



Outline

- Introduction / Motivation
- Programming models on clusters of SMP nodes
- Case Studies / pure MPI vs. hybrid MPI+OpenMP
- Mismatch Problems
- Thread-safety quality of MPI libraries
- Case Studies / pure OpenMP
- Summary



Major Programming models on hybrid systems

- Pure MPI (one MPI process on each CPU)
- Hybrid MPI+OpenMP
 - shared memory OpenMP
 distributed memory MPI
- OpenMP inside of the SMP nodes
 SMP nodes

 MPI between the nodes via node interconnect
 Node Interconnect
- Other: Virtual shared memory systems, HPF, ...
- Often hybrid programming (MPI+OpenMP) slower than pure MPI

 why?















- Performance using Gigabit Ethernet (GE):
- Hybrid implementations outperform pure MPI implementation
- BT Hybrid V1: shows best scalability
- BT Hybrid V1: best performance employing 16 MPI processes and 4 or 5 threads respectively:
- BT Speed-up on Sun Fire E6800 (GE)
- Tight interaction between MPI and OpenMP limits number of threads that can be used efficiently
- BT Hybrid V2 achieves best performance using 4 MPI processes employing 16 threads each:
 - Large messages saturate slow network and limit number of MPI processes that can be used efficiently



Characteristics of Hybrid Codes

• BT Hybrid V1:

Hybrid Parallel Programming

Slide 28 / 122

- Message exchange with 2 neighbour processes
- Many short messages
- Length of messages remain the same
- Increase number of threads:
 - Increase of OpenMP barrier time (threads from different MPI processes have to synchronize)
 - Increase of MPI time (MPI calls within parallel regions are serialized)

Rabenseifner, Hager, Jost, Keller

• BT Hybrid V2:

- Message exchange with 6 neighbour processes
- Few long messages
- Length of messages decreases with increasing number of processes
- Increase number of threads:

 Increase of OpenMP barrier time

Observation on fast networks:

- Single Level MPI:
 - Best performance, best scalability
 - Coarse-grained well balanced distribution and scheduling of work
- Hybrid MPI/OpenMP V2 did not yield performance advantage
- Hybrid MPI/OpenMP V1:
 - Implementation non-typical: pipelined thread execution, communication within parallel regions.
 - Low percentage of useful thread work time:
 - 1D data distribution limits parallelism on coarse grain
 - OpenMP introduces extra synchronization overhead at the end of parallel regions
 - Interaction of OpenMP and MPI yields thread pre-emption and thread migration
 - Performance improves through explicit binding.



Observation on slow networks:

- Hybrid MPI/OpenMP V1 showed better performance than V2 or pure MPI:
 - Message exchange with only 2 neighbours vs 6 neighbours
 - Many short messages vs few longer messages:
 - BT V1 4x16: 14880 send, avg. length 10600 bytes
 - BT V2 4x16: 960 send, avg. length 116360 bytes
 - Long messages sent by many MPI processes potentially saturate a slow network quickly.



The Multi-zone NAS Parallel Benchmarks Nested set up zones MPI/OpenMP MLP OpenMP Time step sequential sequential sequential initialize MPI MLP zones OpenMP inter-zones Processes Processes data copy+ exchange Call MPI OpenMP boundaries sync. exchange boundaries intra-zones OpenMP OpenMP OpenMP timestep zones Multi-zone versions of the NAS Parallel LU/SP/BT Benchmarks LU,SP, and BT • Two hybrid sample implementations verify • Load balance heuristics part of sample codes • Nested OpenMP based on NanosCompiler extensions was developed for this study Hybrid Parallel Programming H L R S Rabenseifner, Hager, Jost, Keller Slide 31 / 122

Using MPI/OpenMP call omp set numthreads (weight) subroutine ssor(u, rsd, ...) do step = 1, itmax call exch_qbc(u, qbc, nx,...) !\$OMP PARALLEL DEFAUL(SHARED) !\$OMP& PRIVATE(m, i, j, k...) do k = 2, nz-1 !\$OMP DO call mpi send/recv do j = 2, ny-1 do i = 2, nx-1do zone = 1, num zones do m = 1, 5if (iam .eq. pzone id(zone)) then rsd(m, i, j, k) =dt*rsd(m,i,j,k-1) call ssor(u,rsd,...) end do end if end do end do end do !\$OMP END DO nowait end do end do !\$OMP END PARALLEL Hybrid Parallel Programming н L R Rabenseifner, Hager, Jost, Keller Slide 32 / 122





