

Performance evaluation of DECK combining multithreading and communication on Myrinet and SCI clusters*

César De Rose[‡] Fábio Oliveira^{¶†} Fausto Blanco[‡] Marcos Barreto[¶]
Philippe Navaux[¶] Rafael Ávila[¶] Tiago Ferreto[‡]

[¶] Institute of Informatics
Federal University of Rio Grande do Sul
PO Box 15064 — 91501-970 Porto Alegre, Brazil
<http://www-gppd.inf.ufrgs.br/projects/mcluster>
E-mail: {fabreu,barreto,navaux,avila}@inf.ufrgs.br

[‡] High Performance Research Centre (CPAD)
Catholic University of Rio Grande do Sul
Av. Ipiranga, 6681 — 90619-900 Porto Alegre, Brazil
<http://www.cpad.pucrs.br>
E-mail: {derose,blanco,ferreto}@cpad.pucrs.br

Abstract

Paper submitted to PDPTA'01

This paper presents a performance evaluation of DECK (Distributed Execution and Communication Kernel), a multithreaded parallel programming environment for clusters of SMPs, with the parallel implementation of the classical Mandelbrot fractal generation and Laplace's Equation algorithms. The applications have been run on Myrinet and SCI clusters and the results are compared to corresponding MPI implementations. The comparison shows that DECK is able to achieve very good performance when multithreading is combined with communication, without the need of multiple processes on a single machine.

Keywords: cluster computing, multithreading, parallel programming environment, Myrinet, SCI.

1. Introduction and context

Current high-speed communication technologies such as Myrinet and SCI are usually interfaced by low-level programming libraries such as GM [4], BIP [10] and SISI [3], which have been introduced in order to efficiently exploit the low latency and high throughput offered by those technologies. However, such libraries present a rather complex API and are not targeted at the end user, but instead they serve as a basis for the implementation of higher level environments such as MPI [8]. Still, one of the most common characteristics of clusters, which is the availability of SMP

nodes, is usually not taken advantage of, and frequently not even supported, given that many MPI implementations are not even thread-safe.

Following this idea, we have designed and implemented DECK (Distributed Execution and Communication Kernel) [2], an environment for parallel programming on SMP clusters. The main goal of DECK is to try to efficiently integrate multithreading and communication in a single environment. In this paper we present a performance evaluation of two flavours of DECK, running on Myrinet and SCI clusters, with the parallel implementation of two algorithms usually present in the HPC community, namely the Mandelbrot fractal generation [7] and the solving of Laplace's Equation [9]. In addition, we compare the obtained results to those of equivalent implementations executed with versions of MPI for both communication technologies.

The paper is structured as follows: in Section 2 we present some background information with an overview of DECK and its implementations for Myrinet and SCI; in Section 3 we describe the applications, their implementation on DECK and MPI and analyse the obtained results; finally, Section 4 brings the authors' conclusions and future directions.

2. The DECK environment

The API of DECK is conceptually object-oriented, providing basic abstractions for parallel applications, such as threads and mail boxes, and more elaborate services like naming and collective communication. Figure 1 shows the internal structure of DECK. The bottom layer is called μ DECK, being responsible for the basic abstractions: threads, semaphores, messages and mail boxes. The upper layer of DECK is a service layer, whose services can be

*Results for Myrinet have been measured at the High Performance Research Centre (CPAD), PUCRS/HP; results for SCI have been measured at the Institute of Informatics, UFRGS

[†]Presenter

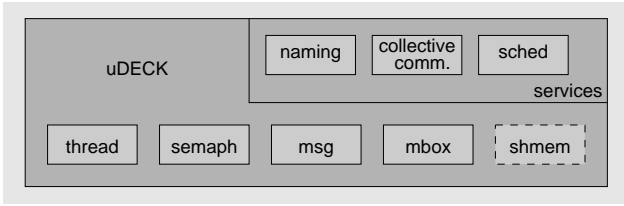


Figure 1. The internal structure of DECK.

chosen at compilation time.

Two implementations of DECK are currently available, for both Myrinet and SCI communication technologies. Multithreading is implemented on both by making use of standard POSIX Threads [6] calls. The communication issues are described next.

