulation object—is invoked upon the finalization of any external library function that has been intercepted by the corresponding wrapper and is then executed in platform mode. If the invocation to **check-and-switch** in line 1 leads to no switch to platform context, then the flag indicating whether we are running in application or platform mode is correctly aligned with the return to application mode.

As for the integration between **check-and-switch** and the timer-interrupt handling function **extra-tick-manager**, we have the following structure:

---

```

void extra-tick-manager()
1. if(platform_mode)
2.   return; //already platform mode running - no control flow variation
3. else
4.   check-and-switch(); // do we need an execution flow variation?

```

---

If the timer-interrupt handler is activated while already running in platform mode, then no control flow variation needs to take place. In fact, if platform mode is currently associated with simulation object context, it means that the check on whether some higher priority task needs to be CPU-dispatched will be carried out right before returning to application mode via the `_try_leave_platform_mode` macro. If the current mode is not the platform one, then the handler triggers **check-and-switch** to verify whether higher priority tasks need to be carried out. This may lead to context-switch the currently processed event (say, the running simulation object) off the CPU.

The check in line 1 and the avoidance of the variation of the current execution flow if the worker thread is already in platform mode guarantee that whichever platform level block of code is executed along any worker thread as a non-preemptable action, which is a fundamental prerequisite. In fact, locks on data structures or memory regions might be acquired by some worker thread once the application has trapped into platform mode along that thread, which might be necessary to correctly manage the triggered

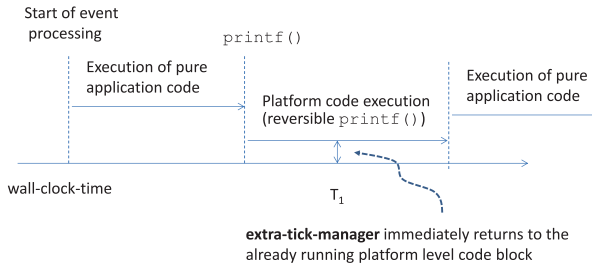


Fig. 5. Management of extra-ticks in the interleave between application and platform code blocks within an event processing wall-clock-time window.

service (see, e.g., Pellegrini and Quaglia [2015]). Hence, the in place critical section cannot be context-switched off the CPU.<sup>5</sup>

A schematization of the behavior of our fine-grain time-sharing Time Warp system in relation to the delivery of timer-interrupts while already running in platform mode is provided in Figure 5, where we show the arrival of an extra-tick at wall-clock-time  $T_1$ , with consequent activation of the **extra-tick-manager**. In this scenario, the interrupt handler simply returns given that at the same time instant the thread was already running a platform-level reversible version of the `printf()` function, via interception. However, the check on whether higher priority tasks would need to be dispatched by preempting the current event is anyhow carried out in our architecture as soon as the interrupted platform level function will attempt to return to application mode, which is done via the `_try_leave_platfrom_mode` macro.

### 2.3. Overall Configuration of the Time-Sharing Support

By the architectural organization of the event preemption support described in Section 2.2, in our Time Warp system the platform level software has two different triggers for context-switching a simulation object off the CPU: (a) timer-interrupts and (b) returns from platform mode. However, while the returns from platform mode are intrinsically related to the activities of the application level software (since they are triggered depending on the interaction between application level modules and the intercepted external libraries), timer-interrupts (and the cost/benefit they induce) depend on the configuration of the extra-tick interval. The shorter the length of the extra-tick period, which we denote as  $\Delta ET$ , the higher the expected overhead caused by timer-interrupts. However, shorter  $\Delta ET$  values can provide more opportunities for event preemption and prompt CPU reassignment to higher priority events/tasks, thus likely improving the effectiveness of the fine-grain time-sharing approach in reducing the amount of rollback.

To cope with the selection of  $\Delta ET$  and to optimize the synergy between the above two triggers for event preemption, we devise the following scheme. We denote with  $\widehat{\Delta ET}$  the minimum length of the extra-tick interval, which still induces negligible overhead due to extra-tick delivery. As we shall discuss in Section 3, the value of  $\widehat{\Delta ET}$  can be determined by running an ad-hoc benchmark in the early phase of the installation of our fine-grain time-sharing Time Warp system. Once determined  $\widehat{\Delta ET}$ , synergistic exploitation of the two different triggers for event preemption is based on the run-time estimation of (i) the average event granularity for the specific application, which we

<sup>5</sup>A way to cope with the interruption of platform level code blocks would be to design the platform level software according to the concept of “safe places,” which characterizes preemptable, for example, real-time, operating system kernels. However, this type of design is aside of the core focus of our time-sharing proposal.

refer to as  $\Delta_e$ , and (ii) the average number of switches to (and then back from) platform mode while processing an individual event. We recall that these switches (if any) take place while running in simulation object context (see the state diagram in Figure 4). We denote such an average number of switches as *APMS* (Application/Platform Mode Switches), hence the expression

$$T = \frac{\Delta_e}{APMS} \quad (1)$$

represents the average wall-clock-time interval after which the application software spontaneously provides control back to the platform software while an event processing phase is in place. In our implementation, the computation of  $\Delta_e$  and *APMS* has been based on the exponential mean of samples. The samples for computing  $\Delta_e$  are taken by monitoring, via `gettimeofday`, the CPU time spent for processing individual events. The samples for computing *APMS* are taken by counting the number of times the macro `_leave_platform_mode` is executed while processing any individual event.<sup>6</sup> Given that, according to the rules specified in Section 2.2, this already provides opportunities for CPU reassignment, the combination of the two different triggers for event preemption is based, in our design, on the relation between  $T$  and  $\widehat{\Delta ET}$ . Specifically, we dynamically switch on/off the extra-tick delivery along a thread depending on whether the following inequality is verified (or not):

$$\frac{T}{\widehat{\Delta ET}} \geq (1 + \alpha), \quad (2)$$

where  $\alpha$  is a tunable parameter whose value falls in the interval  $[-1, \infty]$ . If  $\alpha$  is set to the minimum value  $-1$ , then the Time Warp worker thread registers itself on the `dev_extra_tick` device file, thus being interrupted by the timer each  $\widehat{\Delta ET}$  time units, independently of the simulation objects' run-time behavior (in terms of switches between application and platform modes). For  $\alpha \rightarrow \infty$ , Equation (2) would not be satisfied—except for  $\widehat{\Delta ET}$  tending to zero, which is not realistic—leading the Time Warp worker thread to deregister itself from the `dev_extra_tick` device file, thus fully renouncing to be periodically interrupted for possible CPU reassignment.