Benchmark Characteristics

Aggregate sizes:

- Class W: aggregate 64x64x8 grid points
- Class A: aggregate 128x128x16 grid points
- Class B: aggregate 304x208x17 grid points
- BT-MZ: •
 - #Zones: 16 (Class W), 16 (Class A), 64 (Class B)
 - Size of the zones varies widely:
 - large/small ≈ 20
 - · requires multi-level parallelism to achieve a good load-balance
- LU-MZ:
 - #Zones: 16 (Class W), 16 (Class A), 16 (Class B)
 - Size of the zones identical:
 - no load-balancing required
 - limited parallelism on outer level
- SP-MZ:
 - #Zones: 16 (Class W), 16 (Class A), 64 (Class B)
 - Size of zones identical



Performance of BT-MZ on SGI Origin 3000

Rabenseifner, Hager, Jost, Keller

Slide 36 / 122





















OpenMP only

Intel[®] Compilers with Cluster OpenMP – Consistency Protocol

Basic idea:

- Between OpenMP barriers, data exchange is not necessary, i.e., visibility of data modifications to other threads only after synchronization.
- When a page of sharable memory is not up-to-date, it becomes *protected*.
- Any access then faults (SIGSEGV) into Cluster OpenMP runtime library, which requests info from remote nodes and updates the page.
- Protection is removed from page.
- Instruction causing the fault is re-started, this time successfully accessing the data.



Real consistency protocol is more complicated

- · Diffs are done only when requested
- Several diffs are locally stored and transferred later if a thread first reads a page after several barriers.
- Each write is internally handled as a read followed by a write.
- If too many diffs are stored, a node can force a "reposession" operation, i.e., the page is marked as invalid and fully re-send if needed.
- Another key point:
 - After a page has been made read/write in a process, no more protocol traffic is generated by the process for that page until after the next synchronization (and similarly if only reads are done once the page is present for read).
 - This is key because it's how the large cost of the protocol is averaged over many accesses.

RS

Courtesy of J. Cownie, Intel

ΗL

- I.e., protocol overhead only "once" per barrier
- Examples in the Appendix





Comparison: MPI based parallelization $\leftrightarrow \rightarrow$ DSM

- MPI based:
 - Potential of boundary exchange between two domains in one large message
 - → Dominated by **bandwidth** of the network
- DSM based (e.g. Intel[®] Cluster OpenMP):
 - Additional latency based overhead in each barrier
 - → May be marginal
 - Communication of updated data of pages
 - → Not all of this data may be needed
 - \rightarrow i.e., too much data is transferred
 - → Packages may be to small
 - → Significant latency
 - Communication not oriented on boundaries of a domain decomposition
 - → probably more data must be transferred than necessary

 Hybrid Parallel Programming

 Slide 76 / 122
 Rabenseifner, Hager, Jost, Keller



hybrid MPI+OpenMP

 OpenMP only

by rule of thumb:

Communication

may be

10 times slower

than with MPI





Mismatch Problems

- Topology problem
- **Unnecessary intra-node communication** [with pure MPI] •
- Inter-node bandwidth problem
- Sleeping threads and saturation problem

[with hybrid MPI+OpenMP]

RS

L

н

[with pure MPI]

- [with hybrid MPI+OpenMP]
 - [with masteronly] [with pure MPI]
- Additional OpenMP overhead
 - Thread startup / join
 - Cache flush (data source thread communicating thread sync. \rightarrow flush)
- **Overlapping communication and computation** [with hybrid MPI+OpenMP]
 - an application problem \rightarrow separation of local or halo-based code
 - a programming problem \rightarrow thread-ranks-based vs. OpenMP work-sharing

- a load balancing problem, if only some threads communicate / compute
- Communication overhead with DSM [with pure (Cluster) OpenMP]

no silver bullet, i.e., each parallelization scheme has its problems

Hybrid Parallel Programming Rabenseifner, Hager, Jost, Keller Slide 79 / 122

No silver bullet

- The analyzed programming models do not fit on hybrid architectures
 - whether drawbacks are minor or major
 - depends on applications' needs
 - problems ...

