

MPI for Python: Performance improvements and MPI-2 extensions

Lisandro Dalcín*, Rodrigo Paz, Mario Storti, Jorge D'Elía

Centro Internacional de Métodos Computacionales en Ingeniería (CIMEC), Instituto de Desarrollo Tecnológico para la Industria Química (INTEC), Consejo Nacional de Investigaciones Científicas y Tecnológicas (CONICET), Universidad Nacional del Litoral (UNL), (S3000GLN) Santa Fe, Argentina

Received 16 November 2006; received in revised form 20 July 2007; accepted 12 September 2007

Available online 7 October 2007

Abstract

MPI for Python provides bindings of the message passing interface (MPI) standard for the Python programming language and allows any Python program to exploit multiple processors.

In its first release, MPI for Python was constructed on top of the MPI-1 specification defining an object-oriented interface that closely followed the MPI-2 C++ bindings, and provided support for communications of general Python objects. In the latest release, this package is improved to enable direct blocking/non-blocking communication of numeric arrays, and to support almost all MPI-2 features.

Improvements in communication performance have been tested in a Beowulf class cluster. Results showed a negligible overhead in comparison to compiled C code.

MPI for Python is open source and available for download on the web (<http://mpi4py.scipy.org/>).

© 2007 Elsevier Inc. All rights reserved.

Keywords: Message passing; MPI; High-level languages; Parallel Python

1. Introduction

During the last decade, high-performance computing has become an affordable resource to many more scientists and engineers than ever before. The conjunction of quality open source software and commodity hardware strongly influenced the now widespread popularity of dedicated Beowulf [4] class clusters and cluster of workstations. Message-passing has proven to be an effective computational model, specially suited for (but not limited to) distributed memory architectures. Although portable message-passing parallel programming used to be a nightmare in the past because of the many incompatible options developers were faced with, this situation definitely changed after the MPI Forum [17] released its standard specification, which rapidly gained widespread acceptance.

At the same time, the popularity of scientific computing environments such as MATLAB and IDL has increased considerably. Users simply feel more productive in such interactive environments with tight integration of simulation and

visualization. They are alleviated of low-level details associated to compilation/linking steps, memory management and input/output of traditional programming languages like Fortran, C, and C++. However, native support for parallel processing is absent and motivated different approaches to overcome it [14,15].

Recently, the Python programming language has attracted the attention of many users and developers in the scientific community. Python offers a clean and simple syntax, is a very powerful language, and allows skilled users to build their own computing environment, tailored to their specific needs and based on their favorite high-performance Fortran, C, or C++ codes [16]. Sophisticated but easy to use and well integrated packages are available for interactive work [25,26], visualization [2,12], efficient multidimensional array processing [23], and scientific computing [13].

Following the aforementioned trends, some researchers have taken advantage of Python for writing the high-level parts of large-scale, massively parallel scientific applications and driving simulations in parallel architectures [3,10], while others have tried to make available the benefits of parallel computing to general Python codes using MPI [11,18,22].

In this work, the latest advances in the development of MPI for Python [6] are reported. MPI for Python is a package for

* Corresponding author. Fax: +54 342 4511169.

E-mail addresses: dalcinl@intec.unl.edu.ar (L. Dalcín), rodrigop@intec.unl.edu.ar (R. Paz), mstorti@intec.unl.edu.ar (M. Storti), jdelia@intec.unl.edu.ar (J. D'Elía).

the Python programming language enabling general applications to exploit multiple processors by using any available MPI implementation as a back-end.

The next section presents a brief overview of MPI, Python and MPI for Python. Section 3 describes the most relevant features added to MPI for Python. Section 4 presents some efficiency comparisons between MPI for Python and compiled C code communicating numeric arrays. Finally, Section 6 presents some conclusions and plans for future work.

2. Background

2.1. What is MPI?

MPI [19,20], the *Message Passing Interface*, is a standardized and portable message-passing system designed to function on a wide variety of parallel computers. The standard defines the syntax and semantics of library routines (MPI is not a programming language extension) and allows users to write portable programs in the main scientific programming languages (Fortran, C, and C++).

Since its release, the MPI specification has become the leading standard for message-passing libraries in the world of parallel computers. Implementations are available from vendors of high-performance computers and well known open source projects like MPICH [9,21] and Open MPI [7,24].

MPI defines a high-level abstraction for fast and portable inter-process communication [8,27]. Applications can run in clusters of (possibly heterogeneous) workstations or dedicated nodes, (symmetric) multiprocessors machines, or even a mixture of both. MPI hides all the low-level details, like networking or shared memory management, simplifying development and maintaining portability, without sacrificing performance. However, for communicating data items non-contiguous in memory or corresponding to user-defined structures, MPI still relies on manual packing/unpacking of data buffers, or creation of derived MPI datatypes alongside language datatypes.

2.2. What is Python?

Python [31] is a modern but mature, easy to learn, powerful programming language with a constantly growing community of users. It has efficient high-level data structures and a simple but effective approach to object-oriented programming with dynamic typing and dynamic binding. Python's elegant syntax, together with its interpreted nature, makes it an ideal language for scripting and rapid application development in many areas on most platforms.

