

Hybrid Parallel Programming: Performance Problems and Chances

Rolf Rabenseifner*

Summary

Most HPC systems are clusters of shared memory nodes. Parallel programming must combine the distributed memory parallelization on the node inter-connect with the shared memory parallelization inside of each node. Various hybrid MPI+OpenMP programming models are compared with pure MPI. Benchmark results of several platforms are presented. This paper analyzes the strength and weakness of several parallel programming models on clusters of SMP nodes. Benchmark results show, that the hybrid-masteronly programming model can be used more efficiently on some vector-type systems, although this model suffers from sleeping application threads while the master thread communicates. This paper analyses strategies to overcome typical drawbacks of this easily usable programming scheme on systems with weaker inter-connects. Best performance can be achieved with overlapping communication and computation, but this scheme is lacking in ease of use.

Keywords. OpenMP, MPI, Hybrid Parallel Programming, Threads and MPI, HPC, Performance.

1 Introduction

Most systems in High Performance Computing are clusters of shared memory nodes. Such hybrid systems range from small clusters of dual-CPU PCs up to largest systems like the Earth Simulator consisting of 640 SMP nodes connected by a single-stage cross-bar and with SMP nodes combining 8 vector CPUs on a shared memory [6, 9, 23, 26]. Optimal parallel programming schemes enable the application programmer to use the hybrid hardware in a most efficient way, i.e., without any overhead induced by the programming scheme. On distributed memory systems, message passing, especially with MPI [8, 17, 18], has shown to be the mainly used programming paradigm. One reason of the success

of MPI was the clear separation of the optimization: communication could be improved by the MPI library, while the numerics had to be optimized by the compiler. On shared memory systems, directive-based parallelization was standardized with OpenMP [20], but there is also a long history of proprietary compiler-directives for parallelization. The directives handle mainly the work sharing; there is no data distribution.

On hybrid systems, i.e., on clusters of SMP nodes, parallel programming can be done in several ways: one can use pure MPI, or some schemes combining MPI and OpenMP, e.g., calling MPI routines only outside of parallel regions (which is herein named the *masteronly* style), or using OpenMP on top of a (virtual) distributed shared memory (DSM) system. A classification on MPI and OpenMP based parallel programming schemes on hybrid architectures is given in Sect. 2. Unfortunately, there are several mismatch problems between the (hybrid) programming schemes and the hybrid hardware architectures. Often, one can see in publications, that applications may or may not benefit from hybrid programming depending on some application parameters, e.g., in [10, 12, 15, 28].

Sect. 3 gives a list of major problems often causing a degradation of the speed-up, i.e., causing that the parallel hardware is utilized only partially. Sect. 4 shows, that there isn't a silver bullet to achieve an optimal speed-up. Sect. 5 discusses strategies to overcome typical drawbacks of the hybrid masteronly style. With these modifications, efficiency can be achieved together with the ease of parallel programming on clusters of SMPs. The conclusions are provided in Sect. 6.

2 Parallel programming on hybrid systems, a classification

Often, *hybrid MPI+OpenMP programming* denotes a programming style with OpenMP shared memory parallelization inside the MPI processes (i.e., each MPI process itself has several OpenMP threads) and com-

*High-Performance Computing-Center (HLRS), University of Stuttgart, Allmandring 30, D-70550 Stuttgart, Germany rabenseifner@hlrs.de, www.hlrs.de/people/rabenseifner/

municating with MPI between the MPI processes, but *only outside of parallel regions*. For example, the MPI parallelization is based on a domain decomposition, the MPI communication mainly exchanges the halo information after each iteration of the outer numerical loop, and these numerical iterations itself are parallelized with OpenMP, i.e., (inner) loops inside of the MPI processes are parallelized with OpenMP work-sharing directives. However, this scheme is only one style in a set of different hybrid programming methods. This hybrid programming scheme will be named *masteronly* in the following classification, which is based on the question, when and by which thread(s) the messages are sent between the MPI processes:

1. **Pure MPI:** each CPU of the cluster of SMP nodes is used for one MPI process. The hybrid system is treated as a flat massively parallel processing (MPP) system. The MPI library has to optimize the communication by using shared memory based methods between MPI processes on the same SMP node, and the cluster interconnect for MPI processes on different nodes.
2. Hybrid MPI+OpenMP without overlapping calls to MPI routines with other numerical application code in other threads:
 - (a) **Hybrid masteronly:** MPI is called only outside parallel regions, i.e., by the master thread.
 - (b) **Hybrid multiple/masteronly:** MPI is called outside the parallel regions of the application code, but the MPI communication is done itself by several CPUs: The thread parallelization of the MPI communication can be done
 - automatically by the MPI library routines, or
 - explicitly by the application, using a full thread-safe MPI library.

In this category, the non-communicating threads are sleeping (or executing some other applications, if non-dedicated nodes are used). This problem of idling CPUs is solved in the next category:

3. Overlapping communication and computation: While the communication is done by the master thread (or a few threads), all other non-communicating threads are executing application

code. This category requires, that the application code is separated into two parts: the code that can be overlapped with the communication of the halo data, and that code that must be deferred until the halo data is received. Inside of this category, we can distinguish two types of sub-categories:

- How many threads communicate:
 - (A) **Hybrid funneled:** Only the master thread calls MPI routines, i.e., all communication is funneled to the master thread.
 - (B) **Hybrid multiple:** Each thread handles its own communication needs (B1), or the communication is funneled to more than one thread (B2).
 - Except in case B1, the communication load of the threads is inherently unbalanced. To balance the load between threads that communicate and threads that do not communicate, the following load balancing strategies can be used:
 - (I) **Fixed reservation:** reserving a fixed amount of threads for communication and using a fixed load balance for the application between the communicating and non-communicating threads.
 - (II) **Adaptive.**
4. **Pure OpenMP:** based on virtual distributed shared memory systems (DSM), the total application is parallelized only with shared memory directives.