Hybrid Parallel Programming

Slide 80 / 122

- to utilize the CPUs the whole time
- > to achieve the full inter-node network bandwidth

H L R S

- to minimize inter-node messages
- > to prohibit intra-node
- message transfer,
- synchronization and
- balancing (idle-time) overhead

Rabenseifner, Hager, Jost, Keller

> with the programming effort

Chances for optimization

- with hybrid masteronly (MPI only outside of parallel OpenMP regions), e.g.,
 - > Minimize work of MPI routines, e.g.,
 - application can copy non-contiguous data into contiguous scratch arrays (instead of using derived datatypes)
 - MPI communication parallelized with multiple threads to saturate the inter-node network
 - by internal parallel regions inside of the MPI library
 - by the user application
 - > Use only hardware that can saturate inter-node network with 1 thread
 - > Optimal throughput:
 - reuse of idling CPUs by other applications
- On constellations:
 - Hybrid Masteronly

with several MPI multi-threaded processes on each SMP node

Hybrid Parallel Programming Slide 81 / 122 Rabenseifner, Hager, Jost, Keller



Summary of mismatch problems

Performance and Programming Problems with	Pure MPI	Master- only 1 process per node	Master- only several processes per node	Over- lapping 1 process per node	Over- lapping several processes per node	Pure OpenMP: e.g., Intel Cluster OpenMP
Application topology problem (neighbor domains inside of SMP node)	4		4		4	4
Additional MPI communication inside of SMP nodes	4		4		4	
Do we achieve full inter-node bandwidth on constellations?		444		4		444
Sleeping CPUs while MPI communication	(4)	44	4			4
Additional OpenMP overhead		4	4	4	4	
Separation of (a) halo data and (b) inner data based calculations				44	44	
OpenMP work sharing only partially usable				44	44	
Load balancing problem due to hybrid programming model				4	4	

Outline



- Programming models on clusters of SMP nodes
- Case Studies / pure MPI vs. hybrid MPI+OpenMP
- Mismatch Problems

Thread-safety quality of MPI libraries

Rainer Keller, High Performance Computing Center Stuttgart (HLRS)

- Case Studies / pure OpenMP
- Summary



Thread-safety of MPI Libraries

- Make most powerful usage of hierarchical structure of hardware:
- Efficient programming of clusters of SMP nodes
 SMP nodes:

Rabenseifner, Hager, Jost, Keller

Dual/multi core CPUs

Hybrid Parallel Programming

Slide 84 / 122

- Multi CPU shared memory
- Multi CPU ccNUMA
- · Any mixture with shared memory programming model



- No restriction to the usage of OpenMP for intranode-parallelism:
 - OpenMP does not (yet) offer binding threads to processors
 - OpenMP does not guarantee thread-ids to stay fixed.
- OpenMP is based on the implementation dependant thread-library: LinuxThreads, NPTL, SolarisThreads.

H L R S



Standardized by ISO as ISO/IEC 9945-1:1996.

Rabenseifner, Hager, Jost, Keller

1 H L

RS

Hybrid Parallel Programming

Slide 88 / 122

/* omp end parallel */

SOMP END PARALLEL

Rabenseifner, Hager, Jost, Keller

Hybrid Parallel Programming

Slide 87 / 122

RS





Duplicated

Odd-/Even Split MPI COMM WORLD

Zero-and-Rest

Intercommunicator

Two-dimensional

Cartesian

Hybrid Parallel Programming

Slide 95 / 122

MPI COMM WORLD

Halved

Halved

MPI COMM WORLD

Intercommunicators

RS

Fully-connected

Topology

HL

Reversed

Cartesian

Rabenseifner, Hager, Jost, Keller

MPI COMM WORLD

Intracomm merged of

Halved Intercomms

• Threaded running of already implemented tests:







Examples of failed multi-threaded tests

Standard send in comm. "Reversed MPI_COMM_WORLD": P2P tests Ring, comm Reversed MPI_COMM_WORLD, type MPI_INT mpi_test_suite: ./../../../../ompi/mca/pml/ob1/pml_ob1_sendreq.c:196: mca_pml_ob1_match_completion_free: Assertion `0 == sendreq->req_send.req_base.req_pml_complete' failed. 2-threads Collective (Bcast, Bcast) on different comms wrong data: mpirun -np 4 ./mpi_test_suite -r FULL -j 2 -t "Bcast" -c "MPI_COMM_WORLD, Duplicated MPI_COMM_WORLD" 2-threads Collective (Bcast, Gather) on different comms hangs: mpirun -np 4 ./mpi_test_suite -r FULL -j 2 -t "Bcast, Gather" -c "MPI_COMM_WORLD, Duplicated MPI_COMM_WORLD" Of course a test-suite may contain errors as well -1

📊 🗖 🖬 H L R S 🕷

Of course, a test-suite may contain errors as well ,-]



Thread support in MPI libraries

· The following MPI libraries offer thread support:

Implemenation	Thread support level
MPlch-1.2.7p1	Always announces MPI_THREAD_FUNNELED.
MPIch2-1.0.4	ch:sock3 (default) supports MPI_THREAD_MULTIPLE
Intel MPI 2.0	MPI_THREAD_FUNNELED
Intel MPI 3.0	MPI_THREAD_SERIALIZED
SGI MPT-1.14	Not thread-safe?
IBM MPI	Full MPI_THREAD_MULTIPLE
Nec MPI/SX	Full MPI_THREAD_MULTIPLE

- Examples of failures in MPI libraries uncovered are shown.
- Failure logs are shown only for Open MPI.