The Python interpreter and its extensive standard library are available in source or binary form without charge for all major platforms, and can be freely distributed. It can be easily extended with new functions and data types implemented in C or C++ and is also suitable as an extension language for customizable applications that require a programmable interface.

Python is an ideal candidate for writing the higher-level parts of large-scale scientific applications and driving simulations in parallel architectures. Python codes are quickly developed,

easily maintained, and can achieve a high degree of integration with other libraries written in compiled languages.

2.3. What is MPI for Python?

Python has enough networking capabilities as to develop an MPI implementation in “pure Python”, i.e., without using compiled languages or third-party MPI implementation. This approach is attractive as it assures almost immediate portability to any platform capable of running Python. However, even some basic MPI features (like communication scheduling or datatype management) would require a fair amount of work to be developed. Moreover, achieving reasonable performance and accessing special platform features or dedicated communication hardware would necessarily require the support of compiled languages. Additionally, the availability of a full-featured MPI package for Python implemented on top of a third-party MPI implementation is expected to notably ease the integration of many other MPI-based libraries.

MPI for Python [5] is a Python package providing bindings of the MPI standard, allowing any Python program to exploit multiple processors. This package is constructed on top of the MPI-1/MPI-2 specification and defines an object-oriented interface that closely follows MPI-2 C++ bindings.

The core of MPI for Python is implemented as an extension module written in C in order to access and take advantage of any other MPI implementation providing a C interface. The exposed Python interface is implemented with Python code defining the relevant class hierarchies, functions, and constants. This higher-level code is supported by the C extension module.

In its first release, MPI for Python provided support for blocking point-to-point and collective communications of general Python objects, as well as many facilities for managing process groups and defining new communication domains. Its API was designed with a focus on translating syntax and semantics of standard MPI-2 bindings from C++ to Python. Users with only a basic knowledge of standard C/C++ MPI bindings were able to use this package without having to learn a new interface.

2.4. Related projects

As MPI for Python started and evolved, many ideas were borrowed from other well known open source projects related to MPI and Python.

pyMPI [18] rebuilds the Python interpreter and adds a built-in module for message passing. It permits interactive parallel runs, which are useful for learning and debugging, and provides an environment suitable for basic parallel programming. There is no full support for defining new communicators, process topologies or intercommunicators. General Python objects can be messaged between processors, but there is no support for direct communication of numeric arrays.

Pypar [22] is a rather minimal Python interface to MPI. There is no support for communicators or process topologies. It does not require the Python interpreter to be modified or recompiled. General Python objects of any type can be communicated.

There is also good support for communicating numeric arrays and practically full MPI bandwidth can be achieved.

Scientific Python [11] provides a collection of Python modules that are useful for scientific computing. Among them, there are interfaces to MPI and BSP (*Bulk Synchronous Parallel programming*). The MPI interface is simple but incomplete and does not resemble the MPI specification. However, there is good support for efficiently communicating numeric arrays.

3. New features in MPI for Python

This section presents a survey of MPI capabilities and the new available features in MPI for Python to enhance communication performance and better support classic MPI-1 operations [19] in a Python programming environment. The recent availability of free, high quality, open source MPI-2 implementations strongly motivated the inclusion of another set of features, in order to provide support full for almost all MPI-2 extensions [20].

3.1. Object serialization

The Python standard library supports different mechanisms for data persistence. Many of them rely on disk storage, but *pickling* and *marshaling* can also work with memory buffers.

The `pickle` (slower, written in pure Python) and `cPickle` (faster, written in C) modules provide user-extensible facilities to serialize general Python objects using ASCII or binary formats. The `marshal` module provides facilities to serialize built-in Python objects using a binary format specific to Python, but independent of machine architecture issues.

MPI for Python can communicate any general or built-in Python object taking advantage of the features provided by `cPickle` and `marshal` modules. Their functionalities are wrapped in two classes, `Pickle` and `Marshal`, defining `dump()` and `load()` methods. These are simple extensions, being completely unobtrusive for user-defined classes to participate as they actually use the standard pickle protocol, but carefully optimized for serialization of Python objects on memory streams exposing the standard Python buffer protocol. This approach is also fully extensible; that is, users are allowed to define new, custom serializers implementing the generic `dump()/load()` interface.

Any provided or user-defined serializer can be attached to communicator instances. They will be routinely used to build binary representations of objects to communicate (at sending processes) and restoring them back (at receiving processes).

3.2. Direct communication of memory buffers

Although simple and general, the serialization approach (i.e., *pickling* and *unpickling*) previously discussed imposes important overheads in memory as well as processor usage, especially in the scenario of objects with large memory footprints being communicated. The reasons for this are simple. Pickling general Python objects, ranging from primitive or container built-in types to user-defined classes, necessarily requires

computer resources. Processing is needed for dispatching the appropriate serialization method (that depends on the type of the object) and doing the actual packing. Additional memory is always needed, and if its total amount is not known *a priori*, many reallocations can occur. Indeed, in the case of large numeric arrays, this is certainly unacceptable and precludes communication of objects occupying half or more of the available memory resources.

