

ASSESSING LOAD-SHARING WITHIN OPTIMISTIC SIMULATION PLATFORMS

Roberto Vitali
Alessandro Pellegrini
Francesco Quaglia

High Performance and Dependable Computing Systems (HPDCS) Research Group
DIS - Sapienza, University of Rome

ABSTRACT

The advent of multi-core machines has led to the need for revising the architecture of modern simulation platforms. One recent proposal we made attempted to explore the viability of load-sharing for optimistic simulators run on top of these types of machines. In this article, we provide an extensive experimental study for an assessment of the effects on run-time dynamics by a load-sharing architecture that has been implemented within the ROOT-Sim package, namely an open source simulation platform adhering to the optimistic synchronization paradigm. This experimental study is essentially aimed at evaluating possible sources of overheads when supporting load-sharing. It has been based on differentiated workloads allowing us to generate different execution profiles in terms of, e.g., granularity/locality of the simulation events.

1 INTRODUCTION

Multicore/multiprocessor machines have become a wide-spread reality, and the possibility to get access to these platforms at low costs is growing up to a greater extent, placing the need for a systematic approach concerning computational power usage, which is a key factor in high performance simulations.

Historically, high performance simulation platforms have been based on the partitioning of the simulation model into several distinct objects, handled by Logical Processes (LPs) (Fujimoto 1990), which are allowed to concurrently execute simulation events on, e.g., clusters of machines. On the other hand, the typical organization of the underlying simulation-kernel layer has been based on a multi-process paradigm, where each process runs in single-threaded mode and schedules its locally hosted LPs for event execution according to a time-interleaved scheme. For this type of organization, (dynamic) unbalance of the workload, associated with (dynamic) changes of the computational power demand across the LPs, has been traditionally tackled via load-balancing approaches targeted at migrating LPs from the overloaded simulation-kernel instances to the underloaded ones. Examples of this type of approaches, showing differentiated levels of transparency towards the application-level software, can be found in (D'Angelo and Bracuto 2009; Glazer and Tropper 1993; Carothers and Fujimoto 2000; Peluso, Didona, and Quaglia 2011)

On the other hand, the multicore/multiprocessor organization characterizing modern hardware platforms offers new potentialities that could be fully exploited by abandoning the (simple) single-threaded programming approach in favor of multi-threaded programming paradigms. Along this direction, a recent achievement presented in (Jagtap, Abu-Ghazaleh, and Ponomarev 2012) has shown how the reshuffle of multi-process optimistic simulation platforms to multi-threaded versions, particularly for the case of the ROSS open source platform (Carothers, Bauer, and Pearce 2000), can provide noticeable benefits in terms of performance thanks to the optimization of cross simulation-kernel communication. In particular, the conjunction of (a) the usage of shared-memory to support message passing and (b) the reliance on a coherent view of virtual-addressing across the threads operating within a same process (each one implementing an instance of the simulation-kernel), has led to a significant reduction of the amount of data to be copied when transmitting/receiving messages. However, the architectural organization proposed in (Jagtap,

Abu-Ghazaleh, and Ponomarev 2012) still relies on the traditional view according to which one instance of the simulation-kernel runs as a single thread, thus having (at most) a single CPU-core as the computational power assigned to it. Another attempt has been presented in (Li-li, Ya-shuai, Yi-ping, Shao-liang, and Ling-da 2011) where a global scheduling mechanism (based on a centralized event queue) is exploited in order to assign simulation work (i.e., event processing on different LPs) to different active threads within the same platform. However, this approach may suffer from reduced scalability, thus being suited for contexts where a reduced number of CPU-cores is available (it has been evaluated with no more than 8 CPU-cores).

A more recent advance we have proposed in (Vitali, Pellegrini, and Quaglia 2012) provides a paradigm shift in the design of optimistic simulation-kernels by introducing a reference architecture based on a symmetric multi-threaded approach. Here, each instance of the optimistic simulation-kernel is able to run multiple worker threads, which can take care of the execution of whichever locally hosted LP. According to this organization, the number of worker threads within each kernel instance can be dynamically scaled up/down in a seamless manner, which opens the possibility for a new approach to the usage of computational power. Specifically, for dynamically changing computational requirements across the LPs, re-balanced runs can be achieved by scaling down the amount of worker threads operating within under-loaded kernel instances and scaling up the number of worker threads operating within over-loaded instances. Hence, *load-sharing* policies are actually pursued, with the meaning that the whole computational load is shared across all the available CPU-cores which are dynamically bind to one kernel instance or the other depending on the aforementioned changes in the application requirements (and associated amounts of worker threads per kernel instance). Overall, with this approach LPs' migration facilities are no more mandatorily required in order to optimize resource usage in multicore/multiprocessor contexts, since the load-sharing approach exploits the orthogonal concept of computational power migration at the simulation-kernel level. This can not only provide different (hopefully better) tradeoffs between housekeeping overhead and resource exploitation for productive work, but can also help application transparency since the application programmer will be no longer requested to provide modules for, e.g., relocating the LPs across different simulation-kernel processes (hence across different address spaces), which is generally not transparently supported except for a few advanced state management architectures (Peluso, Didona, and Quaglia 2011).

Early experimental data have shown how the load-sharing architecture depicted in (Vitali, Pellegrini, and Quaglia 2012), and implemented within the ROME OpTimistic Simulator (ROOT-Sim) package (Quaglia, Pellegrini, and Vitali 2011), which is an open source general purpose simulation platform adhering to the optimistic synchronization paradigm, can provide low overhead (e.g. for synchronizing worker threads executed within a same simulation-kernel instance) while allowing performance optimizations in the context of dynamically changing workloads.

