

A Case for Non-Blocking Collective Operations

Torsten Hoefler^{1,3}, Jeffrey M. Squyres², Wolfgang Rehm³, and Andrew Lumsdaine¹

¹ Indiana University, Open Systems Lab, Bloomington IN 47405, USA
{[htor](mailto:htor@cs.indiana.edu),[lums](mailto:lums@cs.indiana.edu)}@cs.indiana.edu,

² Cisco Systems, San Jose, CA 95134 USA
jsquyres@cisco.com

³ Technical University of Chemnitz, 09107 Chemnitz, Germany
{[htor](mailto:htor@cs.tu-chemnitz.de),[rehm](mailto:rehm@cs.tu-chemnitz.de)}@cs.tu-chemnitz.de

Abstract. Non-blocking collective operations for MPI have been in discussion for a long time. We want to contribute to this discussion and to give a rationale for the usage these operations and assess their possible benefits. A LogGP model for the CPU overhead of collective algorithms and a benchmark to measures it are provided and show a large potential to overlap communication and computation. We show that non-blocking collective operations can provide at least the same benefits as non-blocking point to point operations already do. Our claim is that actual CPU overhead for non-blocking collective operations depends on the message size and the communicator size and benefits especially highly scalable applications with huge communicators. We prove that the share of the overhead of the overall communication time of current blocking collective operations gets smaller with bigger communicators and larger messages. We show that the user level CPU overhead is less than 10% for MPICH2 and LAM/MPI using TCP/IP communication, which leads us to the conclusion that, by using non-blocking collective communication, ideally 90% idle CPU time can be freed for the application.

Keywords: Collective communication, Overlap, Non-blocking communication, Message passing (MPI)

1 Introduction

Non-blocking collective operations and their possible benefits have already been discussed at meetings of the MPI standardization committee. The final decision to not include them into the MPI-2 standard fell at March 6, 1997⁴. However, the fact that the decision was extremely marginal (11 yes / 12 no / 2 abstain) suggests that the role of non-blocking collective operations is still debatable. Our contention is that non-blocking collective operations are a natural extension to the MPI-2 standard. We show that non-blocking collective operations can be beneficial for a class of applications to utilize the available CPU time more

⁴ see: <http://www.mpi-forum.org/archives/votes/mpi2.html>

efficiently and decrease the time to solution of these applications significantly. Further, we discuss two main problems of blocking collective communication which limit the scalability of applications.

First, blocking collective operations have a more or less synchronizing effect on applications which leads to unnecessary wait time. Even though the MPI standard does not define other blocking collective operations than `MPI_BARRIER` to be strictly synchronizing, most used algorithms force many processes to wait for other processes due to data dependencies. In this way, synchronization with a single process is enforced for some operations (e.g., a `MPI_BCAST` can not be finished until the root process called it) and the synchronizing behavior of other operations highly depends on the implemented collective algorithm. In either case, pseudo-synchronizing behavior often leads to many lost CPU cycles, a high sensitivity to process skew (e.g., due to daemon processes which “steal” the CPU occasionally and introduce a pseudo-random skew between processes [1,2]), and a high sensitivity to imbalanced programming (e.g., some processes do slightly more computation than others each round).

Second, most blocking collective operations can not take much advantage of modern interconnects which enable communication offload to support efficient overlapping of communication and computation. Abstractly seen, each supercomputer or cluster consists of two entities, the CPU which processes data streams and the network which transports data streams. In many networks, both entities can act mostly independently of each other, but the programmer has no chance to use this parallelism efficiently if blocking communication (point-to-point or collective) is used.

Another rationale to offer non-blocking semantics for collective communication is an analogy between many modern operating systems and the MPI standard. Most modern operating systems offer possibilities to overlap computation on the host CPU with actions of other entities (for example harddisks or the network). Asynchronous I/O and non-blocking TCP/IP sockets are today’s standard features to communicate. The MPI standard offers non-blocking point-to-point communication which can be used to overlap communication and computation. It would be a natural extension to offer also a non-blocking interface to the collective operations.