Thread support within Open MPI

Hybrid Parallel Programming Slide 100 / 122 Rabenseifner, Hager, Jost, Keller

• In order to enable thread support in Open MPI, configure with:

configure --enable-mpi-threads --enable-progress-threads

- This turns on:
 - Support for threaded initialization functions and internal checks to enable locking when run with threads
 - Progress threads to asynchronously transfer/receive data per network BTL.
 - However, some BTLs (mvapi, openib, mx) are still marked non-thread-safe.

HLR S

Outline



- Programming models on clusters of SMP nodes
- Case Studies / pure MPI vs. hybrid MPI+OpenMP
- Mismatch Problems
- Thread-safety quality of MPI libraries

Case Studies / pure OpenMP

- First Experiences with Intel[®] Cluster OpenMP (CLOMP)
 Georg Hager, Regionales Rechenzentrum Erlangen (RRZE)
- Summary

Hybrid Parallel Programming Slide 101 / 122 Rabenseifner, Hager, Jost, Keller HLR S

General Remarks on Intel® Cluster OpenMP (CLOMP)

- CLOMP == "extreme" ccNUMA
 - very long latencies, expensive non-local access
 - page replications can lead to memory problems
 - but: placement is handled "automatically"
- Consequence: A well-optimized, ccNUMA-aware OMP code that scales well on Altix does not necessarily scale well with CLOMP
 - example: boundary code must be optimized for local access
- · Good stability on all systems with latest CLOMP release
- No problems and good performance with IP over IB
 - native IB not working yet (but check latest CLOMP versions!)



Overview

- Cluster OpenMP is part of every 9.1 Intel compiler
 - separate license must be purchased
- Systems used
 - EM64T (dual Nocona) with Gbit Ethernet and Infiniband, Debian 3.1 (Sarge)
 - Itanium2 (HP zx6000) with Gbit Ethernet, SLES9pl3
 - AMD Opteron is supported with latest CLOMP compiler versions
- · Basic numbers: Triad tests on Nocona nodes
- Application: Lattice-Boltzmann code
 - influence of algorithmic details (locality of access, page sharing)
 - data layout considerations
- Odds and ends



General Remarks

- · Problems
 - memory footprint is about 2.5 times larger than expected from serial code (270MB instead of 61MB for vector triad)
 - Partially resolved by Intel (Jim C.)
 - Problem is specific to RRZE kernel and system libs
 - huge core dumps even with small sharable heap and resident memory (2.4GB core with 200MB code)
 - Problem is specific to RRZE kernel and system libs









Filled vs. Half-filled nodes

- 2 ways to "fill the node"
 - 1. Keep unique names in hostfile and use 2 "real" OpenMP threads per node with --process_threads=2
 - 2. Duplicate names in hostfile and use --process_threads=1
- · Observations
 - $-\,$ breakdown of performance compared to the half-filled case for large N $\,$
 - Improvement with OpenMP for medium-sized arrays
 - --process_threads=2: quite erratic performance data
- Breakdown was actually expected (the same happens on single node with pure OpenMP)
- Erratic behaviour
 - influence of "loaded" switch? (improbable)
 - Threads losing CPU affinity?









Conclusions on CLOMP



Outline

- Introduction / Motivation
- Programming models on clusters of SMP nodes
- Case Studies / pure MPI vs. hybrid MPI+OpenMP
- Mismatch Problems
- Thread-safety quality of MPI libraries
- Case Studies / pure OpenMP