2.1. DECK/BIP

For the implementation of DECK on Myrinet [1] we have used the BIP library, mainly because of its performance and simple API. BIP provides basic point-to-point communication primitives (asynchronous *send()* and *recv()*) on Myrinet, with the constraint that the exchange of large messages (in BIP terms) must follow a *rendez-vous* semantics, which means that a given *recv()* must always be posted before the corresponding *send()*. In order to accomplish that in DECK, we have used a handshaking protocol where the sender makes use of small request messages before sending a large message. Such requests are handled, on the receiver side, by a dedicated thread, the *rv-daemon*, created at initialisation time.

This additional thread may have a significant impact on communication performance depending on the communication pattern. Figure 2 shows a raw performance evaluation of DECK and MPI-BIP [12], an MPI implementation on top of BIP. The graph shows the influence of the *rv-daemon* for messages larger than 1K, where the performance of DECK is around only 80% of that of MPI.

2.2. DECK/SCI

The SCI version of DECK makes use of the SISCO API. SISCO—Software Infrastructure for SCI—is a specification of standard primitives for SCI programming, based on the shared segment notion, proposed by a group of partners from both the academia and the industry. One implementation of SISCO is available in the SSP (Scali Software Platform), a software package distributed with Scali Wulfkits [11], upon which our SCI cluster is based.

Three communication protocols were developed for the implementation of DECK on SCI, in order to guarantee low latency and minimal overhead for small messages, as well

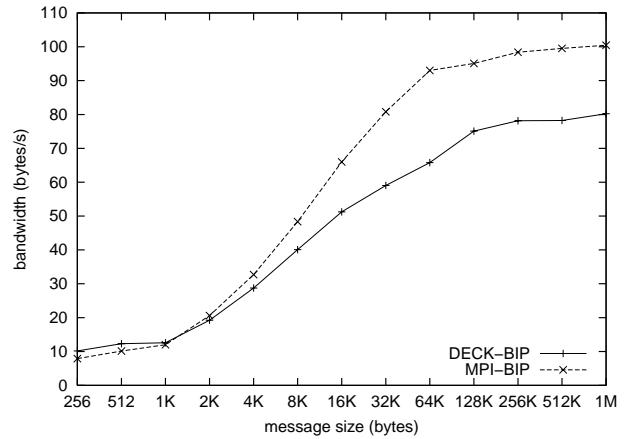


Figure 2. Bandwidth of DECK and MPI-BIP on Myrinet.

as high bandwidth for large ones. A protocol specialised in exchanging small messages achieves low latency by integrating communication and signalling into the same SCI packet, whereas high bandwidth is accomplished through a carefully designed zero-copy communication mechanism.

Figure 3 shows the raw performance evaluation of DECK/SCI and ScaMPI, the Scali implementation of MPI [5]. The graph reveals that the performance of DECK/SCI is clearly better than that of ScaMPI. It is possible, with DECK/SCI, to achieve a minimal latency of $4.66 \mu\text{s}$ and a maximum bandwidth above 84 Mbytes/s. These results are really near the raw performance of the SCI network.

In contrast to DECK/BIP, DECK/SCI is more efficient than MPI for all message sizes. This is, in part, due to the fact that DECK/SCI does not use any additional thread to manage communication.

3. The implemented applications

3.1. Execution environment

In order to investigate the potential and the feasibility of the proposed parallel programming environment for the target system, we have implemented two parallel algorithms that make heavy use of both computation and communication on the target system. These applications have been run on two SMP clusters available at our universities:

- a Myrinet cluster at the High Performance Research Centre (CPAD — PUCRS/HP); the cluster is composed of 16 nodes, each one being a HP E60 NetServer with two Intel Pentium III 550MHz processors and 128M DRAM, resulting in a 32 processor system; the

