

# OIT Technical Report

## MPI Performance of the UMBC *hpc.rs.umbc.edu* Cluster

*Samuel G Trahan*  
[straha1@umbc.edu](mailto:straha1@umbc.edu)  
UMBC OIT

**Abstract:** This technical report examines the performance of OpenMPI 1.2.6, MVAPICH2 1.0.3 and MVAPICH 1.0.1 on the UMBC *hpc.rs.umbc.edu* cluster. This report presents the bandwidth of MPI\_Alltoall, MPI\_Allreduce, MPI\_Bcast, MPI\_Send, MPI\_Recv and MPI\_Reduce, as well as the round trip time of MPI\_Send/MPI\_Recv pairs. Due to recent hardware issues, this report is restricted to 1 to 12 machine cases. A bad Infiniband card on *node032* ended the study before the planned 16, 24 and 32 machine cases. Additionally, this paper examines a bug in MVAPICH 1.0.1 that prevents its use for jobs that require more than seven nodes, as well as other issues in MVAPICH2 and OpenMPI. Lastly, this report explains the need for load balancing on *hpc.rs.umbc.edu*.

### **Introduction**

The Message Passing Interface (MPI) is the industry standard for communication between distributed memory jobs. MPI-based programs execute simultaneously on several machines at once. The various instances of the program (referred to as *processes*) communicate by passing messages between each other. MPI makes this easy by providing subroutines that perform message passing for the program. The MPI implementation handles the details of the particular underlying communication interface, and algorithms for efficient communication, allowing scientists and other programmers to worry about matters more relevant to their work.

MPI supports a number of communication patterns through C, Fortran and C++ application programming interfaces (APIs; essentially sets of callable functions). Most of these communication patterns involve one process sending a message to all other processes (or receiving a message from all other processes), some involve communication between only two processes and others involve all processes communicating with each other simultaneously. Some of the more common MPI functions are the ones benchmarked by this study:

- *MPI\_Send*, *MPI\_Recv* – *MPI\_Send* sends a single message from one process to another specific process. *MPI\_Recv* receives the sent message.
- *MPI\_Bcast* – *MPI\_Bcast* sends a message from one process to all other processes.
- *MPI\_Reduce* – *MPI\_Reduce* takes an array of numbers on one process, splits the array up into several parts and sends each part to a different process. The processes then perform some task on the data (such as summing it or finding the maximum value). They then send the result back to the originating process, which combines the results (such as by summing the sums, or getting the maximum of the maximums).
- *MPI\_Alltoall* – *MPI\_Alltoall* sends a distinct message from every process to every other process. (That is, if you have  $n$  processes, then each one sends  $n-1$  messages – one to each other process.)

- *MPI\_Allreduce* – *MPI\_Allreduce* is somewhat of a combination of *MPI\_Alltoall* and *MPI\_Reduce*. As in *MPI\_Reduce*, one process provides a list of numbers which is then split up evenly among all other processors. The processors then work together to perform a single computation on the data, such as summing the numbers. Unlike *MPI\_Reduce*, the result is not merely sent back to the originating process. It is also sent to all other processes.

In addition to bandwidth measurements, we also will examine unidirectional latency of *MPI\_Send* and round trip time of *MPI\_Send*/*MPI\_Recv* pairs. The unidirectional latency measures how long it takes to tell the MPI implementation to send a message out. The round trip time measures how total time it takes to send a message out (via *MPI\_Send*) and receive one back again (via *MPI\_Recv*).

Throughout this document, *node* or *machine* will refer to one physical computer. *Processor* or *core* will refer to one processor core on one machine. The *hpc.rs.umbc.edu* cluster is made up of thirty-two nodes, each with two processor cores. Due to hardware failures, this study is limited to jobs running on twelve nodes or less.

## Methodology

This study uses the Low Level Characterization Benchmark (*llcbench*) benchmarking suite. It is designed to benchmark the smallest possible subset of the MPI implementation at a time. It does this through the use of tight loops surrounded by timers. For example, here is the loop in the *MPI\_Send*/*MPI\_Recv* bandwidth test:

```
TIMER_START;
for (i=0; i<cnt; i++)
    mp_send(dest_rank, 1, sendbuf, bytes);
mp_recv(dest_rank, 2, destbuf, 4);
TIMER_STOP;
```

Here, *mp\_send* and *mp\_recv* are *llcbench*'s wrappers around *MPI\_Send* and *MPI\_Recv*, respectively. *TIMER\_START* and *TIMER\_STOP* are C Preprocessor macros that calculate the current time via the use of the UNIX *clock\_gettime* function's real time mode, which returns the time in nanoseconds. The single call to *mp\_recv* ensures that all of the *mp\_send* calls finish before the call to *TIMER\_STOP*. The reason that works is that the remote process does not call *mp\_send* until after it has called *mp\_recv cnt* times (and process 0's *mp\_recv* cannot complete until someone has sent data to process 0).