A baseline setting for  $\alpha$  could be represented by the value zero, leading the Time Warp thread to register itself as one to be extra-ticked with period  $\widehat{\Delta ET}$  if the frequency according to which the execution of an event (taking place in simulation object context) switches to/from platform mode does not overstep the one of the extra-tick delivery. However, this setting would lead to reduced opportunities for event preemption if the switches to/from platform mode were not uniformly distributed along the lifetime of a simulation event. More conservative values of  $\alpha$ , say, in the interval  $[-1, -0.1]$ , would likely avoid this phenomenon. With these settings, even if  $T$  is less than  $\widehat{\Delta ET}$ , meaning that the execution in simulation object context spontaneously and frequently switches to platform mode, thus providing opportunity for event preemption, we still retain the possibility to achieve the same objective via the timer-interrupt scheme, which is done to avoid having a portion of the event processing interval uncovered by switches to platform mode.

As a last note, the quantity  $T$  in Equation (1) can be estimated at run-time on a per worker thread basis, so each worker thread can operate its decision on whether to dynamically register or deregister itself as one to be extra-ticked (depending on Equation (2)) independently of the other ones. This would allow coping with scenarios

<sup>6</sup>In the implementation the counter is directly updated by the macro upon its execution.

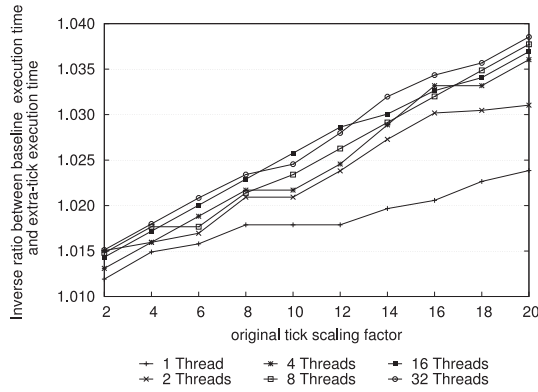


Fig. 6. Extra-tick overhead vs. variations of both the extra-tick period and the number of threads.

with simulation objects exhibiting different (heterogeneous) execution profiles, possibly giving rise to different  $\Delta_e$  and/or  $APMS$  values across the simulation objects' sets managed by the different threads.

### 3. EXPERIMENTAL RESULTS

#### 3.1. Determining $\widehat{\Delta ET}$

To determine the value of  $\widehat{\Delta ET}$ , an ad-hoc benchmark can be run in the early phase of installation of our fine-grain time-sharing Time Warp system on the target computing platform. Since PDES engines based on Time Warp are CPU-bound applications (given the absence of wait/block phases, at least in the presence of tasks/events to be actually carried out), our ad-hoc benchmark is simply made up by a multi-thread application, where each thread executes a busy loop of a given duration. We initially run this benchmark in a baseline configuration with no registration of the threads on `dev_extra_tick`, which leads to no extra-tick delivery. Then we run the benchmark again by registering the threads on `dev_extra_tick`, thus leading them to be periodically interrupted by the extra-tick logic. In our case, the target computing platform is a 32-core 64-bit HP ProLiant NUMA server, equipped with four 2GHz AMD Opteron 6128 processors (each one equipped with 8 CPU-cores) and 64GB of RAM. The operating system is Linux SUSE, kernel version 3.16.7. We used this same platform for all the experiments whose outcomes are reported in this section.

In the original configuration of the Linux kernel, the timer was set to issue an interrupt (a tick) each 1ms. When running our benchmark, we experimented with different values of the extra-tick interval achieved by scaling the original tick by a factor between 2 and 20, leading to experiment with extra-tick periods in the interval between 500 and 50 $\mu$ s. This allowed us to observe how the extra-tick overhead scales versus the length of the extra-tick period. In our benchmark, the extra-tick handler only increments a counter of delivered extra-ticks, and then returns control to the interrupted execution flow. This is aligned with the objectives of this benchmarking phase given that we only aim at evaluating the cost for delivering extra-ticks, independently of the usefulness of the delivery (hence independently of any real action to be taken upon extra-tick arrival). We also varied the number of running threads between 1 and 32, thus studying how the overhead varies versus the amount of threads managed according to the extra-tick logic.

We show in Figure 6 the inverse ratio between the execution time of the baseline configuration (no extra-tick), which we roughly report to be on the order of 30s, and

the execution time achieved with extra-ticks delivered to the application according to the selected scaling factor. Each sample has been computed as the average over five different runs of a same configuration; nonetheless, very minimal variation has been observed among the different sampled values. By the plot we see that the overhead induced by the extra-tick operating mode is less than 4% even when the scaling factor of the original tick interval (which we recall is of 1 millisecond) is set to the value 20, meaning that the extra-tick is delivered with granularity of  $50\mu\text{s}$ . Another interesting trend is that the overhead appears to be slightly higher when running the benchmark with larger number of threads. This is due to the slightly increased interference by common kernel level threads automatically started up by the operating systems (e.g., `kworker` threads), which is naturally induced when the benchmark runs by exploiting more CPU-cores. In fact, kernel level threads, although not being CPU-bound, lead anyhow to periodic operating system context switches, which in turn force our extra-tick management logic to more frequently interact with the timer, in terms of setting the requested interrupt period (depending on whether a `dev_extra_tick` registered thread, or not, is CPU-dispatched by the kernel).

Scaling factors of the original tick lower than the maximum value 20 lead the overhead by extra-tick delivery to be further reduced, at the expense of reduced opportunities for timer-interrupt based preemptions in the time-sharing Time Warp systems in case of adoption of such lower scaling factors. Also, extra-tick interval length of  $50\mu\text{s}$ , beyond still providing minimal overhead, looks a suited value—in terms of opportunities for preempting an event currently being processed—when considering complex PDES workloads characterized by event granularity well above the order of a few (or a few tens of) microseconds. These workloads can be considered as typical targets for Time Warp synchronization, and more generally for classical PDES methodologies. We intrinsically target this category of workloads with our fine-grain time-sharing Time Warp proposal, at least for what concerns timer-interrupt triggered preemptions. On the other hand, discrete event models with (very) fine grain events spontaneously bring control back to platform mode (after the CPU-dispatch of some event) in a prompt manner. Hence, they naturally allow the platform to promptly react to priority variations even when events are processed in non-preemptable manner. Still, for these workloads, our time-sharing Time Warp system offers the possibility to preempt a CPU-dispatched simulation object by relying on switches back from platform mode.

Finally, the overhead determined via this benchmark can be considered a worst-case reference value, since the busy loop run by the threads is not interfered by factors such as memory access latency (and its variation as a function of locality of the accesses), which would tend to reduce the relative per-instruction cost of extra-tick delivery.

### 3.2. Performance Results with the PHOLD Benchmark

In this section, we provide performance data collected by relying on the well-known PHOLD benchmark [Fujimoto 1990b]. The relevance of using PHOLD lies in that it entails events that are loosely coupled with the underlying Time Warp platform. In fact, they are simple CPU busy loops, which lead to no invocation of application-external libraries intercepted by the underlying platform.<sup>7</sup> This kind of workload allows assessing the benefits from our fine-grain time-sharing Time Warp system in scenarios where preemptions can essentially be triggered by fine-grain timer-interrupts.