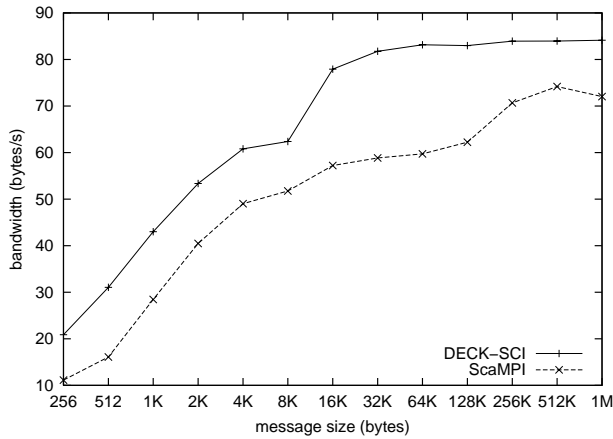


Figure 3. Bandwidth of DECK and ScaMPI on SCI.

nodes are connected by a 1.28 + 1.28Gbit/s Myrinet network based on Lanai7 32-bit cards; the operating system is Slackware Linux 7.1 and BIP is used to drive the Myrinet boards. The MPI implementation used in the tests is MPI-BIP

- an off-the-shelf SCI cluster available at the Institute of Informatics of UFRGS¹, composed of 4 bi-processed Pentium III 500MHz nodes, each with 256M of RAM (totalising 8 processors); the nodes run Conectiva Linux² version 5.0 and Scali Software Platform version 2.0.2, which includes ScaMPI and SISCI.

All the results have been obtained using the maximum number of nodes on each cluster, each figure being a mean of 10 executions.

3.2. Mandelbrot fractal

This application calculates the Mandelbrot set and draws a two-dimensional picture of it. It iterates the Mandelbrot function for every pixel on the picture, and sets the colour of the pixel according to the iteration where the orbit “escapes”, with a maximum iteration value of 17500. There are no dependencies between the calculations done in different regions of the picture but different regions may be easier to calculate than others.

Our parallel version of the application uses a master/slave model where the master is responsible for distributing parts of the picture to be calculated and for drawing the already calculated parts. The picture is partitioned in horizontal slices and slaves keep requesting new slices to

¹Partially financed by FAPERGS, FINEP and CNPq

²A Brazilian derivative of Red Hat Linux.

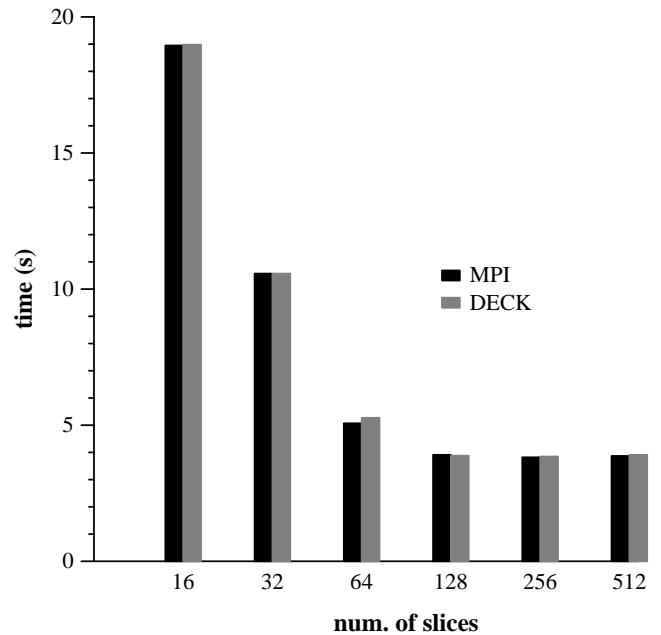


Figure 4. Results for the Mandelbrot algorithm on Myrinet.

calculate until all slices are ready. Slaves communicate only with the master in a 1:n pattern.

Figures 4 and 5 present the results for a 600×600 fractal calculated with different number of slices, varying from 16 to 512. A higher number of slices results in smaller messages and better load balancing but increases the number of messages.

On MPI we have run 2 processes on each node, and on DECK the implementation uses two threads for the calculations and one dedicated to communication.

The results obtained with DECK are practically equivalent to those of MPI, with differences lying mostly on the tenths of seconds. On Myrinet the best results are obtained for a number of slices between 128 and 256. Notice that in this application the influence of the *rv-daemon* does not affect the final results due to the more irregular communication pattern, confirming that the use of threads can be suitable for applications involving high-performance communication.