The next section describes related work in the field of overlap of computation and communication and the avoidance of synchronization. Section 2 gives some information about possible benefits of non-blocking collective communication. Section 3 presents benchmark results for a selected set of operations followed by a conclusion of this work.

1.1 Related Work

The obvious benefits of overlapping communication with computation and leveraging the hardware parallelism efficiently with the usage of non-blocking communication is well documented. Analyses [3,4,5] try to give an assessment of the capabilities of MPI implementations to perform overlapping for point-to-point communications. Many groups analyze the possible performance benefits for real

applications. Liu et al. [6] showed possible speedups up to 1.9 for several parallel programs. Brightwell et al. [7] classifies the source of performance advantage for overlap and Dimitrov [8] uses overlapping as fundamental approach to optimize parallel applications for cluster systems. Other studies, as [9,10,11,12] apply several transformations to parallel codes to enable overlapping. However, little research has been done in the field of non-blocking collectives. Studies like [13,14] mention that non-blocking collective operations would be beneficial but do not provide a measure for it. Kale et al. [15] analyzed the applicability of a non-blocking personalized exchange to a small set of applications in practice. However, many studies show that non-blocking communication and non-blocking collectives *may* be beneficial. Our work contributes to the field because we actually assess the potential performance benefits of a non-blocking collective implementation.

2 Possible Performance Benefits

The most obvious benefits of non-blocking collective operations are the avoidance of explicit pseudo synchronization and the ability to leverage the hardware parallelism stated in Section 1. The pseudo-synchronizing behavior of most algorithms cannot be avoided, but non-blocking collective operations process the operation in the background, which enables the user to ignore most synchronization effects. Common sources for de-synchronization, process skew and load imbalance are not easily measurable. However, results can increase the application running time dramatically, as shown in [16]. Theoretical [17] and practical analyses [18,16] show that operating system noise and resulting process skew is definitively influencing the performance of parallel applications using blocking collective operations. Non-blocking collective operations avoid explicit synchronization unless it is necessary (if the programmer wants to wait for the operation to finish). This enables the programmer to develop applications which are more tolerant of process skew and load imbalance.

Another benefit is to use the parallelism of the network and computation layers. Non-blocking communication (point-to-point and collective) allows the user to issue a communication request to the hardware, perform some useful computation, and ask later if it has been completed. Modern interconnect networks can perform the message transfer mostly independently of the user process. The resulting effect is that, for several algorithms/applications, the user can overlap the communication latency with useful computation and ignore the communication latency up to a certain extent (or totally). This has been well analyzed for point-to-point communication (see Section 1.1). Non-blocking collective operations allow the programmer to combine the benefits of collective communication [19] with all benefits of non-blocking communication. The following subsections analyze the communication behavior of current blocking collective algorithms and implementations and show that only a fraction of the CPU time is involved into communication related computation. In relation to previous studies we show, theoretically and practically, that a similar percentage, in many cases

even more, idle CPU time as with non-blocking point-to-point communication can be gained. We assume that the biggest share of the remaining (idle) CPU time can be leveraged by the user if overlap of communication and computation together with non-blocking collective communication can be applied.

2.1 Modelling CPU and Network Activity

This subsection gives an estimation of the theoretical CPU idle time during a collective operation. The CPU idle time during the communication will be modelled and benchmarked. Precise models for collective operations are presented in [20] and for barrier synchronization in [21]. Both studies show that the LogP [22] or LogGP [23] model is able to predict the communication time sufficiently accurately if the processes enter the collective operation simultaneously.

We analyze the three collective operations `MPI_BARRIER`, `MPI_ALLREDUCE`, and `MPI_BCAST` without loss of generality, in detail. As shown in [24,25,26], these three operations are frequently used in real applications. However, the results can also be applied to all other collective operations.