In order to improve the representativeness of this study, the PHOLD benchmark configuration we selected entails three different execution phases having the same virtual time duration, which we refer to as A, B and C. PHASE-A is lightweight, being it characterized by event duration of the order of  $30\mu\text{s}$ . PHASE-B is middleweight,

<sup>7</sup>The only exception is the usage of pseudo-random generators.

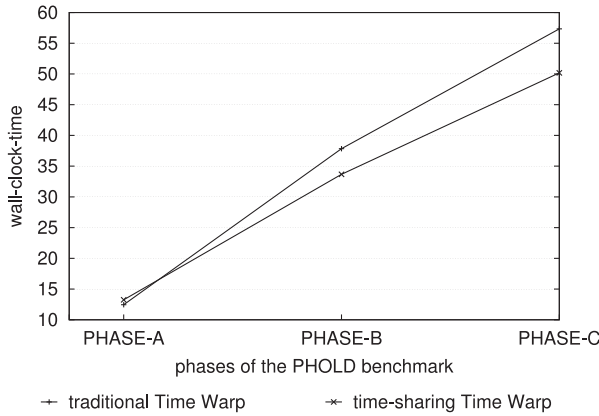


Fig. 7. Results with the PHOLD benchmark.

being it characterized by event granularity of the order of  $150\mu\text{s}$ . Finally, PHASE-C is heavyweight, since the events have granularity of the order of  $300\mu\text{s}$ . We configured the benchmark with 2025 simulation objects connected as a bi-dimensional mesh, which have been equally distributed among 32 worker threads operating in the Time Warp system. Five events (say, jobs) per simulation object have been inserted in the system. The jobs are routed randomly among neighbors by scheduling new events with exponential timestamp increment. This configuration of PHOLD (coupled with the selected level of parallelism in the underlying platform, say, 32 threads) gave rise to a speculative execution pattern characterized by infrequent rollbacks undoing large numbers of events. The overall efficiency (say, the ratio between the number of committed events, and the total number of processed events, say, committed plus rolled back) that has been observed in executions with the traditional configuration of the Time Warp system was on the order of 50%, with minor variations in the different execution phases.

In Figure 7 we show the execution time of the different phases of the PHOLD benchmark for the case of both traditional Time Warp and the time-sharing version.<sup>8</sup> Both these configurations rely on the same core PDES engine, namely ROOT-Sim, within which the time-sharing support has been integrated. The runs have been carried out by setting the extra-tick period to  $50\mu\text{s}$ . Also, the parameter  $\alpha$  in Equation (2) has been set to the baseline value zero. This led time-sharing executions to have worker threads (dynamically) registered on the `dev_extra_tick` device file along PHASE-B and PHASE-C, but not along the lightweight PHASE-A. On the other hand, the two different platform configurations, time-sharing and traditional, were run with either the `dev_extra_tick` device file active or not, respectively. Therefore, the execution time of PHASE-A can help assessing the overhead by the device file logic, in relation to checking whether a thread that has been CPU-dispatched by the operating system kernel is a registered one or not, plus the overhead for managing simulation object contexts.

By the results we see that the time-sharing Time Warp version introduces about 4% overhead during PHASE-A, along which it provides no advantage due to scarce impact of preemptions, since extra-ticks are not delivered in this phase. However, when switching to PHASE-B and then to PHASE-C, the time-sharing version allows growing reduction of the execution time. In particular, the time-sharing Time Warp system

<sup>8</sup>All the reported samples have been computed as the average over 10 runs.

Table I. Ratio Between the Efficiency of Traditional Execution and the Efficiency of Time-sharing Execution

execution phase	$\frac{\text{efficiency of traditional Time Warp}}{\text{efficiency of time-sharing Time Warp}}$
PHASE-A	1.003
PHASE-B	0.935
PHASE-C	0.917

is between 9% and 11% faster than the traditional one. This happens thanks to the delivery of extra-ticks during both PHASE-B and PHASE-C, possibly triggering CPU reassignment to higher priority events/tasks, which leads the time-sharing version to reduce the amount of wasted computation. In particular, we show in Table I the ratio between the efficiency of the traditional execution, and the efficiency of the time-sharing execution in the different phases. By the data, time-sharing Time Warp provides about 7% better efficiency along PHASE-B and about 8% better efficiency along PHASE-C. This result, in combination with the particular rollback pattern with unfrequent but long rollbacks, allows for boosting the final performance gain. Such a gain is not only originated by the reduction of the number of events that are eventually undone but also by the reduction of rollback management costs, such as the cost for managing anti-messages, which may not scale linearly, for example, for locality reasons.<sup>9</sup>

### 3.3. Performance Results with a Data Store Model

As an alternative workload to PHOLD, in this section we consider a real-world discrete event model of an in-memory key/value data store system. This type of models has recently become attractive (e.g., for capacity planning purposes), since real platforms based on this data storing paradigm have become a first-class technology in modern (e.g., cloud-based) infrastructures.

We simulated a distributed data store with 64 nodes, each one modeled by a separate simulation object, where data are partitioned and the partitions are distributed across the nodes. Batches of transactional data access requests are delivered to each node by proper simulation events that are self-generated by the same simulation object modeling the node, which are scheduled following an exponential distribution of their timestamps. The transactions may entail accessing the local partition or remote partitions, and the access to remote data partitions leads to cross-simulation-object exchange of simulation events, carrying as payload the set of transactional requests that require access to the remote partition. The batching factor determines the actual workload to be simulated, hence the resource requirements for executing the simulation. We have considered two different configurations of this model, a lightweight configuration characterized by batching factor set to 10 and a heavy one characterized by batching factor set to 20. The transactional requests within each batch are processed (in the simulation) by having them managed via a round-robin scheme, resembling the assignment of real CPU resources in a multi-thread data management system. For the lightweight configuration, the average CPU requirement for simulating the event delivering the batch of transactions is of the order of  $300\mu s$ . Instead, the heavy configuration has CPU requirement of the order of slightly less than  $500\mu s$ . One primary objective of this type of simulation is the determination of the data store performance (e.g., its

<sup>9</sup>In the configuration of PHOLD we have used, each undone event by a rollback operation requires sending a corresponding anti-message, which leads to costs for both send and receive tasks, including the cost for scanning output/input queues of the simulation objects, operations that lead to reduced locality, especially for longer rollbacks. On the other hand, as the PHOLD benchmark is an application with almost no state, the checkpoint/restore cost for the data structure representing the state of the simulation objects does not influence performance significantly vs. variations of the amount of rollback.

