

Interconnection Network Simulation Using Traces of MPI Applications

J. Miguel-Alonso · J. Navaridas · F. J. Ridruejo

Received: 25 October 2007 / Accepted: 7 November 2008 / Published online: 6 December 2008
© Springer Science+Business Media, LLC 2008

Abstract This paper addresses the utilization of traces taken from MPI applications to do simulation-based performance studies of parallel computing systems. Different mechanisms to capture traces are discussed, pointing out important limitations of some of them. One of these limitations is the invisibility of message interchanges in collective operations, which is circumvented modifying a trace-capturing library. During a simulation, trace records must be simulated in causal order, to fully comply with application semantics. Alternatives to follow this order, and the risks of not following it, are presented and discussed. The techniques introduced in this paper have been implemented in an in-house developed simulation environment, which is used in two example studies to show its usefulness: an evaluation of alternatives for interconnection network design, and a performance prediction study in which traces from one machine are used to estimate the execution times of applications running in a different machine.

Keywords Interconnection network simulation · Traces of parallel applications · Message passing interface

J. Miguel-Alonso (✉) · J. Navaridas · F. J. Ridruejo
Department of Computer Architecture and Technology, The University of the Basque Country,
P.O. Box 649, 20080 San Sebastian, Spain
e-mail: j.miguel@ehu.es

J. Navaridas
e-mail: javier.navaridas@ehu.es

F. J. Ridruejo
e-mail: franciscojavier.ridruejo@ehu.es

1 Introduction

Simulation is one of the most widely used tools for performance evaluation of computing systems, including parallel computers. A simulation-based study requires a model of the system being evaluated (a model that may have very different levels of accuracy) and also a mechanism to supply a *representative* workload.

In the field of interconnection networks, in which we place this paper, many simulation-based studies use synthetic traffic patterns, such as random uniform traffic or permutations of interest, to feed the simulator. This kind of synthetic workload is of great interest because it is easy to implement, sometimes it may support analytical studies, and may be representative of the ways applications use the network. However, a comprehensive evaluation requires actual workloads, otherwise important aspects of parallel applications cannot be understood in detail. For example, many applications pass through different phases, in which the ways of using the network differ widely; pressure on the network may be very intense in some phases, but the inter-dependencies amongst processes may lower the utilization of the interconnection infrastructure in some others.

Actual traffic may be generated using an execution-driven environment, in which applications run on real (or simulated) processors and are connected to a simulated interconnection network. This set-up provides very high levels of evaluation accuracy, but cannot be easily scaled to thousands of processors. It may also fail, victim of unexpected interactions between components, as shown in [16, 24]. For this reason, a frequently used alternative is the utilization of traces of parallel applications.

We can obtain traces of large systems, even using small ones. For example, given an application that uses a static logical topology, a cluster of 10 PCs can be used to generate a trace running on 200 nodes—we only need to run 20 processes per available computer. Timing information would not be representative of a real, 200-CPU computer; however, the (spatial) patterns defined by the sequence of interchanged messages are valid, and the distribution of message sizes is valid too.

The main contribution of this research work is the trace-based simulation tool-set included in INSEE [23], an in-house developed simulation environment for the evaluation of interconnection networks. INSEE is able to accept many kinds of workloads for the simulation: synthetically generated messages, full system simulation and traces. The focus of this paper is, precisely, on simulation with traces. We explain the mechanisms required to capture traces from MPI-based parallel applications, point out some important limitations in the way events are captured and stored in the trace files, and explain the way we circumvent some of those limitations. After that, we describe in detail the way traces must be processed in order to fully comply with application semantics, with particular attention to the order in which events are processed to avoid violations of message causality. A simple extension of the trace processing mechanism allows us to carry out performance prediction studies. Finally, we put INSEE to work and include two simple example studies carried out with traces from actual applications. Note that the problems and solutions discussed in this paper, although presented from the point of view of the INSEE developers and users, are not specific to this tool.

The rest of this paper is organized as follows. Section 2 introduces INSEE, with focus on its utilization with traces. Section 3 describes how traces are obtained. Sections 4 and 5 discuss the mechanisms used by our simulator to process traces, preserving application semantics. In Sect. 6 we use INSEE for two performance studies. Section 7 reviews some related work. Conclusions of the paper are summarized in Sect. 8.

2 Evaluation of Networks Using Simulation: INSEE

In this section we briefly introduce INSEE [23], the Interconnection Network Simulation and Evaluation Environment developed at the University of the Basque Country. The two main elements of INSEE are FSIN, a Functional Simulator of Interconnection Networks, and TrGen, a Traffic Generator [22].

FSIN has been designed to provide a fast simulation engine for interconnection networks, both direct (meshes and tori) and multistage (trees, including fat-trees), with different architectural characteristics. Its small footprint allows us to simulate, on an off-the-shelf desktop computer, large size networks: we have carried out experiments with 64K-node networks using less than 2 GB of RAM. In the case of direct networks, each node represents a router attached to a computing element which, in fact, is the source (and sink) of the traffic managed by the network. In the case of trees, computing elements are attached to switches at the lowest level, using interface cards.

The management of workloads is carried out by TrGen. A workload is a collection of messages that are generated by computing elements (actually, by message sources), then packetized and passed to FSIN (which simulates the way they traverse the interconnection network), then reassembled and, finally, delivered to the computing elements (actually, to the message sinks). When generating a workload, we have to define

- The *spatial* distribution of messages: source nodes and destination nodes.
- The *size* distribution of the messages. These come in many sizes, depending on the application.
- The *temporal* distribution of the generation of messages. In some cases, workload generators simply generate random numbers (following a certain distribution) that determine the inter-generation intervals. In some others traffic is *reactive*, meaning that there are *causal relationships* between them: the arrival of a message to a certain destination node triggers the generation of a new message from that node.

Very simple *synthetic workloads* use statistical distributions to generate destinations, sizes and inter-generation times. A special class of this traffic is what we call application-inspired traffic [17], for which we emulate the behavior of some kernels of scientific applications, including causality. It is important to point out that simulations using traces, and full system simulation, use spatial, size and temporal distributions exactly as defined by the application that was instrumented, or that is being executed. When performing *full system simulations*, an external toolset based on Simics [11] fully simulates a collection of computing nodes (including hardware, drivers, operating system, message passing library, and running application) attached to an interconnection network, which is simulated by FSIN [16]. To perform *trace based*

simulation, we use traces obtained using the mechanisms provided by MPI implementations, although some modifications to these are required in order to get *extended* trace files, because FSIN only deals with point-to-point operations. An extended trace file includes the detailed message passing involved in collective operations, something that is not visible in *regular* trace files.

We stated before that FSIN allows the simulation of very large networks. Unfortunately, this can be done *only* with synthetic traffic. The trace capturing environment, or the ability to fully simulate collections of computers, limits the node count for the other traffic generation arrangements. At any rate, the focus of this work is on trace based simulation.

In the following sections we will explain the way INSEE deals with traces, from the mechanism used to capture them, to the way they are consumed by the simulator. We will explain how traces can be used not only for performance evaluation but also for performance prediction.