We assume the usual LogP/LogGP communication parameters (**L**atency, **o**verhead, **g**ap, **G**ap per byte and **P**rocessors) and γ to assess computation:

We derive simplified LogGP models for networks adhering the properties defined in Section 2.2 in [21] (full bisectional bandwidth; full duplex; unlimited forwarding rate; L, o are constant; $o > L > g$). We model point-to-point message based implementations with logarithmic running time ($O(\log_2 P)$) of all three operations. We assume the dissemination principle to perform `MPI_BARRIER` (1), analyzed in [21]. Our model for `MPI_ALLREDUCE` (2) assumes a simple binomial tree reduce implementation followed by `MPI_BCAST` and our `MPI_BCAST` (3) model assumes a binomial tree implementation (compare proposed models in [20]).

$$t_{barr} = (2o + L) \cdot \lceil \log_2 P \rceil \quad (1)$$

$$t_{allred} = 2 \cdot (2o + L + m \cdot G) \cdot \lceil \log_2 P \rceil + m \cdot \gamma \cdot \lceil \log_2 P \rceil \quad (2)$$

$$t_{bcast} = (2o + L + m \cdot G) \cdot \lceil \log_2 P \rceil \quad (3)$$

If we come back to the two entities, which are the network and the processor, mentioned in Section 1, we realize that each parameter is “accounted” at a specific entity. The processor is only used by o and γ while the network is used to perform the message transmission (L, g, G). Using this information, we can divide the equations presented above up into processing and network parts:

$$t_{barr}^{CPU} = 2o \cdot \lceil \log_2 P \rceil \quad t_{barr}^{NET} = L \cdot \lceil \log_2 P \rceil \quad (4)$$

$$t_{allred}^{CPU} = (4o + m \cdot \gamma) \cdot \lceil \log_2 P \rceil \quad t_{allred}^{NET} = 2 \cdot (L + m \cdot G) \cdot \lceil \log_2 P \rceil \quad (5)$$

$$t_{bcast}^{CPU} = 2o \cdot \lceil \log_2 P \rceil \quad t_{bcast}^{NET} = (L + m \cdot G) \cdot \lceil \log_2 P \rceil \quad (6)$$

We see that both, t^{CPU} and t^{NET} scale logarithmically with P . However, on modern interconnects the parameters can differ significantly. The following section provides an analysis of these parameters for modern interconnect networks.

2.2 Fitting the Model to Modern Architectures

Modern interconnect architectures, like InfiniBandTM, QuadricsTM, or MyrinetTM, which are used for HPC systems, try to offload a huge share of the communication into the network interface card. Traditional networks, like Ethernet (without offloading), still use the CPU extensively to process network protocols like TCP/IP. However, also Ethernet has been optimized for lower host overhead with simplified protocols [27] as well as direct user level access and protocol offloading [28]. All these new networks and approaches aim to reduce the overhead of the main CPU involved in communication (o parameter). The L parameter is usually greater than the o in modern networks, and the gap between t^{CPU} and t^{NET} grows with the message size as $G \cdot m$ is added. This enables efficient overlapping of computation and communication for point-to-point communication which has been described in the related work section. However, this idea can also be applied to collective communication. As one can see in equations (4),(5),(6), the gap between Network and CPU occupancy also grows with the number of involved processors P . This leads us to the prediction that especially blocking collectives which communicate large data chunks with many processors should be mostly utilizing the network (with an idle CPU). The only exception could be reduction operations, like MPI_ALLREDUCE, because they include processing (reduction) of values on the host CPU. However, in most cases, the bandwidth of the CPU should be much higher than the network bandwidth. In the following section, we evaluate these theoretical expectations with a custom benchmark set which measures the CPU usage during blocking collective operations.

3 Benchmark Results

We implemented a benchmark which measures the CPU utilization for different MPI collective operations. The benchmark uses the standard `gettimeofday()` and `getrusage()` functionality of modern operating systems to measure the idle time. It issues collective calls with different message sizes and communicator sizes. The benchmark methodology is described as pseudocode in Listing 1.1. The `getrusage()` call returns system time and user time used by the running process separately. We chose a high number of iterations (10000) in the inner loop (`max_iters`, Line 5) to be able to neglect the overhead and relative imprecision of the system functions. We conducted the benchmark for different MPI implementations shown in Table 1.