MPI for Python was improved to support direct communication of any object exporting the single-segment buffer interface. This interface is a standard Python mechanism provided by some types (e.g. strings and numeric arrays), allowing access in the C side to a contiguous memory buffer (i.e., address and length) containing the relevant data.

This new feature, in conjunction with the capability of constructing user-defined MPI datatypes describing complicated memory layouts, enables the implementation of many algorithms involving multidimensional numeric arrays (e.g. image processing, fast Fourier transforms, finite difference schemes on structured Cartesian grids) directly in Python, with negligible overhead, and almost as fast as compiled Fortran, C, or C++ codes.

3.3. Non-blocking and persistent communications

On many systems, performance can be significantly increased by overlapping communication and computation. This is particularly true on systems where communication can be executed autonomously by an intelligent, dedicated communication controller. Non-blocking communication is a mechanism provided by MPI in order to support such overlap.

The `Isend()` and `Irecv()` methods of the `Comm` class initiate a send and receive operation, respectively. These methods return a `Request` instance, uniquely identifying the started operation. Its completion can be managed using the `Test()`, `Wait()`, and `Cancel()` methods of the `Request` class.

Often a communication with the same argument list is repeatedly executed within an inner loop. In such a case, communication can be further optimized by using persistent communication, a particular case of non-blocking communication allowing the reduction of the overhead between processes and communication controllers. This kind of optimization can also alleviate the extra overheads associated to interpreted, dynamic languages like Python, especially for fine-grained tasks.

The `Send_init()` and `Recv_init()` methods of the `Comm` class create a persistent request for a send and receive operation, respectively. These methods return an instance of the `Prequest` class, a subclass of the `Request` class. The actual communication can be effectively started using the `Start()` method, and its completion can be managed as previously described.

The inherently asynchronous nature of non-blocking communications currently imposes some restrictions in what can be communicated using MPI for Python. Communication of memory buffers, as described in Section 3.2 is fully supported. However, communication of general Python objects using serialization, as described in Section 3.1, is possible but not

transparent since objects must be explicitly serialized at sending processes, while receiving processes must first provide a memory buffer large enough to hold the incoming message and next recover the original object. Additionally, the management of Request objects and associated memory buffers involved in communication requires a careful and rather low-level coordination in order to avoid synchronization problems. Similar issues arise in many other Python packages. The buffer interface as currently defined in Python is rather simple, lacking of a lock/unlock semantics for accessing the exposed memory.

3.4. MPI-2 extensions

3.4.1. Dynamic process management

An MPI-1 application is static; that is, no processes can be added to or deleted from an application after it has been started. This limitation was addressed in MPI-2. The new specification added a process management model providing a basic interface between an application and external resources and process managers. This extension can be really useful, especially for serial applications built on top of parallel modules, or parallel applications with a client/server model. The MPI-2 process model provides a mechanism to create new processes and establish communication between them and the existing MPI application. It also provides a mechanism to establish communication between two existing MPI applications, even when one did not “start” the other.

In MPI for Python, new processes can be created by calling the `Spawn()` method within an intracommunicator (i.e., an `Intracomm` instance). This call returns a new intercommunicator (i.e., an `Intercomm` instance), which can be used to perform point to point and collective communications between the parent and child groups of processes.

Alternatively, disjoint groups of processes can establish communication in a client/server approach. Server applications must first call the `Open_port()` function to open a “port” and the `Publish_name()` function to publish a provided “service”, and next call the `Accept()` method within an `Intracomm` instance. Client applications can first find a published “service” by calling the `Lookup_name()` function, which returns the “port” where a server can be contacted; and next call the `Connect()` method within an `Intracomm` instance. Both `Accept()` and `Connect()` methods return an `Intercomm` instance. When connection between client/server processes is no longer needed, all of them must cooperatively call the `Disconnect()` method of the `Comm` class. Additionally, server applications can release resources by calling the `Unpublish_name()` and `Close_port()` functions.

3.4.2. One-sided communications

One-sided communications (also called *Remote Memory Access, RMA*) supplements the traditional two-sided MPI communication model with a one-sided interface that can take advantage of the capabilities of RMA network hardware. This extension lowers latency and software overhead in applications written using a shared-memory-like paradigm. The semantics

of one-sided communication are fairly complex. The MPI RMA API revolves around the use of objects called *windows*, which intuitively specify regions of a process’s memory that have been made available for remote operations.

Windows are created by calling the `Create()` method of the `Win` class at all processes within a communicator and specifying a memory buffer (i.e., a base address and length). Three one-sided operations for remote write, read and reduction are available using the `Put()`, `Get()`, and `Accumulate()` methods respectively within a `Win` instance. These methods need an offset into the window and an integer rank identifying the remote target. This one-sided operations are implicitly non-blocking and must be synchronized.

Windows are synchronized by using two primary modes. Active target synchronization requires the origin process to call the `Start()/Complete()` methods at the origin process and target process cooperates by calling the `Post()/Wait()` methods. There is also a collective variant provided by the `Fence()` method. Passive target synchronization is more lenient, only the origin process calls the `Lock()/Unlock()` methods.

