

# High-Speed LANs: New Environments for Parallel and Distributed Applications

Patrick Geoffray    Laurent Lefèvre    CongDuc Pham    Loïc Pryllif  
Olivier Reymann†    Bernard Tourancheau    Roland Westrelin

RHDAC, Université Lyon 1, Claude Bernard

†LIP, Ecole Normale Supérieure de Lyon

e-mail: bip-team@lhpc.univ-lyon1.fr

**Abstract.** As the technology for high-speed networks has incredibly evolved this last decade, the interconnection of workstations at gigabits rates and low prices has become a reality. These clusters, based on regular workstations (e.g. PCs), can now be used in place of traditional parallel computers with no possible comparison on the prices! In this article, 3 applications (high performance computing, distributed shared memory system and parallel simulation) that were traditionally executed on expensive parallel machines are ported on a Myrinet-based cluster of PCs. The results show that the performances of these new architectures can be very close to those obtained on state-of-the art parallel computers.

## 1 New technologies for parallel applications

For a long time parallel computers were the solution for people with high computation needs. If the processing units of such massively parallel processors can now be taken from the commodity market, the interconnection networks and the software are still highly customized. While sequential computers have always seen a dramatic cut down in their prices every year, parallel computers took the opposite direction because of the decreasing demand.

However, there are not so many solutions for having more computation power: one has to use several processors. The choice resides on how to make the parallel system: from custom or standard products. If the first choice was before justified by the low-quality of products from the commodity market, this is not the case any more. There is the possibility to go a complete step farther by building parallel systems entirely from standard components. Processing units are just regular workstations or PCs that can be bought in any supermarket for a few thousands of dollars, interconnection networks are also taken from the high-speed regular market of LANs (local area network) such as Gigabits Ethernet, ATM, and Myrinet. These architectures are often referred to as Network Of Workstations (NOW). Several research teams have launched projects dealing with NOWs used as parallel machines. Early experiments with IP-based implementations have shown disappointing performances so the goal of most of these teams is to design the software needed to make clusters built with commodity components and high-speed networks really efficient. The NOW project of UC

Berkeley [1] was one of the first one. Its main contribution is the Active Messages [12] layer that provides high performance access to the network. In this sense, it is very similar to the BIP software we are developing on our cluster.

The objective of the paper is two-folds: (*i*) to show experiments on different types of LANs with several communication softwares and (*ii*) to compare our Myrinet-based test-bed with the BIP software to traditional parallel computers. To do so, 3 applications covering very different areas of distributed and parallel computing are ported on NOWs with a focus on Myrinet-based networks. The paper is organized as follows: Section 2 presents the high-speed environment based on a Myrinet network and the BIP communication software. Section 3 compares state-of-the-art parallel computers and our Myrinet test-bed with the NAS benchmarks. Section 4, 5 and 6 presents the applications (high performance computing, distributed shared memory system and parallel simulations) that were ported on Myrinet. Conclusions are given in Section 7.

## 2 The high-speed LAN environment based on Myrinet

The Myricom LAN[3] target was chosen for its performance over the Gbits/s, its affordable price and its software openness (all the software and specifications are freely available for customers). There are several features that make this kind of technology much more suitable than a traditional commodity network:

- the hardware provides an end-to-end control flow that guarantees a reliable delivery and alleviates the problem of implementing in software the reliability on top on an lossy channel. As message losses are exceptional, it is possible to use algorithms that focus on very low overheads in the normal case,
- the interface card has a general purpose processor which is powerful enough to handle most of the communication activity without interrupting the main processor. In particular, it provides an efficient overlapping of communication and computation.
- the interface card has up to one megabyte memory for buffers. As the network being as fast as the computer bus, this memory isolates transfers on the network from transfer on the bus.

### 2.1 BIP

BIP stands for Basic Interface for Parallelism. The idea was to build it with a library interface accessible from applications that will implement a high speed protocol on the Myrinet network. BIP provides only a protocol with low level functionalities. Although specialized parallel applications could interface directly with it, it is intended as the base of other protocol layers like IP, and higher level APIs like the well established MPI and PVM. All the applications described later on in this paper actually use our MPI implementation based on top of BIP. Highly optimized, the raw communication performance for BIP is about  $5\mu s$  latency one-way. With large messages, the bandwidth goes up to 125MByte/s. The

MPI layer adds about  $5\mu s$  for small messages. As will show the other examples in the paper, this allows applications to scale more by allowing the use of a finer grain of computation. BIP is described in more details in [11].

