

Optimization of MPI_Allreduce and MPI_Reduce

Rolf Rabenseifner, Panos Adamidis
{rabenseifner, adamidis}@hls.de

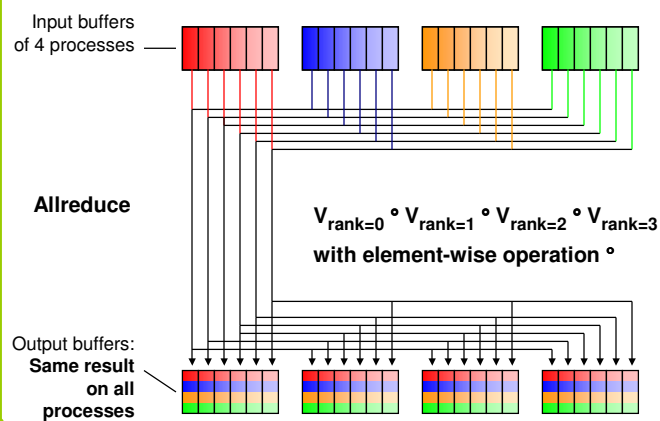
University of Stuttgart
High-Performance Computing-Center Stuttgart (HLRS)
www.hls.de



MPI_Allreduce & MPI_Reduce Optim.
Slide 1
Hochleistungsrechenzentrum Stuttgart

HLRS

What do we want to do



MPI_Allreduce & MPI_Reduce Optim. Rolf Rabenseifner
Slide 2 / 20
Hochleistungsrechenzentrum Stuttgart

HLRS

Basic Principles

Principle I

- Different optimizations for latency and bandwidth
 - Latency optimization, e.g.,
 - sending the full input buffers to all processors
 - executing the reduction on all processors
 - Bandwidth optimization:
 - splitting the input buffers
 - transferring cross-wise between processes
 - reduction operation only on partial buffers
 - allgather step at the end
- } i.e., reduce_scatter



MPI_Allreduce & MPI_Reduce Optim. Rolf Rabenseifner
Slide 3 / 20 Höchstleistungsrechenzentrum Stuttgart

H L R | S

Basic Principles

Principle II

- In case where the number of processors is a power-of-two, then optimization is possible by buffer halving and distance doubling
- In case where the number of processors is non-power-of-two, various algorithms are shown.

Background

- **37%** of MPI time in **MPI_Allreduce**
- **25%** of user time with **non-power-of-two** number of processes
 - data from automatic profiling of all customers on HLRS CRAY T3E



MPI_Allreduce & MPI_Reduce Optim. Rolf Rabenseifner
Slide 4 / 20 Höchstleistungsrechenzentrum Stuttgart

H L R | S

Rabenseifner's Algo., Nov. 1997

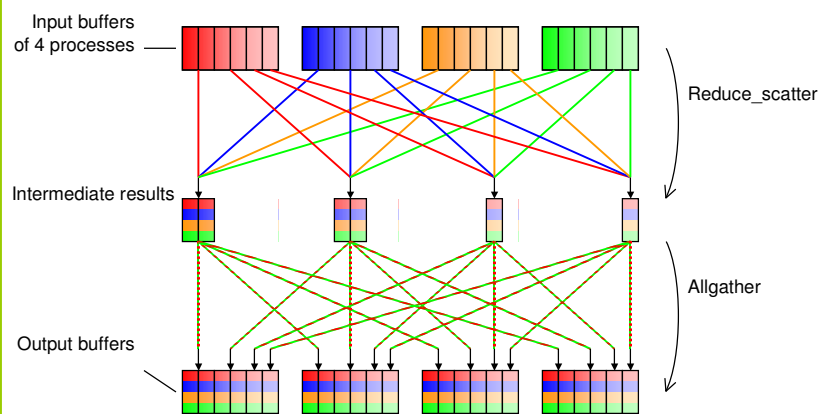
- Standard algorithm used in mpich1:
 - MPI_Reduce = binomial tree
 - MPI_Allreduce = binomial tree + MPI_Bcast
 - Binomial tree is inefficient
 - logarithmic behavior but in each iteration, **half of the processes gets inactive** → **bad load balancing**
- Better algorithms (butterfly-algorithms):
 - MPI_Reduce = Reduce_scatter + Gather
 - MPI_Allreduce = Reduce_scatter + Allgather



MPI_Allreduce & MPI_Reduce Optim. Rolf Rabenseifner
Slide 5 / 20 Höchstleistungsrechenzentrum Stuttgart