3.4.3. Extended collective operations

In the MPI-1 specification, collective communications were only defined for intracommunicators. The MPI-2 specification introduces extensions generalizing many of the collective routines to intercommunicators. They can be really useful for collective interaction between disjoint group of processes created or connected as described in Section 3.4.1.

MPI for Python was enhanced in order to support these extensions. The `Barrier()`, `Bcast()`, `Gather()`, `Scatter()`, `Allgather()`, `Alltoall()`, `Reduce()`, and `Allreduce()` methods are defined for both `Intracomm` and `Intercomm` classes. They are able to collectively communicate general Python objects, as discussed in Section 3.1, or memory buffers, as discussed in Section 3.2. Notably, scan and exclusive scan operations as defined in MPI do not apply to intercommunicators; that is, the `Scan()` and `Exscan()` methods are only available for `Intracomm` instances.

3.4.4. Parallel I/O

POSIX provides a model of a widely portable file system. However, the optimization needed for parallel I/O cannot be achieved with this interface, and can only be if the parallel I/O system provides a high-level interface supporting partitioning of file data among processes and a collective interface supporting complete transfers of global data structures between process memories and files. Additionally, further efficiencies can be gained via support for asynchronous I/O, strided accesses, and control over physical file layout on storage devices.

The common patterns for accessing a shared file (broadcast, reduction, scatter, gather) is expressed using user-defined MPI datatypes. Compared to communication patterns of point to point and collective communications, this approach has the advantage of added flexibility and expressiveness. Data access operations (read and write) are defined for different kinds of

positioning (using explicit offsets, individual file pointers, and shared file pointers), coordination (non-collective and collective), and synchronism (blocking, non-blocking, and split collective).

All these features are available in MPI for Python by using instances of the `File` class. Parallel files are created by calling method `Open()` at all processes within a communicator; they can be closed or even destroyed by calling `Close()` and `Delete()` methods, respectively. The data layout in the file can be set and queried with the `Set_view()` and `Get_view()` methods, respectively. Data access is provided by many methods related to read and write operations, but with different behavior regarding positioning, coordination, and synchronism.

4. Testing

Some efficiency tests were run on the Beowulf class cluster *Aquiles* [29] at CIMEC. Its hardware consists of 80 diskless single processor computing nodes with Intel Pentium 4 Prescott 3.0 GHz 2 MB cache processors, Intel Desktop Board D915PGN motherboards, Kingston Value RAM 2 GB DDR 400 MHz memory, and 3Com 2000ct Gigabit LAN network cards, interconnected with a 3Com SuperStack 3 Switch 3870 48-ports Gigabit Ethernet. MPI for Python was compiled on a *Linux 2.6.17* box using *GCC 3.4.6* with *Python 2.5.1*. The chosen MPI implementation was *MPICH2 1.0.5p4*. Communications between processes involved numeric arrays, they were provided by *NumPy 1.0.3*.

The first test consisted in blocking send and receive operations (`MPI_SEND` and `MPI_RECV`) between a pair of nodes. Messages were numeric arrays of double precision (64 bits) floating-point values. The two supported communications mechanisms, serialization and memory buffers, were compared against compiled C code. A basic implementation of this test using MPI for Python with direct communication of memory buffers (translation to C or C++ is straightforward) is shown below.

```
from mpi4py import MPI
from numpy import empty, float64

comm = MPI.COMM_WORLD
rank = comm.Get_rank()
array1 = empty(2**16, dtype=float64)
array2 = empty(2**16, dtype=float64)
sendbuf = [array1, 2**16, MPI.DOUBLE]
recvbuf = [array2, 2**16, MPI.DOUBLE]

wt = MPI.Wtime()
if rank == 0:
    comm.Send(sendbuf, 1, tag=0)
    comm.Recv(recvbuf, 1, tag=0)
elif rank == 1:
    comm.Recv(recvbuf, 0, tag=0)
    comm.Send(sendbuf, 0, tag=0)
wt = MPI.Wtime() - wt
```

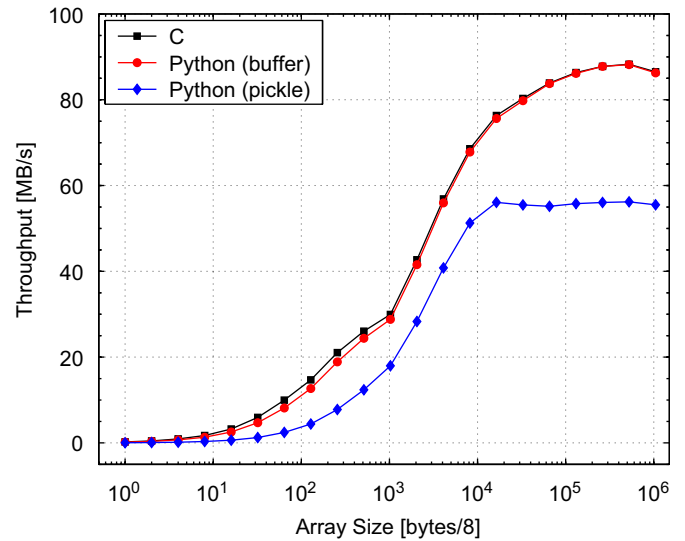


Fig. 1. Throughput in blocking send and receive.