In this article we complement such an early study by reporting an extensive experimental characterization of the load-sharing architecture when considering differentiated application-level settings. This characterization is essentially aimed at determining potential sources of overheads and/or scalability limitations. In particular, we report a set of measures that, in a complementary manner to the coarse grain data reported in (Vitali, Pellegrini, and Quaglia 2012), allow a fine grain analysis of the actual run-time dynamics of the load-sharing architecture in comparison to those achievable with a traditional multi-process organization of the simulation platform.

The remainder of this paper is organized as follows. In Section 2 we provide the reader with an overview of the load-sharing architecture. A discussion on the structure of the experimental assessment, and on related target parameters to be observed, is provided in Section 3. The actual experimental results are presented in Section 4.

2 OVERVIEW OF THE LOAD-SHARING ARCHITECTURE

As hinted, the load-sharing architecture presented in (Vitali, Pellegrini, and Quaglia 2012) is based on a symmetric multi-threaded approach where each worker thread running within each single simulation-kernel

instance has the ability to execute in both application- and kernel-mode, and can control and take care of the execution of whichever locally hosted LP. This type of approach has relations with what happens in operating systems targeted at multicore/multiprocessor machines, where the CPU-scheduler controlling a specific CPU-core is generally allowed to dispatch whichever ready-to-run thread. On the other hand, different LPs may mutually issue interactions, e.g., via event scheduling services, that are actually supported by entering the kernel-mode along the thread taking care of running the LP that issues the interaction.

However, while the execution in user-mode intrinsically relies on data partitioning, since, in compliance with the original specification of the optimistic synchronization protocol, namely the Time Warp protocol (Jefferson 1985), each LP operates on private (per-LP) data structures implementing its current state, this is not the same when worker threads operate in kernel-mode. In particular, kernel-mode execution may require that a worker thread taking care of running LP_a needs to access the kernel level meta-data associated with whichever locally hosted LP_b , e.g., the event-queue associated with LP_b , given that LP_b may figure out as the recipient of an issued interaction, such as a newly scheduled event from LP_a . If a different worker thread is allowed to concurrently operate in kernel-mode exactly on that queue, which is the case for the symmetric organization, a critical section must be guaranteed, requiring some form of synchronization. We note that, depending on proper dynamics related to the specific simulation model, such a situation may occur frequently, thus imposing the need for a properly tailored management of synchronization aspects within the symmetric multi-threaded organization.

In order to prevent kernel-level synchronization phases from becoming a bottleneck, in our load-sharing proposal we have devised that each cross-LP interaction logically represents an interrupt, which does not get atomically finalized upon its acceptance, hence avoiding the need for acquiring locks on target data structures according to a possibly adverse timing (e.g. when the data structure is already locked for the finalization of a concurrently issued interaction). Instead, the finalization of the interaction takes place by adhering to a top/bottom-half scheme, resembling the scheme typically used for the implementation of interrupt-drivers in operating systems targeted at multicore/multiprocessor machines.

A graphical representation of the outcoming architecture is provided in Figure 1. When an interaction is issued, a lightweight top-half module is executed which only inserts a bottom-half task into a queue associated with the destination LP, in order to allow finalization at a later instant of time, more conveniently wrt synchronization. We note that an interaction to be treated according to the top/bottom-half scheme may be issued in three different scenarios:

- When an LP runs in forward mode and produces events to be destined to other locally-hosted LPs.
- When an LP runs in rollback mode (i.e., the kernel layer is currently recovering its state due to the occurrence of a causality violation). In this case an interaction associated with an anti-event might be destined to some locally hosted LP (An anti-event, also known as anti-message, is a *negative* copy used to annihilate a previously issued interaction, namely an already scheduled event. Anti-events are used to propagate the effects of causality errors across the LPs by retracting events scheduled during the causal inconsistent portion of the simulation.).
- When the messaging layer locally notifies an interaction, namely an event or an anti-event, whose source LP is hosted by a different simulation-kernel instance.

Basically, our approach can be supported by relying on a spin-lock array, named `LP_LOCKS`, having one entry for each LP hosted by the multi-threaded simulation-kernel. `LP_LOCKS[i]` is used to implement a fast critical section for the access to the bottom-half queue associated with the i -th LP hosted by the kernel, either for inserting a new bottom-half task to be eventually flushed, or for taking care of unlinking the current chain, in order to flush the pending bottom-halves.

With this organization, each worker thread can (in principle) take care of flushing the pending interactions currently recorded within the bottom-half queue of any LP, which can be done at convenient time instants, namely when no other worker thread is already doing this same job. In this way, wait-phases for exclusive accesses/updates to the kernel-level data structures associated with whichever LP get eliminated. The

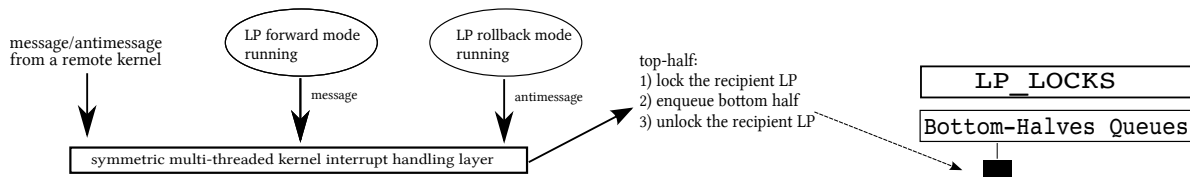


Figure 1: Outline of the top/bottom-halves architecture.

support for this type of operation would simply consist in an additional per-LP spin-lock, accessible in non-blocking mode. Successful non-blocking access would give rise to the possibility for the worker thread to operate the flush of the bottom-half queue for that LP. On the other hand, unsuccessful non-blocking access would simply tell the worker thread to skip taking care of that LP.