I compiled *llcbench* using OpenMPI 1.2.6, MVAPICH 1.0.1 and MVAPICH 1.0.3, each with the PGI compiler and the GCC compiler (six different combinations). I ran the *llcbench* benchmark suite for all six combinations on 1, 2, 3, 4, 6, 8 and 12 machines, using 1, 2 and 4 processors per machine.

## Results

The results are split into three sections. The first section discusses the latency and two process communication bandwidth (the good news). The second section discusses the multiprocess communication bandwidth (the bad news). The last section discusses the critical bug in MVAPICH 1.0.3 and other MPI-related issues.

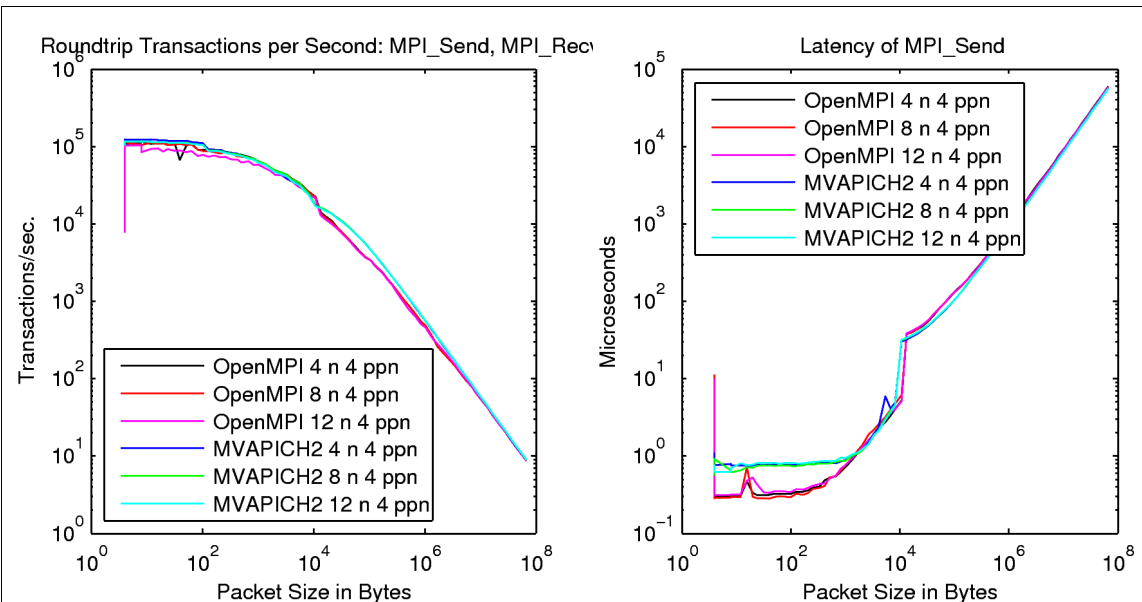
### Two Process Communication

The two process communication methods analyzed here use the *MPI\_Send* and *MPI\_Recv* commands to send messages between the rank 0 process and the rank *n*-1 process, where *n* is the number of processes.

The first result I will present here is the latency and round trip time – two aspects of the same benchmark. The latency benchmark calls MPI\_Send five thousand times and calculates the average time taken by an MPI\_Send call. That benchmark can be deceptive since the MPI implementation has the freedom to store multiple MPI\_Send messages and send them all at once in a single communication.

