

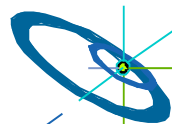
# Communication Bandwidth of Parallel Programming Models on Hybrid Architectures

Rolf Rabenseifner  
rabenseifner@hlrs.de

*This talk is given by Matthias Müller on the WOMPEI 2002*

University of Stuttgart  
High-Performance Computing-Center Stuttgart (HLRS)  
www.hlrs.de

WOMPEI 2002 at ISHPC-IV, Kansai Science City, Japan, May 15-17, 2002



Cluster Programming Models

Slide 1

Höchstleistungsrechenzentrum Stuttgart



## Motivation

- HPC systems
  - often clusters of SMP nodes
  - i.e., hybrid architectures
- Hybrid programming models, e.g.,
  - MPI & OpenMP
  - MPI & automatic loop parallelization
- Often hybrid programming slower than pure MPI programming
  - why?
  - bandwidth problems?



## Goals

- Using the communication bandwidth of the hardware (that I have bought)
- Appropriate parallel programming models
- Pros & Cons of existing programming models
- Work horses for legacy & new parallel applications
- Optimization of the middleware



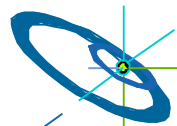
# Programming Large Scale Hybrid Systems

- Interconnect
  - ccNUMA
  - Remote Direct Memory Access (RDMA)
  - only message passing
- Programming models
  - OpenMP on ccNUMA
  - OpenMP cluster extensions
  - MLP (Multi Level Parallelism)
  - Co-Array Fortran & UPC
  - One-sided communication (MPI-2 & shmem)



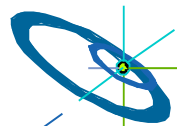
## On ccNUMA

- OpenMP on large partitions
  - single level of parallelism
  - Nested parallelism
    - **not yet implemented in most OpenMP compilers**
- OpenMP cluster extensions
  - first touch mechanism
  - data distribution extensions
  - **may optimize the data locality**
  - **may reduce communication on interconnect**
- Multi Level Parallelism (MLP) from NASA/Ames
  - a Fortran wrapper to System V shared memory (shm)
  - multi-threaded processes access variables in shared memory segments
  - cheap load balancing:  
changing the number of threads of each process



# RDMA (Remote Direct Memory Access)

- Co-Array Fortran / Unified Parallel C (UPC)
  - access to variables in other threads/processes is done via additional array subscript in brackets, addressing the thread/process by its rank
  - multi-dimensional ranking is possible
  - not tailored for clusters of SMP, but usable
  - without overhead for
    - additional message passing
    - additional temporary data copies
  - Lack of portable / public compiling system
  - Architecture may allow **separation** of optimization
    - **Communication:**
      - **Compiled into** ›› RDMA+synchronization and  
›› normal sequential code
    - **Computation:**
      - **Optimizing compiler for sequential code**
  - This **separation** was the beginning of the **success of MPI**
  - But **not yet done** for Co-Array Fortran / UPC



## RDMA platforms (continued)

- MPI-2 one-sided operations
- shmem
- all RDMA programming models can be used also on ccNUMA platforms



## Neither NUMA nor RDMA required

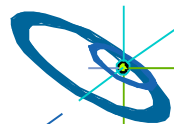
- MPI / PVM on node interconnect
- Inside of the SMP node:
  - MPI → pure MPI (the MPP-MPI model)
  - OpenMP → hybrid programming MPI+OpenMP
- Other models **focus of this report**
  - HPF
  - OpenMP and DSM





## MPI + OpenMP

- MPI\_Init\_threads(required, &provided) MPI 2.0: provided=
- categories of thread-safety: MPI\_THREAD\_...
  - no thread-support by MPI ...\_SINGLE
  - MPI process may be *sometimes* multi-threaded, ~~...\_SINGLE~~  
(parallel regions) and MPI is called only if only the master-thread exists
  - Same, but the other threads may sleep ~~...\_SINGLE~~
  - MPI may be called only outside of OpenMP parallel regions ~~...\_SINGLE~~ **...\_MASTERONLY** { proposal for MPI 2.1
  - Same, but all other threads may compute ...\_FUNNELED
  - Multiple threads may call MPI, but only one thread may execute an MPI routine at a time ...\_SERIALIZED
  - MPI may be called from any thread ...\_MULTIPLE ■



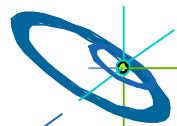
## MPI + OpenMP

- using OMP MASTER → MPI\_THREAD\_FUNNELED needed
  - no implied barrier !
  - no implied cache flush !
- using OMP SINGLE → MPI\_THREAD\_SERIALIZED needed
- A[i] = ...  
OMP MASTER / SINGLE  
MPI\_Send(A, ...)  
OMP END MASTER / SINGLE  
A[i] = new value
  - OMP BARRIER
  - OMP BARRIER
- Same problem as with MPI\_THREAD\_MASTERONLY:
  - **all application threads are sleeping while MPI is executing** ■