H L R I S

Reduce_scatter and Allgather

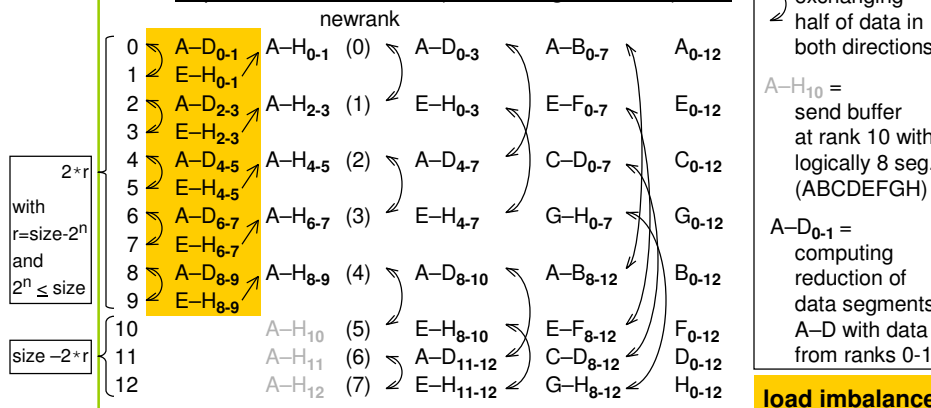


MPI_Allreduce & MPI_Reduce Optim. Rolf Rabenseifner
Slide 6 / 20 Höchstleistungsrechenzentrum Stuttgart

H L R I S

Scheme with Rabenseifner's Algo., Nov. 1997 (1st part)

Rank 1st part: Reduce scatter ... (with **halving** the buffers)



Always computing: $\{ [(0+1)+(2+3)] + [(4+5)+(6+7)] \} + \{ [(8+9)+(10)] + [(11)+(12)] \}$

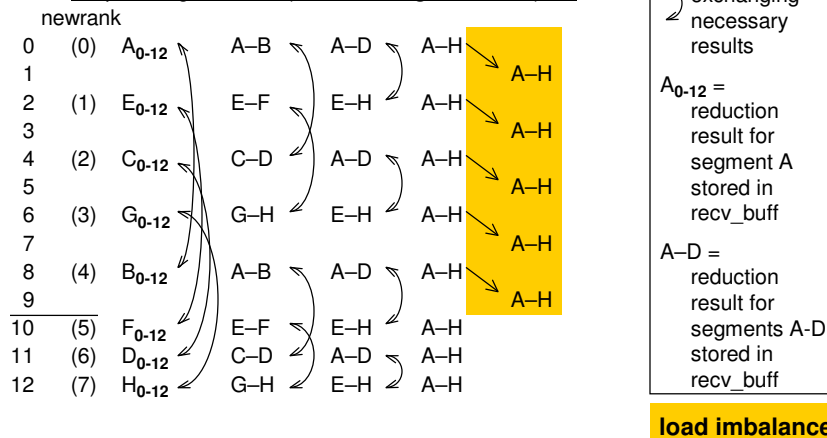


MPI_Allreduce & MPI_Reduce Optim. Rolf Rabenseifner
 Slide 7 / 20 Höchstleistungsrechenzentrum Stuttgart

H L R I S

Scheme with Rabenseifner's Algo., Nov. 1997 (2nd part)

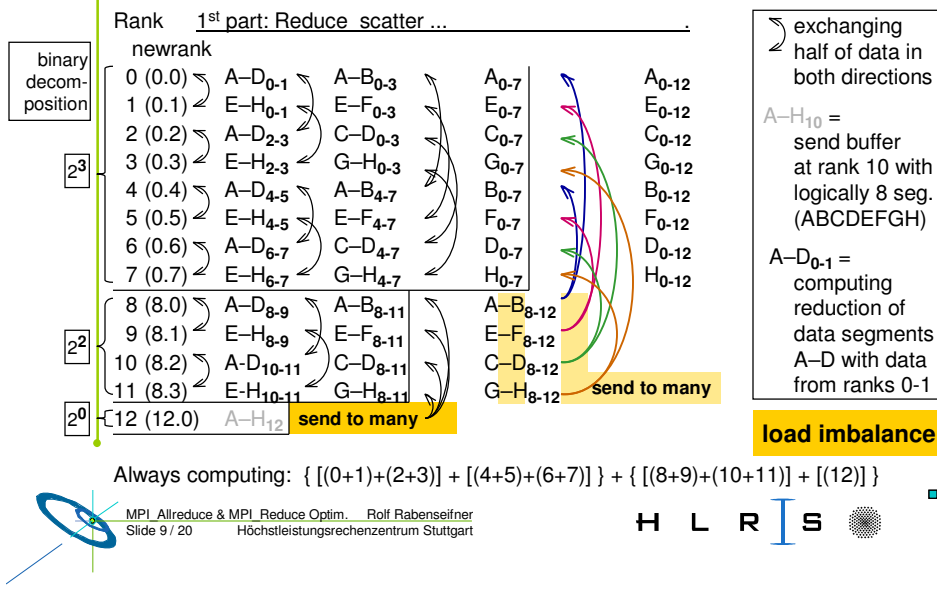
Rank 2nd part: Allgather ... (with **doubling** the buffers)