3 Generation of Traces from MPI Applications

MPICH [9] is one of the most widely used implementations of the Message Passing Interface (MPI) [12], a standard programming interface for parallel applications based on processes that communicate and synchronize explicitly, interchanging messages. The MPI standard defines a profiling mechanism called PMPI (“P” from Profiling) that allows programmers to intercept all calls to MPI functions. This mechanism is often used to implement libraries to generate traces of applications, or to obtain profiling information.

3.1 Generating Trace Files with MPE

The MPICH distribution includes MPE (Multi-Processing Environment) [10], a set of libraries and tools to generate and analyze traces of parallel applications. The tracing ability is based on PMPI, so MPE can be used with any MPI implementation, not only with MPICH. However, in our discussion we will only consider the MPICH/MPE combination. To trace-enable an application, we just need to compile it using the compiler wrappers offered by MPICH (mpicc, mpicxx, mpif77, mpif90) with the “-mpilog” flag activated—no change in the source code is required.

Trace-enabled applications run as normal but, when finished, write a trace file that consists of a set of time-stamped records (events) that describe the dynamic behavior of the application during its execution. Records in a trace file include:

- Information about MPI functions invoked by the application processes. Each function invocation generates two records: one when a process calls the function and another one when the function returns. A pair of these records represents a “state”: the first one indicates when the process enters a given state, the second when the process exits from it. States are defined on a process-by-process basis. There are no “global” states. The most relevant fields of the state records are:

- Process identifier
 - Timestamp
 - Record type (state start/state end, MPI function)
- Information about message interchanges, only for point-to-point operations. Two records are generated per message, one when generated and another one when received. The basic information of the message records includes:
- Process identifier
 - Timestamp
 - Record type, or operation (send/receive)
 - Identifier of the “other” party (destination for a send, source for a receive)
 - A message tag
 - Message size

From now on, we will use these abbreviations: SS means State-Start, SE means State-End, MS means Message-Send, and MR means Message-Receive.

A trace file can be analyzed using tools such as Jumpshot [29], distributed with MPE. Figure 1 shows a screenshot of this tool analyzing a CG.W.8 benchmark (Conjugate Gradient with 8 tasks, class W, included in the well-known NAS Parallel Benchmarks, NPB [2]). The legend (left) indicates the color codes used in the bars that represent states. Messages are represented by arrows.

This way of generating trace files has some limitations. For our purposes, two are the most relevant:

1. *Collective operations.* This class of operations, that involve synchronization and communication among multiple processes, are only represented as states. The trace file does not include any record that reflects the way messages are interchanged to

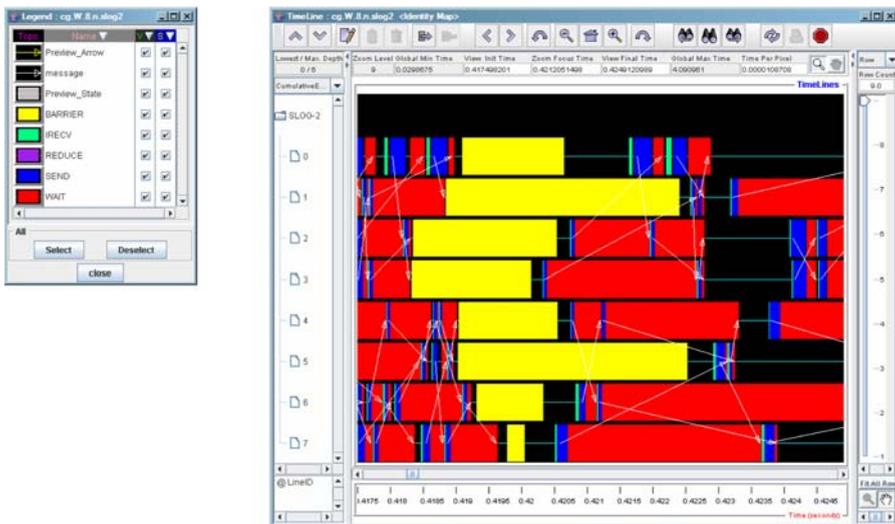


Fig. 1 Screenshot of Jumpshot visualizing a trace file generated by CG.W.8

implement collective operations. For example, in Fig. 1 it is possible to see that an MPI_Send state at a given node is related to a message that departs from that node. In contrast, the MPI_Barrier states are not related to any message.

2. *States vs. messages.* In order to fully understand a point-to-point state, you need to consider the information provided in the corresponding state records as well as information contained in separate message records. For example, an $\langle n \ ts1 \ SS \ MPI_Send \rangle$ state start record indicates that the node identified as $\#n$ tries, at time $ts1$, to send a message, but details about the message are found in a separate $\langle n \ ts2 \ MS \rangle$ message send record.

Collective operations can be implemented in many different ways. Often, the implementation is done at the MPI library level, using point-to-point operations to perform broadcast, reduce, gather/scatter, etc. This approach is very flexible, because collectives will work on top of any network, and is the one of choice in the case of popular MPI implementations, including MPICH. However, some networks provide support for collectives, or have topological properties that make some implementations more efficient than others, and the generic libraries cannot take advantage of these characteristics. A tailor-made implementation of collectives would be much more efficient. For example, [1] discusses the implementation of collective operations in an MPI library specifically designed for the IBM BlueGene/L system.

In the following subsection we will discuss how to overcome the first of these limitations and how, for simulation purposes, the state records (those that register when processes enter in and exit from MPI operations) can be safely ignored.

3.2 Generating Extended Trace Files

We have explained how trace files do not include detailed information about collective operations, because the details of how they are implemented are invisible to the application. A study of the internals of the MPICH implementation of MPI showed that collectives are, by default, carried out using point-to-point messages. The MPICH designers could have chosen to use some internal message-passing functions; fortunately for us, they decided instead to use the standard MPI point-to-point passing functions. For example, MPI_Broadcast is implemented using MPI_Send and MPI_Recv (the most basic message interchange functions), and MPI_Barrier is implemented using MPI_Sendrecv (a combination of MPI_Send and MPI_Recv in a single operation). The details of the default implementation of collectives are not accessible via the PMPI profiling interface, but this limitation is intentional—and makes sense, because other implementations are possible. We have modified the sources of MPICH to change this behavior, making the hidden operations visible through the profiling interface. This has no consequence in terms of communication semantics/timings.