Each of these categories of hybrid programming has different reasons, why it is not appropriate for some classes of applications or classes of hybrid hardware architectures. The paper focuses on the first three methods; for pure OpenMP approaches, the reader is referred to [2, 11, 13, 16, 22, 24, 25]. Other parallel programming models that can also be used on clusters of SMPs are, e.g., UPC [4, 7], Co-Array Fortran [19], MLP [5], and HPF [1]. Different SMP parallelization strategies in the hybrid model are studied in [27] and in [3] for the NAS parallel benchmarks. The following section shows major problems of mismatches between programming and hardware architecture.

0	1	4	5	8	9	12	13
2	3	6	7	10	11	14	15
16	17	20	21	24	25	28	29
18	19	22	23	26	27	30	31
32	33	36	37	40	41	44	45
34	48	49	38	42	52	53	46
35	50	51	39	43	54	55	47

Figure 1: Example for optimized ranking on a hybrid system with 4 CPUs per SMP node

3 Mismatch problems

All these programming styles on clusters of SMP nodes have advantages, but also serious disadvantages based on mismatch problems between the (hybrid) programming scheme and the hybrid architecture. With pure MPI, minimizing of the inter-node communication requires that the application-domain’s neighborhood-topology matches with the hardware topology. Pure MPI also introduces intra-node communication on the SMP nodes that can be omitted with hybrid programming. But, such MPI+OpenMP programming is not able to achieve full inter-node bandwidth on all platforms for any subset of inter-communicating threads. With masteronly style, all non-communicating threads are idling. CPU time is also wasted, if all CPUs of an SMP node communicate, although a few CPUs are already able to saturate the inter-node bandwidth. With hybrid masteronly programming, additional overhead is induced by all OpenMP synchronization, but also by additional cache flushing between the generation of data in parallel regions and the consumption in subsequent message passing routines and calculations in subsequent parallel sections.

Overlapping of communication and computation is a chance for an optimal usage of the hardware, but causes serious programming effort in the application itself, in the OpenMP parallelization and in the load balancing.

3.1 The topology problem with pure MPI

To minimize the inter-node communication needs, it is important that the MPI processes hosted by each SMP node are neighbors in the domain decomposition. Otherwise more data is shipped over the inter-node network instead of using the shared memory based communication inside of the nodes.

With structured Cartesian grids, sometimes, the de-

fault ranking of the MPI processes used in the application’s domain decomposition does not minimize the inter-node traffic. Some vendors provide options to choose several ranking schemes in `MPI_COMM_WORLD`, e.g., (a) sequential ranking inside of the SMP nodes, or (b) round robin ranking, numbering serially all first CPUs of the nodes, then the second ones, and so on. The application programmer may use the Cartesian topology interface in MPI-1 to find an optimal topology map (application topology to clustered hardware topology), but this optimization feature is still not implemented in many MPI libraries. As an alternative, the user can implement his/herself a process ranking that fits to the sequential ranking in `MPI_COMM_WORLD`. Fig. 1 shows an example ranking of 8x7 Cartesian grid, distributed on 14 4-way SMP nodes. In principal, the MPI processes are combined to SMP nodes in a first step, and afterwards, all MPI processes can be numbered sequentially. With `MPI_Comm_create`, such an optimal hybrid numbering scheme can be mapped onto the natural Cartesian numbering (0-11 on the first row, ...). With a subsequent creation of an Cartesian communicator (without reordering of the ranks), one can achieve an optimal mapping of natural Cartesian ranking on the hybrid hardware architecture, and allowing the usage of all MPI-1 Cartesian functions.

The topology problem with unstructured grids on hybrid architectures can be solved by implementing the load balancing with two steps: First, the total grid is mapped on the SMP nodes (e.g., by using ParMETIS [14] with the number of nodes), and second, on each node, the grid cells for the node are distributed to the MPI processes belonging to the node (e.g., again with ParMETIS on the sub-communicator combining these MPI processes). This method is fast, but not optimal. Nevertheless, it can be used as long as stable products are not available for hybrid platforms.

In both cases (structured and unstructured grids), a modification of the application is often necessary to achieve an optimal mapping of the pure MPI processes onto the CPUs of a cluster of SMP nodes.

3.2 Unnecessary intra-node communication

With pure MPI, additional intra-node communication is necessary between the MPI processes. With hybrid MPI+OpenMP programming, this intra-node communication is substituted by direct accesses to the application data structures in the shared memory.

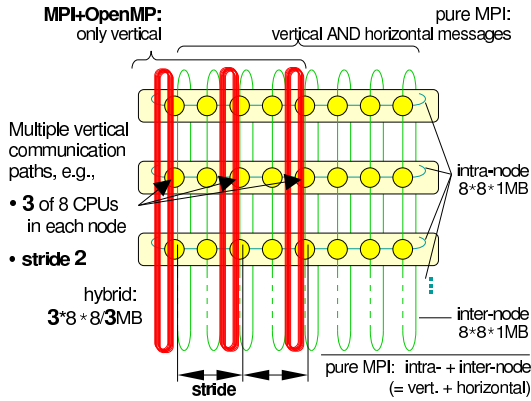


Figure 2: Communication pattern with *hybrid MPI+OpenMP* style and with *pure MPI* style.