Additionally, in order to improve the efficiency of the caching hierarchy, we have proposed a temporary binding mechanism that assigns each worker thread the responsibility to flush the bottom-half queues of only a subset of the locally hosted LPs. Such a thread also has the responsibility to dispatch its bound LPs for event execution in time-interleaved mode. In this way, direct cache contention while handling user-level or kernel-level execution of the LPs may only occur for LPs bound to the same worker thread (hence operating on top of the same CPU-core), depending on the interleave of the operations performed by the worker thread (e.g., forward execution operations) related to one or the other of these LPs. In other words, cache contention across different worker threads (e.g., mutual cache invalidation) may only occur in relation to the access/manipulation of kernel-level synchronization data structures, such as the bottom-half queues. In (Vitali, Pellegrini, and Quaglia 2012), rules for periodically selecting the LPs to be bind to a specific worker thread are defined, together with some rules adopted to reassign the computational power (i.e., the CPU-cores) to the different simulation-kernel instances by scaling up/down the amount of worker threads per instance. These are particularly suited for the case of dynamic workloads.

3 OVERVIEW OF THE EXPERIMENTAL ASSESSMENT

In this section we provide an overview of our experimental assessment of the load-sharing architecture. We will present data related to a specific implementation of this architecture which has been integrated within the ROOT-Sim package (Quaglia, Pellegrini, and Vitali 2011). As a consequence, beyond the evaluation of the general mechanisms at the base of the load-sharing approach, we will also focus on specific effects due to the integration of load-sharing within ROOT-Sim since this imposes some constraints on run-time dynamics properly related to specific ROOT-Sim's subsystems.

All the parameters that will be object of the experimental assessment will be discussed in this section, by also motivating why they have been selected in the analysis. To this end, we briefly recall the structure and capabilities of ROOT-Sim, so to provide the basis for easing the comprehension of the discussion. As a final preliminary note, our reference architecture for the assessment will be represented by the original multi-process version of ROOT-Sim.

3.1 The Architecture of ROOT-Sim

ROOT-Sim is an open source C/MPI-based simulation package targeted at POSIX systems, which implements a general-purpose parallel/distributed simulation environment relying on the optimistic (i.e., rollback-based) synchronization paradigm. It offers a very simple programming model relying on the classical notion of simulation-event handlers (both for processing events and for accessing a committed and globally consistent state image upon GVT (GVT - Global Virtual Time - represents the commitment horizon. No causality violation can even occur for processed events whose timestamp falls before the current GVT value. GVT updates typically trigger memory recovery procedures, e.g., of obsolete state logs.) calculations), to be implemented according to the ANSI-C standard, and transparently supports all the services required to parallelize the execution.

As for management and recoverability of LPs' state, which are crucial aspects for the design of effective optimistically synchronized environments, two main architectural approaches have been adopted. First, dynamic memory allocation/release requests by the application, performed via the standard `malloc` library, are *hooked* by the kernel and redirected to a wrapper. Second, the simulation platform is "*context-aware*", i.e., it has an internal state which distinguishes whether the current execution flow belongs to the application-level code or the platform's internals. In the former case, the hooked calls are redirected via the wrapper to an internal Memory Map Manager (called DyMeLoR), which handles per-LP allocation/deallocation operations by maximizing memory locality of the state layout for each single LP, and by maintaining meta-data which identify the state memory map and make it correctly recoverable to past values (Tocaceli and Quaglia 2008).

Concerning GVT calculation, ROOT-Sim relies on an optimized asynchronous approach based on a message acknowledgment scheme used to solve the well-known transient message problem. Within this scheme, each kernel instance keeps track of all the messages sent to the other instances in an aggregate manner (i.e., via counters). Also, to reduce the communication overhead, each instance sends cumulative acknowledgment messages according to a window-based approach. Finally, to overcome the simultaneous reporting problem, each kernel instance temporarily stops sending acknowledgment messages during the execution of the GVT protocol.

ROOT-Sim also supports a very peculiar service that, once a new GVT value is available, transparently rebuilds a Committed and Consistent Global Snapshot (CCGS), formed by a collection of individual LPs' states (Cucuzzo, D'Alessio, Quaglia, and Romano 2007). This occurs via update operations applied to local committed checkpoints of individual LPs so to eliminate mutual dependencies among the final-achieved state values. Once the CCGS is built, each LP gains control via an ad-hoc callback within the API, by also having access to the copy of its state image belonging to the CCGS. Such a service can support, e.g., termination detection schemes based on global predicates evaluated on a committed and consistent global snapshot.

3.2 Evaluating the Effects on Caching and Memory Accesses

It is clear that the internal organization of the load-sharing architecture can impact locality, which may give rise to variations of the effectiveness of the caching hierarchy. This may occur, e.g., due to the presence of kernel-level data structures shared across multiple threads, which are instead avoided in traditional multi-process platforms. In addition, we note that concurrent accesses can produce an impact on bus contention, due to locking operations needed for synchronizing threads' execution in critical sections. In order to provide quantitative data related to potential variations of the execution locality and its effects, we have decided to focus on three parameters:

- The latency for taking a checkpoint of the LP state.
- The latency for reloading a previously taken checkpoint in case of rollback.
- The event execution latency.

The first two parameters are associated with memory intensive operations, since each log or restore operation entails spanning across the LP's state or the log buffer in read mode. They represent therefore good metrics for determining how efficiently these read operations are supported thanks to the effects of the caching hierarchy. On the other hand, the event execution latency is a reflection of the locality expressed by the application, and of how well such a locality is supported via the caching system.

In addition, we have decided to measure scheduling operations' latency in order to assess the effects of multi-threading on data structures which are accessed sparsely. We have explicitly decided to rely on a $O(n)$ Shortest Timestamp First (STF) scheduler, which determines the next event to be processed by going over LPs' input queues for identifying the pending event associated with the minimum timestamp. We

consider this to be a significant measure, representative of the dynamics proper of a large set of operations which are essential in a simulation platform, such as queues scanning and log-chains traversing.