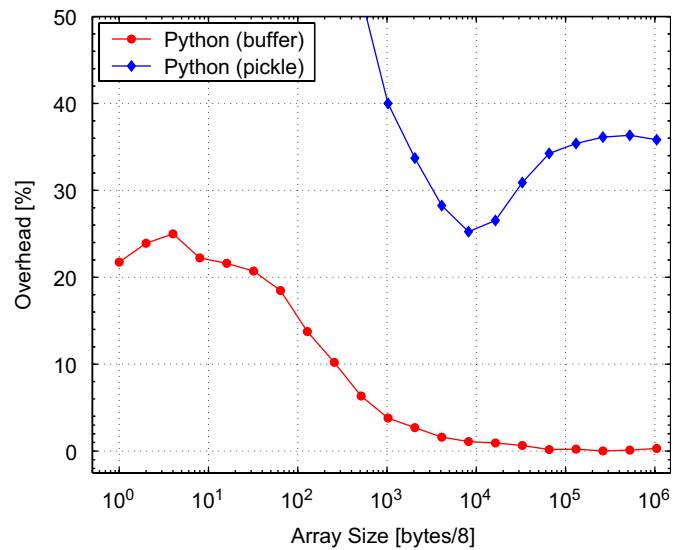


Fig. 2. Relative overhead in blocking send and receive.

Results are shown in Figs. 1 and 2. Throughput is computed as $2S/\Delta t$, where S is the basic message size (in megabytes), and Δt is the measured wall-clock time. Clearly, the overhead introduced by object serialization degrades overall efficiency; the maximum throughput in Python is about 60% of the one in C. However, the direct communication of memory buffers introduces a negligible overhead for medium-sized to long arrays.

The second test consisted in a small variation of the first one. The interchange of messages consisted in a bidirectional send/receive operation (`MPI_SENDRECV`). Results are shown in Figs. 3 and 4. In comparison to the previous test, the overhead introduced by object serialization is lower (the maximum throughput in Python is about 75% of the one in C) and the overhead communicating memory buffers is similar (and again, it is negligible for medium-sized to long arrays).

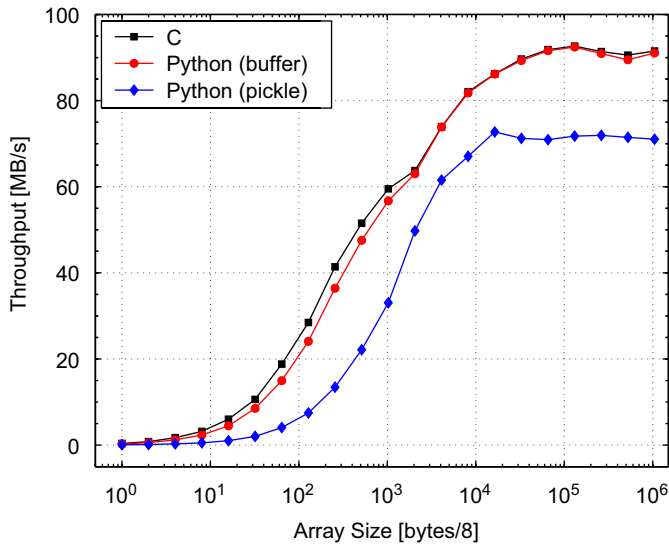


Fig. 3. Throughput in bidirectional send/receive.

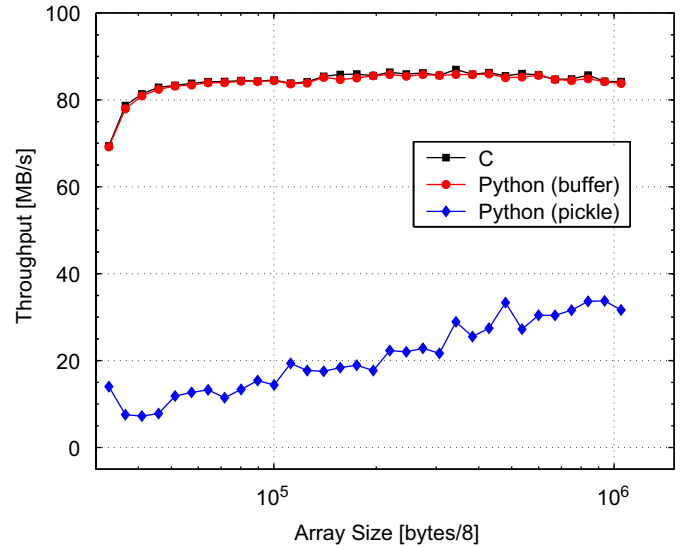


Fig. 5. Throughput in all-to-all.