This drawback of pure MPI programming is less important, if the intra-node communication is significantly faster than the inter-node communication.

In this case, the major question is, whether the topology problem described in the previous section is solved, i.e., pure MPI and hybrid MPI+OpenMP versions of the application exchange the same amount of data on the inter-node network.

3.3 The inter-node bandwidth problem

With hybrid masteronly or funneled style, all communication must be done by the master thread. The benchmark measurements in Fig. 4 and the inter-node results in Tab. 1 show, that on several platforms, the available aggregated inter-node bandwidth can be achieved only, if more than one thread is used for the communication with other nodes.

In this benchmark, all SMP nodes are located in a logical ring. Each CPU sends messages to the corresponding CPU in the next node and receives from the previous node in the ring. The benchmark is done with pure MPI, i.e., one MPI process on each CPU, except for Cray X1, where we used as smallest entity an MSP (which itself has 4 SSPs [=CPUs]). Fig. 2 shows the communication patterns. The aggregated bandwidth per node is defined as the number of all bytes of all messages on the inter-node network divided by the time needed for the communication and divided by the number of nodes. Note that in this definition, each message is counted only once, and not twice.¹ Fig. 3 shows the

¹The hardware specification typically presents the duplex

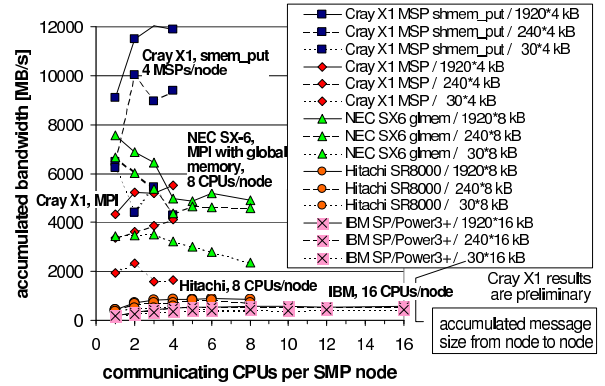


Figure 3: Aggregated bandwidth per SMP node.

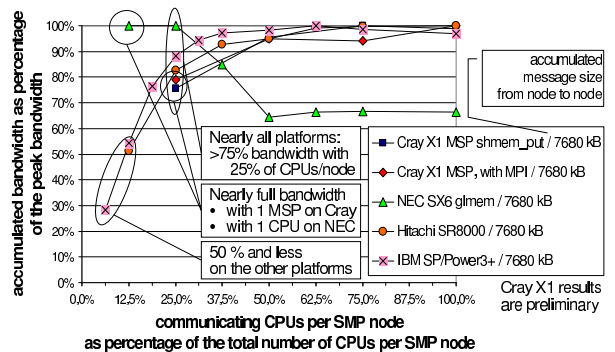


Figure 4: Aggregated bandwidth per SMP node.

absolute bandwidth over the number of CPUs (or MSPs at Cray X1), Fig. 4 shows relative values, i.e., the percentage of the achieved peak bandwidth in each system over the percentage of CPUs of a node.

One can see, that only on the NEC SX-6 and Cray X1 systems, one can achieve more than 75 % of the peak bandwidth already with **one** CPU (or MSP on Cray X1) per node.

On the other systems, the hybrid masteronly or funneled programming scheme can achieve only a small percentage of the peak inter-node bandwidth. [21] has compared the pure MPI with the masteronly scheme. For this comparison, each MPI process in the pure MPI scheme has also to exchange messages between the pro-

inter-node bandwidth by counting each message twice, i.e. as incoming and outgoing message at a node, e.g., on a Cray X1, 25.6 GB/s = 2*12.8 GB/s; the measured 12 GB/s (shmem.put) must be compared with the 12.8 GB/s of the hardware specification.

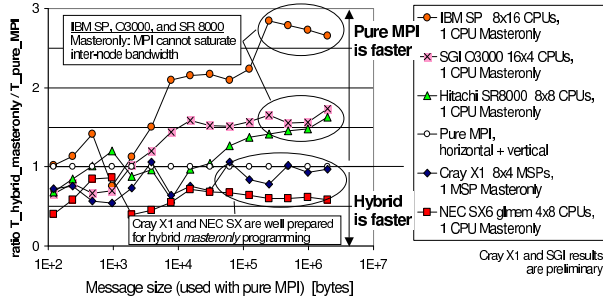


Figure 5: Ratio of hybrid communication time to pure MPI communication time.

cesses in the same node. These intra-node messages have the same size as the inter-node messages, (c.f. Fig. 2). Fig. 5 shows the ratio of the inter-node communication time of hybrid MPI+OpenMP masteronly style divided by the time needed for the inter- and intra-node communication with pure MPI. In case of hybrid masteronly style, the messages must transfer the accumulated data of the parallel inter-node messages in pure MPI style, i.e., the message size is multiplied with the number of CPUs of an SMP node. In Fig. 5, one can see a significantly better communication time with pure MPI, on those platforms, on which the inter-node network cannot be saturated by the master thread. In this benchmark, the master thread in the masteronly scheme was emulated by an MPI process (and the other threads by MPI process waiting in a barrier). On most platforms, the measurements were verified with a benchmark using hybrid compilation and hybrid application start-up. The diagram compares experiments with the same aggregated message-size in the inter-node communication; on the x-axis, the corresponding number of bytes in the pure-MPI experiment is shown. This means, e.g., that the message size in the hybrid-masteronly experiment on a 16-CPU-per-node system is 16 times larger than in the experiments with pure MPI.