With the modified MPICH, the generated, extended trace files follow the scheme described above, but they are longer because they include more detail. When visualized, using Jumpshot, a regular trace and an extended trace present different pictures. In Fig. 2 we can see a screenshot of a visualization of the extended version of a trace file, which corresponds to the same CG.W.8 benchmark used in Fig. 1. Note how the boxes inside the MPI_Barrier state represent the way this collective is implemented

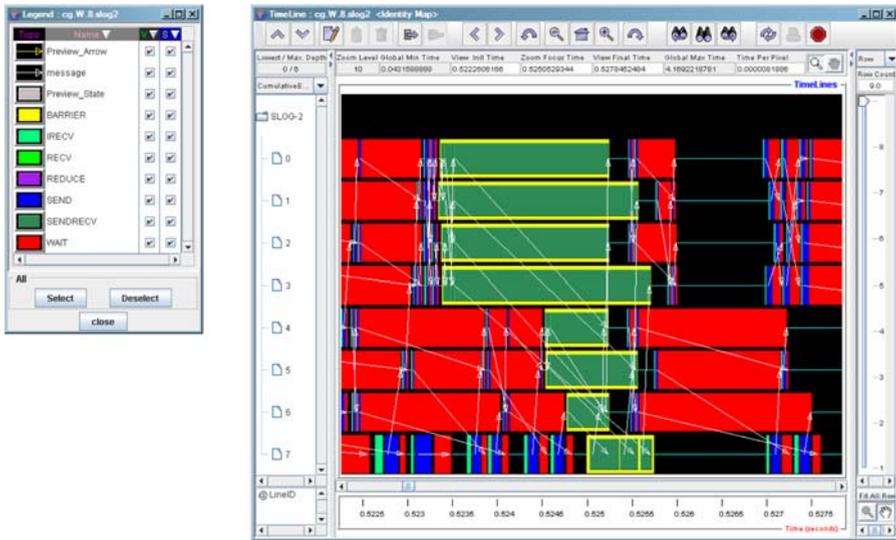


Fig. 2 Screenshot of Jumpshot visualizing an extended trace file generated by CG.W.8

using MPI_Sendrecv, and how messages interchanged in those operations are clearly visible.

From this point onwards, when discussing trace files we actually mean extended trace files.

4 Using Traces to Feed Simulations

We now describe how (extended) trace files can be used to provide realistic communication workloads to simulators of interconnection networks. We make the following assumptions:

- If we simulate a network with N nodes, the trace includes information about exactly N communicating processes; that is, there is a one-to-one relationship between application processes and network-attached nodes. For simplicity, we place processes into the nodes in consecutive order, that is, process n goes to node n .
- Simulators deal with the exchange of packets. Applications generate/consume messages of variable sizes, which need to be split into fixed or variable-sized packets (depending on the network technology). We discuss message interchanges as if the network delivered them directly, although implicitly we are considering message segmentation into fixed-size packets at origin, as well as message reassembly at destination.

If the first assumption is not met, simulation would be possible too, but the injection mechanisms and the placement policies would be more complex. There could be nodes (CPUs) without assigned processes (that would not participate in the simulation), or

nodes housing several concurrent processes. In this case a table mapping processes to nodes would be enough to deliver messages to the appropriate destination.

4.1 First Approach: Inject as Fast as You Can

An initial, and rather unrefined, approach to feeding a simulator with events taken from a trace file is as follows:

1. Ignore all the state (SS, SE) and MR records—in other words, use only the MS records with information about messages sent.
2. Split the trace file in one list per simulated node, and arrange the lists in timestamp order.
3. Make each node inject messages from its list into the network as fast as the network can accept them. Network backpressure is used to modulate the injection of load into the network.

As we can see, timestamps are used only to impose an order. The main justification for this decision is that we focus on network performance: we want to measure how fast a network can deal with a given workload, so we want to stress it, making it our bottleneck. The timing information included in a trace file is affected by issues that fall outside our control: the actual network used in the instrumented experiment, the processors and their speeds, the MPI implementation, the overhead of the instrumentation system, the number of processes that share a CPU, etc. We want to isolate the simulation from these facts. If we want to carry out performance predictions at a system level, we should take into considerations all these issues—which would make an already complex problem close to unworkable.

This approach to simulation accurately reproduces, using the information captured in the trace, the *spatial* communication pattern of the application (sources and destinations), and also the *message sizes*. However, it fails to reproduce the *temporal* pattern, which should respect message causality and reflect the actual way message exchanges are interleaved, as required by the application.

4.2 Second Approach: Follow Causal Order

The previous approach does not take into account the causal relationships between messages. In an actual execution of a parallel application, it often happens that a process stalls while waiting to receive a new message. Process execution is only resumed when the expected message arrives. We may emulate this behavior as follows:

1. Ignore all state records—in other words, use only MS and MR records.
2. Split the trace file into timestamp-ordered lists, one list per simulated node. Note that each node $\#n$ has an ordered record list (an “event queue”) of $\langle n \text{ timestamp MS destination size} \rangle$ and $\langle n \text{ timestamp MR origin size} \rangle$ records. In the following steps timestamps are used only to order records.
3. Create, at each node, a receipt list, initially empty, that will store messages delivered by the network.

4. At each node, do the following:
 - a. If the first record in the event queue is an MS, remove it and inject the corresponding message into the network.
 - b. If it is an MR record, check if a corresponding message (matching origin, destination, tag, and size) is in the receipt list. If it is there, remove both entries. Otherwise, do nothing.
 - c. When the simulator delivers a message, put it in the receipt list.

This procedure is depicted in Fig. 3. Its main implication is that an MR record puts the injection process on hold until the corresponding message is actually received from the network. In the figure, node #0 cannot advance, because it is waiting for a message from node #1, even if a message from node #2 has been received already. In contrast, node #1 can advance because the required message from node #0 has been delivered. This mechanism reproduces the actual way messages were interleaved when running the application, complying with the causal order between a receipt and the subsequent sends it may trigger.

The main drawback of this approach is that it *may* be excessively conservative. If we look again at Fig. 3, we see that the event queue of node #0 says that, after receiving a message from node #2, it is possible to send a message to node #3, and it happens that the message has been received already. However, node #0 is stalled (waiting for a message from node #1). We may wonder if application semantics is adequately emulated. Is it *really* necessary to receive the message from node #1 before advancing?

There is not a single answer to this question: it has to be discussed on a case-by-case (application-by-application) basis. For example, in [21] the causal ordering enforced by our simulator is considered valid in the context of cc-NUMA machines, because

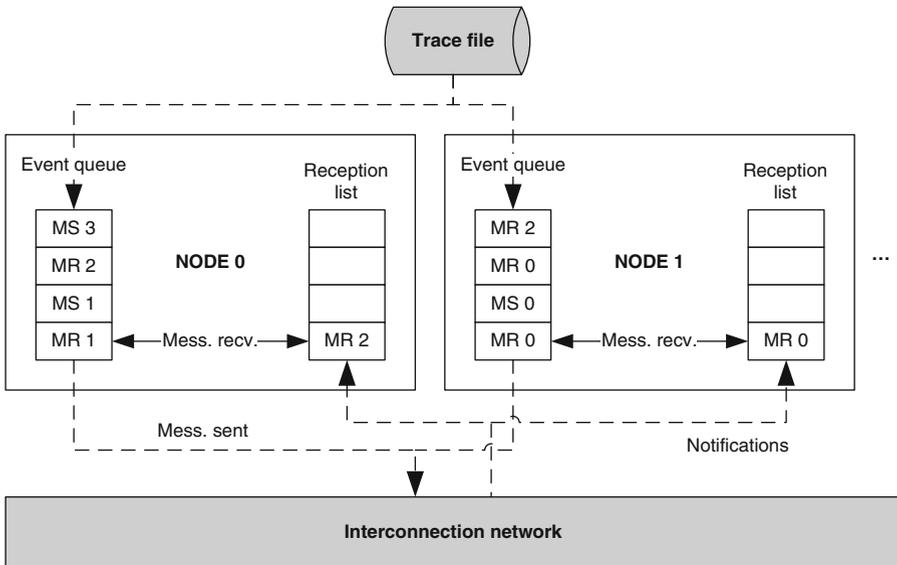


Fig. 3 Interface between the trace file and the network simulator. In the boxes, “MS *n*” means that a message is sent to node #*n*, and “MR *n*” means that a message is expected from node #*n*

message interchange is reactive: a message sent requires a response before allowing the process to advance. However, in general terms “*The parallel execution semantics, as reflected in the message communication operations and how the message data is used, determines process dependencies and message event ordering relationships, but only partially. Non-deterministic execution allows for alternative message event orderings.*” [28]. We further discuss this issue in Sect. 5.