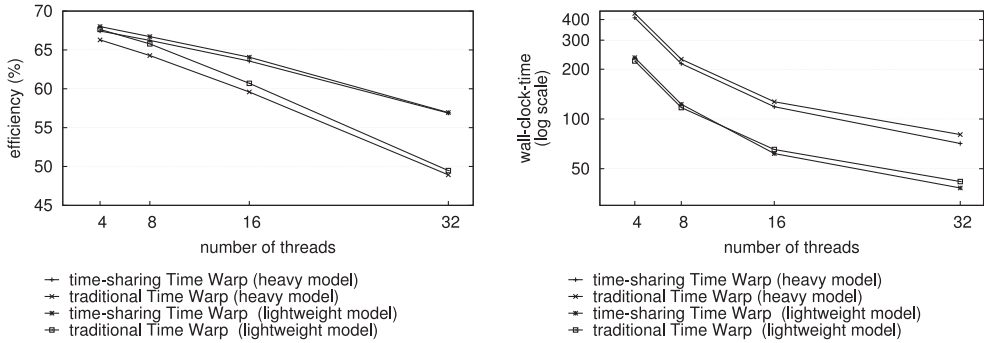


Fig. 8. Results with the data store model.

throughput), while varying the size of data partitions, the transaction access pattern (which may give rise to aborts depending on the materialization of data conflicts), and the locality of the accesses to the partitions.

In this study, we still set the extra-tick interval to the value  $50\mu s$ , while selecting the most conservative value for the parameter  $\alpha$ , which has been set to  $-1$ . This choice is motivated by the fact that this simulation model makes use of dynamic memory allocation/deallocation services for keeping the meta-data representing the transactional requests. Therefore, it generates a non-negligible amount of interactions with the underlying Time Warp platform, since the dynamic memory services are intercepted to make them reversible (as hinted, via the DyMeLoR library [Toccaceli and Quaglia 2008]). Since this already gives rise to switches to/from platform mode while processing the events, which provide opportunities for event preemptions, setting  $\alpha$  to such a conservative value allows us to still exploit opportunities for preemptions thanks to timer-interrupts. In fact, according to Equation (2), with such a conservative value the worker threads operating in the time-sharing configuration of the Time Warp system register themselves on the `dev_extra_tick` device file.

Beyond being focused on a real-world simulation model, this study complements the one with PHOLD for a few additional aspects. First, this time we varied the degree of execution parallelism by varying the number of worker threads between 4 and 32. This has been done to compare traditional and time-sharing executions of the Time Warp system with different concurrency degrees, which in turn give rise to different amounts of rollback. Second, the type of interactions among the simulation objects in the data store model gives rise to a rollback pattern that is opposite to the one observed with PHOLD. In fact, it is made up by frequent rollback occurrences, each one undoing a reduced number of events.

We report in Figure 8 (left side) the variation of the efficiency of the simulation run<sup>10</sup> while varying the number of worker threads for both traditional and time-sharing Time Warp systems. When running with low parallelism degree (say, four worker threads), both the systems show relatively high efficiency, which is slightly less than 70%. The traditional version gives rise to a bit reduced efficiency limited to the case of the heavy model configuration. However, while scaling up the degree of parallelism, the efficiency provided by the traditional Time Warp system rapidly degrades, falling just below 50% when running with 32 worker threads in both heavy and lightweight model configurations. Instead, the time-sharing version allows keeping the efficiency value significantly higher, leading to the order of 57% efficiency for the case of 32

<sup>10</sup>Also for this study we report average values over 10 different samples.



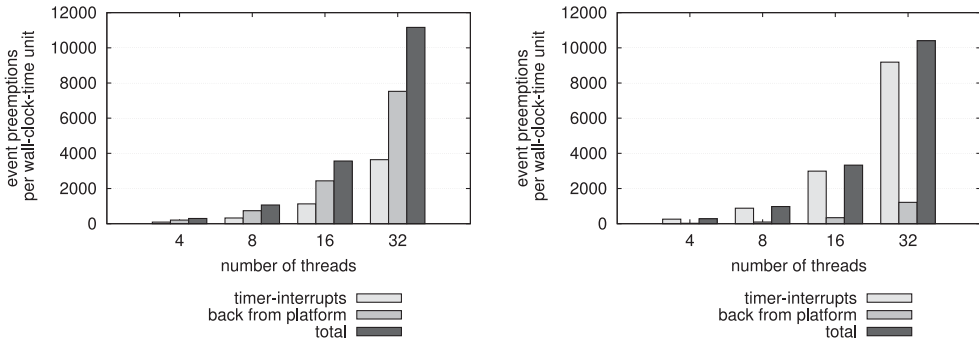


Fig. 9. Frequency of event preemptions for the heavy data store model—original (left) and varied (right) memory allocation patterns.

worker threads in both model configurations. Overall, the efficiency provided by the time-sharing Time Warp system stands up to 14% better than the one provided by the traditional Time Warp system. This advantage is reflected into a reduction of the simulation model execution time, as shown in Figure 8 (right side), which improves when increasing the degree of parallelism. Particularly, when relying on 32 worker threads the performance gain by the time-sharing Time Warp system is of the order of 15% for the heavy model configuration and of the order of 11% for the lightweight model configuration. As compared to PHOLD, this time the advantage provided by the time-sharing configuration on the side of efficiency does not further boost in terms of final performance, which is due to the different rollback pattern. With more frequent rollback occurrences, the time-sharing Time Warp system has improved chances for early detection of (potential) causality violations. However, the gain by reducing the overall amount of rollback is essentially due to the avoidance of processing events that would be eventually undone, rather than to significant reductions of the cost for managing rollback phases (given that the rollback length is very short). Still, the gain by the time-sharing version, at the point of maximum parallelism, is significant. Also, the maximum speedup by the time-sharing Time Warp system compared to sequential executions of the same data store models has been observed to be of the order of 20.

In Figure 9 we analyze the run-time dynamics of the time-sharing Time Warp system from a finer grain perspective. In particular, we report the number of event preemptions per wall-clock-time unit triggered either by timer-interrupts or by switches back from platform mode and the sum of the two (marked as “total” in the plots). On the left side, we show these data for the case of the original version of the data store model (heavy configuration), while on the right side we show the values achieved by running a modified version that is functionally equivalent, but where the instantiation of the meta-data for simulating the transactional requests does not take place following the round-robin scheme according to which the advancement of transactions execution is simulated. Rather, we instantiate these meta-data right upon starting the simulation phase of the whole batch of transactions. This variant only anticipates the instantiation operation at the begin of the event that delivers the batch of transactions to be simulated. Hence, differently from the original version, the interaction with the underlying platform, which intercepts the dynamic memory management requests for instantiating the meta-data, are much more clustered along the execution phase of simulation events. By the data in Figure 9, we can see that the amount of preemptions triggered by timer-interrupts definitely increases in the second configuration, which also shows a reduction of the incidence of preemptions originated by switches back from platform mode. However, the important indication that is conveyed by these data is that the total