Benchmark platforms were: a Cray X1 with 16 nodes at Cray; the NEC SX-6 with 24 nodes and IXS interconnect at the DKRZ, Hamburg, Germany; the Hitachi SR8000 with 16 nodes at HLRZ, Stuttgart, Germany; the IBM SP-Power3 at NERSC, USA; the SGI Origin 3000 *Lomax* with 512 CPUs at NASA/Ames Research Center, NAS, USA; the SUN Starfire 9500 cluster at the RWTH Aachen, Germany.

3.4 The sleeping-threads problem and the saturation problem

The two most simple programming models on hybrid systems have both the same problem although they look quite different: With hybrid masteronly style the non-master threads are sleeping while the master communicates, and with pure MPI, all threads try to communicate while only a few (or one) threads already can saturate the inter-node network bandwidth (expecting that the application is organized in communicating and computing epochs). If one thread is able to achieve the full inter-node bandwidth (e.g., NEC SX-6, see Fig. 3), then both problems are equivalent. If one thread can only achieve a small percentage (e.g., 28% on IBM SP), then the problem with masteronly style is significantly higher.

As example on the IBM system, if an application communicates 1 sec in the pure MPI style (i.e. $1 \times 16 = 16$ CPUsec), then this program would need about $16/0.28 = 57$ CPUsec in masteronly style, and if one would use 4 CPUs for the inter-node communication (4 CPUs achieve 88.3%) and the other 12 threads for overlapping computation, then only $4/0.883 = 4.5$ CPUsec would be necessary.

If the inter-node bandwidth cannot be achieved by one thread, then it may be a better choice to split each SMP node into several MPI processes that are itself multithreaded. Then, the inter-node bandwidth in the pure MPI and hybrid masteronly model are similar and mainly the topology, intra-node communication, and OpenMP-overhead problems determine which of both programming styles are more effective. When overlapping communication and computation, this splitting can also solve the inter-node bandwidth problem described in the previous section.

3.5 Additional OpenMP overhead

The OpenMP parallelization inside of the MPI processes may cause an additional overhead. The creation of parallel regions and the synchronization at their end induces additional work, mainly, if a fine-grained OpenMP parallelization is used. This overhead can be reduced with a coarse-grained OpenMP parallelization: The parallel regions are started only once at the beginning of the application, and OMP MASTER and BARRIER directives are used for synchronizing before and after the MPI communication. If it is not possible to parallelize the total work of the MPI processes, or if the

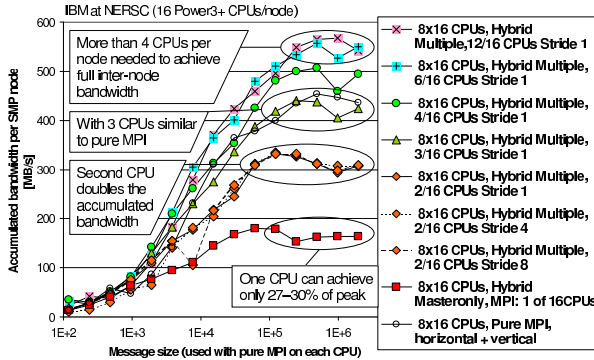


Figure 6: Aggregated bandwidth per SMP node on IBM SP with 16 Power3+ CPUs per node.

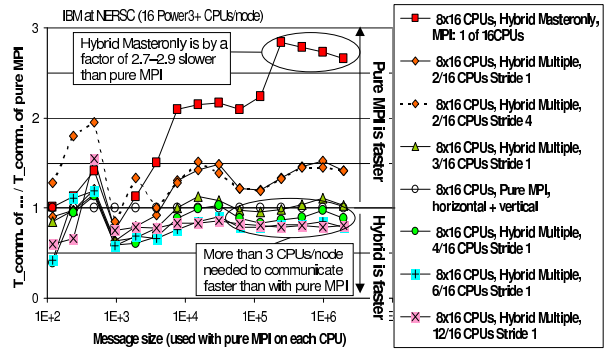


Figure 7: Ratio of communication time in hybrid models to pure MPI programming on IBM SP.

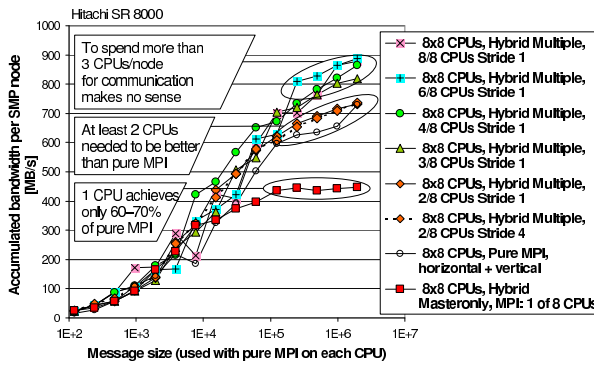


Figure 8: Aggregated bandwidth per SMP node on Hitachi SR 8000.