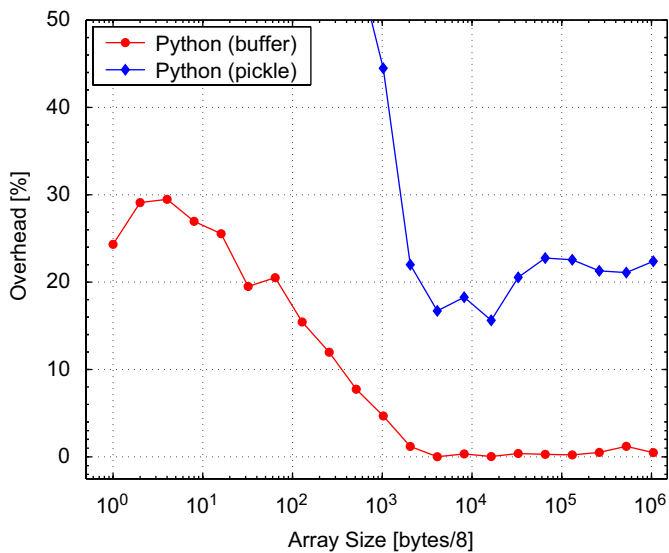


Fig. 4. Relative overhead in bidirectional send/receive.

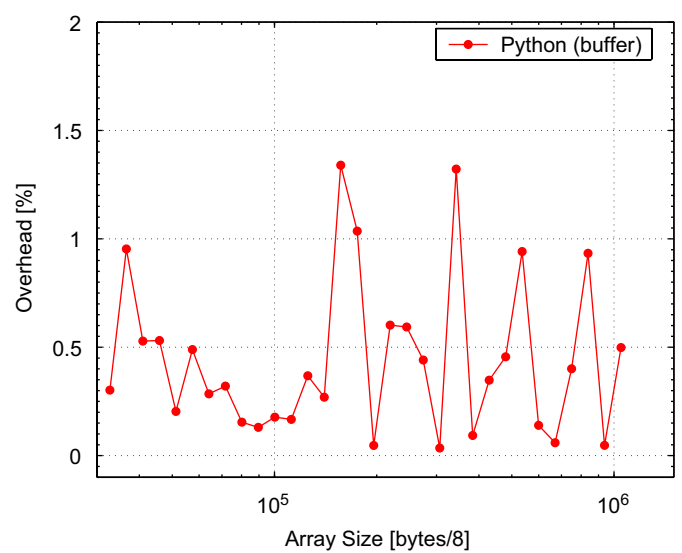


Fig. 6. Relative overhead in all-to-all.

The third test consisted in an all-to-all collective operation (MPI_ALLTOALL) on 16 nodes. As in previous tests, messages were numeric arrays of double precision floating-point values. Results are shown in Figs. 5 and 6. Throughput is computed as $2(N - 1)S/\Delta t$, where N is the number of nodes, S is the basic message size (in megabytes), and Δt is the measured wall-clock time. The overhead introduced by object serialization is notably more significant than in previous tests; the maximum throughput in Python is about 40% of the one in C. However, the overhead communicating memory buffers is always below 1.5%.

Interested readers should review previous results from a similar set of tests [6], but obtained with older hardware components and software distributions.

Finally, some remarks are worth to be done on the implementation of the previous tests and the obtained efficiency results.

The snippet of code shown in the discussion above was included just for reference. The actual implementation took into account memory preallocation (in order to avoid paging effects) and parallel synchronization (in order to avoid asynchronous skew in the start-up phase). Timings were measured many times inside a loop over each single run and the minimum were taken.

All the tests involved communications of one-dimensional, contiguous *NumPy* arrays. For those kind of objects, serialization and deserialization with the *pickle* protocol is implemented quite efficiently. The serialization step is accomplished with memory copying into a raw string object. The deserialization step reuses this raw string object, thus saving from extra memory allocations and copies. These optimizations are possible because the total amount of memory required for serialization is known in advance and all items have a common data type corresponding to a C primitive type. For more general Python

objects, the serialization approach is expected to achieve lower performance than the previously reported.

5. Use cases

MPI for Python is a general purpose package. Currently, it does not provide other functionalities than those defined in the MPI-2 standard. However, MPI for Python is able to effectively support parallel applications written in Python, especially in those making use of external parallel libraries and routines written in compiled languages. This was the main motivation for starting its development and guides its evolution. Some projects taking advantage of MPI for Python are commented below.

PETSc-FEM [28,30] is a parallel multi-physics finite elements code in active development at CIMEC. PETSc-FEM is built around MPI and PETSc libraries, targeting computational fluid dynamics simulations on dedicated Beowulf class clusters. MPI for Python is used in conjunction with Python wrappers to PETSc libraries and PETSc-FEM core functionalities to support Python-driven large-scale simulations.

The Synergia project [1] is developed by the Advanced Accelerator Modelling team at the Fermilab Computing Division. The software simulates the behavior of particles in an accelerator, emphasizing three-dimensional models of collective beam effects. This project uses Python-driven Fortran 90 and C++ libraries to create beam dynamics simulations. The individual beam dynamics libraries were designed to use MPI to support parallel computation. MPI for Python is used from the Python-driven simulations to support communication within MPI at the Python level.