 Hybrid Parallel Programming

 Slide 118 / 122
 Rabenseifner, Hager, Jost, Keller

Summary

Slide 120 / 122

Acknowledgements

- I want to thank
 - Gerhard Wellein, RRZE
 - Monika Wierse, Wilfried Oed, and Tom Goozen, CRAY
 - Holger Berger, NEC
 - Reiner Vogelsang, SGI
 - Gabriele Jost, NASA
 - Dieter an Mey, RZ Aachen
 - Horst Simon, NERSC
 - Matthias Müller, HLRS
 - my colleges at HLRS



erformance and Programming roblems with	Pure MPI	Master- only 1 process per node	Master- only several processes per node	Over- lapping 1 process per node	Over- lapping several processes per node	Pure OpenMP: e.g., Intel Cluster OpenMP
Application topology problem neighbor domains inside of SMP node)	4		4		4	4
Additional MPI communication insid	4		4		4	
Do we achieve full inter-node bandwidth on constellations?		likh		K		444
Sleeping CPUs while MPI communication	(4)	(1/4/4)	4			4
Additional OpenMP overhead		4	4	4	4	
Separation of (a) halo data and b) inner data based calculations				44	44	
DpenMP work sharing only partially usable				44	44	
oad balancing problem due to				4	4	

Rabenseifner, Hager, Jost, Keller

к

5



Abstract

Half-Day Tutorial (Level: 25% Introductory, 50% Intermediate, 25% Advanced)

Rolf Rabenseifner, HLRS, Germany Gabriele Jost, Sun Microsystems, Germany Rainer Keller, HLRS, Germany

Georg Hager, University of Erlangen-Nuremberg, Germany Bainer Keller, HLRS, Germany

Abstract. Most HPC systems are clusters of shared memory nodes. Such systems can be PC clusters with dual or quad boards, but also "constellation" type systems with large SMP nodes. Parallel programming must combine the distributed memory parallelization on the node inter-connect with the shared memory parallelization inside of each node.

This tutorial analyzes the strength and weakness of several parallel programming models on clusters of SMP nodes. Various hybrid MPI-OpenMP programming models are compared with pure MPI. Benchmark results of several platforms are presented. A hybrid-masteronly programming model can be used more efficiently on some vector-type systems, but also on clusters of dual-CPUs. On other systems, one CPU is not able to saturate the inter-node network and the commonly used masteronly programming model suffers from insufficient inter-node bandwidth. The thread-safety quality of several existing MPI libraries is also discussed. Case studies from the fields of CFD (NAS Parallel Benchmarks, in detail), Climate Modelling (POP2, maybe) and Particle Simulation (GTC, maybe) will be provided to demonstrate various aspect of hybrid MPI/OpenMP programming.

Another option is the use of distributed virtual shared-memory technologies which enable the utilization of "near-standard" OpenMP on distributed memory architectures. The performance issues of this approach and its impact on existing applications are discussed. This tutorial analyzes strategies to overcome typical drawbacks of easily usable programming schemes on clusters of SMP nodes.



Intel® Compilers with Cluster OpenMP – Consistency Protocol – Examples

Notation

- ..=A[i] Start/End Start/end a read on element i on page A
- A[i]=.. Start/End Start/end a write on element i on page A, trap to library
- Twin(A) Create a twin copy of page A
- WriteNotice(A) Send write notice for page A to other processors
- DiffReq_A_n(s:f) Request diffs for page A from node n between s and f
- Diff_A_n(s:f) Generate a diff for page A in writer n between s and where s and f are barrier times. This also frees the twin for page A.



Node 0	Node 1	
Barrier 0	Barrier 0	
A[1]= Start		
Twin(A)		
A[2]= End		
	A[5]= Start	
	Twin(A)	
	A[5]= End	
Barrier 1	Barrier 1	
WriteNotice(A)	Writenotice(A)	
A[5]= Start		
Diffreq_A_1(0:1)->		
•	<-Diff_A_1(0:1)	
Apply diffs		
A[5]= End		
Barrier 2	Barrier 2	
WriteNotice(A)		

Hybrid Parallel Programming Slide 127 <u>Rabenseifner,</u> Hager, Jost, Keller



Courtesy of J. Cownie, Intel

Node 0	Node 1	Node 2
Barrier 0	Barrier 0	Barrier 0
A[1]= Start		
Twin(A)		
A[1]= End		
Barrier 1	Barrier 1	Barrier 1
WriteNotice(A)		
A[2]= (no trap to library)		
Barrier 2	Barrier 2	Barrier 2
(No WriteNotice(A) required)		
A[3]= (no trap to lib)		
	=A[1] Start	
-	<-Diffreq_A_0(0:2)	
Diff_A_0(0:2)->		
	Apply diffs	
	=A[1] End	
Barrier 3	Barrier 3	Barrier 3
(no WriteNotice(A) required because diffs		
were sent after the A[3]=)		
A[1]= Start		
Twin(A)		
Barrier 4	Barrier 4	Barrier 4
WriteNotice(A)		
		=A[1] Start
• •		Diffreq_A_0(0:4)
Create Diff_A_0(2:4) send Diff_A_0(0:4)->		
		Apply diffs
-		=A[1] End
		0