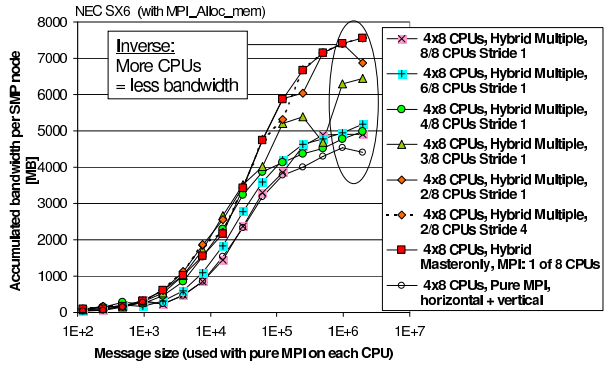


Figure 9: Aggregated bandwidth per SMP node on NEC SX-6.

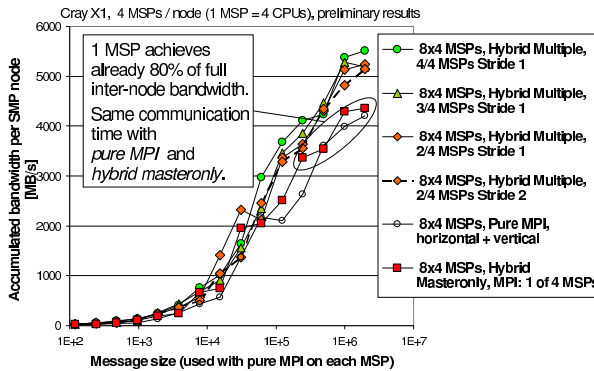


Figure 10: Aggregated bandwidth per SMP node on Cray X1, MSP-based MPI-parallelization.

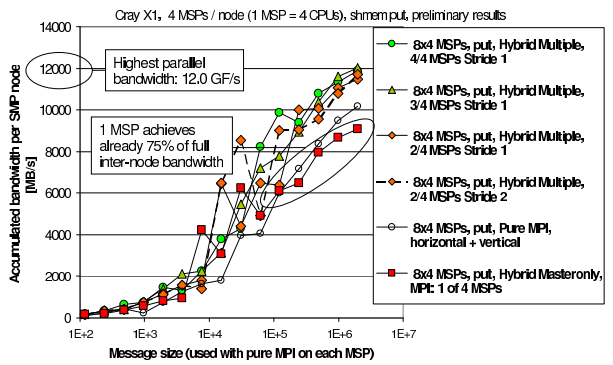


Figure 11: Aggregated bandwidth per SMP node on Cray X1. MSP-based and MPI_Sendrecv is substituted by shmem_put.

OpenMP-parallelization cannot be balanced perfectly, then the speedup will be also reduced according to Amdahl's law.

With pure MPI programming, the communication can be accelerated if the data is sent from the cache. Hybrid programming typically needs a barrier synchronization between the generation of data in one thread and sending the data by another thread (usually the master thread). Such barrier synchronization implies automatically an OMP FLUSH operation that may internally causes a cache flush.

3.6 Overlapping communication and computation

Overlapping of communication and computation is necessary to fully exploit the hardware. But there are three major drawbacks bound up with this programming model: First, most applications are not prepared to distinguish between numerics that can be computed without halo data (can be done immediately), and numerics that needs halo data, i.e., that must be deferred upon halo communication is done. Second, a coarse-grained OpenMP programming style is necessary. All work-sharing must be done manually based on the OpenMP-ranks of the threads (each thread has an MPI-rank and an OpenMP-rank); the OpenMP work-sharing directives cannot be further used. And third, load balancing must be implemented to compensate the different loads on communicating and computing threads.

4 Bite the bullet

Each parallel programming scheme on hybrid architectures has one or more significant drawbacks. Depending on the needed resources of an applications, the drawbacks may be major or only minor.

Programming without overlap of communication and computation

One of the two problems, *sleeping-threads* and *saturation problem* is indispensable. The major design criterion may be the topology problem:

- If it cannot be solved, pure MPI may cause too much inter-node traffic, but the masteronly scheme

implies on some platforms a slow inter-node communication due to the inter-node bandwidth problem described above.

- If the topology problem can be solved, then we can compare hybrid masteronly with pure MPI: On some platforms, wasting inter-node bandwidth with masteronly style is the major problem; it causes more CPUs longer idling than with pure MPI. For example on an IBM SP system with 16 Power3+ CPUs on each SMP node, Fig. 6 shows the aggregated bandwidth per node with the experiment described in Sect. 3.3. The *pure MPI horizontal+vertical* bandwidth is defined in this diagram by dividing the amount of inter-node message bytes (without counting the intra-node messages) by the time needed for inter- and intra-node communication, i.e., the intra-node communication is treated as overhead. One can see, that more than 4 CPUs per node must communicate in parallel to achieve full inter-node bandwidth. At least 3 CPU per node must communicate in the hybrid model to beat the pure MPI model. Fig. 7 shows the ratio of the execution time in the hybrid models to the pure MPI model. A ratio greater than 1 shows that the hybrid model is slower than the pure MPI model.