Many MPI libraries are implemented in a non-blocking manner which means that the CPU overhead is, due to polling, 100% regardless of other factors. Only LAM/MPI with TCP/IP and MPICH2 with TCP/IP used blocking communication to perform the collective operations. However, it is totally correct to use polling to perform blocking MPI collective operations because, at least for single threaded MPI applications, the CPU is unusable anyways and polling has usually slightly lower overhead than interrupt based (blocking) methods.

```

1 for( proc=1; proc<nproc; proc=proc*2) {
  create_communicator(nproc, comm);
  for( size=1; size<maxsize; proc=proc*2) {
5   gettimeofday(t1); getrusage(r1);
   for( i=0; i<max_iters; i++)
     MPI_Coll(comm, size, MPI_BYTE, ...)
   getrusage(r2); gettimeofday(t2);
  }}

```

Listing 1.1. Benchmark Methodology (pseudocode)

Implementation	Networks
LAM/MPI 7.1.2	InfiniBand, TCP/IP
MPICH2 1.0.3	TCP/IP
Open MPI 1.1a3	InfiniBand and TCP/IP
OSU MVAPICH 0.9.4	InfiniBand

Table 1. Tested MPI Implementations

We investigated all collective operations for LAM/MPI and MPICH2 and want to discuss the frequently used `MPI_ALLREDUCE` (cf. [26]) in detail. Both `MPI_ALLREDUCE` implementations exhibit a similar behavior and use only a fraction of the available CPU power for communicators with more than 8 nodes. MPICH2 causes in the average of all measurement points less than 30% CPU load while LAM/MPI consumes less than 10%. We see also that the data size plays an important role because there may be switching points in the collective implementation where the collective algorithms or underlying point-to-point operations are changed (e.g., 128kb for MPICH2). However, this overhead includes the TCP/IP packet processing time spent in the kernel to transmit the messages which is measured with the `getrusage()` function as system time. User level, kernel-bypass, and offloading communication hardware like InfiniBand, Quadrics or Myrinet does not use the host CPU to process packets and does not enter the kernel during message transmission. Figure 1 shows the user level CPU usage (without TCP/IP processing) for both examples from above. It shows that the CPU overhead for `MPI_ALLREDUCE`, which implies a user level reduction operation in our case, is below 10% in the average for MPICH2 and below 3% for LAM/MPI. These figures show also that the share of CPU idle time grows with communicator and data size. Other collective operations which are not shown here due to space restrictions exhibit a similar behavior.

However, generally speaking, the time to perform a collective operation grows also with communicator and data size. This means that the overall (multiplicative) CPU waste is even higher. Figure 2 shows the absolute CPU idle time of both implementations, several collective operations, and a fixed communicated data size with varying communicator sizes. The effect of growing CPU waste

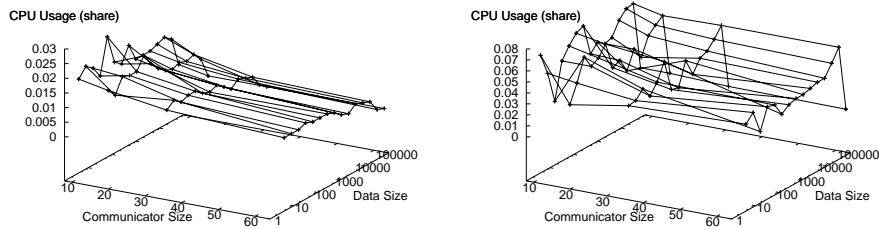


Fig. 1. MPI_ALLREDUCE (user time) overheads for LAM/MPI (left) and MPICH2 (right).

during blocking collectives is clearly visible. Especially the MPI_ALLTOALL operation, which usually scales worst, shows high CPU idle times with a growing number of participating processes.