3.3 Evaluating the Impact on MPI Operations

In case of interactions between LPs hosted by different kernel instances, instead of relying on the top/bottom-half scheme, messages are directly provided in input to the MPI layer. Given that MPI does not support multi-threading, accesses have been serialized by exploiting again critical sections supported via spin-locking. The same has been done for probing MPI and issuing message receive operations by the worker threads, which are ultimately reflected in the execution of a top-half module (see again Figure 1).

Clearly this approach may induce delays on the worker threads when compared to the multi-process scheme. However, once fixed the total amount of threads (and hence of CPU-cores) running the simulation platform, in either multi-process or multi-threaded mode, there is a non-zero likelihood that two LPs hosted by different kernel instances within the multi-process organization are hosted by the same kernel instance when running in multi-threaded mode. Hence the mutual interactions between these LPs, if any, will not require passing via MPI. This likely leads to a reduced amount of interactions to be handled via MPI, since some interactions will be locally treated at the level of the top/bottom-half subsystem. Overall, to account for the above effects we have decided to evaluate:

- The time spent while interacting with the MPI layer.
- The time spent while managing the data structures supporting interactions via the top/bottom-half architecture, which, we recall, might include the time spent while synchronizing concurrent worker threads within the access to bottom-half queues.

A joint analysis of the two above parameters would allow understanding dynamics related to the actual handling of the interactions across the LPs involved within the simulation model.

We note that a possible approach to reduce the synchronization costs in the load-sharing architecture while interacting with the MPI layer would be represented by message aggregation. In fact, messages (namely events and anti-events) destined to remote multi-threaded kernel instances could be aggregate into local buffers and only periodically sent towards the destination. This can reduce the frequency of interactions with the MPI layer, thus favoring a reduction of the overhead when considering the case of synchronized accesses to MPI by multiple worker threads. Given that we have not yet embedded a similar optimization within ROOT-Sim, for what concerns the interaction with MPI, the experimental assessment can be related to a kind of worst case architectural configuration.

3.4 Evaluating the Impact on GVT and Global Snapshot Operations

In the symmetric multi-threaded version of ROOT-Sim, the GVT subsystem has been modified in order to account, within the global reduction determining the new GVT value, for the timestamps of events/anti-events that have not yet been reflected into the event queues of the recipient LPs due to the fact that they are still pending within bottom-half queues. These events/anti-events represent a sort of in-transit information, exhibiting similarities (and hence requiring similar management approaches) with traditional in-transit messages traveling via the messaging subsystem (MPI in our case) across different kernel instances.

Beyond the above issue, another relatively significant intervention while integrating the load-sharing approach within ROOT-Sim is related to the CCGS subsystem. As hinted, this subsystem is in charge of reconstructing, upon GVT calculations, committed and consistent global states, formed by collections of individual LPs' states. These individual states are then passed in input to an application level callback where the programmer is allowed to inspect the committed computation results. In the original multi-process version of ROOT-Sim, each active thread, individually representing an active kernel instance, is allowed to process those callbacks since they are intrinsically sequentialized along the execution of that same thread. Instead, for the symmetric multi-threaded organization, the active worker threads are not all allowed to do

this same job since this would lead to inconsistencies on the content of the (default) file used for tracing the output on each kernel instance. As a consequence, we have decided to synchronize all the worker threads operating within the same kernel instance in such a way to allow a single worker thread to run CCGS facilities. This reduces the computational power of the load-sharing architecture during the phases where the CCGS protocol is run. Hence we have decided to report in the assessment the latency observed when running GVT plus CCGS protocols upon committing a new portion of the simulation in order to quantify this phenomenon.

3.5 Evaluating the Effects on the Overall Rollback Pattern

Since a rollback happens upon receiving an out-of-order event to be executed, this can more likely arise for larger gaps between different LPs' local clocks possibly caused by a different workload being processed across different simulation-kernels. Therefore, if the computational power is dynamically redistributed among the various simulation-kernel instances in order to achieve more balanced runs, as it occurs in the load-sharing architecture, local clocks are expected to diverge less, and in case a rollback operation must be performed, the rollback length (i.e., the amount of executed events which must be undone in order to reach the correct Local-Virtual-Time to restart the execution from) is expected to be reduced. On the other hand, even for balanced workloads, the employment of the top/bottom-half architecture generates a different timing of actions, in terms of information reflection within the LPs' input queues, which can secondarily impact the rollback pattern. In order to evaluate these secondary effects, we have explicitly measured the following parameters:

- Rollback probability, evaluated as the ratio between the amount rollback operations and the amount of executed events.
- Rollback length, expressed as the average number of undone events per rollback operation.
- Efficiency, which is measured as the ratio between the amount of committed and executed events.

4 EXPERIMENTAL RESULTS

4.1 Benchmark Applications and Setting

The experimental assessment has been based on the PCS (*Personal Communication Service*) simulation model, which has been thoroughly described in (Vitali, Pellegrini, and Quaglia 2012). This is a parameterizable GSM communication model — explicitly modeling cells' response upon different call arrival rates — that has been configured to produce a workload relatively uniform across the various LPs (each one modeling a different wireless cell), with a communication pattern which involves sending messages only to each LP's neighbors (due to mobile devices hand-offs).

The call inter-arrival time (τ_A) is exponentially distributed, and the average call duration is set to 2 minutes. Three different configurations of the model have been executed, namely with τ_A set to 0.4, 0.8, and 1.2 respectively, to achieve channel utilization factors on the order of 35%, 15%, and 10% respectively, while the residence time of an active device within a cell has been set to a mean value of 5 min and follows the exponential distribution. The variations of τ_A determine model instances with different profiles in terms of both event granularity and memory requirements. Specifically, the lower the value of τ_A , the larger CPU/memory requirements. Also, lower values for τ_A imply greater computation to communication ratios.