4.3 Traces for Performance Prediction

In the previous subsections we have stated that timestamp information is used to arrange events in temporal (or causal) order, but is otherwise ignored. This decision makes sense if we want to study the performance of the interconnection network. However, the network is only part of a system. Applications run on a collection of compute nodes, whose behavior is also visible through traces.

How can we use trace-driven simulation to estimate the time needed to run an application? The starting point is a *real* system, in which the application is run and traces are obtained. Those traces contain communication events, as can be seen in Figs. 1 and 2. In the figures we can also see “empty” spaces between MPI states. These spaces represent the time spent by processes outside MPI calls. If, during that time, processes are doing useful CPU work (of course, this is not always the case), then a simulator can be fed with the traces and a set of parameters that define the CPU characteristics, as well as the network characteristics, and used to predict overall application performance.

The simulator runs as described in Sect. 4.2 (following causal relationships). Times between communication states (between the SE record that corresponds to the end point of an MPI operation and the SS record that corresponds to the starting point of the following one) are converted into CPU states, with their SS and SE events. In this procedure we could scale the duration of these states, to simulate faster or slower CPUs. During the simulation run, when an “SS CPU” record for a given node is processed, injections from that node are stopped, and will be resumed only after processing the corresponding “SE CPU” event. In other words, the node is kept “busy” for the time required by the CPU state.

With this set-up, the simulation takes into consideration the characteristics of the CPUs, as well as the characteristics of the interconnection network, to estimate the time required to execute a message-passing application. Any change in the CPUs, or in the network, will be reflected as a change in the time required to consume the trace file.

5 Reproducing Application Semantics Accurately

5.1 Record Order in MPE Traces

A naïve user that analyzes an MPE trace file may think that events are recorded exactly when they happen during the application run, and that timestamps are accurate. This is not true, for several reasons. An obvious one is that the program file has been

instrumented in order to generate the trace, so the program is not running as fast as it would when not instrumented. The second reason is more subtle, and to understand it we need to explain how MPI traces are generated in the MPE environment. Our discussion is also valid for other tracing tools based on PMPI. Note that we assume that we deal only with point-to-point operations, because collective operations, if present, are also included in terms of the underlying point-to-point primitives that implement them. We focus on a subset of the MPI point-to-point operations, in order to discuss the relevant characteristics of the way they are logged without introducing unnecessary details. We consider that this subset is still sufficient, because it includes most (if not all) operations necessary to run the applications included in the NPB.

As we explained before, MPE logs are generated using instrumented versions of all the MPI functions. An $\langle n \text{ ts SS MPI_X} \rangle$ record is generated when process $\#n$ invokes the MPI_X function; an $\langle n \text{ ts SE MPI_X} \rangle$ record is generated when this function returns. MS and MR records (messages sent/received) are also generated by these instrumented routines, and only inside them. Note the implication of this way of working: *message send and receipt are not logged when they happen.*

- An MS record is logged *after* the process has entered into a state in which it requests sending a message (MPI_Send, MPI_Isend, MPI_Sendrecv), and *before* the process exits from that state. The message may be injected into the network much later due to a variety of reasons: semantics of MPI operations (immediate operations), previous messages queued, network congestion, decision of the kernel's scheduler, etc.
- An MR record is logged *after* the process has entered into a state in which it is actively waiting for a message or collection of messages (MPI_Recv, MPI_Wait, MPI_Sendrecv, MPI_Waitall), and *just before* the process exits from that state. It may happen that a message has been received from the network interface long before, but this receipt is not logged until the receiving process has entered a waiting state.

This behavior is clearly visible in Figs. 1 and 2. Note that MPI_Irecv states can be safely ignored, because the actual receipt of a message is recorded in a subsequent MPI_Wait or MPI_Waitall; there are no arrows arriving to the MPI_Irecv states. At any rate, when the trace file includes an MR record, it is there because the process really needs it to advance. So, in the simulation, it is necessary to receive that message before allowing the process to proceed—we will further discuss this issue in the following subsection.

The main conclusion here is that application-generated logs are not accurate because they do not reflect the exact moments in which messages are actually sent or received. The actual injection of a message may have happened later than indicated in the trace file, and the actual receipt of a message may have taken place before the time indicated by the record timestamp—sometimes, long before.

5.2 Receipts from MPI_ANY_SOURCE

A skilled MPI programmer knows that it is possible to indicate a wildcard, instead of a source process, in point-to-point receive operations: MPI_Recv(..., 3, ...) executed

at process #0 forces this process to pause until a message from process #3 is received. In contrast, `MPI_Recv(..., MPI_ANY_SOURCE, ...)` pauses the process until a message *from any source* is received. In terms of records in a trace file, the first call and the second one are indistinguishable. Both generate three records for process #0:

```
<0 t0 SS MPI_Recv>
<0 t1 MR ...>
<0 t2 SE MPI_Recv>
```

The only difference *could* be in the second record: the first call guarantees that a message was received from process #3, while the other one may contain any (valid) process identifier.

The programmer may use `MPI_ANY_SOURCE` just for convenience: it may happen that the sender is known beforehand, so that it is not necessary to make it explicit. However, its main purpose is to allow processes to wait for messages that could arrive from any source, when the origin of the next useful message cannot be known a-priori. Let us explore this issue by means of a simplistic scenario of a master-slave application implemented using three application processes. Process #0 is the master, and processes #1 and #2 act as slaves. The protocol is as follows. A slave, when free to perform some work, sends a job request to the master. Then, the master replies with a task to perform.

A beginner in MPI programming could code the application as shown in Fig. 4, Version A. However, a more experienced programmer would use Version B of the code (Fig. 4, right).

Code Version A forces an unnecessary receipt order, which may delay program progress. For example it may happen that a message from slave #2 is already buffered, but one from slave #1 has not been received yet. The process is stalled in the first sentence of the loop, even when the third (receipt of a job request from slave #2) and fourth (sending a task to slave #2) could be executed without risk. This would not happen with Version B of the program, where the utilization of `MPI_ANY_SOURCE` at the receipt side would allow the process to make progress as soon as possible.