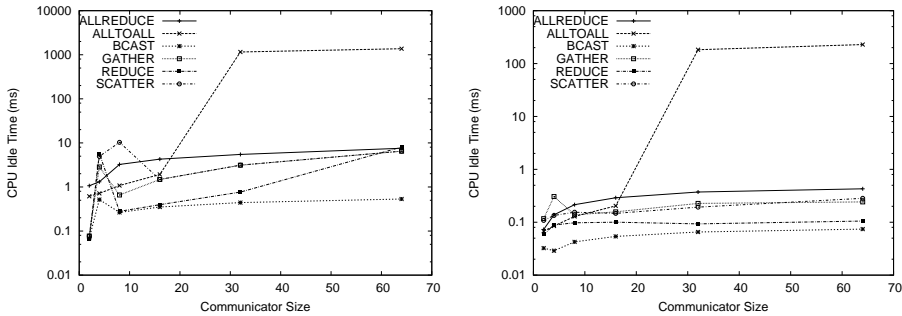


Fig. 2. CPU idle time for some collective functions with varying communicator sizes for a constant data size of 1kB (left: LAM/MPI, right: MPICH2).

Figure 3 shows that absolute CPU idle time of both implementations, for a fixed communicator size, and varying data sizes. The CPU waste is even higher and scales worse than for the varying communicator size, nearly linearly with the data size (the figures are plotted with a logarithmic scale).

4 Conclusions

We show that the addition of non-blocking collective operations to the MPI-2 standard would be a natural extension to the existing interface. We model the potential performance benefit of overlapping communication with computation

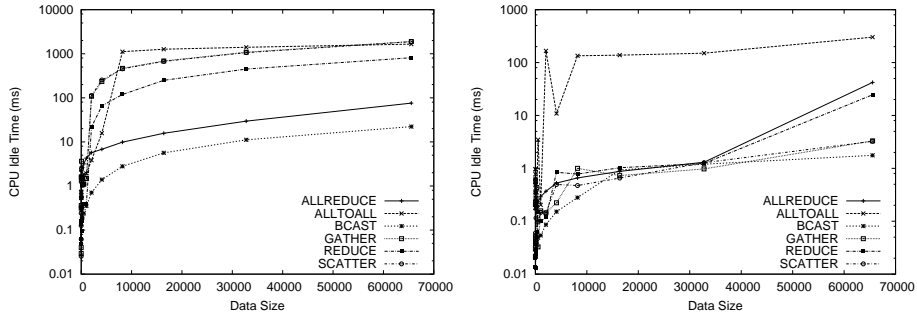


Fig. 3. CPU idle time for some collective functions with varying data size for a constant communicator size of 16 processes (left: LAM/MPI, right: MPICH2).

during collective operations. The model is proven and quantified with an extensive analysis of the CPU overhead for TCP/IP based networks. The results show clearly that, using TCP/IP, more than 70% of the CPU time is wasted in average during blocking collective operations. We assume that the gap is more than 90% for offloading based networks such as InfiniBand, Quadrics or Myrinet which do not process messages on the host CPU. Absolute measurements show the wasted time per collective which can easily be converted into wasted CPU cycles. These considerations lead to possible optimizations using non-blocking collective operations.

We propose a simple double buffering scheme to enable the use of non-blocking collective communication for existing parallel applications or algorithms. Examples include [29] and can be found at the LibNBC webpage [30]. Other double-buffering based schemes to optimize parallel implementations of more algorithms (e.g. Gaussian elimination) can be easily derived.

We implemented a portable library (LibNBC, [31]) supporting non-blocking collective operations on top of MPI-1 and port scientific applications to use the new semantics. However, implementing collective semantics on top of MPI-1 cannot easily take advantage of special hardware features to support collective communication (e.g., a hardware barrier [32]). We are planning to move the non-blocking collective implementation into the extensible Open MPI collective framework [33] to enable hardware optimized non-blocking collectives. We do also propose a MPI-2 extension [34] to support non-blocking collective operations in the MPI standard.

The NBC library is available at: <http://www.unixer.de/NBC/>.