## 2.2 Experimentation platform

Our test-bed for the experiments presented in the next sections consists of 8 nodes, each with an Intel Pentium Pro 200 MHz, 64MBytes of RAM, a 440FX chipsets and a myrinet board with LANai 4.1 and 256Kbytes of memory. The test-bed is provided by the LHPC (Laboratoire pour les Hautes Performances en Calcul), a cooperation between ENS-Lyon and Matra Système Information.

## 2.3 Running the NAS benchmarks

The NAS parallel benchmarks are a set of parallel programs designed to compare the performance of supercomputers. Each program try to test a given aspect of parallel computation that could be found in real applications. These benchmarks are provided both as single processor versions and parallel codes using MPI. We selected 3 different benchmarks and compiled them with MPI-BIP. They are then run on 1, 4 and 8 nodes of our Myrinet cluster described previously. Table 1 gives our measurements and comparisons with several parallel computers. Data for parallel computers come from <http://science.nas.nasa.gov/Software/NPB>.

**IS (Integer Sort)** sorts 8388608 keys distributed on the processors. The test uses a lot of small message communications and needs few processing power.

**LU** solves a finite difference discretization of the 3D compressible Navier-Stokes equations. A 2-D partitioning of the  $64 \times 64 \times 64$  data grid is done. Communication of partition boundary data occurs after completion of computation on all diagonals that contact an adjacent partition. We have a relatively large number of small communications of 5 words each.

**SP** solves three sets of uncoupled systems of equations of size  $64 \times 64 \times 64$ , first in the x, then in the y, and finally in the z direction. This algorithm uses a multi-partition scheme. The granularity of communications is kept large and a few messages are sent.

The speedup values summarized in table 1 show the very good performances of MPI-BIP and Myrinet. Super-linear speedups in our case can be explained by the large cache size and the super-scalar architecture of the PPro. Reported super-linear speedups for the Sun Enterprise show that a super-linear speedup is no exception. The favorable super-linear speedup for the LU test may be related to the fact that this test has been compiled with `f2c` that apparently produce a slower code than a commercial Fortran software. The behavior of MPI-BIP is excellent both for small messages and large messages. The speedup for the IS test with MPI-BIP is always better than those obtained on parallel machines. However, when a lot of computational power is needed (SP), a gap appears between traditional parallel machines and our cluster: the weak floating-point unit of the PPro 200 shows its limits!

	MPI-BIP on PPro 200		IBM SP (66/WN)		Cray T3E-900		SGI Origin 2000-195		Sun Enter- prise 4000	
	Mop/s	Speedup	Mop/s	Speedup	Mop/s	Speedup	Mop/s	Speedup	Mop/s	Speedup
IS @ 4 proc	2.44	4.5	2.2	3.1	3.2	N/A	2.1	3.8	1.6	3.8
IS @ 8 proc	2.11	8.8	2.0	5.6	3.8	N/A	2.4	8.6	1.5	6.9
LU @ 4 proc	23.68	6	59.0	3.7	67.6	N/A	96.8	N/A	36.9	4.1
LU @ 8 proc	22.73	11.6	57.2	7.1	66.4	N/A	103.3	N/A	37.3	8.4
SP @ 4 proc	10.10	3.4	42.1	3.6	43.0	N/A	60.3	N/A	25.0	3.7

**Table 1.** Performance for NAS Benchmarks on various platforms.

### 3 High performance computing: the Thesee application

*Thesee* is a 3D panel method code, which calculates the characteristic of a wing in an inviscid, incompressible, irrotational, and steady airflow, in order to design new paragliders and sails. Starting from a sequential version [6], *Thesee* has been parallelized using the ScaLAPACK[2] library routines to be run on NOWs. The parallelization process has been done in a systematic manner to keep the development cost low. 3 parts have been identified:

- **Part P1** is the fill-in of the element influence matrix from the 3D mesh. Its complexity is  $O(n^2)$  with  $n$  the mesh size (number of nodes). Each matrix element gives the contribution of a double layer (source + vortex) singularity distribution on facet  $i$  at the center of facet  $j$ .
- **Part P2** is the LU decomposition of the element influence matrix and the resolution of the associated linear system ( $O(n^3)$ ), in order to calculate the strength of each element singularity distribution.
- **Part P3** is the speed field computation. Its complexity is  $O(n^2)$  because the contribution of every node has to be taken into account for the speed calculation at each node. Pressure is obtained using Bernoulli equations.