Exa. 3 (start)

Node 0	Node 1	Node 2	Node 3
Barrier 0	Barrier 0	Barrier 0	Barrier 0
A[1]= Start	A[5]= Start		
Twin(A)	Twin(A)		
A[1]= End	A[5]= End		
Barrier 1	Barrier 1	Barrier 1	Barrier 1
WriteNotice(A)	WriteNotice(A)		
A[2]= Start	A[1]= Start		
Diffreq_A_1(0:1)->	<-Diffreq_A_0(0:1)		
Diff_A_0(0:1)->	<-Diff_A_1_(0:1)		
Apply diff	Apply diff		
Twin(A)	Twin(A)		
A[2]= End	A[1]= End		
Barrier 2	Barrier 2	Barrier 2	Barrier 2
WriteNotice(A)	WriteNotice(A)		
A[3]= Start	A[6]= Start		
Diffreq_A_1(1:2)->	<-Diffreq_A_A(1:2)		
Diffs_A_0(1:2)	<-Diffs_A_1(1:2)		
Apply diffs	Apply diffs		
Twin(A)	Twin(A)		
A[3]= End	A[6]= End		
		=A[1] Start	
		-Diffreq_A_0(0:2)	
		<-Diffreq_A_1(0:2)	
Create Diff_A_0(1:2)	Create Diff_A_1(1:2)		
Send Diff_A_0(0:2)->	Send Diff_A_1(0:2)->		
		Apply all diffs	
		=A[1] End	

Courtesy of J. Cownie, Intel

Rolf	Raben	seifner



Dr. Rolf Rabenseifner studied mathematics and physics at the University of Stuttgart. Since 1984, he has worked at the High-Performance Computing-Center Stuttgart (HLRS). He led the projects DFN-RPC, a remote procedure call tool, and MPI-GLUE, the first metacomputing 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 Parallel Computing - Training and Application Services at HLRS. He is involved in MPI profiling and benchmarking, e.g., in the HPC Challenge Benchmark Suite. In recent projects, he studied parallel I/O, parallel programming models for clusters of SMP nodes, and optimization of MPI collective routines. In workshops and summer schools, he teaches parallel programming models in many universities and labs in Germany.



Barrier 3	Barrier 3	Barrier 3
Writenotice(A)		
-Diffs_A_1_(2:3)		
Barrier 4	Barrier 4	Barrier 4
		=A[1] Start
		<-Diffreq_A_0(0:4
		<-Diffreq_A_1(0:4
Create Diff_A_1(2:4)		
Send Diff_A_1(0:4)->		
		Apply diffs
		=A[1] End
	1	·
	Writenotice(A) Create Diff_A_1(2:3) Create Diff_A_1(2:4) Send Diff_A_1(0:4)->	Writenotice(A) Create Diff A_1(2:4) Create Diff A_1(0:4)->

These examples may give an impression of the overhead induced by the Cluster OpenMP consistency protocol.



Georg Hager



Dr. Georg Hager studied theoretical physics at the University of Bayreuth, specializing in nonlinear dynamics. Since 2000 he is a member of the HPC Services group at the Regional Computing Center Erlangen (RRZE), which is part of the University of Erlangen-Nürnberg. His daily work encompasses all aspects of user support in High Performance Computing like tutorials and training, code parallelization, profiling and optimization and the assessment of novel computer architectures and tools.

In his dissertation he developed a shared-memory parallel density-matrix renormalization group algorithm for ground-state calculations in strongly correlated electron systems. Recent work includes architecture-specific optimization strategies for current microprocessors and special topics in shared memory programming.

HLRS





Gabriele Jost



Gabriele Jost obtained her doctorate in Applied Mathematics from the University of Göttingen, Germany. For more than a decade she worked for various vendors (Suprenum GmbH, Thinking Machines Corporation, and NEC) of high performance parallel computers in the areas of vectorization, parallelization, performance analysis and optimization of scientific and engineering applications.

In 1998 she joined the NASA Ames Research Center in Moffett Field, California, USA as a Research Scientist. Here her work focused on evaluating and enhancing tools for parallel program development and investigating the usefulness of different parallel programming paradigms. In 2005 she moved from California to the Pacific Northwest and joined Sun Microsystems as a staff engineer in the Compiler Performance Engineering team. Her task is the analysis of compiler generated code and providing feedback and suggestions for improvement to the compiler group. Her research interest remains in area of performance analysis and evaluation of programming paradigms for high performance computing.