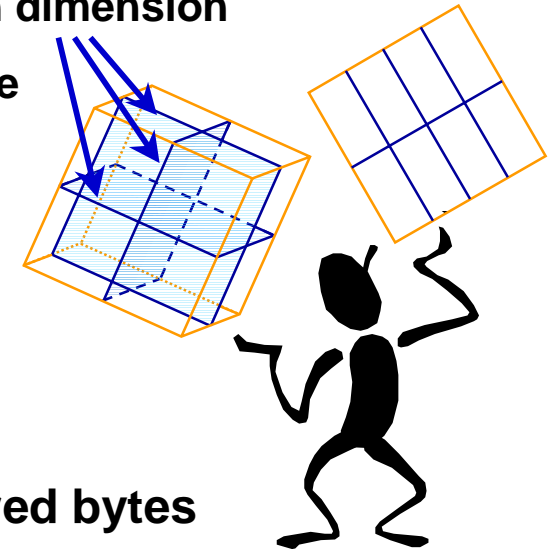
## Pure MPI on hybrid architectures

- Optimizing the communication
    - best ranking of MPI processes on the cluster
      - based on MPI virtual topology
      - sequential ranking: **0, 1, 2, 3, ... 7**   **8, 9, 10, ... 15**   **16, 17, 18 ... 23**
      - round robin: **0, N, 2N ... 7N**   **1, N+1, ... 7N+1**   **2, N+2, ... 7N+2**
- (Example: N = number of used nodes, 8 threads per node)



## Pure MPI on hybrid architectures (continued)

- Additional message transfer inside of each node
  - compared with MPI+OpenMP
  - Example: 3-D (or 2-D) domain decomposition
    - e.g. on 8-way SMP nodes
    - one (or 1–3) additional cutting plane in each dimension
    - **expecting same message size on each plane**
      - outer boundary (pure MPI)
      - inner plane (pure MPI)
      - outer boundary (MPI+OpenMP)
  - *pure MPI* compared with *MPI+OpenMP*:  
**only doubling the total amount of transferred bytes**

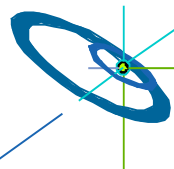


## Benchmark results

- On Hitachi SR8000, b\_eff<sup>1)</sup> benchmark on 12 nodes

	b_eff	b_eff Lmax <sup>2)</sup>	3-d-cyclic average	3-d-cyclic Lmax <sup>2)</sup>
aggregated bandwidth – <b>hybrid</b> [MB/s] (per node)	<b>1535</b> (128)	<b>5565</b> (464)	<b>1604</b> (134)	<b>5638</b> (470)
aggregated bandwidth – <b>pure MPI</b> [MB/s] (per process)	<b>5299</b> (55)	<b>16624</b> (173)	<b>5000</b> (52)	<b>18458</b> (192)
$\frac{bw_{\text{pure MPI}}}{bw_{\text{hybrid}}}$ (measured)	<b>3.45</b>	<b>2.99</b>	<b>3.12</b>	<b>3.27</b>
$\frac{size_{\text{pure MPI}}}{size_{\text{hybrid}}}$ (assumed)	2 (based on last slide)			
$\frac{T_{\text{hybrid}}}{T_{\text{pure MPI}}}$ (concluding)	<b>1.73</b>	<b>1.49</b>	<b>1.56</b>	<b>1.64</b>

→ communication in this hybrid model is about 60% slower than with pure MPI



## Interpretation of the benchmark results

- If the inter-node bandwidth cannot be consumed by using only one processor of each node  
→ the pure MPI model can achieve a better aggregate bandwidth
- If  $\text{bw}_{\text{pure MPI}} / \text{bw}_{\text{hybrid}} > 2$  &  $\text{size}_{\text{pure MPI}} / \text{size}_{\text{hybrid}} < 2$   
→ faster communication with pure MPI
- If  $\text{bw}_{\text{pure MPI}} = \text{bw}_{\text{hybrid}}$  &  $\text{size}_{\text{pure MPI}} > \text{size}_{\text{hybrid}}$   
→ faster communication with hybrid MPI+OpenMP



## Other Advantages of Hybrid MPI+OpenMP

- No communication overhead inside of the SMP nodes
- Larger message sizes on the boundary
  - reduced latency-based overheads
- Reduced number of MPI processes
  - better speedup (Amdahl's law)
  - faster convergence,  
e.g., if multigrid numeric is computed only on a partial grid



## Workaround for single-threaded MPI implementations in the hybrid MPI+OpenMP model

- Work of MPI routines is done by single thread
  - other processors of the SMP nodes may sleep

Communication aspects on last slides  
Now, local work of the MPI routines, executed by the CPUs:

- **Workaround for single-threaded MPI implementations**
  - concatenation of strided message data:
    - not by MPI with derived datatypes
    - by multi-threaded user code
  - reduction operations (MPI\_reduce / MPI\_Allreduce):
    - numerical operations by user-defined multi-threaded call-back routines
    - no rules in the MPI standard about multi-threading of such call-back routine





## Overlapping computation & communication

The model:

- Hybrid MPI+OpenMP
- At least MPI\_THREAD\_FUNNELED
- While master thread calls MPI routines:
  - all other threads are computing!