IPython [26,25] is a well established open source project providing an interactive shell far superior to Python's default one. The next-generation IPython, currently being developed, will provide facilities for interactive parallel and distributed computing. MPI for Python will support IPython to be used in a fully interactive manner.

6. Conclusions

MPI for Python provides a base layer for applying the message-passing paradigm in parallel applications written in Python. It makes use of any available MPI implementation retaining the syntax and semantics of the standard MPI specification. It is specially suited for integration of other MPI-based libraries and assist in the development of complex Python-driven applications targeted to run on parallel environments.

MPI for Python can communicate general Python objects as well as any Python object exposing a memory buffer. In the later case, efficiency tests have shown that performance degradation is negligible, even for medium-sized numeric array objects. In fact, the introduced overhead is far smaller than the normal one associated to the use of interpreted versus compiled languages.

Future work will be directed towards the improvement of MPI for Python by adding some currently unsupported MPI functionalities like datatype decoding, attribute catching and interoperability with Fortran libraries. Additionally, an automatic mapping between MPI datatypes and NumPy datatypes will be

provided in order to simplify the parallelization of demanding applications involving multidimensional array processing.

Acknowledgments

The authors gratefully acknowledge many helpful discussions with Brian Granger. They also thank Fernando Pérez and unknown reviewers for the careful reading of the first version of the paper and many useful comments and recommendations.

This work has received financial support from *Consejo Nacional de Investigaciones Científicas y Técnicas* (CONICET, Argentina, Grants PIP 02552/00, PIP 5271/05), *Universidad Nacional del Litoral* (UNL, Argentina, Grant CAI + D 2005-10-64), and *Agencia Nacional de Promoción Científica y Tecnológica* (ANPCyT, Argentina, Grants PICT 12-14573/2003, PME 209/2003).

References

- [1] J. Amundson, P. Spentzouris, J. Qiang, R. Ryne, Synergia: a 3D accelerator modelling tool with 3D space charge, *J. Comput. Phys.* 211 (1) (2006) 229–248.
- [2] P. Barrett, J. Hunter, P. Greenfield, Matplotlib—a portable Python plotting package, *Astronomical Data Analysis Software & Systems XIV*, 2004.
- [3] D.M. Beazley, P.S. Lomdahl, Feeding a large scale physics application to Python, in: *Proceedings of 6th International Python Conference*, San Jose, California, 1997, pp. 21–29.
- [4] Beowulf.org, The Beowulf cluster site, (<http://www.beowulf.org/>), 2007.
- [5] L. Dalcín, MPI for Python, (<http://www.mpi4py.scipy.org/>), 2007.
- [6] L. Dalcín, R. Paz, M. Storti, MPI for Python, *J. Parallel Distributed Comput.* 65 (9) (2005) 1108–1115.
- [7] E. Gabriel, G.E. Fagg, G. Bosilca, T. Angskun, J.J. Dongarra, J.M. Squyres, V. Sahay, P. Kambadur, B. Barrett, A. Lumsdaine, R.H. Castain, D.J. Daniel, R.L. Graham, T.S. Woodall, Open MPI: goals, concept, and design of a next generation MPI implementation, in: *Proceedings of the 11th European PVM/MPI Users' Group Meeting*, Budapest, Hungary, 2004, pp. 97–104.
- [8] W. Gropp, S. Huss-Lederman, A. Lumsdaine, E. Lusk, B. Nitzberg, W. Saphir, M. Snir, MPI—The Complete Reference: Volume 2, The MPI-2 Extensions, second ed., vol. 2, The MPI-2 Extensions, MIT Press, Cambridge, MA, USA, 1998.
- [9] W. Gropp, E. Lusk, N. Doss, A. Skjellum, A high-performance, portable implementation of the MPI message passing interface standard, *Parallel Comput.* 22 (6) (1996) 789–828.
- [10] K. Hinsen, The Molecular Modelling Toolkit: a new approach to molecular simulations, *J. Comput. Chem.* 21 (2) (2000) 79–85.
- [11] K. Hinsen, ScientificPython: Home page, (<http://starship.net/~hinsen/ScientificPython/>), 2007.
- [12] J.D. Hunter, Matplotlib, (<http://matplotlib.sourceforge.net/>), 2003–2007.
- [13] E. Jones, T. Oliphant, P. Peterson, et al., SciPy: open source scientific tools for Python, (<http://www.scipy.org/>), 2001–2007.
- [14] J. Kepner, A multi-threaded fast convolver for dynamically parallel image filtering, *J. Parallel Distributed Comput.* 63 (3) (2003) 360–372.
- [15] J. Kepner, S. Ahalt, MatlabMPI, *J. Parallel Distributed Comput.* 64 (8) (2004) 997–1005.
- [16] H.P. Langtangen, Python Scripting for Computational Science, Simula Research Lab, Department of Informatics, University of Oslo, 2006.
- [17] Message Passing Interface Forum, Message Passing Interface (MPI) Forum Home Page, (<http://www.mpi-forum.org/>), 1994.
- [18] P. Miller, pyMPI: putting the py in MPI, (<http://pympi.sourceforge.net/>), 2000–2007.
- [19] MPI Forum, MPI: a message passing interface standard, *Internat. J. Supercomputer Appl.* 8 (3/4) (1994) 159–416.