One can notice a slight increase in the execution time when 512 slices are used; we believe this is a tendency from this point on, since slices become too small to compensate the time spent to send them between nodes. In other words, the grain of parallelism in this case is too small.

On SCI the behaviour is the same, with execution times also starting to increase at 512 slices (the absolute values are higher since the SCI cluster has only 4 nodes). DECK is slightly better than MPI for 64 and 128 slices.

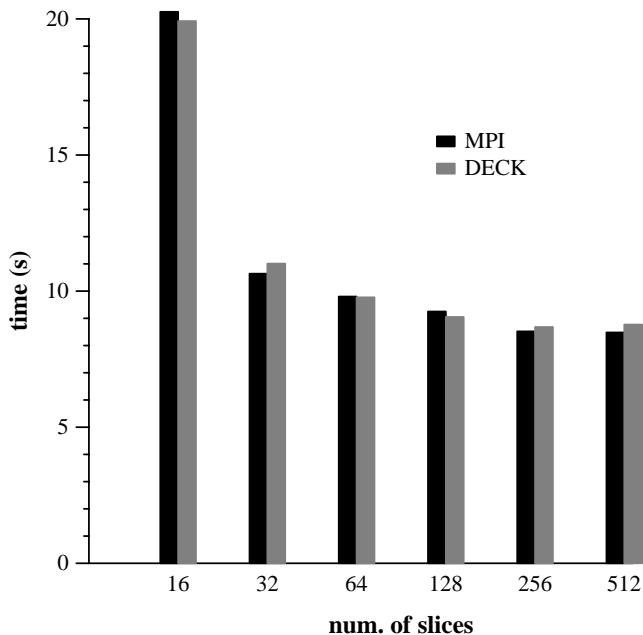


Figure 5. Results for the Mandelbrot algorithm on SCI.

3.3. Laplace equation

This application calculates the temperature distribution in a slab of a hypothetical homogeneous material completely insulated on the edges. Initially, the slab is at one uniform temperature throughout and a heat source is applied to one of the borders. The Laplace equation is used to solve this problem, being applied by means of an iterative method. The surface of the slab is divided in square sections, where each intersection is a point in a grid. The finer the grid, the more accurate the approximation, and the larger the problem.

From the adaption of the Laplace equation to a computational method, in this case the Gauss-Seidel approach [9], the temperature of a given grid point at a given iteration is taken to be the average of the temperatures of the four surrounding grid points at the previous iteration. Border temperatures are always kept constant.

Our parallel version of the application partitions the image in rectangular regions, assigning each region to a node. Figure 6 presents the results for DECK and MPI for a slab of size 600×600 .

Execution times are shown for both Myrinet and SCI distributing the regions equally among the available nodes. The same approach for achieving multiprocessing has been used, i.e., on MPI each machine executes two MPI processes, and on DECK two calculating threads are created.

In this application DECK has shown more encouraging

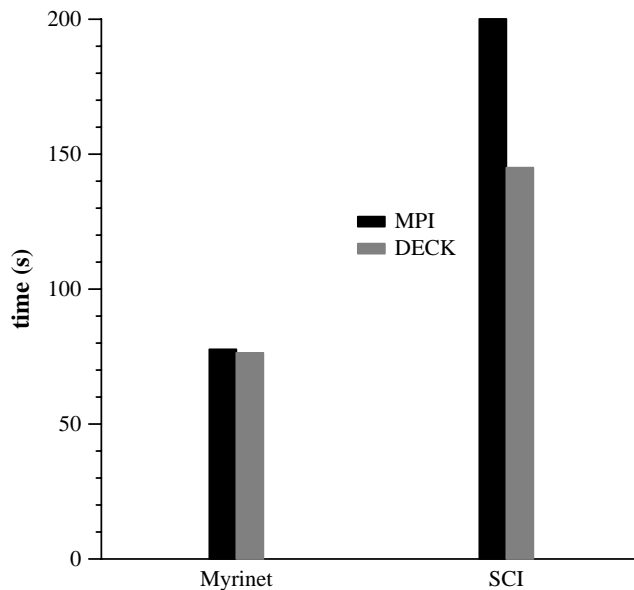


Figure 6. Results for Laplace's Equation.