For the above scenario, we have run experiments with 1024 wireless cells, modeled as hexagons covering a square region, each one hosting 1000 wireless channels. The checkpointing interval χ has been set to the fixed value of 20, in order to avoid run-time dynamics fluctuations potentially caused by self-adjusting checkpointing policies. In order to clearly show the actual overhead due to the load-sharing architecture, we have run our experiments in a static fashion, i.e., by forcing the power reassignment procedure within the multi-threaded kernel not to modify the initial even allocation of worker threads to kernel instances.

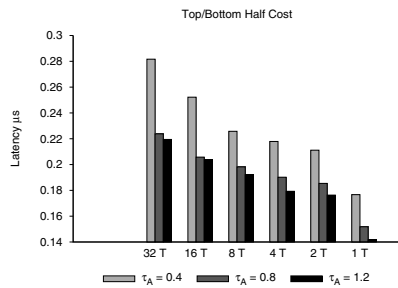


Figure 2: Top/bottom-halves.

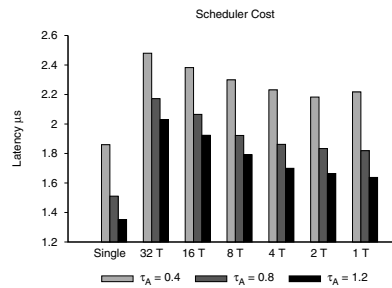


Figure 3: Scheduling

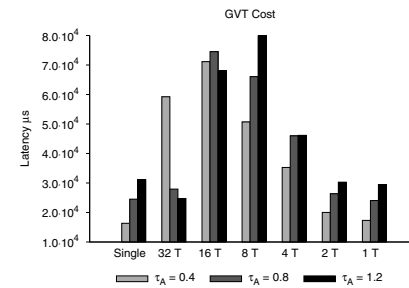


Figure 4: GVT and CCGS

This allows us to check what is the overhead associated with monitoring, managing, and reassignment operations without any benefit from the actual load-sharing approach.

We have run our set of experiments on a HP ProLiant 64-bit NUMA server equipped with four 2GHz AMD Opteron 6128 processors and 64GB of RAM. Each processor has 8 cores (for a total of 32 cores) that share a 10MB L3 cache (5 MB per each 4-cores set), and each core has a 512KB private L2 cache. The operating system is 64-bit Debian 6, with Linux Kernel version 2.6.32.5.

4.2 Results

4.2.1 Assessing Operation Costs

Top/Bottom-Halves Processing. In order to ensure correctness, whenever a top/bottom-half operation must be performed by some worker thread, a lock on the bottom-half queue associated with the destination LP must be taken (although in the case of bottom-halves processing, the lock is only needed for de-queueing the events' chain currently registered within the queue). If the number of concurrent worker threads grows, the contention on the queues is increased, since the worker threads synchronizing on this resource must wait for the lock to be granted to them. At the same time, since a higher number of available worker threads entails a higher number of handled LPs per-kernel instance, this statistically reduces contention on per-LP queues, so that this latency is expected to grow, but up to a certain (not large) extent. Additionally, in our implementation we explicitly relied on pre-allocation for reserving buffers used to keep track of bottom-halves. This choice is guided by the fact that relying on the `malloc` library to allocate nodes can result in a costly operation when executed in a multi-threaded environment, since its internal synchronization relies on `futexes`. If the size of the pre-allocated buffer is well-tuned, the contention on top-half registration is reduced to the minimum.

In Figure 2 we show the per-event latency related to the management of top/bottom-halves. By the plots we see that when the number of per-kernel worker threads increases, the related cost increases just linearly and moderately, given the above considerations.

Scheduling. In Figure 3, the per-event latency related to scheduling operations is provided. By the plots, we can see that the non-multi-threaded implementation (which we refer to as “Single”) shows a latency which is on the order of 15% smaller than the load-sharing one. This small overhead is related to the fact that the load-sharing architecture implements a mapping between LPs handled by a certain worker thread and the actual thread. Therefore, in order to perform CPU scheduling operations, a worker thread must first check which are the LPs it is currently handling. This is an operation which is not executed in the non-multi-threaded implementation. Nevertheless, this difference is not enough to justify 15% latency increase. In fact, a significant additional difference between the two implementations relies on the fact that the load-sharing architecture makes large use of locking primitives for ensuring correctness. This entails a higher number of in-memory accesses for trying to acquire spin-locks, which in turn increases memory bus contention (We note that this result is expected to be different on hardware architectures which rely on *cache locking*, i.e., a cache-coherence protocol which ensures atomicity of in-cache operations by

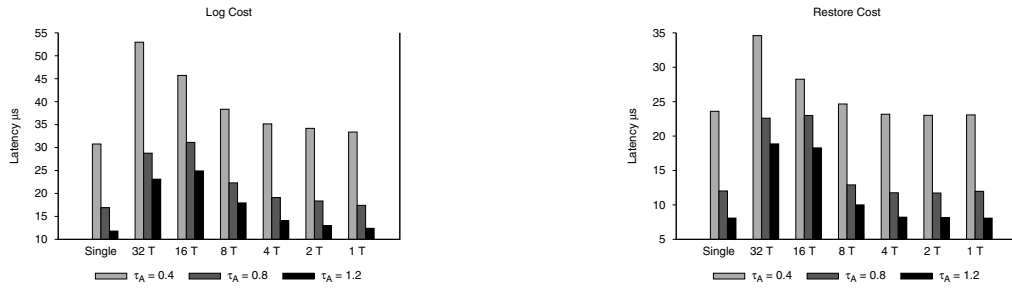


Figure 5: Log/restore operations.

relegating accesses to the highest available levels, therefore avoiding bus contention if data accessed by other threads is not related at all.) and can affect procedures which access (large) data structures sparsely, as the scheduling operation does.