References

1. Wagner, A., Buntinas, D., Panda, D.K., Brightwell, R.: Application-bypass reduction for large-scale clusters. In: 2003 IEEE International Conference on Cluster Computing (CLUSTER 2003), IEEE Computer Society (2003) 404–411

2. Terry, P., Shan, A., Huttunen, P.: Improving application performance on hpc systems with process synchronization. *Linux J.* **2004**(127) (2004) 3
3. Iancu, C., Husbands, P., Hargrove, P.: Hunting the overlap. In: PACT '05: Proceedings of the 14th International Conference on Parallel Architectures and Compilation Techniques (PACT'05), Washington, DC, USA, IEEE Computer Society (2005) 279–290
4. III, J.W., Bova, S.: Where's the Overlap? - An Analysis of Popular MPI Implementations (1999)
5. Lawry, W., Wilson, C., Maccabe, A.B., Brightwell, R.: Comb: A portable benchmark suite for assessing mpi overlap. In: CLUSTER, IEEE Computer Society (2002) 472–475
6. Liu, G., Abdelrahman, T.: Computation-communication overlap on network-of-workstation multiprocessors. In: Proc. of the Int'l Conference on Parallel and Distributed Processing Techniques and Applications. (1998) 1635–1642
7. Brightwell, R., Underwood, K.D.: An analysis of the impact of mpi overlap and independent progress. In: ICS '04: Proceedings of the 18th annual international conference on Supercomputing, New York, NY, USA, ACM Press (2004) 298–305
8. Dimitrov, R.: Overlapping of Communication and Computation and Early Binding: Fundamental Mechanisms for Improving Parallel Performance on Clusters of Workstations. PhD thesis, Mississippi State University (2001)
9. Calland, P.Y., Dongarra, J., Robert, Y.: Tiling on systems with communication/computation overlap. *Concurrency - Practice and Experience* **11**(3) (1999) 139–153
10. Baude, F., Caromel, D., Furmento, N., Sagnol, D.: Optimizing metacomputing with communication-computation overlap. In: PaCT '01: Proceedings of the 6th International Conference on Parallel Computing Technologies, London, UK, Springer-Verlag (2001) 190–204
11. Danalis, A., Kim, K.Y., Pollock, L., Swamy, M.: Transformations to parallel codes for communication-computation overlap. In: SC '05: Proceedings of the 2005 ACM/IEEE conference on Supercomputing, Washington, DC, USA, IEEE Computer Society (2005) 58
12. Abdelrahman, T.S., Liu, G.: Overlap of computation and communication on shared-memory networks-of-workstations. (2001) 35–45
13. Dubey, A., Tessera, D.: Redistribution strategies for portable parallel FFT: a case study. *Concurrency and Computation: Practice and Experience* **13**(3) (2001) 209–220
14. Brightwell, R., Riesen, R., Underwood, K.D.: Analyzing the impact of overlap, offload, and independent progress for message passing interface applications. *Int. J. High Perform. Comput. Appl.* **19**(2) (2005) 103–117
15. Kale, L.V., Kumar, S., Vardarajan, K.: A Framework for Collective Personalized Communication. In: Proceedings of IPDPS'03, Nice, France (2003)
16. Petrini, F., Kerbyson, D.J., Pakin, S.: The case of the missing supercomputer performance: Achieving optimal performance on the 8, 192 processors of *asci q*. In: Proceedings of the ACM/IEEE SC2003 Conference on High Performance Networking and Computing, 15-21 November 2003, Phoenix, AZ, USA, CD-Rom, ACM (2003) 55
17. Agarwal, S., Garg, R., Vishnoi, N.: The impact of noise on the scaling of collectives: A theoretical approach. In: 12th Annual IEEE International Conference on High Performance Computing, Goa, India (2005)
18. Jones, T., Dawson, S., Neely, R., Jr., W.G.T., Brenner, L., Fier, J., Blackmore, R., Caffrey, P., Maskell, B., Tomlinson, P., Roberts, M.: Improving the scalability of