## Latency and Round Trip Time

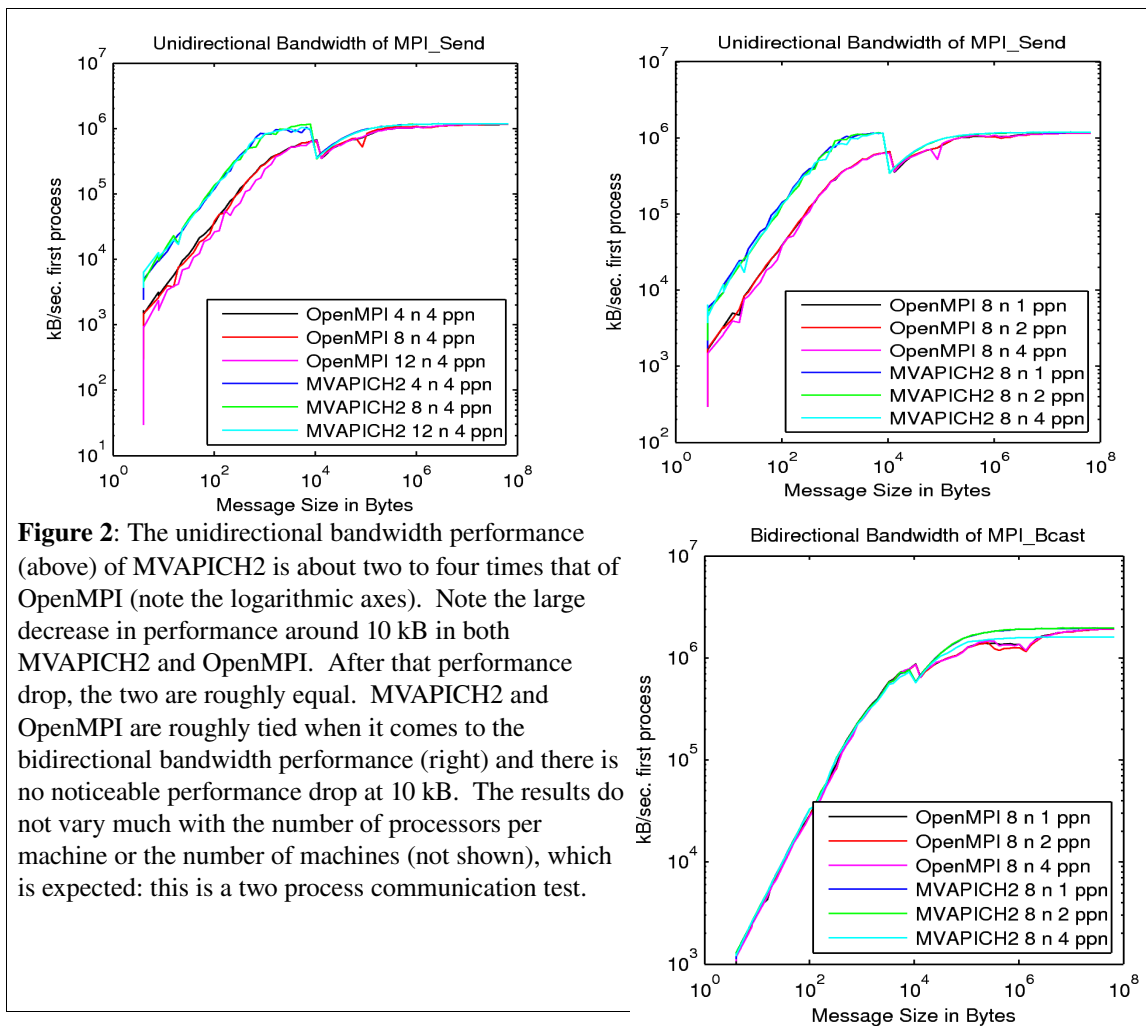
The roundtrip benchmark does not have that problem. It calls MPI\_Send followed by MPI\_Recv, repeating that pair of calls five thousand times. It then calculates the average time taken by a pair of calls. In order for the MPI\_Recv to succeed, the remote process must have received and acknowledged the MPI\_Send. That ensures that the MPI implementation cannot simply cache five thousand MPI\_Send calls and send them all at once. The results of the round trip benchmark do not vary noticeably between MPI implementations, machine counts or processor per node variations. OpenMPI displays slightly better performance in the latency benchmark for small message sizes, indicating that it does a better job of caching multiple messages before sending them (see Figure 1).



**Figure 1:** The round trip performance (left) is measured in transactions per second. It starts at a rate of 90,000 to 115,000 transactions per second, which corresponds to a round trip time of about 10 microseconds – equivalent to a unidirectional latency of only 5 microseconds. The round trip time remains quite fast until around 1 kB messages, when it begins to drop off roughly linearly with the message size. There are no noticeable differences between the different MPI implementations. There are also no differences between different numbers of processors per node (not shown). The latency (right) displays similar behavior, except that OpenMPI has about half of the latency of MVAPICH2 for small message sizes. Note that the latency benchmark, unlike the round trip benchmark, allows the MPI implementation to cache several MPI\_Send calls before sending the messages. Thus the latency benchmark indicates that OpenMPI does a better job of caching MPI\_Sends than MVAPICH2 does.

## Bandwidth: Unidirectional and Bidirectional

The next type of benchmark is two process bandwidth benchmarks. The *llcbench* suite contains unidirectional and bidirectional two process bandwidth benchmarks. The unidirectional bandwidth benchmark calls `MPI_Send` many times, followed by a single `MPI_Recv` (to ensure that the MPI implementation actually sends all of the `MPI_Send` messages). It then reports the average bandwidth in kB/s. The bidirectional bandwidth benchmark calls `MPI_Send` once, then `MPI_Recv` once, and repeats the pair thousands of times. It then reports the average bandwidth. In the unidirectional bandwidth benchmarks, MVAPICH2 consistently outperforms OpenMPI by a factor of two to four for message sizes of less than 10 kB. In the bidirectional bandwidth benchmark, the two have approximately equal performance (Figure 2).