Each of these parts are parallelized independently and are linked together by the redistribution of the matrix data. For each part, the data distribution is chosen so as to insure the best possible efficiency of the parallel computation. The ScaLAPACK library uses a block cyclic data distribution on a virtual grid of processors. This solution provides a good load-balance because the processors receive matrix elements from different locations of the original matrix (as opposed to a classical full block decomposition). Communication overheads are reduced to a minimum by preserving the row and the column shape of the matrix (most of communication of 1D arrays can then happen without a complex index computation). Tests are run beforehand to choose the best grid shape and the best block size of the data distribution for our problem on each of the platforms. In the parallel version of *Thesee*, we used the following parameters for the data distribution of the LU factorization: the block size is  $32 \times 32$  and the processor grid shape is a 1D-grid that gives the best overall computation timings.

### 3.1 Performance results

We ran tests with 3 different interconnection networks on 2 different platforms: (i) Ethernet and ATM network of SUN Sparc5 85MHz with Solaris and (ii) Ethernet and Myrinet network of Pentium-Pro 200MHz under Linux. On both systems, PVM and MPI were used. For IP over Myrinet, the LAM implementation of MPI is used. Otherwise it is the MPI-BIP user level implementation based on MPICH. The overhead for MPI or PVM is similar regarding the small set of primitives involved.

The efficiency of the fill-in of the influence matrix (part 1) and the speed field computation (part 3) is roughly the same on each configuration. The code for these parts is “embarrassingly” well suited for parallel execution and thus the speed-up is almost linear. On the other hand, the LU factorization that involves a lot more communications is highly dependent on the network software and hardware performances.

For space consideration, only results for Ethernet vs. Myrinet are shown in the paper. For Ethernet and ATM, the test shows that high-speed networks designed for long-distance communication are not well-suited to system-oriented communication. Although ATM provides more throughput, the gain obtained is small because the initialization time for each communication on this network is similar to the one on Ethernet. Unfortunately, this initialization time represents the larger part of the communication delay.

System size		Timings (seconds) [speedups]		
		PVM/IP/Ethernet (on SUN Sparc)	PVM/IP/Myrinet (on Pentium Pro)	MPI/BIP/Myrinet (on Pentium Pro)
Sequential	902 × 902	10.0		
2 Proc	902 × 902	9.6 [1.03]	6.7 [1.49]	5.0 [1.97]
4 Proc	902 × 902	10.2 [0.97]	4.7 [2.10]	3.1 [3.24]
Sequential	1722 × 1722	87.4		
2 Proc	1722 × 1722	56.4 [1.54]	44.3 [1.97]	38.1 [ <b>2.29</b> ]
4 Proc	1722 × 1722	45.9 [1.90]	28.1 [3.10]	21.3 [ <b>4.09</b> ]

**Table 2.** Timings and speedups with Myrinet (using the best block size).

**Comparison between Ethernet and Myrinet** Table 2 presents the timings results obtained on the Myrinet test-bed along with the BIP software. Since the communication/computation ratio is rather high, we expect better performances than those obtained with ATM. The best results are achieved with MPI-BIP and we can see that IP-based implementation can not fully exploit the low latencies of the Myrinet hardware. With MPI-BIP, the low latency for a basic send communication ( $9\mu s$  for this case) has an incredible impact on the speedup when compared to the Ethernet run. Super-linear speedup with 4 processors can be explained by a better cache hit ratio in the parallel version of the code. As the matrix is distributed cyclically on the processors, the computation occurs on blocked data that fits better in the cache during the LU decomposition, leading to a better use of the processor’s pipeline units. One more reason is an

increase in the overlapping of computation over communications in the parallel LU decomposition since the communication cost is greatly reduced with the BIP/Myrinet platform.