On systems with 8 CPUs per node, the problem may be reduced, e.g., as one can see on a Hitachi SR 8000 in Fig. 8. On some vector type systems, one CPU may already be able to saturate the inter-node network, as shown in Fig. 9–11. Note: the aggregated inter-node bandwidth on the SX-6 is reduced, if more than one CPU per node tries to communicate at the same time over the IXS. Fig. 10 and 11 show preliminary results on a Cray X1 system with 16 nodes. Each SMP node consists of 4 MSPs (multi streaming processors). Each MSP itself consists of 4 SSPs (single streaming processors). With MSP-based programming, each MSP is treated as a CPU, i.e., each SMP node has 4 CPUs (=MSPs) that internally use an (automatic) thread-based parallelization (= *streaming*). With SSP-based programming, each SMP node has 16 CPUs (=SSPs). Preliminary results with the SSP-mode have shown, that the inter-node bandwidth is partially bound to the CPUs, i.e., that the behavior is similar to the 16-way IBM system.

Similar to the multi-threaded implementation of MPI on the Cray MSPs, it is also possible on all other platforms to use multiple threads inside of the MPI communication routines if the application uses the hybrid

	Master-only, inter-node bandw. [GB/s]	pure MPI, inter-node bandw. [GB/s]	Master-only bw / max. inter-node bw [%]	pure MPI, intra-node bandw. [GB/s]	memory bandwidth [GB/s]	Peak performance [GFLOP/s]	max. inter-node bw / peak perf. [B/FLOP]	#nodes * #CPUs per SMP node
Cray X1, shmem_put preliminary results	9.27	12.34	75 %	33.0	136	51.2	0.241	8 * 4 MSPs
Cray X1, MPI preliminary results	4.52	5.52	82 %	19.5	136	51.2	0.108	8 * 4 MSPs
NEC SX-6, MPI with global memory	7.56	4.98	100 %	78.7 (93.7+)	256	64	0.118	4 * 8 CPUs
NEC SX-5Be local memory	2.27	2.50 _{a)}	91 %	35.1	512	64	0.039	2 * 16 CPUs a) only 8 CPUs
Hitachi SR8000	0.45	0.91	49 %	5.0	32+32	8	0.114	8 * 8 CPUs
IBM SP Power3+	0.16	0.57+	28 %	2.0	16	24	0.023	8 * 16 CPUs
SGI Origin 3000 preliminary results	0.10	0.30+	33 %	0.39+	3.2	4.8	0.063	16 * 4 CPUs
SUN-fire (preliminary)	0.15	0.85	18 %	1.68				4 * 24 CPUs

Table 1: Inter- and Intra-node bandwidth for large messages compared with memory bandwidth and peak performance. All values are aggregated over one SMP node. Each message counts only once for the bandwidth calculation. Message size is 16 MB, except +) with 2 MB.

masteronly scheme. The MPI library can easily detect whether the application is inside or outside of a parallel region. With this optimization (described in more detail in Sect. 5), the communication time of the hybrid masteronly model should always be shorter than the communication time in the pure MPI scheme.

Programming with overlap of communication and computation

Although overlapping communication with computation is the chance to achieve fastest execution, this parallel programming style isn't widely used due to the lack of ease of use. It requires a coarse-grained and thread-rank-based OpenMP parallelization, the separation of halo-based computation from the computation that can be overlapped with communication, and the threads with different tasks must be load balanced.

Advantages of the overlapping scheme are: (a) the problem that one CPU may not achieve the inter-node bandwidth is no longer relevant as long as there is enough computational work that can be overlapped with the communication; (b) the saturation problem is solved as long as not more CPUs communicate in parallel than necessary to achieve the inter-node bandwidth; (c) the sleeping threads problem is solved as long as all computation and communication is load balanced among the threads.

A detailed analysis of the performance benefits of overlapping communication and computation can be found in [21].

5 Optimization Chance

On Cray X1 with MSP-based programming and on NEC SX-6, the *hybrid masteronly* communication pattern is faster than the *pure MPI*. Although both systems have vector-type CPUs, the reasons for these performance results are quite different: On the NEC SX-6, the hardware of one CPU is really able to saturate the inter-node network if the user data resides in global memory. On the Cray X1, each MSP consist of 4 SSPs (=CPUs). MPI communication issued by one MSP seems internally to be multi-streamed by all 4 SSPs. With this multi-threaded implementation of the communication, Cray can achieve 75–80% of the full inter-node bandwidth, i.e., of the bandwidth that can be achieved if all MSPs (or all SSPs) communicate in parallel.

This approach can be generalized for the masteronly style. Depending on whether the application itself is translated for pure MPI approach, hybrid MPI + automatic SMP-parallelization, or hybrid MPI+OpenMP, the linked MPI library itself can also be parallelized with OpenMP directives or vendor-specific directives.

Often, the major basic capabilities of an MPI library are to put data into a shared memory region of the

destination process (RDMA put), or to get data from the source process (RDMA get), or to locally calculate reduction operations on a vector, or to handle derived datatypes and data. All these operations (and not the envelop handling of the message passing interface) can be implemented multi-threaded, e.g., inside of a parallel region. In the case, that the application calls the MPI routines outside of parallel *application* regions, the parallel region inside of the MPI routines will allow a thread-parallel handling of these basic capabilities. In the case, the application overlaps communication and computation, the parallel region inside of the MPI library is a nested region and will get only the (one) thread on which it is already running. Of course, the parallel region inside of MPI should only be launched, if the amount of data that must be transferred (or reduced) exceeds a given threshold.

This method optimizes the bandwidth without a significant penalty to the latency. On the Cray X1, currently only 4 SSPs are used to stream the communication in MSP mode achieving only 75–80% of peak. It may be possible to achieve full inter-node bandwidth, if the SSPs of an additional MSP would also be applied. With such a multi-threaded implementation of the MPI communication for masteronly-style applications, there is no further need (with respect to the communication time) to split large SMP nodes into several MPI processes each with a reduced number of threads (as proposed in Sect. 3.4).