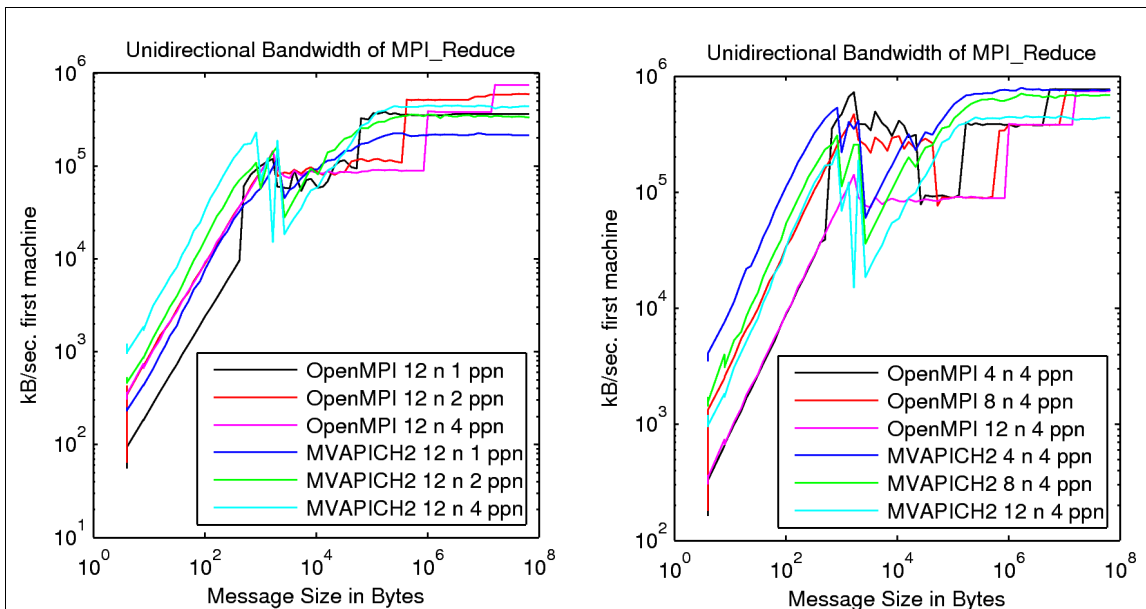
## Multiprocess Benchmarks

The next class of benchmarks examine the bandwidth of MPI\_Allreduce, MPI\_Alltoall, MPI\_Bcast, MPI\_Reduce. All four of those MPI calls communicate between all processes in the MPI program, unlike MPI\_Send and MPI\_Recv. As with MPI\_Send and MPI\_Recv, *llcbench* calculates the bandwidth in terms of the number of bytes per second sent out from process 0. However, in these benchmarks, I convert the bandwidth to kB/s for the whole machine by dividing *llcbench*'s bandwidth result by the number of processors used per node.

The performance of MPI\_Reduce (Figure 3) and MPI\_Allreduce (Figure 4) is reasonable. MPI\_Alltoall (Figure 5) suffers serious performance problems, while MPI\_Bcast (Figure 6) displays similar performance to MPI\_Reduce and MPI\_Allreduce.