A trace file generated by Code Version A would always have message records in the same order: the one shown in the **Trace a** column of Fig. 5. It is easy to understand that a record reflecting that slave #1 sends a message to master #0 (`<1 ts1 MS 0>`) must be recorded in the trace file somewhere before record timestamped `ta2`; this record would match the one timestamped `ta1`. In the same way, a record `<2 ts2 MS 0>` must be somewhere before record `ta4`, to match with record `ta3`. Note, again, that the actual

CODE VERSION A:	CODE VERSION B:
...	...
do {	do {
MPI_Recv(..., 1, ...);	MPI_Recv(..., MPI_ANY_SOURCE, ..., &sender);
MPI_Send(..., 1, ...);	MPI_Send(..., sender, ...);
MPI_Recv(..., 2, ...);	} while pending_tasks;
MPI_Send(..., 2, ...);	...
} while pending_tasks;	
...	

Fig. 4 Excerpts of sample codes for a master-slave application

Trace a:	Trace b:	Trace c:	Trace d:	Trace a' :
...
<0 ta1 MR 1>	<0 tb1 MR 2>	<0 tc1 MR 1>	<0 td1 MR 2>	<0 ta1 MR ANY>
...
<0 ta2 MS 1>	<0 tb2 MS 2>	<0 tc1 MS 1>	<0 td2 MS 2>	<0 ta2 MS 1>
...
<0 ta3 MR 2>	<0 tb3 MR 1>	<0 tc2 MR 1>	<0 td3 MR 2>	<0 ta3 MR ANY>
...
<0 ta4 MS 2>	<0 tb4 MS 1>	<0 tc3 MS 1>	<0 td4 MS 2>	<0 ta4 MS 2>
...

Fig. 5 Excerpts of trace files for Version A (Trace a) and Version B (Traces a, b, c and d) of the master-slave application. Trace a' is a modification of Trace a using wildcard receives

receipt from node #2 could have happened before the receipt from node #1, but the trace file would not reflect this circumstance.

Now, let us suppose we used Code Version B. A sequence of events equal to that generated by Version A (probably with different timestamps) would be valid, but **Trace b**, **Trace c** and **Trace d**, also shown in Fig. 5, are equally valid. Only one of them would be actually recorded, depending on aspects such as workload assigned to those processors, relative CPU speeds, the characteristics and status of the network, etc. Let us further suppose that the trace actually recorded looks like **Trace a**.

When doing a simulation we do not want to force that particular order, because it may introduce unnecessary delays. We could use wildcard receives, because this information is in the trace file (not in the MR records, but in the corresponding state records). After a small manipulation of the MR records, **Trace a** can be modified to look like **Trace a'**—with which we feed the simulator. Immediately, we must suspend the master process (#0) at event timestamped *ta1* while waiting for a matching receipt. These are two possible scenarios:

1. The simulator delivers a message (job request) from slave #1 to the master. This unblocks the master, and the simulation continues. The master sends a message (containing a task to perform) to slave #1 and blocks again, waiting for a message from slave #2 that, eventually, will be delivered.
2. The simulator delivers a message from slave #2 to the master. This unblocks the master and the simulation continues. Then, as directed by the trace, master sends a message to #1—something that is not consistent with the event order in the trace file. We interpret this as a violation of application’s semantics.

In order to obey causality relationships among messages we should not use wildcards in the trace files. The sequence of events actually stored in the trace may not be the only valid one, but at least we know that it is semantically valid.

6 Experimental Work Using Traces

First of all, we want to lay emphasis on this point: the purpose of this section is *only* to illustrate the kind of research work that can be done with a trace-driven simulation toolset, such as INSEE. We have used this environment to perform many performance studies, published elsewhere. To cite some examples, in [5] INSEE was used to

evaluate the performance impact of using twisted wrap-around links in mixed-radix twisted tori; in [16] we cross-validated the trace-driven simulation with an execution-driven environment based on Simics; and in [25] we studied the performance of a congestion control mechanism, comparing results obtained with synthetic workloads with those obtained with traces.

We will use the trace-processing abilities of INSEE to carry out two, very different, example performance studies. First we evaluate the impact on performance of different strategies of routing and virtual channel management. Then we estimate the time to execute an application (Conjugate Gradient) on three different target multicomputers. In the experiments we use a network with an 8-ary 2-cube topology—in other words, a 2D torus with 64 nodes. We have several trace files that can be used on networks of this size, which were obtained running the applications included in the NPB suite [2], class W, on 64 nodes.

6.1 Experimenting with Virtual Channel Management and Routing

Our first example study consists of an evaluation of the effects of using several virtual channels per physical link, and the impact of using adaptive routing. Figure 6 represents the network we model, and the details of each router. Each router has 4 bidirectional links (X+, X−, Y+ and Y−), each one connecting it to a different, neighboring router. In simple routers each link has associated input and output ports, with some buffer space for in-transit traffic. However, it is a common practice to associate several virtual channels (VC) to each link; each VC manages its own buffers. In the figure, there are 3 VC per link—we can see that in detail for link X+, shared by virtual channels X0+, X1+, X2+.

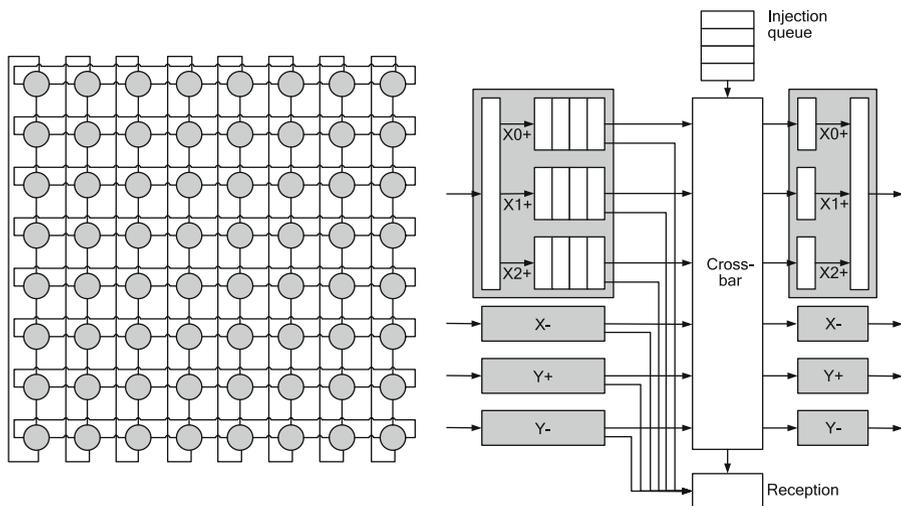


Fig. 6 Left: an 8-ary 2-cube (2D torus with 64 nodes). Right: model of the router simulated by FSIN

Routers, perform routing decisions in order to make packets advance from source to destination. There are multiple variants of routing algorithms, but we will only consider these:

1. Dimension-order, oblivious routing (DOR). A packet must traverse first as many hops as necessary in the X axis (in the row where it was injected) to reach its destination column. Then, it has to move in the Y axis (up or down) until reaching the destination.
2. Adaptive routing, using minimal paths. A packet can jump freely from a given VC to any other VC, continuing in the same axis (row or column) or switching. However, the jump must move the packet closer to the destination.