6 Conclusions

Different programming schemes on clusters of SMPs show different performance benefits or penalties on the hardware platforms benchmarked in this paper. Table 1 summarizes the results. Cray X1 with MSP-based programming and NEC SX-6 are well designed for the hybrid MPI+OpenMP masteronly scheme. On the other platforms, as well as on the Cray X1 with SSP-based programming, the master thread cannot saturate the inter-node network which is a significant performance bottleneck for the masteronly style.

To overcome this disadvantage, a multi-threaded implementation of the basic device capabilities in the MPI libraries is proposed in Sect. 5. Partially, this method is already implemented in the Cray X1 MSP-based MPI-library.

This enhancement of current MPI implementations implies that the hybrid masteronly communication

should be always faster than pure MPI communication. Both methods still include the sleeping threads or saturated network problem, i.e., that more CPUs are used for communicating than really needed to saturate the network. This drawback can be solved with overlapping of communication and computation, but this programming style needs extreme programming effort.

To achieve an optimal usage of the hardware, one can also try to use the idling CPUs for other applications, especially low-priority single-threaded or multi-threaded non-MPI applications if the parallel high-priority hybrid application does not use the total memory of the SMP nodes.

Acknowledgments

The author would like to acknowledge his colleagues and all the people that supported this project with suggestions and helpful discussions. He would especially like to thank Gerhard Wellein at RRZE, Monika Wierse, Wilfried Oed, and Tom Goozen at CRAY, Holger Berger at NEC, Gabriele Jost at NASA, Dieter an Mey at RWTH Aachen, and Horst Simon at NERSC for their assistance in executing the benchmark on their platforms. This research used resources of the HLRS Stuttgart, LRZ Munich, RWTH Aachen, Cray Inc., NEC, NASA/AMES, and resources of the National Energy Research Scientific Computing Center, which is supported by the Office of Science of the U.S. Department of Energy.

Biographies

Rolf Rabenseifner studied mathematics and physics at the University of Stuttgart. Since 1984, he is working at the High Performance Computing Center Stuttgart (HLRS). He led the projects DFN-RPC, a remote procedure call tool, and MPI-GLUE, the first meta-computing MPI combining different vendor's MPIs without loosing the full MPI interface. In his dissertation, he developed a controlled logical clock as global time for trace-based profiling of parallel and distributed applications. Since 1996, he has been a member of the MPI-2 Forum. From January to April 1999, he was an invited researcher at the Center for High-Performance Computing at Dresden University of Technology. Currently, he is head of the Department for Parallel Computing of the HLRS. His research interests include hybrid parallel programming models, MPI profiling, and

benchmarking of parallel communication and I/O. He is also involved in teaching projects.

References

- [1] Siegfried Benkner, Thomas Brandes, *High-Level Data Mapping for Clusters of SMPs*, in proceedings of the 6th International Workshop on High-Level Parallel Programming Models and Supportive Environments, HIPS 2001, San Francisco, USA, April 2001, Springer LNCS 2026, pp 1–15.
- [2] Rudolf Berrendorf, Michael Gerndt, Wolfgang E. Nagel and Joachim Prumerr, *SVM Fortran*, Technical Report IB-9322, KFA Jlich, Germany, 1993.
www.fz-juelich.de/zam/docs/printable/ib/ib-93/ib-9322.ps
- [3] Frank Cappello and Daniel Etiemble, *MPI versus MPI+OpenMP on the IBM SP for the NAS benchmarks*, in Proc. Supercomputing'00, Dallas, TX, 2000. <http://citeseer.nj.nec.com/cappello00mpi.html>
www.sc2000.org/techpaper/papers/pap.pap214.pdf
- [4] William W. Carlson, Jesse M. Draper, David E. Culler, Kathy Yelick, Eugene Brooks, and Karen Warren, *Introduction to UPC and Language Specification*, CCS-TR-99-157, May 13, 1999, <http://www.super.org/upc/>, www.gwu.edu and <http://projects.seas.gwu.edu/~hpc1/upcdev/upctr.pdf>.
- [5] Robert B. Ciotti, James R. Taft, and Jens Petersohn, *Early Experiences with the 512 Processor Single System Image Origin2000*, proceedings of the 42nd International Cray User Group Conference, SUMMIT 2000, Noordwijk, The Netherlands, May 22–26, 2000, www.cug.org.
- [6] The Earth Simulator. www.es.jamstec.go.jp
- [7] Tarek El-Ghazawi and Sébastien Chauvin, *UPC Benchmarking Issues*, proceedings of the International Conference on Parallel Processing, 2001, pp 365–372.
<http://projects.seas.gwu.edu/~hpc1/upcdev/UPCbench.pdf>
- [8] William Gropp, Ewing Lusk, Nathan Doss, and Anthony Skjellum, *A high-performance, portable implementation of the MPI message passing interface standard*, in Parallel Computing 22–6, Sep. 1996, pp 789–828.
<http://citeseer.nj.nec.com/gropp96highperformance.html>
- [9] Shinichi Habataa, Mitsuo Yokokawa, and Shigemune Kitawaki, *The Earth Simulator System*, in NEC Research & Development, Vol. 44, No. 1, Jan. 2003, Special Issue on High Performance Computing.
- [10] Georg Hager, Frank Deserno, and Gerhard Wellein, *Pseudo-Vectorization and RISC Optimization Techniques for the Hitachi SR8000 Architecture*, in High Performance Computing in Science and Engineering in Munich '02, Springer-Verlag Berlin Heidelberg, 2003.
- [11] Jonathan Harris, *Extending OpenMP for NUMA Architectures*, in proceedings of the Second European Workshop on OpenMP, EWOMP 2000. www.epcc.ed.ac.uk/ewomp2000/proceedings.html
- [12] D. S. Henty, *Performance of hybrid message-passing and shared-memory parallelism for discrete element modeling*, in Proc. Supercomputing'00, Dallas, TX, 2000.
<http://citeseer.nj.nec.com/henty00performance.html>
www.sc2000.org/techpaper/papers/pap.pap154.pdf
- [13] Matthias Hess, Gabriele Jost, Matthias Müller, and Roland Rühle, *Experiences using OpenMP based on Compiler Directed Software DSM on a PC Cluster*, in WOMPAT2002: Workshop on OpenMP Applications and Tools, Arctic Region Supercomputing Center, University of Alaska, Fairbanks, Aug. 5–7, 2002. <http://www.hlrs.de/people/mueller/papers/wompat2002/wompat2002.pdf>
- [14] Georg Karypis and Vipin Kumar. *A parallel algorithm for multilevel graph partitioning and sparse matrix ordering*, Journal of Parallel and Distributed Computing, 48(1): 71–95, 1998.
<http://citeseer.nj.nec.com/karypis98parallel.html>
<http://www-users.cs.umn.edu/~karypis/metis/>
- [15] Richard D. Loft, Stephen J. Thomas, and John M. Dennis, *Terascale spectral element dynamical core for atmospheric general circulation models*, in proceedings, SC 2001, Nov. 2001, Denver, USA.
www.sc2001.org/papers/pap.pap189.pdf
- [16] John Merlin, *Distributed OpenMP: Extensions to OpenMP for SMP Clusters*, in proceedings of the Second European Workshop on OpenMP, EWOMP 2000. www.epcc.ed.ac.uk/ewomp2000/proceedings.html