### MPI\_Reduce Bandwidth Per Machine

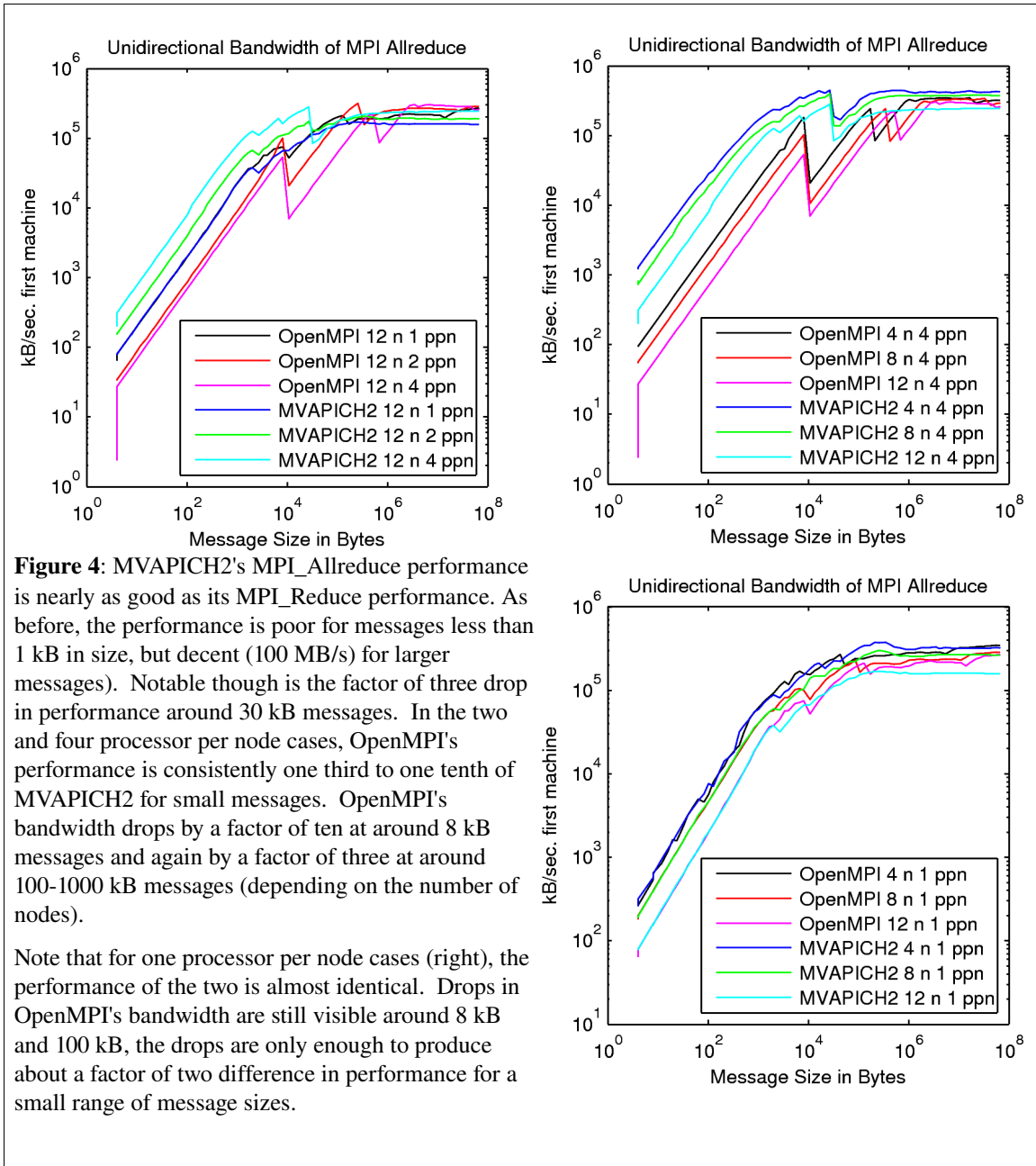
In the MPI\_Reduce benchmark, process 0 assembles an array of MPI\_INT values and instructs the processes in the MPI program to collectively sum the array. The message size reported in Figure 3 is the total size of that array. There are no guarantees made about the algorithm the MPI implementation will use to divide up the array among the available processes – the only guarantee is that the processes will collectively sum the data.



**Figure 3:** MVAPICH2 outperforms OpenMPI by about a factor of three for small message sizes. MVAPICH2's bandwidth drops abruptly by a factor of two to five at messages around 3 kB in size. A similar drop is visible in the OpenMPI performance when using four or eight nodes (right) though that feature is not present in the one and two processor per machine. cases (not shown). One important result to take away from this is that MPI\_Reduce's performance is very poor for messages less than 1 kB in size for both OpenMPI and MVAPICH2.