Log/Restore. In Figure 5(a) we present the per-log cost. By the plots, we can see that, independently of the workload, a higher number of worker threads (i.e., a smaller number of concurrent kernel instances) presents a higher latency. In particular, the cost starts growing when running with 4 kernel instances, each one handling 8 worker threads, and the greatest difference is on the order of 30% (15 μs). This latency increase wrt the number of worker threads is also influenced by the aforementioned internal synchronization at the level of the `malloc` library.

In Figure 5(b) the per-restore latency is shown. Fluctuations between different workloads are on the order of 25%, which is due to the fact that a higher number of threads entails a higher number of in-cache buffers invalidations.

GVT and CCGS Computation. Figure 4 shows the plots related to the GVT and CCGS execution latency. As hinted in Section 3.4, the load-sharing version of our simulation-kernel allows a single worker thread to perform CCGS operations, due to critical races on the output.

At the same time, during GVT operations, a procedure for computing the actual workload which the various kernel instances are following through is executed. This can be seen as a distributed agreement among the kernels to determine which is the best number of worker threads per-kernel instance to evenly share the current workload. Again, this procedure is based on MPI message exchanges.

By the plots, we can see that, when running with a small number of worker threads, the latency is comparable with the one achieved by the single-threaded kernel. On the other hand, a larger number of worker threads entails a higher latency. This is related to the serialization required for consistently accessing event queues which are needed to compute the GVT reduction and to evaluate the future (expected) LPs' workload. Additionally, inter-kernel (MPI-based) communication is exploited to correctly follow through the distributed agreement on the best-suited number of worker threads. We additionally note that when running with 32 worker threads, the latency is reduced, since in this configuration there is no actual need to rely on MPI for executing the agreement procedure, since data structures are already locally available.

Inter-Kernel Communication. By the plots in Figure 6, the executions relying on 32 and 16 worker threads exhibit inter-kernel communication latency which is almost two and one orders of magnitude greater than other configurations, respectively, thus giving rise to a smaller event throughput, as depicted in Figure 9. This is related to the fact that these configurations process a larger number of uncommitted events (as reported in Figure 14), since most of the processed events are rolled back. In fact, in Figure 12 we can see that these configurations show, among the others, the higher rollback probability, along with a non-minimal rollback length. This gives rise to low efficiency (as reported in Figure 13).

The high inter-kernel communication latency exhibited by these configurations is related to the higher contention on the MPI layer due to the large amounts of message exchanges which is related to a larger number of events and anti-events generated (for the 16-threads configuration), to a higher number of GVT phases as described in Section 4.2.1 (for both configurations), and to the higher number of MPI probe

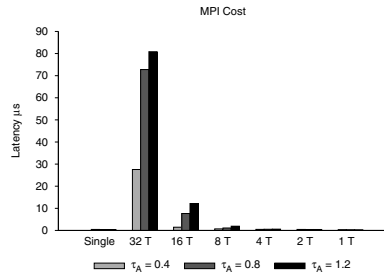


Figure 6: Inter-kernel communication.

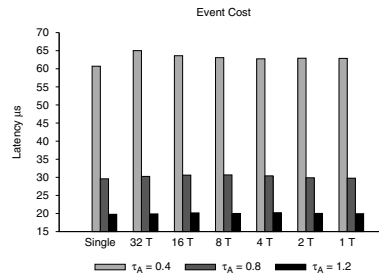


Figure 7: Event execution.

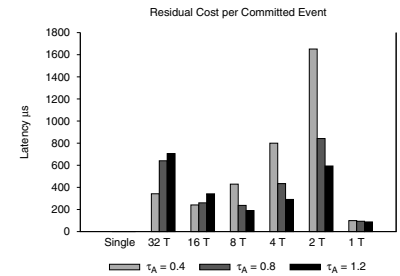


Figure 8: Residual cost.

operations (for both configurations). As for the latter aspect, in the 32-threads configuration, we left the simulation-kernel to perform probe operations towards the MPI even though no message will ever income (given that the run relies on a single kernel instance). This has been done just to observe the effects of the interactions with the MPI layer when scaling up the number of worker threads to the maximum value admitted in relation to the amount of CPU-cores available from the underlying computing platform.

Events' Execution. In Figure 7 we show the per-event execution latency. Although different values of the τ_A parameter produce different events granularity values (due to the different simulation model's load, which produces a higher amount of data to be processed for SIR regulation), different configurations do not affect significantly the event's latency, with small differences that are in the order of 3-5%. This emphasizes the fact that the load-sharing architecture does not affects locality more significantly than the non-multi-threaded one, as far as events' execution is concerned.

Residual Cost. To complete the punctual assessment of our load-sharing architecture, in Figure 8 we present the plots for the residual cost. This includes all the per-committed-event costs which do not appear in the above measurements. Essentially, the residual cost shows which is the time spent for scanning/processing input/output-queues, and for all the other housekeeping operations needed to let the simulation correctly advance. These operations are performed by worker threads on a per-LP basis, but entail accessing memory sparsely, being subject to secondary effects related to memory contention, as depicted in the previous analysis.

As for housekeeping operations, the input queue management and the ack-handling subsystems have a great importance. The former entails calling the `malloc` library for reserving memory buffers which are used to store messages destined to locally handled LPs. Concurrently requesting memory buffers to the `malloc` library involves synchronization mechanisms based on `futexes`, which are likely to increase the overhead related to the registration of messages. The latter is a subsystem in charge of facing the well-known transient message issue for GVT computation, for which a window-based ack mechanism has been adopted. The implementation relies on a lock for each time window (one per kernel), which is acquired during an update operation.

As it can be seen by the plots, the highest residual cost is associated with the two-worker-threads configuration. In fact, in this case there is a small contention wrt the number of threads, but since there are 16 kernel instances running, there is a higher inter-kernel message exchange volume which entails trying to acquire the lock more frequently. The one-thread configuration does not show this contention effect, while increasing the number of threads reduces the need to update the window, thus reducing contention as well.