- parallel jobs by adding parallel awareness to the operating system. In: Proceedings of the ACM/IEEE SC2003 Conference on High Performance Networking and Computing. (2003) 10
19. Gorlatch, S.: Send-recv considered harmful: Myths and realities of message passing. *ACM Trans. Program. Lang. Syst.* **26**(1) (2004) 47–56
 20. Pjesivac-Grbovic, J., Angskun, T., Bosilca, G., Fagg, G.E., Gabriel, E., Dongarra, J.J.: Performance Analysis of MPI Collective Operations. In: Proceedings of the 19th International Parallel and Distributed Processing Symposium, 4th International Workshop on Performance Modeling, Evaluation, and Optimization of Parallel and Distributed Systems (PMEO-PDS 05), Denver, CO (2005)
 21. Hoefler, T., Cerquetti, L., Mehlan, T., Mietke, F., Rehm, W.: A practical Approach to the Rating of Barrier Algorithms using the LogP Model and Open MPI. In: Proceedings of the 2005 International Conference on Parallel Processing Workshops (ICPP'05). (2005) 562–569
 22. Culler, D., Karp, R., Patterson, D., Sahay, A., Schauser, K.E., Santos, E., Subramanian, R., von Eicken, T.: LogP: towards a realistic model of parallel computation. In: Principles Practice of Parallel Programming. (1993) 1–12
 23. Alexandrov, A., Ionescu, M.F., Schauser, K.E., Scheiman, C.: LogGP: Incorporating Long Messages into the LogP Model. *Journal of Parallel and Distributed Computing* **44**(1) (1995) 71–79
 24. Vetter, J.S., Mueller, F.: Communication characteristics of large-scale scientific applications for contemporary cluster architectures. In: IPDPS '02: Proceedings of the 16th International Parallel and Distributed Processing Symposium, Washington, DC, USA, IEEE Computer Society (2002) 96
 25. Brightwell, R., Goudy, S., Rodrigues, A., Underwood, K.: Implications of application usage characteristics for collective communication offload. *International Journal of High-Performance Computing and Networking* **4**(2) (2006)
 26. Rabenseifner, R.: Automatic mpi counter profiling. In: 42nd CUG Conference, CUG Summit 2000. (2000)
 27. Hoefler, T., Reinhardt, M., Mehlan, T., Mietke, F., Rehm, W.: Low overhead ethernet communication for open mpi on linux clusters. In: Submitted to EuroPVM'06. (2006)
 28. Shivam, P., Wyckoff, P., Panda, D.: Emp: zero-copy os-bypass nic-driven gigabit ethernet message passing. In: Supercomputing '01: Proceedings of the 2001 ACM/IEEE conference on Supercomputing (CDROM), New York, NY, USA, ACM Press (2001) 57–57
 29. Hoefler, T., Gottschling, P., Rehm, W., Lumsdaine, A.: Optimizing a Conjugate Gradient Solver with Non-Blocking Collective Operations. (2006) Accepted for publication at the ParSim 2006 Workshop.
 30. LibNBC: <http://www.unixer.de/NBC> (2006)
 31. Hoefler, T., Lumsdaine, A.: Design, Implementation, and Usage of LibNBC. Technical report, Open Systems Lab, Indiana University (2006)
 32. Hoefler, T., Mehlan, T., Mietke, F., Rehm, W.: Adding Low-Cost Hardware Barrier Support to Small Commodity Clusters. In: 19th International Conference on Architecture and Computing Systems - ARCS'06. (2006) 343–350
 33. Squyres, J.M., Lumsdaine, A.: The Component Architecture of Open MPI: Enabling Third-Party Collective Algorithms. In: Proceedings, 18th ACM International Conference on Supercomputing, Workshop on Component Models and Systems for Grid Applications, St. Malo, France (2004)
 34. Hoefler, T., Squyres, J., Bosilca, G., Fagg, G., Lumsdaine, A., Rehm, W.: Non-Blocking Collective Operations for MPI-2. (2006)