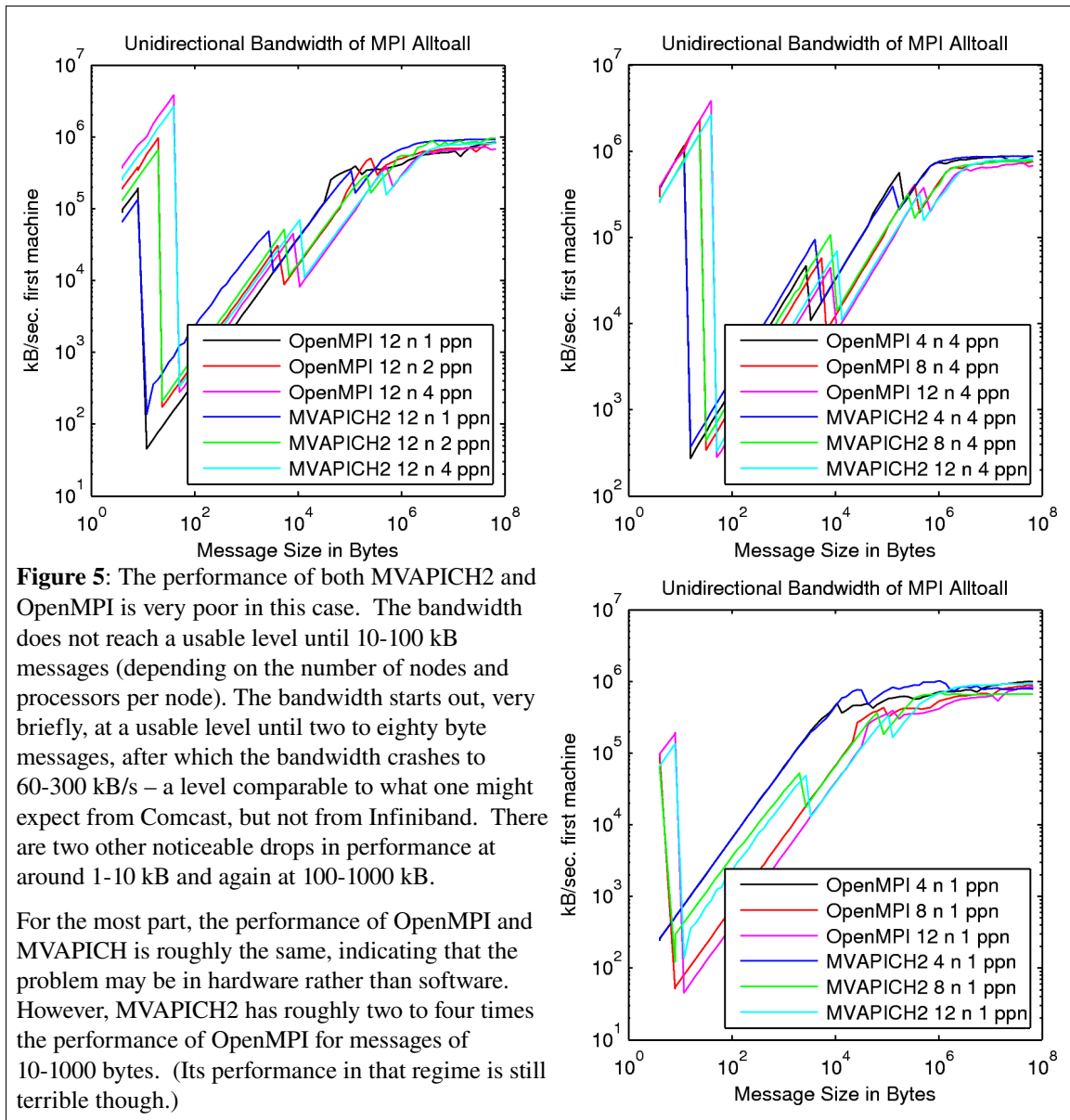
## MPI\_Allreduce Bandwidth Per Machine

MPI\_Allreduce works much like MPI\_Reduce – process 0 creates an array of MPI\_INT values and orders the collection of processors to sum the results. Unlike MPI\_Reduce, MPI\_Allreduce forwards the results of that sum to all processors. The performance of MPI\_Allreduce is decent, though there are some ranges of message sizes over which OpenMPI performs far more poorly than MVAPICH2.