Table II. Ratio Between the Execution Times of the Two Tested Configurations of the Heavy Data Store Model with Time-sharing Time Warp

number of used threads	$\frac{\text{execution time of the original configuration}}{\text{execution time of the varied configuration}}$
4	0.989
8	0.985
16	0.979
32	0.976

amount of preemptions per wall-clock-time unit is very similar in the two scenarios, independently of the number of used worker threads. This points out the robustness according to which timer-interrupts and switches to/from platform mode can be combined in complex workloads especially when relying on conservative values of the parameter  $\alpha$ . In fact, by these results we see that the fine-grain time-sharing Time Warp system does not degrade its ability to early detect priority variations independently of the actual pattern of interaction between application and platform level software. Also, the version of the data store model with clustered allocation of meta-data for simulating the transactional requests has shown execution times very close to the ones observed for the case of the original version with either traditional or time-sharing configurations of the Time Warp system, which is somehow expected given that no relevant change in the actual run-time dynamics were induced. Relative performance values of the two versions of the data store model for the case of time-sharing executions are shown in Table II. Overall, the ability of our fine-grain time-sharing Time Warp system to robustly provide opportunities for event preemptions (as shown in Figure 9) is reflected into performance improvements independently of the interaction pattern between application and platform software.

### 3.4. Performance Results with a Personal Communication System Model

As a third alternative workload, in this section we consider a personal communication system model, namely a real-world application that has already been used in a number of studies assessing optimizations in PDES platforms (see, e.g., Cingolani et al. [2015]). In this application, each simulation object models a wireless cell. We selected a total number of 1,024 cells organized into a hexagonal grid, each one managing 1,000 wireless channels, which provide coverage to mobile devices in a squared region. Cells are modeled in high fidelity, taking into account both interference across different channels within a same cell and power management upon call setup/handoff according to the results in Kandukuri and Boyd [2002]. Depending on the selected setup, this application allows recalculating fading coefficients and actual Signal-to-Interference Ratio (SIR) on the occurrence of specific events (e.g., the startup of a call) and also periodically (to account for, for example, changes of weather conditions in the coverage area). Also, the inter-arrival of calls to mobile devices residing in the coverage area can be configured, thus leading to different values for the wireless channels' utilization factor. This, in its turn, affects both memory and CPU demand by the simulation given that higher utilization factors lead to the need for keeping more records (stored on dynamically allocated buffers) for simulating the concurrently active calls in any cell, and also to more costly operations for scanning and (possibly) updating these records. As a final preliminary note, the interaction across the different simulation objects takes place upon the occurrence of a handoff of a mobile device involved in an ongoing communication, in which case the wireless channel at the source cell is released, and a new one is attempted to be reserved at the destination cell.

In our experimentation we set the average residual residence time in the current cell for a mobile device involved in an on-going call to 5min, while the average call

duration was set to 2min. Both these parameters have been set to follow exponential distributions. Also, we run this model with three different settings for the channel utilization factor, namely 25%, 50%, and 75%, determined by different call inter-arrival rates, with balanced workload on all the simulation objects, and with periodic recalculations of the fading coefficients of active channels. These settings give rise to variations of the average CPU requirement for simulation events from about  $70/80\mu\text{s}$  to about  $150\mu\text{s}$ .

For this study we still set the extra-tick interval to  $50\mu\text{s}$  and  $\alpha$  to the value  $-1$ , say, to its lower bound. With these settings the worker threads operating within the time-sharing Time Warp system keep on being registered on the `dev_extra_tick` device file for the whole lifetime of the simulation in all the tested configurations (say, for any value of the channel utilization factor). This leads to maximal exploitation of timer-interrupts for event preemptions. This choice is motivated by the fact that, unlike the data store model, the processing of an event in the personal communication system model leads to reduced interactions with platform level reversible implementations of memory allocation/deallocation services (since the number of buffer allocations/deallocations per event is much lower than the one characterizing the data store model). Hence, returns from platform mode can play a reduced role in triggering preemptions. On the other hand, compared to the data store model, the event processing routine shows a very different profile, much more based on floating point operations.

For this test-bed application we initially run a modified version, with the aim to assess the overhead imposed by the core facilities offered by the fine-grain time-sharing Time-Warp system. These facilities are (i) the support for managing contexts, and (ii) the delivery of extra-ticks to the PDES platform. This overhead study is somewhat complementary to the one associated with PHASE-A of the execution of the PHOLD benchmark, since in that phase we only assessed the cost for managing contexts. In fact, during that phase of the execution of PHOLD, the worker threads did not register themselves on the `dev_extra_tick` device file, hence no extra-ticks were delivered.

In order to assess the overhead by the aforementioned two facilities we run the personal communication system model by always enforcing a call to complete with the same wireless cell where it was originated. In this way, no interaction at all by the different simulation objects is ever generated, and events are processed by the worker threads always according to non-decreasing values of their timestamps. In such a scenario, the delivery of extra-ticks provides no revenue (since the event being processed will always represent the one with the highest priority bound to a given worker thread), just like the management of separate simulation object contexts (since no simulation object will be ever context-switched off the CPU while processing an event). Also, the absence of rollback in such scenarios allows us to assess the overhead by the two facilities with no interference by rollback management operations (which might lead to, e.g., changes in the locality of the execution due to the access to both checkpoints of the simulation object states and already processed event buffers). In Figure 10 we report the execution time values<sup>11</sup>—achieved with 32 worker threads—when excluding both the management of contexts and the delivery of extra-ticks, or when including these facilities. The former settings represent the common ones for Time Warp systems not embedding the support for fine-grain time-sharing execution of the simulation objects. The reported data show how the overhead by the core facilities enabling fine-grain time-sharing is very limited, except for 25% channel utilization factor, case in which it reaches 7%. In fact, as soon as the event granularity (say, the channel utilization factor) increases, we observe a decrease of the overhead, especially in relation to the management of contexts. This is somehow expected given that longer

---

<sup>11</sup>Still based on the average over 10 runs.

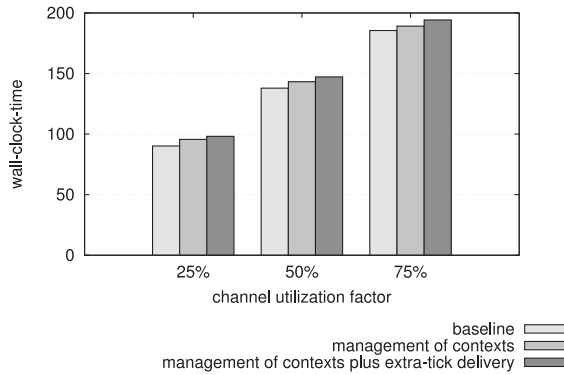


Fig. 10. Execution time with no cross-simulation-object event scheduling.