- [17] Message Passing Interface Forum. *MPI: A Message-Passing Interface Standard*, Rel. 1.1, June 1995, www.mpi-forum.org.
- [18] Message Passing Interface Forum. *MPI-2: Extensions to the Message-Passing Interface*, July 1997, www.mpi-forum.org.
- [19] R. W. Numrich, and J. K. Reid, *Co-Array Fortran for Parallel Programming*, ACM Fortran Forum, volume 17, no 2, 1998, pp 1–31, www.co-array.org, <http://citeseer.nj.nec.com/numrich98coarray.html> <ftp://matisa.cc.rl.ac.uk/pub/reports/nrRAL98060.ps.gz>
- [20] OpenMP Group, www.openmp.org.
- [21] Rolf Rabenseifner and Gerhard Wellein, *Communication and Optimization Aspects of Parallel Programming Models on Hybrid Architectures*, International Journal of High Performance Computing Applications, Sage Science Press, Vol. 17, No. 1, 2003, pp 49–62.
- [22] Mitsuhisa Sato, Shigehisa Satoh, Kazuhiro Kusano, and Yoshio Tanaka, *Design of OpenMP Compiler for an SMP Cluster*, in proceedings of the 1st European Workshop on OpenMP (EWOMP'99), Lund, Sweden, Sep. 1999, pp 32–39. <http://citeseer.nj.nec.com/sato99design.html>
- [23] Tetsuya Sato, *The Earth Simulator*, SC2002, Baltimore, Nov. 16-22, 2002. www.sc2002.org
- [24] Alex Scherer, Honghui Lu, Thomas Gross, and Willy Zwaenepoel, *Transparent Adaptive Parallelism on NOWs using OpenMP*, in proceedings of the Seventh Conference on Principles and Practice of Parallel Programming (PPoPP '99), May 1999, pp 96–106.
- [25] Weisong Shi, Weiwu Hu, and Zhimin Tang, *Shared Virtual Memory: A Survey*, Technical report No. 980005, Center for High Performance Computing, Institute of Computing Technology, Chinese Academy of Sciences, 1998, www.ict.ac.cn/chpc/dsm/tr980005.ps.
- [26] Satoru Shingu, Hiroshi Takahara, Hiromitsu Fuchigami, Masayuki Yamada, Yoshinori Tsuda, Wataru Ohfuchi, Yuji Sasaki, Kazuo Kobayashi, Takashi Hagiwara, Shin-ichi Habata, Mitsuo Yokokawa, Hiroyuki Itoh and Kiyoshi Otsuka, *A 26.58 Tflops Global Atmospheric Simulation with the Spectral Transform Method on the Earth Simulator*, SC2002, Baltimore, Nov. 16-22, 2002. www.sc2002.org/paperpdfs/pap.pap331.pdf
- [27] Lorna Smith and Mark Bull, *Development of Mixed Mode MPI / OpenMP Applications*, in proceedings of Workshop on OpenMP Applications and Tools (WOMPAT 2000), San Diego, July 2000. www.cs.uh.edu/wompat2000/
- [28] Gerhard Wellein, Georg Hager, Achim Basermann, and Holger Fehske, *Fast sparse matrix-vector multiplication for TeraFlop/s computers*, in proceedings of VECPAR'2002, 5th Int'l Conference on High Performance Computing and Computational Science, Porto, Portugal, June 26–28, 2002, part I, pp 57–70. <http://vecpar.fe.up.pt/>