## MPI\_Alltoall Bandwidth Per Machine

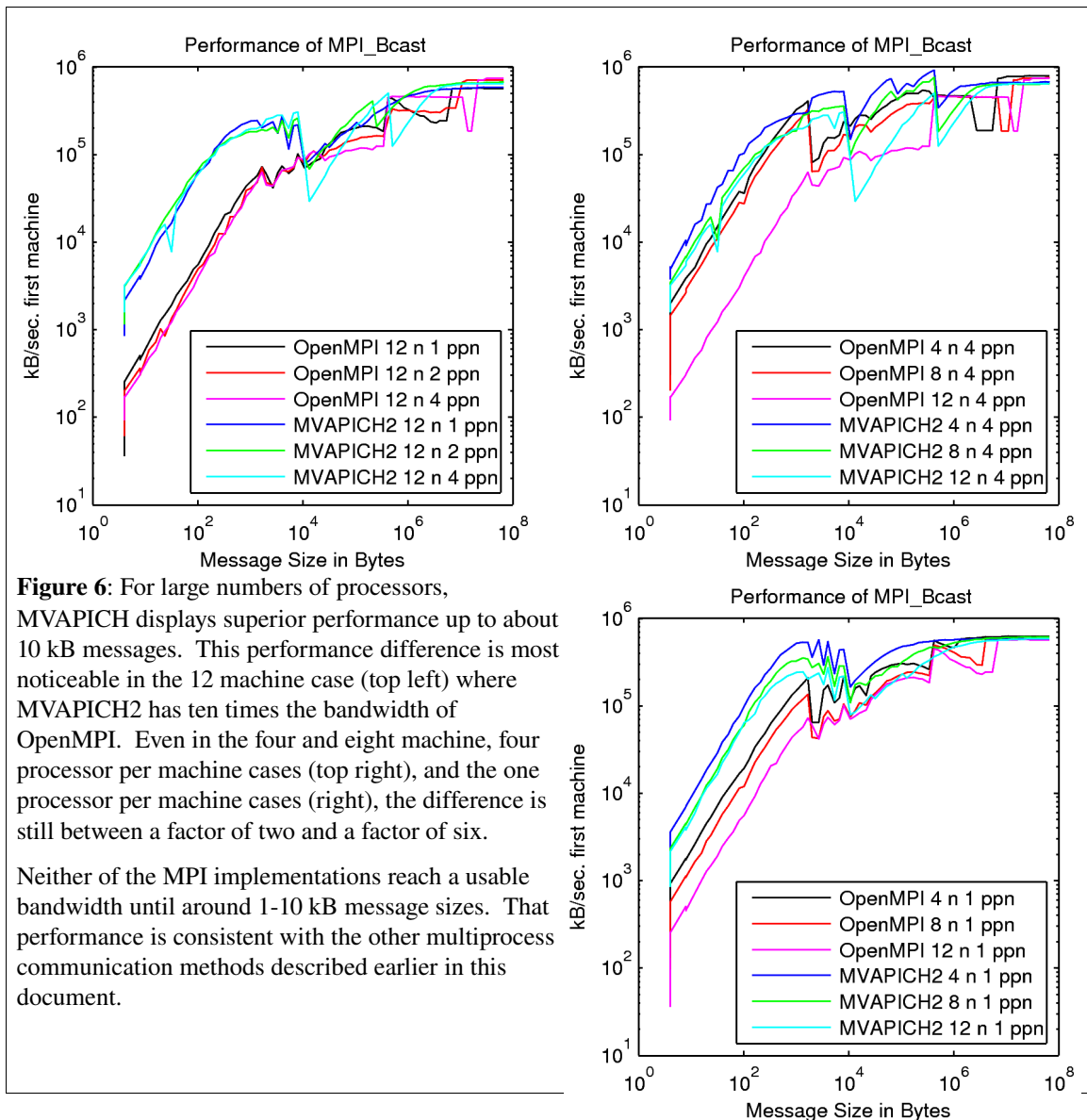
MPI\_Alltoall sends a distinct message from each process to each other process. The MPI implementation is free to distribute these messages by any algorithm it wishes. Unlike the past benchmarks, the results here are very disappointing. Over a large range of message sizes, MVAPICH2 and OpenMPI MPI\_Alltoall bandwidth is atrocious nearly to the point of being useless.





## MPI\_Bcast Bandwidth Per Machine

MPI\_Bcast broadcasts a message from process 0 to all other processes. As with the other calls, there are no guarantees about how the message will get to its destinations. The *llcbench* suite calculates MPI\_Bcast bandwidth by running MPI\_Bcast five thousand times, and then MPI\_Recv once for each process whose rank is greater than zero (to ensure that the broadcasts complete rather than being cached somewhere). I then multiply the bandwidth by the number of processors used per machine to calculate the total bandwidth of the first machine.



**Figure 6:** For large numbers of processors, MVAPICH displays superior performance up to about 10 kB messages. This performance difference is most noticeable in the 12 machine case (top left) where MVAPICH2 has ten times the bandwidth of OpenMPI. Even in the four and eight machine, four processor per machine cases (top right), and the one processor per machine cases (right), the difference is still between a factor of two and a factor of six.

Neither of the MPI implementations reach a usable bandwidth until around 1-10 kB message sizes. That performance is consistent with the other multiprocess communication methods described earlier in this document.



## Other MPI Issues

### **MVAPICH 1.0.1**

You may have wondered about the conspicuous absence of MVAPICH 1.0.1, which was prominently featured in the abstract of this technical report. MVAPICH 1.0.1 has a critical bug causing it to crash the user's program if he, she or it uses MPI\_Bcast to communicate between eight or more machines (cluster nodes – not processes). This bug is a known bug that has already been fixed in a developmental branch of MVAPICH 1.0.1 as of August 21, 2008 (possibly earlier). Due to this bug, I did not test MVAPICH 1.0.1.

### **MVAPICH2 1.0.3**

This MPI implementation displays reasonable performance in most cases and it does not have the critical bug that MVAPICH 1.0.1 has. It appears at first glance to be the best choice for an MPI implementation, however there are significant issues with MVAPICH2:

1. MVAPICH has no *mpirun* script. Thus users must call the low-level *mpdboot*, *mpiexec* and *mpdallexit* functions directly. The *mpiexec* function has an especially confusing argument parsing routine that requires certain arguments to appear in a certain order. In addition, the user must set up a *~/mpd.conf* file in their home directory. This process is overly complex and unnecessarily so. It will likely confuse many users.
2. Dr. Gobbert has noticed that MVAPICH2 does not recover gracefully from a crashed program, unlike OpenMPI which happily exits when your program crashes. This has led to MVAPICH2 jobs getting stuck in the queue.
3. As with OpenMPI, the bandwidth of MVAPICH2's MPI\_Alltoall is horrible.