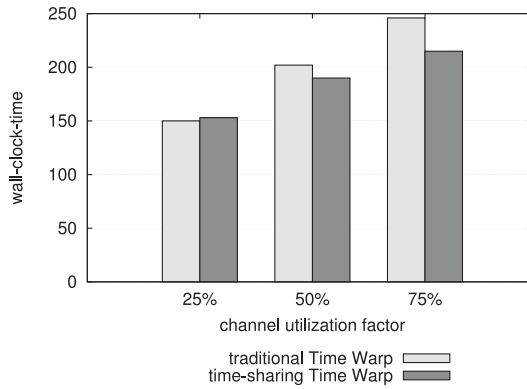


Fig. 11. Results with the personal communication system model.

running events lead to reduced frequency of context-switch operations across different simulation objects over time in scenarios with no preemptions.

In Figure 11 we show execution time results when reintroducing handoff events across cells, say, cross-simulation-object scheduling of events. With these settings, the performance gain provided by the fine-grain time-sharing Time Warp system, compared to the traditional Time Warp execution, increases when increasing the channel utilization factor. The gain is of the order of 7% for the case of utilization factor set to 50%, and of the order of 13% when the utilization factor is further increased up to 75%. For channel utilization factor set to 25% we observe no relevant gain from time-sharing, just because of the reduced potentiality of extra-ticks exploitation (given the reduced wall-clock-time required for processing events in this configuration). This trend is confirmed by data we report in Figure 12, showing the variation of the amount of event preemptions per execution time unit achieved while running in time-sharing mode for the different configurations of the channel utilization factor.

For completeness, we also report in Table III the corresponding execution times for the case of a serial execution of the same identical application code on top of a sequential scheduler based on the Calendar-Queue data structure [Brown 1988], which allows determining the speedup of the parallel runs—hence whether the reported data refer to competitive parallel performance.

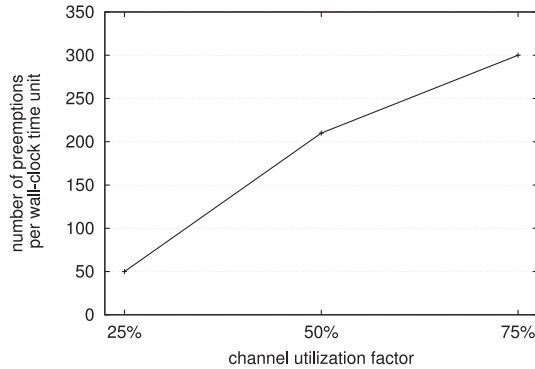


Fig. 12. Number of preemptions per wall-clock-time unit.

Table III. Performance of the Serial Simulator

	Channel Utilization Factor		
	25%	50%	75%
execution time (sec)	3500	5145	7610

#### 4. RELATED WORK

*Approaches Focusing on Performance versus Timer-interrupt Trade-offs.* In the wide area of High Performance Computing (HPC), some literature studies exist on the relation between performance and timer-interrupt frequency. The common idea underlying most of the performance optimization proposals is that the lower the timer-interrupt frequency, the better the final delivered performance [De et al. 2007; Ferreira et al. 2008; Seelam et al. 2010]. The extremism of this approach led to defining *tick-less* operating systems—characterized by extremely reduced timer-interrupt frequency—as the best configuration for hosting HPC applications. However, these studies have been tailored to the case of non-speculative processing, where the work carried out by any thread running on whichever CPU-core is ever useful. In this context there is no need to change the thread execution flow (e.g., periodically) to optimize synchronization dynamics in terms of reduction of wasted computation, which is instead a major objective when dealing with speculative Time Warp systems. Also, the above studies have been tailored to evaluate the effects of the variation of the timer-interrupt frequency in scenarios where the management of the timer-interrupt is still based on the native rules applied by the operating system kernel. In other words, the above proposals have been aimed at simply configuring the timer-interrupt behavior—limited to its frequency—in HPC applications, not at introducing ad-hoc software modules for exploiting timer events, which is instead the approach we followed. In fact, our proposal puts in place a special (and lightweight) mechanism for handling timer-interrupts. Overall, our approach is completely different from the one dealt with by those literature studies, in terms of both reference scenario (speculative versus non-speculative processing) and architectural impact on the system organization.

In the context of speculative PDES systems, the only work we are aware of that deals with the relation between performance and timer-interrupt configuration is the one by Carothers [2002]. Here, the author proposes an approach, which is opposite to ours, where Time Warp threads are allowed to take CPU control for longer periods (thus being not interrupted for a while) to be able to fully execute a simulation model with no interference by other workloads, and to deliver the output in real-time. This

solution is still along the path of tick-less operating systems, with the difference that the tick-less behavior is triggered on-demand (namely whenever a time-critical parallel simulation needs to be executed), hence it is not a static configuration of the underlying operating system. Our approach is fully orthogonal to this one, because our target is the reduction of wasted computation, thanks to an appropriate periodic variation of the control flow along Time Warp threads. Also, while the proposal by Carothers [2002] is based on reserving the computing capacity for Time Warp programs—thus excluding the possibility for other tasks to be run on the system for a while—in our approach we do not create any bias in the usage of the computing system by Time Warp threads and other kinds of threads. We only allow the Time Warp threads to see their own ticks as partitioned into sub-intervals (with proper control flow management at the end of each sub-interval).

*Approaches Targeting Preemptive Rollback.* Our time-sharing Time Warp proposal supports preemptive rollback, a topic that has been somehow studied in literature, mainly in Das et al. [1994] and Santoro and Quaglia [2005]. The solution in Das et al. [1994] targets parallel simulation on shared-memory machines, and is based on direct manipulation of the event list of the recipient simulation object by the thread along which the generation of a new event is handled. With this solution, the sender thread is able to determine the current simulation time of the recipient simulation object and whether any message/anti-message being sent to that object violates causality. If this is the case, then the sender thread notifies the violation to the thread handling the recipient object, which is done to timely interrupt any in-progress activity to execute rollback operations. Our solution is different, since it does not rely on cross-thread signaling. Also, in our approach, any Time Warp thread is allowed to change its current flow (and dynamically dispatch a different simulation event, or simulation object, after preempting the last dispatched one) independently of the materialization of a causality violation, but rather when any need arises to process a higher priority task, bound to a simulation object possibly different from the currently running one. This is done to reduce the likelihood of future rollback generation, not only to react via preemption to an already materialized causality violation. This is basically due to the fact that our fine-grain time-sharing Time Warp system is not limited to the support for preemptive rollback.

As for the preemptive rollback approach in Santoro and Quaglia [2005], it is suited for distributed memory systems while we deal with shared memory multi-core machines. Also, it is based on polling, and the polling code to periodically verify causal consistency of the current event needs to be nested in the application code by the programmer. Instead, our proposal is fully transparent, and exploits back from platform mode and timer-interrupt events, rather than polling. Finally, similar to Das et al. [1994], the solution in Santoro and Quaglia [2005] does not cope with control flow variations associated with the dynamic generation of higher priority events (namely with timestamps lower than that of the event being executed along the thread) that do not directly give rise to a causality violation.