## 4 Distributed Shared Memory System (DSM)

The purpose of DSM is to implement, on top of a distributed memory architecture, a programming model allowing a transparent manipulation of virtually shared data. Thus, in practice, a DSM system has to handle all the communications and to maintain the shared data coherence. We have developed an object-based DSM system called DOSMOS (**D**istributed **O**bjects **S**hared **M**em**O**ry **S**ystem) that allows processes to share in a transparent way a set of objects distributed and replicated over distant processors. DOSMOS integrates novel features:

- **DOSMOS Processes:** a DOSMOS application is composed of two types of processes: **Application processes (A.P.)** contain and execute application code; **Memory processes (M.P.)** manage the whole DSM system, i.e. they provide A.P. with the objects they request and maintain data coherence.
- **Array allocation:** DOSMOS allows to manipulate both basic type variables (integer, float, char...) and distributed arrays which can be split into several “system objects”, replicated among the processors. Various splittings are provided: by row, by column, by block and by cyclic block.
- **Weak consistency protocols:** for efficiency and scalability purposes, DOSMOS gives the opportunity to duplicate shared objects. These replicas have to be kept coherent. DOSMOS implements a weak protocol: the release consistency which provides two synchronization operators: *acquire* and *release*.
- **Hierarchical structuring of the application processes:** processes can be grouped into groups and sub-groups in order to optimize the management of the data coherence.

Previously developed on top of PVM and experimented on Ethernet networks [8], DOSMOS is now available on top of MPI and experimented on MPI-BIP on top of Myrinet networks.

### 4.1 Gram-Schmidt Application

We based our experiments on the Gram-Schmidt application which has been completely studied on Ethernet in [4]. We propose four parallel implementations of the Gram-Schmidt application (Fig. 1) :

- Trivial version: with no splitting of matrix and use of strong consistency protocols. Accesses to shared object are sequentialized.
- Release Consistency and Object Splitting: this version takes the benefit of splitting the large matrix in small adapted shared objects.

- Synchronization 1: also based on Release Consistency protocols and matrix splitting but the matrix object is also used for synchronization of processes.
- Synchronization 2: this optimized version requires an independent object to synchronize processes.

### 4.2 Experiments

Figure 1 demonstrates the benefit of using a high-speed LAN for a DSM system like DOSMOS. Improvements added to the original algorithm are not linked with the application performances. Due to low latencies provided by the Myrinet network combined with MPI-BIP, applications with small improvements provide better speedups than highly optimized versions (synchro2). However, fast communications are not enough for trivial applications (*trivial version* in Fig. 1). Figure 2 shows the scalability provided by DOSMOS on top of MPI-BIP. By increasing the problem size, we also increase the efficiency of applications which benefit from large bandwidth provided by the Myrinet network.

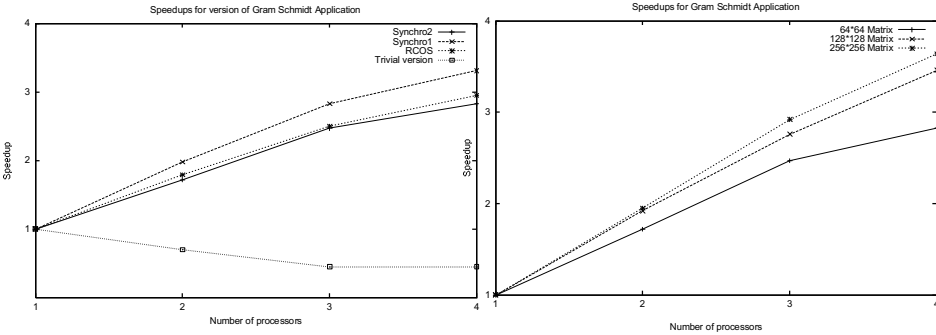
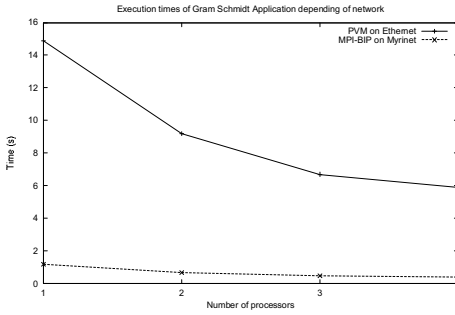


Fig. 1. Various improvements of Fig. 2. Speedup of Gram-Schmidt by increasing the problem size.