These two algorithms can lead to undesirable deadlock situations. To avoid these, we use the bubble routing mechanism described in [20] and used in the IBM BlueGene/L, so that we can state that the network is deadlock-free.

Regarding the utilization of several VC, and combining that with the routing algorithms, we compare routers built with the following designs:

1. *Oblivious 1 VC*. A single VC per physical channel. To ensure deadlock-freedom, bubble-restricted DOR is used.
2. *Oblivious 3 VC*. Three parallel VCs per physical link. Routing is bubble-restricted DOR. This arrangement reduces the effects of head-of-line blocking in the transit queues, so that when several packets are competing to use the same links they can advance faster.
3. *Adaptive 3 VC*. Three VCs per physical link. One of them, the Escape VC, uses bubble-restricted DOR, and the other two are adaptive. Packets can switch VCs, but access to the Escape VC has to follow the bubble restrictions. This arrangement provides the same advantage of the previous one. Furthermore, adaptive (but deadlock-free) routing allows a more efficient utilization of links, especially when packets have to travel long distances.

We configure FSIN to simulate networks built with these three routers. TrGen generates the workload, using the traces from class W of the NPB applications. The most interesting results are those obtained with traces from benchmarks Block-Tridiagonal (BT), Conjugate Gradient (CG) and Integer Sort (IS). Traces contain records of message interchanges, which need to be “packetized” in small blocks (packets) of 64 bytes. The queues in the routers are configured to hold up to 4 of these packets. Bandwidth of the links is 32 bits per cycle.

The simulator reports (among many other things) the number of cycles that the network needs to deliver all the applied workload. As each application is different, the numbers differ widely from one to another. For this reason we present relative values in Fig. 7; the base case (value 1) corresponds to the simpler router architecture (Oblivious 1 VC). The figure represents the average values of 10 simulation runs, as well as the 99% confidence intervals.

From the obtained results, we see that, for the traffic pattern used by BT, the utilization of several VCs per link does not offer any advantage in terms of performance. However, improvements for IS and CG are quite good. For IS, most of the improvement comes from the use of several VCs; adaptivity provides minor additional gains.

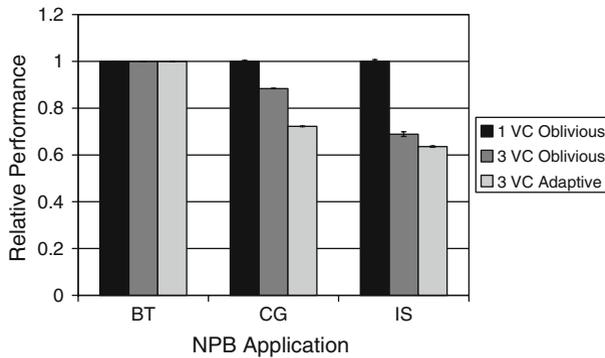


Fig. 7 Effects of using 3 VC per physical link, and of adding adaptivity. Results relative to the base case (Oblivious 1 VC). Average of 10 simulation runs, and 99% confidence intervals

CG benefits less than IS from using several VCs, but is capable of taking advantage of adaptivity. The reasons for these results have to be found in the different characteristics of the traffic patterns generated by the applications. The detailed explanation goes beyond the scope of this paper; however, we can give some clues. Local communications neither benefit from using many VCs, nor from adaptivity; this explains the behavior of BT. Patterns with intense, non-local interchanges can take advantage of many VCs, because its use reduces head of line blocking; this explains the behavior of IS. Adaptivity is useful when the utilization of network resources is not homogeneous, because it helps balancing the workload; this explains why CG improves significantly with adaptivity, and why IS does not. Remember that the improvements reported are for communication clearing time; they do not reflect actual whole-application level improvement.

6.2 Estimating Execution Times

In Sect. 4.3 we described a methodology to estimate the time an application would spend when executed in a “target” architecture different from the one used to capture the traces. As an example, we will estimate the time to execute the CG benchmark (class W.64) in three different scenarios. The original trace file was generated in the MareNostrum Supercomputer, [3] whose interconnection network is a fat-tree implemented with Myrinet-2000 adapters and switches. This network operates at link speed of 2 Gb/s. The compute nodes are PowerPC 970 at 2.3 GHz. Each MPI task runs in a different processor.

The target architectures are three 2D, 8×8 tori, whose network links work, respectively, at 100 Mb/s, 1 Gb/s and 10 Gb/s. We use the “Adaptive 3VC” configuration of virtual channels, as described in the previous experiment. We have not applied any scale to the CPU times, so these target architectures are supposed to use 2.3 GHz PowerPC CPUs. Results of the simulations, reporting estimated execution times for the three link speeds, are summarized in Table 1, along with the actual execution time in the MareNostrum.

Table 1 Actual execution time of CG.W.64 in the MareNostrum, and estimated times for three different target architectures. Times in seconds

MareNostrum	8 × 8 torus 100 Mb/s	8 × 8 torus 1 Gb/s	8 × 8 torus 10 Gb/s
0.54	3.49	0.55	0.25

We can see how CG, a very communication-intensive application, can take advantage of network improvements. The speedup when changing from a 100 Mb/s network to one running at 1 Gb/s is 6.35; this is because at low speeds most of the execution time is due to communication. The improvement when using the 10 Gb/s network, instead of the 1 Gb/s one, is not that great (2.20), but still notable; this is because at high speeds the execution time is more computation-bound.

Note that results obtained in the MareNostrum are just *indicators* of the peak performance reachable with this machine. In this computer we are not using a 64-node network, but a full Myrinet-2000 fat-tree capable of linking the more than 10 K nodes of this computer. The 64 CPUs used in the experiment were not necessarily consecutive; in fact, we don't know the exact placement used by the scheduler, which means that the CPUs could be located in different leaves of the fat-tree. Measurements were taken when the machine was in production, and other applications were running (and using the network) at the same time, so some degree of interference was present—at least in the network, because the compute nodes were used exclusively by our applications. Therefore, it should not look strange that the predictions for our 1 Gb/s network are so close to the measurements with the 2 Gb/s Myrinet-2000 network.

7 Related Work

In the literature we can find many papers discussing different aspects of trace capturing mechanisms and utilization of traces for performance evaluation—see for example [10]. For the specific topic of interconnection networks, Chapters 23 to 25 of [7] are of particular relevance. In addition to performance evaluation and prediction, other common use of traces of parallel applications is visualization, often as a help for debugging and detection of bottlenecks. The list of references would be very long, because this is a well-established area of work, so we focus on the main topics addressed in this paper.

7.1 Alternative Ways of Obtaining Trace Files

MPE is not the only way of acquiring traces of MPI applications. There are many other options available (see [13] for a review), but most of them are based on the instrumentation of application source code, or on the substitution of standard MPI functions by instrumented versions at compile time using the PMPI interface. As they work at the application level, they cannot be totally accurate regarding message send/receive times. An alternative way of getting traces would be to capture information at a lower level. The operating system, or a set of middleware daemons providing services to

running MPI applications, should be capable of recording the actual timestamps of communication events thus providing better timing information. An example of this approach is Sun's Dtrace [18].