4.2.2 Operation Weights

To show how the so-far described costs impact on the overall performance of the load-sharing platform's execution, in Figure 11 we present an aggregation showing the percentage of time spent in the various operations, normalized on committed event's execution. By the plots we can see that the non-multi-threaded

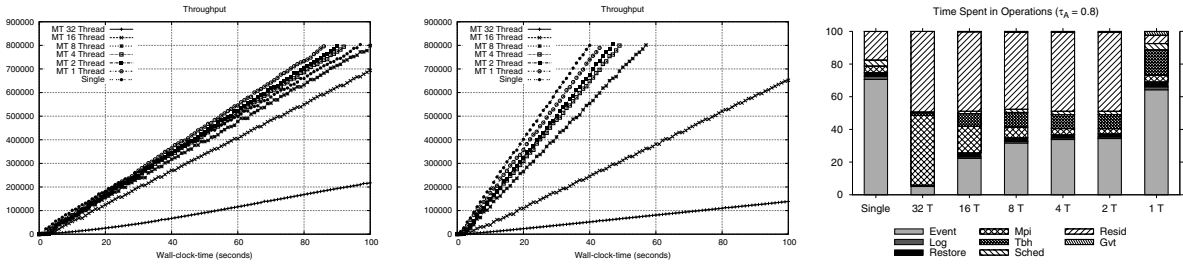


Figure 9: Throughput ($\tau_A = 0.4$). Figure 10: Throughput ($\tau_A = 0.8$). Figure 11: Operation weights.

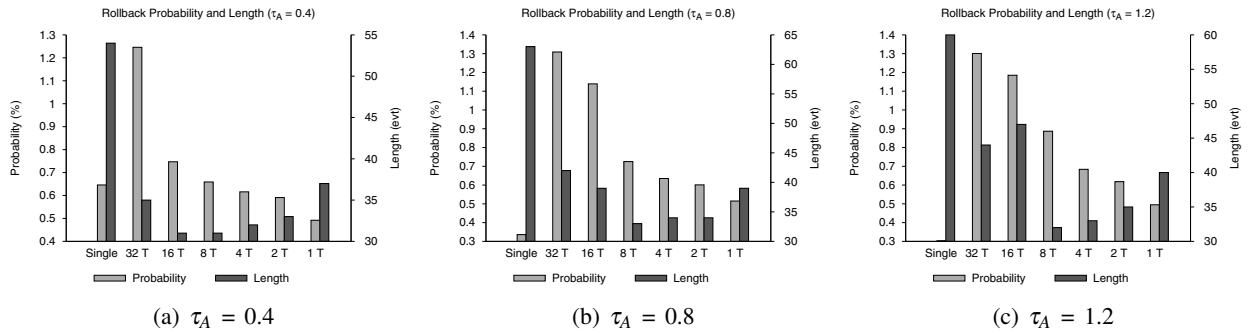


Figure 12: Rollback probability and length.

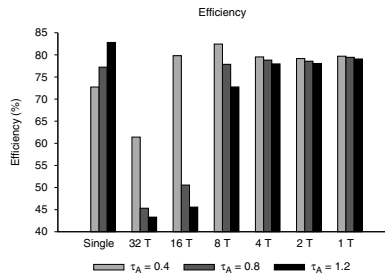


Figure 13: Efficiency.

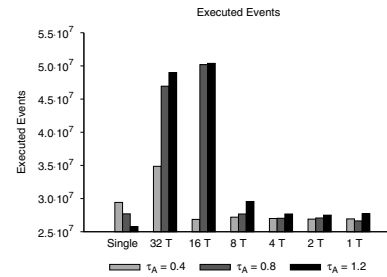


Figure 14: Executed events.

execution, independently of the workload, spends almost 70% of the time in events' processing, similarly to what the load-sharing architecture configured with one worker thread does.

In addition, we can see that as the number of worker threads increases, the amount of time spent in the top/bottom-halves processing decreases, due to the fact that a larger number of LPs is handled by a single kernel instance. At the same time, inter-kernel communication decreases since a higher number of events can be delivered to local LPs. The 32-worker-threads configuration shows an amount of time spent in MPI operations which reduces to the minimum the time spent for event processing, as it was already clearly illustrated in Section 4.2.1.

The high residual time's relevance in the load-sharing architecture running with more than one thread is again related to the window-based ack mechanism, as it was depicted in the previous section.

4.2.3 Global Assessment

To assess our proposed architecture globally, we have measured the cumulated event rate (expressed as the amount of cumulated committed events per Wall-Clock-Time unit), which is a classical indicator of the speed of the optimistic simulation run.

In particular, Figures 9 and 10 show the corresponding throughput values for the benchmark configurations related to the τ_A parameter set to 0.4 and 0.8, respectively. By Figure 10 we can see how the so-far discussed overheads produce a decreasing throughput when the number of worker threads is increased (we remind that in this configuration, the workload is constant and evenly distributed, and the rebalancing procedure is forced not to reassign worker threads to kernel instances, in order to evaluate which are the intrinsic costs of the presented architecture). As it was explained before, the configurations associated with 16 and 32 threads do not scale well, particularly because of the high costs related to MPI operations, and due to a large amount of executed events which get rolled back.

Figure 9 shows the configuration with a different (higher) workload. We note that some load-sharing configurations present a throughput slightly higher than the non-multi-threaded version. This is related to benefits derived from a better exploitation of the caching architecture. Additionally, this is reflected in a higher efficiency, as depicted in Figure 13. As hinted, in this paper our focus is on the overhead analysis for the load-sharing architecture in balanced contexts. Some experimental data related to the benefits from load-sharing with dynamically varying workloads can be found in (Vitali, Pellegrini, and Quaglia 2012).

5 CONCLUSIONS

In this work we have presented a deep study on runtime dynamics related to a symmetric multi-threaded architecture for optimistic simulation-kernels, which is aimed at supporting load-sharing. In particular, the most typical issues related to concurrency have emerged, along with some secondary effects which can significantly effect overall performance.