### **OpenMPI**

OpenMPI is filled with features, and is quite easy to use but is substantially slower than MVAPICH2 in a number of cases. The MPI\_Bcast performance is two to four times slower for messages less than about 10 kB in size. The MPI\_Allreduce performance is ten to thirty times slower until around 30 kB messages. The MPI\_Reduce performance is up to four times slower for messages less than 300 bytes in size. The unidirectional MPI\_Send bandwidth is around one third that of MVAPICH for messages less than 8 kB in size. OpenMPI's bandwidth is littered with sudden performance drops at various message sizes – performance drops that are not present in the MVAPICH2 results (with the notable exception of MPI\_Alltoall). This suggests that OpenMPI may simply be misconfigured.

Fortunately, OpenMPI is configurable at runtime by the user through various flags to *mpirun*. The OpenMPI project FAQ provides a number of suggestions about how to improve bandwidth:

<http://www.open-mpi.org/faq/?category=openfabrics>

That link also provides suggestions as to how system administrators can configure switches and other hardware to improve Infiniband performance.

### **Lack of Load Balancing**

The *hpc.rs.umbc.edu* cluster is currently set up to launch jobs on the highest-numbered cluster nodes first. The first node listed in the user's job list is the highest-numbered node. That is going to wear down the

higher-numbered nodes sooner than the lower-numbered ones. It will have an especially large impact on *node032*'s Infiniband card, due to the nature of MPI programs. The first machine in the machine file gets MPI process 0. That process performs more message passing than the other processes since it tends to be the root of multiprocess communications such as MPI\_Bcast, MPI\_Alltoall, MPI\_Reduce and MPI\_Allreduce. Indeed, twice now *node032* has crashed during benchmarking, and the second time the Infiniband card failed before the rest of the machine. To make matters worse, the PBS queuing system is unable to end a job cleanly if it cannot contact the job launching demon on the job's initial node (the node on which the user's qsub script runs). That means any 32-node jobs that are running when *node032* has crashed will tie up the entire cluster until a system administrator can restart the relevant demons.

It appears that *hpc.rs.umbc.edu* is using MAUI, which includes the MOAB scheduler. According to the administrator's documentation for the MOAB scheduler, there is no load balancing. The "load balancing" page in the manual states that the scheduler only lets you mark nodes as available or unavailable – you cannot prioritize them. Marking a node as available or unavailable is considered a "load balancing" feature since, on machines that run several jobs at a time, you can mark a machine as unavailable if it has too many jobs running. That doesn't help load-balance a cluster – it only keeps a node from running more jobs than it has processors. Unfortunately, standard OpenPBS doesn't include any real load balancing either.

One possible solution to this problem is to simply reorder the node list once a week or so. There is no way to tell the *pbs\_server* demon to re-read its node list file – one can only achieve that effect by killing and restarting the *pbs\_server* demon. A possible solution to this problem is to create a script that will watch the PBS queue and the terminals on the head node, wait for a time when the queue and terminals are all idle, and then restart the *pbs\_server* demon with a new node list. The ideal solution would be a new scheduling algorithm, possibly as part of a commercial queuing system such as the commercial version of PBS or IBM's LoadLeveler.

## Conclusions and Future Work

Unfortunately, there are no good MPI implementations to choose from on the *hpc.rs.umbc.edu* cluster at this point. The best performer is MVAPICH2 1.0.3. The easiest to use is OpenMPI, but it performs much more poorly than MVAPICH2 1.0.3. MVAPICH 1.0.1 is not usable at this time due to a bug that causes the user's program to crash if he or she uses MPI\_Bcast between eight or more machines. Both MVAPICH2 and OpenMPI have extremely poor, bordering on useless, MPI\_Alltoall performance between message sizes of 10 to 10,000 bytes.

Due to a failure of the Infiniband card on *node032*, I have temporarily suspended these benchmarks. The reason *node032* is critical is that PBS refuses to kill a job if the job's initial node (the one running one's qsub script) goes down. Due to the lack of load balancing, *node032* is nearly always the initial node. Thus, if any of my jobs crash, they tie up most or all of the cluster until a system administrator can restart the pbs server. I hope to continue later on if the hardware becomes more reliable.

In the future, if MVAPICH 1.0.1 is fixed, I may re-run the benchmarks with that MPI implementation. The limited benchmarks I've run for MVAPICH 1.0.1 (using one to six machines) have shown that its performance is nearly the same as MVAPICH2 1.0.3. I also intend to test OpenMPI and MVAPICH2 all the way up to 32 nodes, and I will re-test OpenMPI while tuning its performance.