Rainer Keller



Rainer Keller is a scientific employee at the High Performance Computing Center Stuttgart (HLRS) since 2001. He earned his diploma in Computer Science at the University of Stuttgart. Currently, he is the head of the group Applications, Models and Tools at the HLRS.

His professional interest are Parallel Computation using and working on MPI with Open MPI and shared memory parallelization with OpenMP, as well as distributed computing using the Meta-Computing Library PACX-MPI.

His work includes performance analysis and optimization of parallel applications, as well as the assessment of and porting to new hardware technologies, including the training of HLRS users in parallel application development. He is involved in several European projects, such as HPC-Europa.



References (with direct relation to the content of this tutorial)

- NAS Parallel Benchmarks: http://www.nas.nasa.gov/Resources/Software/npb.html
- R.v.d. Wijngaart and H. Jin, NAS Parallel Benchmarks, Multi-Zone Versions, NAS Technical Report NAS-03-010, 2003
- H. Jin and R. v.d.Wijngaart, Performance Characteristics of the multi-zone NAS Parallel Benchmarks, Proceedings IPDPS 2004
- G. Jost, H. Jin, D. an Mey and F. Hatay, Comparing OpenMP, MPI, and Hybrid Programming, Proc. Of the 5th European Workshop on OpenMP, 2003
- E. Ayguade, M. Gonzalez, X. Martorell, and G. Jost, Employing Nested OpenMP for the Parallelization of Multi-Zone CFD Applications, Proc. Of IPDPS 2004



References

Rolf Rabenseifner, Hybrid Parallel Programming on HPC Platforms. In proceedings of the Fifth European Workshop on OpenMP, EWOMP '03, Aachen, Germany, Sept. 22-26, 2003, pp 185-194, www.compunity.org.

- Rolf Rabenseifner, Comparison of Parallel Programming Models on Clusters of SMP Nodes. In proceedings of the 45nd Cray User Group Conference, CUG SUMMIT 2003, May 12-16, Columbus, Ohio, USA.
- Rolf Rabenseifner and Gerhard Wellein, Comparison of Parallel Programming Models on Clusters of SMP Nodes. In Modelling, Simulation and Optimization of Complex Processes (Proceedings of the International Conference on High Performance Scientific Computing, March 10-14, 2003, Hanoi, Vietnam) Bock, H.G.; Kostina, E.; Phu, H.X.; Rannacher, R. (Eds.), pp 409-426, Springer, 2004.
- Rolf Rabenseifner and Gerhard Wellein, Communication and Optimization Aspects of Parallel Programming Models on Hybrid Architectures.

In the International Journal of High Performance Computing Applications, Vol. 17, No. 1, 2003, pp 49-62. Sage Science Press.





References

· Rolf Rabenseifner,

Communication and Optimization Aspects on Hybrid Architectures. In Recent Advances in Parallel Virtual Machine and Message Passing Interface, J. Dongarra and D. Kranzlmüller (Eds.), Proceedings of the 9th European PVM/MPI Users' Group Meeting, EuroPVM/MPI 2002, Sep. 29 - Oct. 2, Linz, Austria, LNCS, 2474, pp 410-420, Springer, 2002.

 Rolf Rabenseifner and Gerhard Wellein, Communication and Optimization Aspects of Parallel Programming Models on Hybrid Architectures.
 If the second sec

In proceedings of the Fourth European Workshop on OpenMP (EWOMP 2002), Roma, Italy, Sep. 18-20th, 2002.

Rolf Rabenseifner,
 Communication Bandwidth of Parallel Programming Models on Hybrid

Architectures. Proceedings of WOMPEI 2002, International Workshop on OpenMP: Experiences and Implementations, part of ISHPC-IV, International Symposium on High Performance Computing, May, 15-17., 2002, Kansai Science City, Japan, LNCS



References

Barbara Chapman et al.:

Toward Enhancing OpenMP's Work-Sharing Directives. In proceedings, W.E. Nagel et al. (Eds.): Euro-Par 2006, LNCS 4128, pp. 645-654, 2006.