By our experimental results, we proved that the overhead associated with the proposed symmetric architecture scales well wrt the number of worker threads used during the simulation, except for particular aspects like MPI communication. Nevertheless, this particular overhead could be faced with complementary techniques, like the one proposed in (Jagtap, Abu-Ghazaleh, and Ponomarev 2012).

REFERENCES

- Carothers, C. D., D. W. Bauer, and S. Pearce. 2000, May. "ROSS: a High Performance Modular Time Warp System". In *Proceedings of the 14th Workshop on Parallel and Distributed Simulation*, 53–60: IEEE Computer Society.
- Carothers, C. D., and R. Fujimoto. 2000. "Efficient Execution of Time Warp Programs on Heterogeneous, NOW Platforms". *IEEE Trans. Parallel Distrib. Syst.* 11 (3): 299–317.
- Cucuzzo, D., S. D'Alessio, F. Quaglia, and P. Romano. 2007. "A Lightweight Heuristic-based Mechanism for Collecting Committed Consistent Global States in Optimistic Simulation". In *Proceedings of the 11th IEEE International Symposium on Distributed Simulation and Real-Time Applications*, 227–234.
- D'Angelo, G., and M. Bracuto. 2009. "Distributed Simulation of Large-Scale and Detailed Models.". *International Journal of Simulation and Process Modelling (IJSPM)* 5 (2): 120–131.
- Fujimoto, R. M. 1990, October. "Parallel Discrete Event Simulation". *Communications of the ACM* 33 (10): 30–53.
- Glazer, D. W., and C. Tropper. 1993. "On Process Migration and Load Balancing in Time Warp". *IEEE Trans. Parallel Distrib. Syst.* 4 (3): 318–327.
- Jagtap, D. A., N. Abu-Ghazaleh, and D. Ponomarev. 2012. "Optimization of Parallel Discrete Event Simulator for Multi-core Systems". In *Proceedings of the 26th Parallel and Distributed Processing Symposium*: IEEE Computer Society.
- Jefferson, D. R. 1985, July. "Virtual Time". *ACM Transactions on Programming Languages and System*:404–425.
- Li-li, C., L. Ya-shuai, Y. Yi-ping, P. Shao-liang, and W. Ling-da. 2011. "A Well-Balanced Time Warp System on Multi-Core Environments". In *Proceedings of the 25th Workshop on Principles of Advanced and Distributed Simulation*, 154–162.

- Peluso, S., D. Didona, and F. Quaglia. 2011. "Application Transparent Migration of Simulation Objects with Generic memory Layout". In *Proceedings of the 25th Workshop on Principles of Advanced and Distributed Simulation*, 169–177: IEEE Computer Society.
- Quaglia, F. and Pellegrini, A. and Vitali, R. 2011, October. "ROOT-Sim: The ROME OpTimistic Simulator:<http://www.dis.uniroma1.it/~hpdc/ROOT-Sim/>".
- Toccaceli, R., and F. Quaglia. 2008. "DyMeLoR: Dynamic Memory Logger and Restorer Library for Optimistic Simulation Objects with Generic Memory Layout". In *Proceedings of the 22nd Workshop on Principles of Advanced and Distributed Simulation*, 163–172. Washington, DC, USA: IEEE Computer Society.
- Vitali, R., A. Pellegrini, and F. Quaglia. 2012. "Towards Symmetric Multi-Threaded Optimistic Simulation Kernels". In *Proceedings of the 26th Workshop on Principles of Advanced and Distributed Simulation*: IEEE Computer Society.

AUTHOR BIOGRAPHIES

ROBERTO VITALI is currently a PhD student in *Dipartimento di Informatica e Sistemistica* at *Sapienza Università di Roma*. He is a member of the *High Performance and Dependable Computing Systems* research group. He achieved the bachelor's degree in Computer Engineering in 2007 and the master's degree in Distributed Systems and Computer Architectures in 2009. His research interests primarily include Distributed Simulation, Computer Architectures and Operating Systems. His email address is vitali@dis.uniroma1.it.

ALESSANDRO PELLEGRINI is currently a PhD student in *Dipartimento di Informatica e Sistemistica* at *Sapienza Università di Roma*, collaborating with the *High Performance and Dependable Computing Systems* research group, where he is working in the area of Distributed Systems and Parallel Simulation. He achieved the bachelor's degree in Computer Engineering in 2008 and the master's degree in Distributed Systems and Computer Architectures in 2010. His other research interests include Code Parallelization, Autonomic Computing, and Machine Learning. His email address is pellegrini@dis.uniroma1.it.

FRANCESCO QUAGLIA received the Laurea degree (MS level) in Electronic Engineering in 1995 and the PhD degree in Computer Engineering in 1999 from the University of Rome "La Sapienza". From summer 1999 to summer 2000 he held an appointment as a Researcher at the Italian National Research Council (CNR). Since January 2005 he works as an Associate Professor at the School of Engineering of the University of Rome "La Sapienza", where he has previously worked as an Assistant Professor since September 2000 to December 2004. His research interests span from theoretical to practical aspects concerning distributed systems and applications, distributed protocols, middleware platforms, parallel discrete event simulation, federated simulation systems, parallel computing applications, fault-tolerant programming, transactional systems, Web-based systems and performance evaluation of software/hardware systems. He has served as Program Co-Chair of PADS 2002 and PADS 2010, as Program Co-Chair of NCA 2007, as Program Co-Chair of SIMUTools 2012, as General Chair of PADS 2008, as General Co-Chair of SIMUTools 2011, and as Tutorial Co-Chair of HPCS 2011. Since 2004, he is an Editorial Board Member of the *International Journal of Simulation and Process Modelling (IJSPM)*. His email address is quaglia@dis.uniroma1.it.