Yet another way of generating a trace file in environments such as networks of workstations could be using a network *sniffer* that captures packet interchanges between computers. A programmable sniffer such as Wireshark [27] could be used to this purpose. A trace file generated this way would contain records with this information: $\langle \text{timestamp}, \text{origin}, \text{destination}, \text{data} \rangle$, with accurate timestamps. However, these tools capture network-level frames, so we do not know whether a burst of packets belongs to the same long message, or are a sequence of smaller ones. Also, we know when a packet has been delivered by the network to the receiving node, but not when it is actually available to the corresponding application process. We have a temporal order of records, but this is not a causal order, because inter-dependencies are not captured. An additional shortcoming of this approach appears when several processes share a single computer (a common scenario when we need traces for large systems): since message interchange between the processes that share a machine is done internally, the network is not used, so we do not have the associated trace records.

7.2 Performance Prediction Using Trace Files

The Dimemas tool, developed at the Technical University of Catalonia [8], can be used to carry out performance prediction studies using trace files (as well as machine descriptions) as its input. Note that the way we model the CPUs within INSEE is utterly simplistic, but the network is simulated with great detail. The Dimemas approach is the opposite: it accepts detailed descriptions of the compute nodes, so that a change in the system architecture is not simulated by simply scaling CPU states; however, the network model is very simple: a collection of parallel buses. Dimemas uses its own trace format and trace-capturing tools, that gather more information than that included in MPE's CLOG traces (for example, thread-level information, and hardware counters)—but that require kernel-level support, not always available. Trace records include not only MPI operations and communications, but also the states of each task. This means that Dimemas traces log when tasks are using the CPU, when they are blocked by other tasks (when sharing CPU) and when they are stalled for I/O operations.

Note that Dimemas was conceived as an application analysis tool that allows to inspect an application in order to locate undesirable behaviors (as bottlenecks or unbalancing between threads), so the high level of detail in CPU and application/task/thread modeling. In contrast, INSEE is used to evaluate interconnection networks, and obviously, the IN is modeled with high level of detail.

7.3 Extended Traces and Collective Support

In Sect. 3, we discussed the generation of extended trace files, as required by our FSIN simulator. It is important to remark that *we do not modify the default implementation of collectives included in MPICH*. These primitives are good for general use, but not optimized for any particular underlying communication fabric. Therefore, when we use

the extended trace files for evaluation purposes, we are testing a target machine with this particular implementation of collectives. The availability of this implementation is of great interest for us, because with it we can focus on the design and evaluation of the point-to-point abilities of the network. However, we know that a good portion of the design effort for a parallel computer should go to the supporting library, including an MPI library with customized collectives [1,4]. A fair assessment of a computer with support for collectives should be done using regular trace files. A well-known machine *with* collective support is the BlueGene [1]. In contrast, clusters built around Myrinet [14] networks *do not* include collective support—BSC’s MareNostrum is a remarkable example [3]. In Myricom’s implementations of MPICH (MPICH-GM on top of the older GM library, and MPICH-MX on top of the MX library [15]) implementation of collectives is not changed, using the default one provided in the original MPICH. They plan to include support for collectives in future releases of MX.

The Dimemas simulator is able to work with regular trace files, and deals with collectives using different models (simplifications). It assumes that the time to complete a collective is

$$time = (latency + (size/bandwidth)) * model_factor$$

where the *model_factor* can be 0 (null time), 1 (constant), N (linear) or an expression to model logarithmic times. These simplifications do not take into consideration the possibility of network contention, and are not universally valid for all hardware and/or software implementation of collectives. This is another proof of the differences in objectives between Dimemas and INSEE.

7.4 Other Simulation Tools

We end this section with a review of simulation tools for interconnection networks. SICOSYS [19], developed at the University of Cantabria, performs simulations of switching components with a high level of detail, providing timing information similar to that achieved using hardware simulators. Its large footprint does not allow it to simulate very large networks, but it is extremely useful for on-chip and on-board networks. SICOSYS can be fed with synthetic workloads and application traces, and can also be integrated with other tools to perform full system simulation.

The Flexim 1.2 simulator [26], developed at the University of Southern California, shares many design principles with FSIN. A main difference is that Flexim is designed for routers using wormhole switching, while FSIN uses virtual cut-through switching. Flexim supports synthesized traffic patterns or trace-driven traffic, although no details of the mechanism involved can be found in the documentation.

The Parallel Programming Laboratory at the University of Illinois at Urbana-Champaign maintains BigNetSim [6]. Its design is very different from INSEE, SICOSYS, or Flexim. It works alongside BigSim, a system emulator able to run applications specifically compiled for it. The emulator captures a collection of tasks (blocks of computation and communication) on a number of processors along with their dependencies and writes these tasks to trace files. BigNetSim reads the traces and

simulates the execution of the original tasks by elapsing time, satisfying dependencies, and spawning additional tasks by passing messages through a detailed network contention model. This generates corrected times for each event (advancing or delaying the simulation clock) which can be used to analyze its performance on the target machine. As happens with Dimemas, focus is more on application performance than on network analysis.

8 Conclusions

In this paper we have introduced a collection of tools and techniques to carry out evaluations, via simulation, of interconnection networks, using realistic workloads, provided by trace files obtained from the execution of actual applications. These techniques have been incorporated in the INSEE toolset.

Firstly, we needed to deal with the information contained in the trace files. As regular traces do not include the details of collective operations, we modified the trace capturing mechanism. Extended trace files allow us to deal only with simple, point-to-point message interchange records.

Then, we have pointed out the two main limitations of traces, namely the lack of accuracy in timing information, and the fact that a trace file includes only a possible valid outcome (record order) of a parallel program execution, but not the only valid one. We have explained that, for the evaluation of interconnection networks, we may ignore timing information, but not the causal relationships implicit in event order. We know that different orderings in event processing *may* be valid (respecting application semantics), but deciding about whether or not altering the order recorded in the trace file requires application-dependent knowledge, and our choice may lead to non-valid sequences of events. So the safest approach, for simulation purposes, is to follow, without exception, the exact event ordering of the trace file.

The trace processing ability integrated into INSEE offers a very flexible tool to evaluate different aspects of interconnection networks for parallel systems. As a way of showing the kind of work that can be carried out with this tool, we provide a simple but illustrative performance study: for some applications, the utilization of multiple virtual circuits per physical channel produces important performance gains, which can be even larger if using adaptive routing; the study shows the extent of the achievable gains. Additionally, we have shown how to estimate the execution time of an application running on a target architecture at different link speeds using traces captured in a real machine.

For future work, we plan to investigate how to integrate application semantics into simulation to allow different (but valid) event orderings.

Acknowledgements This work has been supported by the Spanish Ministry of Education and Science (TIN2007-68023-C02-02) and by the Basque Government (IT-242-07). Mr. Javier Navaridas is supported by a doctoral grant of the UPV/EHU. We gratefully acknowledge the utilization of resources and tools of the Barcelona Supercomputing Center/Centro Nacional de Supercomputación, Spain. We also want to thank the anonymous reviewers for their help improving this paper.