*Approaches Based on Operating Systems Concepts.* Dual-mode execution in Time Warp systems, which we exploit in our proposal, has been also investigated in Pellegrini and Quaglia [2014]. In this proposal, when the worker thread runs in application mode, only a sub-portion of the whole address space is made accessible, namely the sub-portion keeping the memory layout of the CPU-dispatched simulation object. Any access to the state of another object generates a trap that gives control back to the platform code, which actuates proper thread synchronization mechanisms to allow cross-state processing of the events. Unlike Pellegrini and Quaglia [2014], the solution

we present is tailored to variations of the control flow to react to the generation of higher priority simulation events or tasks (such as rollbacks to be processed). Still, we retain application transparency just like [Pellegrini and Quaglia 2014].

Our proposal is clearly related to the work in Jefferson et al. [1987], where Time Warp is instantiated as a special-purpose operating system destined to host discrete event applications to be executed according to the speculative processing paradigm. The core difference between what we are presenting and the proposal in Jefferson et al. [1987] lies in that such an approach uses preemption only in case of causality errors affecting the currently-dispatched simulation event. Rather, we exploit preemption anytime a higher priority task needs to be processed, independently of the actual materialization of causality errors. Thus our solution also tends to anticipate the generation of causality errors.

Our proposal has also relations with recent approaches based on operating system scheduling to support virtual time synchronized advancement in emulation/simulation scenarios (see, e.g., Lamps et al. [2015, 2014] and Yoginath et al. [2012]). These solutions provide scheduling policies of Linux Containers (LXCs) or Virtual Machines (VMs) allowing the emulated components to adjust their speed of operation to align it to the advancement of simulation time. This is typically achieved by scaling up/down the CPU capacity assigned to the different LXCs or VMs. Our solution is orthogonal to these approaches, since we do not work at the level of the operating system scheduling policy. Rather, we customize the operating system management of timer-interrupts—to deliver them with fine granularity and at low cost to the speculative PDES platform—in order to enable optimized CPU assignment to multiple simulation objects run on top of a same thread within a fine-grain time-sharing scheme.

*Approaches Directly Targeting the Reduction of Rollback.* Given that our core target is the reduction of the incidence of causality errors, our proposal is naturally related to literature solutions directly targeting rollback reduction in speculative PDES. We can roughly classify these works in two main categories: (a) the ones based on balanced resource usage (see, e.g., Carothers and Fujimoto [2000], Choe and Tropper [1999], Glazer and Tropper [1993], and Vitali et al. [2012]) and (b) the ones based on bounded optimism (see, e.g., Dickens et al. [1996] and Srinivasan and Reynolds [1998]). In the former case, the target is the one of reducing the skew in the advancement of simulation time at the different simulation objects, which is typically achieved via simulation objects' periodical migration (for balanced execution) across the Time Warp worker threads. These proposals act as long term planners for CPU usage, and do not entail capabilities of reacting to punctual variations of the priority of the events, as instead we do via preemptive CPU reassignment. Overall, we can see the two approaches as orthogonal to each other, hence being ideally combinable. Finally, the solutions based on bounded optimism opt for artificially delaying the execution of events within the speculative processing scheme with the aim at increasing the likelihood of performing useful (not eventually rolled back) work. Some proposal (see, e.g., Srinivasan and Reynolds [1998]) can even dynamically select per-event delays, thus attempting to reduce the incidence of rollback on a fine grain basis. We retain this same capability, but we still fully exploit the available computing power, since we admit truly speculative preemptive event processing, with no artificial delay. Although different in spirit, we can still think of these two approaches as orthogonal and potentially usable in synergy.

## 5. CONCLUSIONS

It is typical that PDES platforms process simulation events in non-preemptive manner. For the case of Time Warp PDES systems, which exploit speculative processing and rollback techniques for causality maintenance, a preemptive approach would provide

the possibility to dynamically reassign the CPU to events (or tasks, such as rollback operations) standing in the past logical time of the currently processed event. This would allow for reducing the incidence of causality errors along the speculated execution path and to more promptly react to the actual generation of the errors. To cope with this issue, we have presented a fine-grain time-sharing version of a Time Warp system, which makes systematic use of event preemption just for the purpose of making the system run, at any time, those events/tasks that are dynamically determined to have the highest priority—they refer to past logical time values compared to the last CPU-dispatched ones. Our proposal is targeted at multi-core machines and Linux/x86-64 platforms. We integrated the fine-grain time-sharing Time Warp architecture, including the ad-hoc Linux module supporting timer-interrupt based preemptions, within an open source speculative PDES platform. Further, we have reported experimental data supporting the effectiveness of our proposal. Indications on how to configure the presented fine-grain time-sharing Time Warp system, in relation to core parameters driving its internal logic, have also been provided, which should favor fruitful usage of our solution with workloads aside the ones used in our experiments.

## REFERENCES

- F. Antonacci, A. Pellegrini, and F. Quaglia. 2013. Consistent and efficient output-streams management in optimistic simulation platforms. In *Proceedings of the ACM SIGSIM Conference on Principles of Advanced Discrete Simulation*. ACM, 315–326.
- P. D. Barnes-Jr., C. D. Carothers, D. R. Jefferson, and J. M. LaPre. 2013. Warp speed: Executing time warp on 1,966,080 cores. In *Proceedings of the 1st ACM SIGSIM Conference on Principles of Advanced Discrete Simulation*. 327–336.
- R. Brown. 1988. Calendar queues: A fast  $O(1)$  priority queue implementation for the simulation event set problem. *Commun. ACM* 31, 1220–1227.
- C. D. Carothers. 2002. Xsim: Real-time analytic parallel simulations. In *Proceedings of the 16th Workshop on Parallel and Distributed Simulation (PADS'02)*. 27–34.
- C. D. Carothers and R. M. Fujimoto. 2000. Efficient Execution of Time Warp Programs on Heterogeneous, NOW Platforms. *IEEE Trans. Parallel Distrib. Syst.* 11, 3, 299–317.
- M. Choe and C. Tropper 1999. On learning algorithms and balancing loads in time warp. In *Proceedings of the 13th Workshop on Parallel and Distributed Simulation*. Springer Verlag, 101–108.
- D. Cingolani, A. Pellegrini, and F. Quaglia. 2015. Transparently mixing undo logs and software reversibility for state recovery in optimistic PDES. In *Proceedings of the 3rd ACM Conference on SIGSIM—Principles of Advanced Discrete Simulation*. 211–222.
- S. R. Das, R. M. Fujimoto, K. Panesar, D. Allison, and M. Hybinette. 1994. GTW: A time warp system for shared memory multiprocessors. In *Proceedings of the 26th conference on Winter simulation (WSC'94)*. Society for Computer Simulation International, 1332–1339.
- P. De, R. Kothari, and V. Mann. 2007. Identifying sources of operating system jitter through fine-grained kernel instrumentation. In *Proceedings of the 2007 IEEE International Conference on Cluster Computing*. 331–340.
- P. M. Dickens, D. M. Nicol, P. F. R. Jr., and J. M. Duva. 1996. Analysis of bounded time warp and comparison with YAWNS. *ACM Trans. Model. Comput. Simul.* 6, 4, 297–320.
- K. B. Ferreira, P. Bridges, and R. Brightwell. 2008. Characterizing application sensitivity to os interference using kernel-level noise injection. In *Proceedings of the 2008 ACM/IEEE Conference on Supercomputing (SC'08)*. IEEE Press, Piscataway, NJ. 19:1–19:12.
- R. M. Fujimoto. 1990a. Parallel discrete event simulation. *Commun. ACM* 33, 10, 30–53.
- R. M. Fujimoto. 1990b. Performance of Time Warp under synthetic workloads. In *Proceedings of the Multi-conference on Distributed Simulation*. Society for Computer Simulation, 23–28.
- D. W. Glazer and C. Tropper. 1993. On process migration and load balancing in time warp. *IEEE Trans. Parallel Distr. Syst.* 4, 3, 318–327.
- HPDCS Research Group. 2012. ROOT-Sim: The ROME OpTimistic Simulator—v 1.0. Retrieved from <http://www.dis.uniroma1.it/hpdc/ROOT-Sim/>.
- S. Jafer, Q. Liu, and G. A. Wainer. 2013. Synchronization methods in parallel and distributed discrete-event simulation. *Simul. Model. Pract. Theory* 30, 54–73.