The implications:

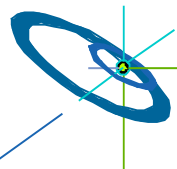
- no communication overhead inside of the SMP nodes
- better CPU usage
  - although inter-node bandwidth may be used only partially
- 2 levels of parallelism:
  - additional synchronization overhead
- Major drawback: load balancing necessary



## Comparing other methods

Memory copies from remote memory to local CPU register and vice versa

Access method	Copies	Remarks	bandwidth $b(\text{message size})$
2-sided MPI	2	internal MPI buffer + application receive buf.	$b(\text{size}) = b_{\infty} / (1 + b_{\infty} T_{\text{latency}} / \text{size})$
1-sided MPI	1	application receive buffer	same formula, but probably better $b_{\infty}$ and $T_{\text{latency}}$
<b>Compiler based:</b> UPC, Co-Array Fortran, HPF, OpenMP on DSM or with cluster extensions	1	page based transfer	extremely poor, if only parts are needed
	0	word based access	8 byte / $T_{\text{latency}}$ , e.g, 8 byte / $0.33\mu\text{s} = 24\text{MB/s}$
	<b>0</b>	<b>latency hiding with pre-fetch</b>	<b><math>b_{\infty}</math></b>
	1	latency hiding with buffering	see 1-sided communication



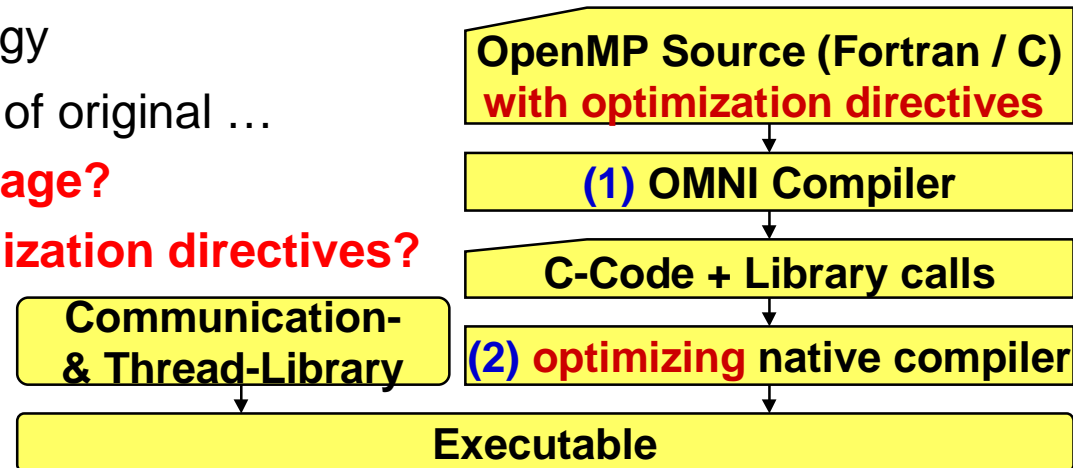
# Compilation and Optimization

- Library based communication (e.g., MPI)
  - clearly separated optimization of
    - (1) communication → MPI library
    - (2) computation → Compiler

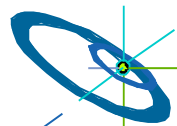
**essential for  
success of MPI**

- Compiler based parallelization (including the communication):
  - similar strategy
  - preservation of original ...

- ... **language?**
- ... **optimization directives?**



- **Optimization of the computation more important than optimization of the communication** ■



## Summary

- **Programming models on hybrid architectures** (clusters of SMP nodes)
  - Pure MPI / hybrid MPI+OpenMP  
/ compiler based parallelization (e.g. OpenMP on clusters)
- **Communication**
  - difficulties with hybrid programming
    - multi-threading with MPI
    - bandwidth of inter-node network
    - overlapping of computation and communication
  - latency hiding with compiler based parallelization
- **Optimization and compilation**
  - separation of optimization
  - we must not lose **optimization of computation**
- **Conclusion:**
  - Pure MPI → MPI+OpenMP → OpenMP on clusters
  - a roadmap full of stones!