References

1. Almási, G., et al.: Optimization of MPI collective communication on BlueGene/L systems. In: Proceedings of the 19th Annual International Conference on Supercomputing, Cambridge, Massachusetts, June 20–22, ICS '05, pp. 253–262. ACM Press, New York, NY (2005)
2. Bailey, D.H., Harris, T., Van der Wignagaart, R., Saphir, W., Woo, A., Yarrow, M.: The NAS Parallel Benchmarks 2.0. Technical Report NAS-95-010, NASA Ames Research Center (1995)
3. Barcelona Supercomputing Center Home Page: Available at <http://www.bsc.es>. Accessed 1 Dec 2008
4. Brightwell, R.: A comparison of three MPI implementations for red storm. In: 12th European PVM/MPI Conference, September 2005. Lecture Notes in Computer Science, vol. 3666, pp. 425–432 (2005). doi:[10.1007/11557265_54](https://doi.org/10.1007/11557265_54)
5. Cámara, J., Moretó, M., Vallejo, E., Beivide, R., Martínez, C., Miguel, J., Navaridas, J.: Mixed-radix twisted torus interconnection networks. In: Proceedings of the 21st IEEE International Parallel & Distributed Processing Symposium—IPDPS '07, Long Beach, CA, March 26–30, 2007
6. Choudhury, N., Mehta, Y., Wilmarth, T.L., Bohm, E.J., Kalé, L.V.: Scaling an optimistic parallel simulation of large-scale interconnection networks. In: Proceedings of the 37th Conference on Winter Simulation, Orlando, Florida, December 04–07, 2005
7. Dally, W.J., Towles, B.: Principles and Practices of Interconnection Networks. Morgan-Kaufmann, Los Altos (2004)
8. Barcelona Supercomputing Center. DIMEMAS. http://www.bsc.es/plantillaA.php?cat_id=475. Accessed 1 Dec 2008
9. Gropp, W., Lusk, E., Doss, N., Skjellum, A.: A high-performance, portable implementation of the MPI Message-Passing Interface standard. *Parallel Comput.* **22**(6), 789–828 (1996). doi:[10.1016/0167-8191\(96\)00024-5](https://doi.org/10.1016/0167-8191(96)00024-5)
10. Karrels, E., Lusk, E.: Performance analysis of MPI programs. In: Proceedings of the Workshop on Environments and Tools for Parallel Scientific Computing (1994)
11. Magnusson, P.S., Christensson, M., Eskilson, J., Forsgren, D., Hillberg, G., Hgberg, J., Larsson, F., Moestedt, A., Werner, B.: Simics: a full system simulation platform. *Computer* **35**(2), 50–58 (2002). doi:[10.1109/2.982916](https://doi.org/10.1109/2.982916)
12. Message Passing Interface Forum: MPI: A Message-Passing Interface Standard. University of Tennessee. Available at <http://www.mpi-forum.org/>. Accessed 1 Dec 2008
13. Moore, S., Cronk, D., London, K., Dongarra, J.: Review of performance analysis tools for MPI parallel programs. In: Cotronis, Y., Dongarra, J. (eds.) 8th European PVM/MPI Users' Group Meeting, Lecture Notes in Computer Science, vol. 2131, pp. 241–248. Springer Verlag, Berlin (2001)
14. Myricom Home Page: Available at <http://www.myri.com>. Accessed 1 Dec 2008
15. Myricom Inc.: Myrinet Express (MX): A High-Performance, Low-Level, Message-Passing Interface for Myrinet, Version 1.1 January 01, 2006. Available at <http://www.myri.com/scs/MX/doc/mx.pdf>
16. Navaridas, J., Ridruejo, F.J., Miguel-Alonso, J.: Evaluation of interconnection networks using full-system simulators: lessons learned. In: Proceedings of the 40th Annual Simulation Symposium, Norfolk, VA, March 26–28, 2007
17. Navaridas, J., Ridruejo, F.J., Miguel-Alonso, J.: On synthesizing workloads emulating MPI applications. In: The 9th IEEE International Workshop on Parallel and Distributed Scientific and Engineering Computing (PDSEC-08), April 14–18, 2008, Miami, Florida, USA
18. OpenSolaris Community: Dtrace. Available at <http://opensolaris.org/os/community/dtrace/>. Accessed 1 Dec 2008
19. Puente, V., Gregorio, J.A., Beivide, R.: SICOSYS: an integrated framework for studying interconnection network performance in multiprocessor systems. In: Proceedings of the 10th Euromicro Workshop on Parallel, Distributed and Network-based Processing (EUROMICRO-PDP 2002) (2002)
20. Puente, V., Izu, C., Beivide, R., Gregorio, J.A., Vallejo, F., Prellezo, J.M.: The adaptative bubble router. *J. Parallel Distr. Comput.* **61**(9), 1180–1208 (2001)
21. Puente, V., Prellezo, J.M., Izu, C., Gregorio, J.A., Beivide, R.: A case study of trace-driven simulation for analyzing interconnection networks: cc-NUMAs with ILP processors. In: Proceedings of the IEEE 8th Euromicro Workshop on Parallel and Distributed Processing, Rhodes, Greece, January 2000
22. Ridruejo, F.J., Gonzalez, A., Miguel-Alonso, J.: TrGen: a traffic generation system for interconnection network simulators. In: 1st International Workshop on Performance Evaluation of Networks for

- Parallel, Cluster and Grid Computing Systems (PEN-PCGCS'05). ICPP 2005 Workshops, 14–17 June 2005
23. Ridruejo, F.J., Miguel-Alonso, J.: INSEE: an interconnection network simulation and evaluation environment. In: *Lecture Notes in Computer Science*, vol. 3648 (Proceedings of the Euro-Par 2005), pp. 1014–1023 (2005). doi:[10.1007/11549468_111](https://doi.org/10.1007/11549468_111)
 24. Ridruejo, F.J., Miguel-Alonso, J., Navaridas, J.: Concepts and components of full-system simulation of distributed memory parallel computers. In: *Proceedings of the HPDC '07*, June 25–29, 2007, Monterey, California, USA
 25. Ridruejo, F.J., Navaridas, J., Miguel-Alonso, J., Izu, C.: Realistic evaluation of interconnection network performance at high loads. In: *Proceedings of the 8th International Conference on Parallel and Distributed Computing Applications and Technologies – PDCAT 2007*, Adelaide, Australia, 3–6 December 2007
 26. SMART group at the U. of Southern California: Information on FlexSim1.2. Available at <http://ceng.usc.edu/smart/FlexSim/flexsim.html>. Accessed 1 Dec 2008
 27. Wireshark Home Page: Available at <http://www.wireshark.org/>. Accessed 1 Dec 2008
 28. Wolf, F., Malony, A., Shende, S., Morris, A.: Trace-based parallel performance overhead compensation. In: *Proceedings of the International Conference on High Performance Computing and Communications (HPCC)*, Sorrento, Italy, Sept. 2005
 29. Zaki, O., Lusk, E., Gropp, W., Swider, D.: Toward scalable performance visualization with Jumpshot. *High Perform. Comput. Appl.* **13**(2), 277–288 (1999). doi:[10.1177/109434209901300310](https://doi.org/10.1177/109434209901300310)