results, presenting the same results as MPI on Myrinet and significantly better results on SCI. To our understanding, this behaviour is caused by different synchronisation semantics which depend on the message sizes. DECK messages are always asynchronous up to 8Kbytes, and the zero-copy protocol is used from this point on. As can be observed in Figure 3, ScaMPI keeps a rather constant evolution in performance throughout the curve, while DECK boosts bandwidth starting at 8Kbytes. Besides, to our knowledge, ScaMPI does not make use of a zero-copy mechanism, employing a *rendez-vous* (synchronous) protocol instead.

4. Conclusions and future directions

The experiments described in this paper have shown a performance evaluation of two implementations of the DECK environment when compared to equivalent MPI versions. In general, the performance presented by DECK with multiple threads is at least as good as that presented by the analysed MPI implementations, which represents an interesting alternative for the exploitation of SMPs in parallel applications.

In addition, in some cases, depending on the communication pattern presented by the application, DECK is able to outperform MPI running with multiple processes per node, as shown in Laplace's Equation. On the other hand, the raw performance evaluation of DECK on both technologies does not correspond to the results measured for the Mandelbrot algorithm.

Our next directions are the fine-tuning of DECK/BIP, with the goal of achieving better performance, and the im-

plementation of additional applications to make a more complete validation of our programming environment. The group has the intention, in the future, of allowing the user to join both Myrinet and SCI technologies in a single application.

5. Acknowledgements

This work has been partially supported by grants from CAPES and CNPq.

References

- [1] M. Barreto, R. Ávila, R. Cassali, A. Carissimi, and P. Navaux. Implementation of the DECK environment with BIP. In *Proc. of the First Myrinet User Group Conference*, pages 82–88, Lyon, France, 2000. Lyon, INRIA Rocquencourt.
- [2] Marcos Ennes Barreto. DECK: Um ambiente para programação paralela em agregados de multiprocessadores. Master's thesis, PPGC da UFRGS, Porto Alegre, 2000.
- [3] F. Giacomini, T. Amundsen, A. Bogaerts, R. Hauser, B. D. Johnsen, H. Kohmann, R. Nordstrøm, and P. Werner. Low-level SCI software requirements, analysis and predesign. Technical report, ESPRIT Project 23174 — Software Infrastructure for SCI (SISCI), May 1998.
- [4] GM. Available at <http://www.myri.com/GM>, December 1999.
- [5] L. P. Huse, K. Omang, H. Bugge, H. Ry, A. T. Haugsdal, and E. Rustad. ScaMPI—design and implementation. In Hermann Hellwagner and Alexander Reinefeld, editors, *SCI: Scalable Coherent Interface: Architecture and Software for High-Performance Compute Clusters*, volume 1734 of *Lecture Notes in Computer Science*, pages 249–261. Springer, Berlin, 1999.
- [6] IEEE. Information technology—portable operating system interface (POSIX), threads extension [C language]. IEEE 1003.1c-1995, 1995.
- [7] Benoit B. Mandelbrot. *The Fractal Geometry of Nature*. W. E. Freeman and Company, New York, 1982.
- [8] MPI Forum. The MPI message passing interface standard. Technical report, University of Tennessee, Knoxville, April 1994.
- [9] W. H. Press et al. *Numerical Recipes in C: The Art of Scientific Computing*. Cambridge University, Melbourne, second edition, 1994.
- [10] Loïc Prylli and Bernard Tourancheau. BIP: A new protocol designed for high performance networking on Myrinet. In José Rolim, editor, *Proc. of PC-NOW'98*, volume 1388 of *Lecture Notes in Computer Science*, pages 472–485. Berlin, Springer, 1998.
- [11] Scali homepage—scalable Linux systems—affordable supercomputing. Available at <http://www.scali.com>, April 2000.
- [12] Roland Westrelin. Une implémentation de MPI pour réseaux locaux à très haut débit: MPI-BIP. In *Proc. of the 11th RENPAR (Rencontres Francophones du Parallélisme)*, Rennes, 1999.