- D. R. Jefferson. 1985. Virtual Time. *ACM Trans. Program. Lang. Syst.* 7, 3, 404–425.
- D. R. Jefferson, B. Beckman, F. Wieland, L. Blume, M. D. Loreto, P. Hontalas, P. Laroche, K. Sturdevant, J. Tupman, L. V. Warren, J. J. Wedel, H. Younger, and S. Bellenot. 1987. Distributed simulation and the time warp operating system. In *Proceedings of the Symposium on Operating Systems (SOSP'87)*. 77–93.
- S. Kandukuri and S. Boyd. 2002. Optimal power control in interference-limited fading wireless channels with outage-probability specifications. *IEEE Trans. Wireless Commun.* 1, 1, 46–55.
- J. Lamps, V. Adam, D. M. Nicol, and , M. Caesar. 2015. Conjoining emulation and network simulators on linux multiprocessors. In *Proceedings of the 3rd ACM-SIGSIM Conference on Principles of Advanced Discrete Simulation*. 113–124.
- J. Lamps, D. M. Nicol, and M. Caesar. 2014. Timekeeper: A lightweight virtual time system for linux. In *Proceedings of the 2nd ACM-SIGSIM Conference on Principles of Advanced Discrete Simulation*. 179–186.
- Y.-B. Lin and E. D. Lazowska. 1991. Processor scheduling for Time Warp parallel simulation. In *Proceedings of the 23rd SCS Multiconference on Advances in Parallel and Distributed Simulation*. IEEE Computer Society, 11–14.
- D. M. Nicol and X. Liu. 1997. The dark side of risk (what your mother never told you about time warp). In *Proceedings of the 11th Workshop on Parallel and Distributed Simulation (PADS'97)*. 188–195.
- A. Pellegrini and F. Quaglia. 2014. Transparent multi-core speculative parallelization of DES models with event and cross-state dependencies. In *SIGSIM Principles of Advanced Discrete Simulation (SIGSIM-PADS'14)*. 105–116.
- A. Pellegrini and F. Quaglia. 2015. NUMA time warp. In *Proceedings of the 3rd ACM Conference on SIGSIM-Principles of Advanced Discrete Simulation*. 59–70.
- A. Pellegrini, R. Vitali, and F. Quaglia. 2011. The ROme OpTimistic Simulator: Core internals and programming model. In *Proceedings of the 4th International ICST Conference on Simulation Tools and Techniques*. *Proceedings of the 4th ICST Conference of Simulation Tools and Techniques (SIMUTools'11)*.
- A. Pellegrini, R. Vitali, and F. Quaglia. 2015. Autonomic state management for optimistic simulation platforms. *IEEE Trans. Parallel Distrib. Syst.* 26, 6, 1560–1569.
- F. Quaglia and V. Cortellessa. 2002. On the processor scheduling problem in time warp synchronization. *ACM Trans. Model. Comput. Simul.* 12.
- R. Rönngren, M. Liljenstam, R. Ayani, and J. Montagnat. 1996. Transparent incremental state saving in Time Warp parallel discrete event simulation. In *Proceedings of the 10th Workshop on Parallel and Distributed Simulation*. IEEE Computer Society, 70–77.
- A. Santoro and F. Quaglia. 2005. Software supports for event preemptive rollback in optimistic parallel simulation on myrinet clusters. *J. Interconnect. Netw.* 6, 4, 435–457.
- S. Seelam, L. L. Fong, A. N. Tantawi, J. Lewars, J. Divirgilio, and K. Gildea. 2010. Extreme scale computing: Modeling the impact of system noise in multicore clustered systems. In *Proceedings of the 24th IEEE International Symposium on Parallel and Distributed Processing (IPDPS'10)*. 1–12.
- S. Srinivasan and Jr., P. Reynolds. 1998. Elastic time. *ACM Trans. Model. Comput. Simul.* 8, 2, 103–139.
- R. Toccaceli and F. Quaglia. 2008. DyMeLoR: Dynamic Memory Logger and Restorer library for optimistic simulation objects with generic memory layout. In *Proceedings of the Workshop on Principles of Advanced and Distributed Simulation*. IEEE Computer Society, 163–172.
- R. Vitali, A. Pellegrini, and F. Quaglia. 2012. Load sharing for optimistic parallel simulations on multi core machines. *SIGMETRICS Perform. Eval. Rev.* 40, 3, 2–11.
- J. Wang, D. Jagtap, N. B. Abu-Ghazaleh, and D. Ponomarev. 2014. Parallel discrete event simulation for multi-core systems: Analysis and optimization. *IEEE Trans. Parallel Distrib. Syst.* 25, 6, 1574–1584.
- D. West and K. Panesar. 1996. Automatic incremental state saving. In *Proceedings of the Workshop on Parallel and Distributed Simulation*. IEEE Computer Society, 78–85.
- S. B. Yoginath, K. S. Perumalla, and B. J. Henz. 2012. Taming wild horses: The need for virtual time-based scheduling of vms in network simulations. In *Proceedings of the 20th IEEE International Symposium on Modeling, Analysis and Simulation of Computer and Telecommunication Systems*. 68–77.

Received November 2015; revised May 2016; accepted October 2016