DSM systems require low-latency networks to provide high performance to applications. Figure 3 shows that DOSMOS can provide high performance distributed objects when combined to low-latency communication protocols like MPI-BIP (the sequential time is computed by using one DOSMOS AP and one MP). However, management of distributed objects (consistency, synchronization...) adds a latency to each distant access combined with the latency of the network (Figure 4). Distant operations concern A.P. which access a shared object managed by their M.P. while Long Distant operations involve one A.P. and two M.P. to access distant objects. The access to a distant object requires a ping-pong communication between two nodes (which takes around 25µs with MPI-BIP). An additional latency is added by DOSMOS to manage complex objects (like arrays) to provide splitting and a transparent access to distributed arrays. As can be seen, the latency added by DOSMOS is kept reasonable and still compatible with the need of high performances by the applications.



**Fig. 3.** Gram-Schmidt with Ethernet and Myrinet Networks.

Operation	Latency
Distant Read Simple Obj.	55
Distant Read Complex Obj.	140
Long Dist. Read Simple Obj.	103
Long Dist. Read Complex Obj.	340
Long Distant Acquire	95

**Fig. 4.** Latency in  $\mu\text{s}$  added by DOSMOS to shared object access.

By using a DSM model, programming a parallel application is made easier but the added overheads usually penalizes the users. These experiments on a high-speed LAN show that designing a parallel application with a DSM system like DOSMOS can merge together ease of programming with high performances.

## 5 Parallel simulation of communication networks

With the emerging of Asynchronous Transfer Mode (ATM) technology, the performance evaluation of future large scale B-ISDN networks based on ATM is a great challenge. Parallel simulation techniques have been proposed during the last 2 decades to reduce the simulation time. These methods fall in two categories: the conservative approach [5] and the optimistic approach [7].

The parallelization process of the network model consists in spatially partitioning the entire network in  $n$  regions and to assign each region to a processor. Cells transfer from one region to another are represented by timestamped messages exchanged between the processors. The messages synchronization is performed conservatively. The system under study includes a routing algorithm with the introduction of link cost and dynamic routing functions that provide load-balancing routing features. The simulation of a routing algorithm implies to simulate (i) the mechanism that consists in constructing and updating the routing tables and (ii) the flow of cells that is handled by the network [10].

From the parallel simulation perspective, this application presents a very small granularity and requires a lot of small messages to be exchanged between processors (as opposed to more traditional high-performance computing applications where the computation part is larger). These messages are either model messages, e.g. cell exchanges between ATM switches, or synchronization messages, e.g. null-messages from the conservative kernel. Therefore the application typically requires very low latencies from the communication system. Previous tests on LANs such as Ethernet were quite disappointing and so far only parallel computers were capable of showing interesting speedups.



## 5.1 Experimental results

In order to compare parallel computers and NOWs, experimental results on a Cray T3E and a Myrinet cluster are presented. The Cray is provided by the *Institut du Développement et des Ressources en Informatique Scientifique* (IDRIS) and consists of 256 processing nodes interconnected by a very low latency 3D torus. Each node on the Cray is a DEC Alpha EV5 with 128 Mo memory running at 300Mhz. The communication libraries are the native SHMEM (SHared MEMory) and MPI. The first one provides a latency of approximately  $7\mu s$  while MPI shows a latency of  $13\mu s$ . The experimental NOW consists of Pentium Pro 200MHz interconnected by the Myrinet network described previously. The communication stack we used is native BIP and MPI-BIP. The latency is about  $10\mu s$  with BIP and  $17\mu s$  with MPI-BIP (for the message size of our application, 70 bytes). Table 3 summarizes the results for a 78-switch network model. The simulation time has been set to 500,000 time slots that represent 0.31s of the real system. More than 50 millions of events are simulated.