Further references

- Sergio Briguglio, Beniamino Di Martino, Giuliana Fogaccia and Gregorio Vlad, Hierarchical MPI+OpenMP implementation of parallel PIC applications on clusters of Symmetric MultiProcessors, 10th European PVM/MPI Users' Group Conference (EuroPVM/MPI'03), Venice, Italy, 29 Sep - 2 Oct, 2003
 Barbara Chapman, Parallel Application Development with the Hybrid MPI+OpenMP Programming Model, Tutorial, 9th EuroPVM/MPI & 4th DAPSYS Conference, Johannes Kepler University Linz, Austria September 29-October 02, 2002
 Luis F. Romero, Eva M. Ortigosa, Sergio Romero, Emilio L. Zapata, Nesting OpenMP and MPI in the Conjugate Gradient Method for Band Systems, 11th European PVM/MPI Users' Group Meeting in conjunction with DAPSYS'04, Budapest, Hungary, September 19-22, 2004
- Nikolaos Drosinos and Nectarios Koziris, Advanced Hybrid MPI/OpenMP Parallelization Paradigms for Nested Loop Algorithms onto Clusters of SMPs, 10th European PVM/MPI Users' Group Conference (EuroPVM/MPI'03), Venice, Italy, 29 Sep - 2 Oct. 2003



Further references

 Holger Brunst and Bernd Mohr, Performance Analysis of Large-scale OpenMP and Hybrid MPI/OpenMP Applications with VampirNG Proceedings for IWOMP 2005, Eugene, OR, June 2005. http://www.fz-juelich.de/zam/kojak/documentation/publications/
 Felix Wolf and Bernd Mohr, Automatic performance analysis of hybrid MPI/OpenMP applications Journal of Systems Architecture, Special Issue "Evolutions in parallel distributed and network-based processing", Volume 49, Issues 10-11, Pages 421-439, November 2003. http://www.fz-juelich.de/zam/kojak/documentation/publications/
 Felix Wolf and Bernd Mohr, Automatic Performance Analysis of Hybrid MPI/OpenMP Applications short version: Proceedings of the 11-th Euromicro Conference on Parallel, Divident deviced base of the second (PDP) Ocena on Parallel, Divident deviced base of the second (PDP) Ocena on Parallel, Divident deviced base of the second (PDP) Ocena on Parallel, Divident deviced base of the second (PDP) Ocena on Parallel, Divident deviced base of the paral (PDP) Ocena on Parallel, Divident deviced base of the paral (PDP) Ocena on Parallel, Divident deviced base of the paral (PDP) Ocena on Parallel, Divident deviced base of the paral (PDP) Ocena on Parallel, Divident deviced base of the paral (PDP) Ocena on Parallel, Divident deviced base of the paral (PDP) Ocena on Parallel, Divident deviced base of the paral (PDP) Ocena on Parallel, Parallel (PDP) Ocena on Parallel, Divident deviced base of the paral (PDP) Ocena on Parallel, Divident deviced base of the paral (PDP) Ocena on Parallel, Parallel (PDP) Ocena on Parallel (PDP) Ocena on Parallel, Parallel (PDP) Ocena on Parallel (P

Distributed and Network based Processing (PDP 2003), Genoa, Italy, February 2003. long version: Technical Report FZJ-ZAM-IB-2001-05.

http://www.fz-juelich.de/zam/kojak/documentation/publications/





Further references

- 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/techpapr/papers/pap.pap214.pdf
- 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
- 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/techpapr/papers/pap.pap154.pdf



Further references

- 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.
- 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/
- 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/

Further references

- Agnieszka Debudaj-Grabysz and Rolf Rabenseifner, Load Balanced Parallel Simulated Annealing on a Cluster of SMP Nodes. In proceedings, W. E. Nagel, W. V. Walter, and W. Lehner (Eds.): Euro-Par 2006, Parallel Processing, 12th International Euro-Par Conference, Aug. 29 - Sep. 1, Dresden, Germany, LNCS 4128, Springer, 2006.
- Agnieszka Debudaj-Grabysz and Rolf Rabenseifner, Nesting OpenMP in MPI to Implement a Hybrid Communication Method of Parallel Simulated Annealing on a Cluster of SMP Nodes. In Recent Advances in Parallel Virtual Machine and Message Passing Interface, Beniamino Di Martino, Dieter Kranzlmüller, and Jack Dongarra (Eds.), Proceedings of the 12th European PVM/MPI Users' Group Meeting, EuroPVM/MPI 2005, Sep. 18-21, Sorrento, Italy, LNCS 3666, pp 18-27, Springer, 2005

H L R S

Hybrid Parallel Programming

Slide 144

Rabenseifner, Hager, Jost, Keller