- [20] MPI Forum, MPI2: a message passing interface standard, High Performance Comput. Appl. 12 (1–2) (1998) 1–299.
- [21] MPICH Team, MPICH: a portable implementation of MPI, (<http://www-unix.mcs.anl.gov/mpi/mpich/>), 2007.
- [22] O. Nielsen, PyPar Home page, (<http://datamining.anu.edu.au/~ole/pypar/>), 2002–2007.
- [23] T. Oliphant, NumPy: Numerical Python, (<http://numpy.scipy.org/>), 2005–2007.
- [24] Open MPI Team, Open MPI: open source high performance computing, (<http://www.open-mpi.org/>), 2007.
- [25] F. Pérez, IPython: an Enhanced Python Shell (<http://ipython.scipy.org/>), 2001–2007.
- [26] F. Pérez, B. Granger, IPython: a system for interactive scientific computing, Comput. Sci. Eng. 9 (3) (2007) 21–29.
- [27] M. Snir, S. Otto, S. Huss-Lederman, D. Walker, J. Dongarra, MPI—The Complete Reference: Volume 1, The MPI Core, second ed., vol. 1, The MPI Core, MIT Press, Cambridge, MA, USA, 1998.
- [28] V.E. Sonzogni, A.M. Yommi, N.M. Nigro, M.A. Storti, A parallel finite element program on a Beowulf cluster, Adv. in Eng. Software 33 (7–10) (2002) 427–443.
- [29] M.A. Storti, Aquiles cluster at CIMEC, (<http://www.cimec.org.ar/aquiles>), 2005–2007.
- [30] M.A. Storti, N. Nigro, R. Paz, L. Dalcín, E. López, PETSc-FEM: a general purpose, parallel, multi-physics FEM program, (<http://www.cimec.org.ar/petscfem>), 1999–2007.
- [31] G. van Rossum, Python programming language, (<http://www.python.org/>), 1990–2007.



Lisandro D. Dalcín is a Fellow Researcher at CONICET (INTEC-UNL, <http://www.cimec.org.ar>). He received his degree in Electromechanical Engineering from “Universidad Tecnológica Nacional, Facultad Regional Concepción del Uruguay” (Concepción del Uruguay, Argentina) in 2001. He is to receive a Ph.D. degree in Engineering and Computational Mechanics from “Universidad Nacional del Litoral”, Santa Fe, Argentina. Currently he is an Assistant Professor of Data Structures and Algorithms, and Numerical Calculus at “Universidad Nacional

del Litoral”, UNL, Santa Fe, Argentina. He is also member of the development team of PETSc libraries (<http://www-unix.mcs.anl.gov/petsc>). His research interests include numerical methods in computational fluid mechanics, scalable solution methods for incompressible fluid flow problems, and parallel computing on distributed memory architectures.



Rodrigo R. Paz is a Researcher at CONICET (INTEC-UNL, <http://www.cimec.org.ar>).

He received his degree in Aeronautical and Mechanical Engineering from “Universidad Nacional de Córdoba”, Córdoba, Argentina, in 1999. In 2006, after postgraduate studies in France (Université Pierre et Marie Curie PARIS VI and Faculté des Sciences of the Université Jean Monnet in Saint Etienne) he received the Ph.D. degree in Engineering and Computational Mechanics from “Universidad Nacional del Litoral”, Santa Fe, Argentina, where he is currently a professor in postgraduate and undergraduate programs.

His research interests include domain decomposition methods, parallel computing, weak/strong fluid/structure coupling, absorbing boundary conditions and multiphase flows.



Mario A. Storti is a Senior Researcher at CONICET (INTEC-UNL, <http://www.cimec.org.ar>).

He received his degree in Physics from “Instituto Balseiro”, Bariloche, Argentina in 1983, and his Ph.D. degree in Chemistry Technology from “Universidad Nacional del Litoral”, Santa Fe, Argentina, in 1990.

His research interests include computational fluid dynamics with finite element methods, free surface and multiphase flows, object oriented programming for finite elements in C++, development of PETSc-FEM

(<http://www.cimec.org.ar/petscfem>) a C/C++ FEM code based on MPI and PETSc, and parallel computing on Beowulf clusters.



Jorge L. D'Elia is a Senior Researcher at CONICET (INTEC-UNL, <http://www.cimec.org.ar>).

He received his degree in Electromechanical Engineering from “Universidad Tecnológica Nacional, Facultad Regional Mendoza”, Mendoza, Argentina, and his Ph.D. degree in Engineering and Computational Mechanics from “Universidad Nacional del Litoral” (Santa Fe, Argentina).

Currently he is an Associate Professor of Computational Theory and Parallel Algorithms and Assistant Professor of Mathematics at “Universidad Nacional del Litoral”, Santa Fe, Argentina.

His research interests include numerical methods, computational mechanics, discrete mathematics, serial and parallel computing.