<i>np.</i>	<i>time SHMEM (s)</i>	<i>speedup SHMEM</i>	<i>time MPI (s)</i>	<i>speedup MPI</i>
1	488	-	488	-
4	131	3.72	253	2.07
8	72	6.77	133	3.66
<i>np.</i>	<i>time BIP (s)</i>	<i>speedup BIP</i>	<i>time MPI (s)</i>	<i>speedup MPI</i>
1	321	-	321	-
4	139	2.30	210	1.52
8	59	5.44	88	3.64

**Table 3.** Comparison between Cray T3E and Myrinet cluster.

As can be seen, the sequential version run faster on a Pentium Pro 200MHz than on a DEC Alpha EV5 300MHz. The main advantage of using standard products is the availability of the most recent processors. As we stated before, parallel computers were interesting because they typically use customized high performance interconnection networks. The low latency on the Cray has a direct impact on the performances of the simulator. Using BIP on Myrinet exploits the small latency provided by the network and therefore the parallel simulation shows very interesting speedups compared to the cost of the hardware used. We have no doubt that this kind of architecture is very promising because of its excellent performance/price ratio.

## 6 Conclusions

Networks of workstations represent a serious alternative to expensive parallel computers. In this paper, we mainly focus on a Myrinet-based cluster with the BIP software for optimized low-level communications. Although the technology has been ready for some time, it is still used only in a confidential manner outside of the high-speed network research community. Also, even if a lot of people use

clusters, a vast majority of them stay with the traditional IP-based implementations of PVM or MPI. Experiments have shown that if the hardware technology has the potential for high performance communications, the communication software layer must be well designed to fully deliver the maximum of performances to the application. One of the principal aims of this paper is to show through a number of applications the maturity of the more recent technologies, and to provide data to help them widespread to the end-users. Generalizing the use of the faster systems available will provide to the end-user community a way to reach a higher level of scalability, and to make parallel solutions usable for a wider range of applications, especially those that require a fine grain decomposition.

## References

1. Anderson, T. E., Culler, D. E., Patterson, D. A., the NOW Team: The Case for Networks of Workstations. *IEEE Micro Magazine*, February 1995.
2. Blackford, Choi, Cleary, d'Azevedo, Demmel, Dhillon, Dongarra, Hammarling, Henry, Petitet, Stanley, Walker, and Whaley: *ScaLAPACK Users' Guide*. SIAM, 1997. <http://www.netlib.org/scalapack/>.
3. Boden, Cohen, Feldermann, Kulwik, Seitz, Seizovic, and Su. MYRINET: A Gigabit per second Local Area Network. *IEEE-Micro*, 15:29–36, February 1995.
4. Brunie, L., Restivo, N., Reymann, O.: Programmation d'applications parallèles sur systèmes à mémoire distribuée partagée et expérimentations sur réseaux hautes performances. In *Calculateurs Parallèles*, **9**(4), pages 417–433.
5. Chandy, K., Misra, J.: Distribution Simulation: A Case Study in Design and Verification of Distributed Programs. *Trans. on Soft. Eng.*, **5**(5) 440–452.
6. Giraudeau L., Perrot G., Petit S., Tourancheau B.: 3-d air flow simulation software for paragliders. TR 96-35, LIP-ENS Lyon, 69364 Lyon, France, 1996.
7. Jefferson, D. R. : Virtual Time. *ACM Trans. on Prog. Lang. and Sys.*, **7**(3) (July 1985) 405–425.
8. Lefèvre, L.: Parallel programming on top of DSM Systems : An Experimental Study. In *Parallel Computing - Environments and Tools for Parallel Scientific Computing III*, **23**(1-2), pages 235–249. April 1997.
9. Mainwaring, A. M., Culler, D. E.: Active messages: Organization and applications programming interface. <http://now.cs.berkeley.edu/Papers/Papers/am-spec.ps>, 1995.
10. Pham, C. D., Brunst, H., Fdida, S.: Conservative Simulation of Load-Balanced Routing in a Large ATM Network Model. In *Proceedings of PADS'98*, May 26-29 1998, Banff, Canada, pp142-149.
11. Prylli, L., Tourancheau B.: BIP: a new protocol designed for high performance networking on myrinet. In *Parallel and Distributed Processing, IPPS/SPDP'98*, bf 1388 of *LNCS*, pages 472–485. Springer-Verlag, April 1998.
12. von Eicken, T., Culler, D. E., Goldstein S. C., Schauser, K. E.: Active Messages: a Mechanism for Integrated Communication and Computation. *Proc. of the 19th Int'l Symp. on Comp. Architecture*, May 1992.