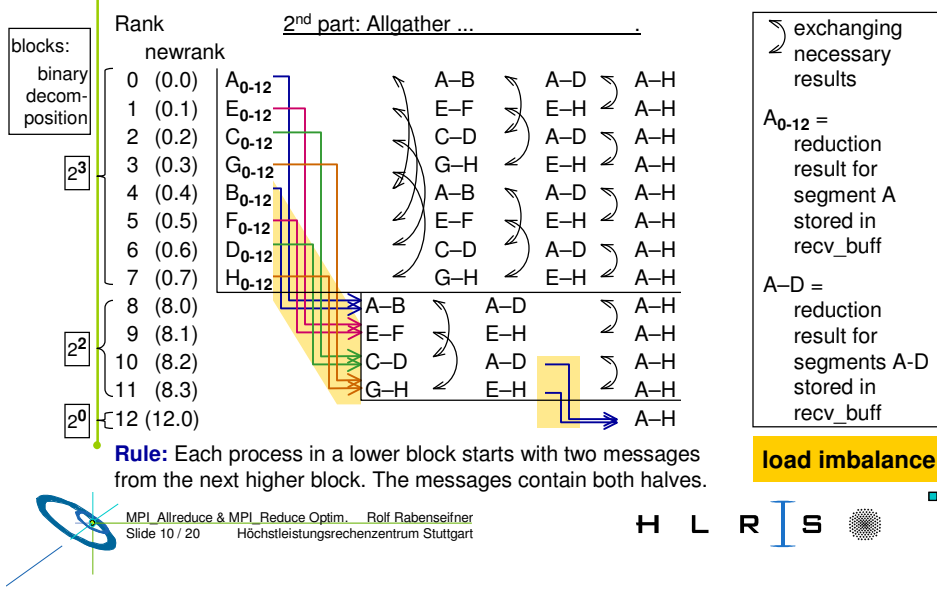
MPI_Allreduce & MPI_Reduce Optim. Rolf Rabenseifner
 Slide 8 / 20 Höchstleistungsrechenzentrum Stuttgart

H L R I S

New Binary Blocks Halving+Doubling, July 2003 (1st part)



New Binary Blocks Halving+Doubling, July 2003 (2nd part)



Compared Protocols

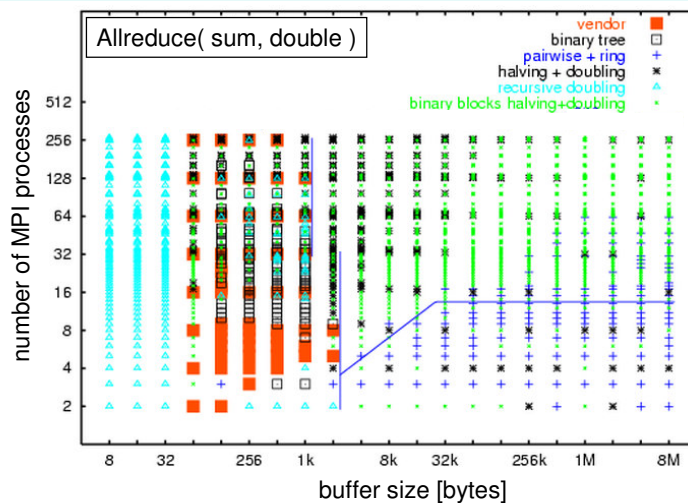
- **Vendor** (MPI_Allreduce and MPI_Reduce of the used MPI library)
- **Binomial tree** + Bcast (i.e., without latency optimization)
- **Recursive doubling** with full buffers (i.e., with latency optimization)
- *Reduce_scatter + Allgather (or Gather)*
 - **Pairwise & Ring**
 - input buffer is divided into (#proc.) pieces of same size
 - optimal load balance but high latency
 - $O(2x \text{ #processes}) + O(\underline{2x \text{ vector size}})$
 - **Halving & Doubling**
 - $O(2x \lceil \log(\text{#processes}) \rceil) + O(4x \text{ vector size})$
 - **Binary Blocks Based Halving & Doubling**
 - normally better than halving & doubling
 - except for special #processes, e.g. 17, 33, 65,....



MPI_Allreduce & MPI_Reduce Optim. Rolf Rabenseifner
Slide 11 / 20 Höchstleistungsrechenzentrum Stuttgart

HLRS

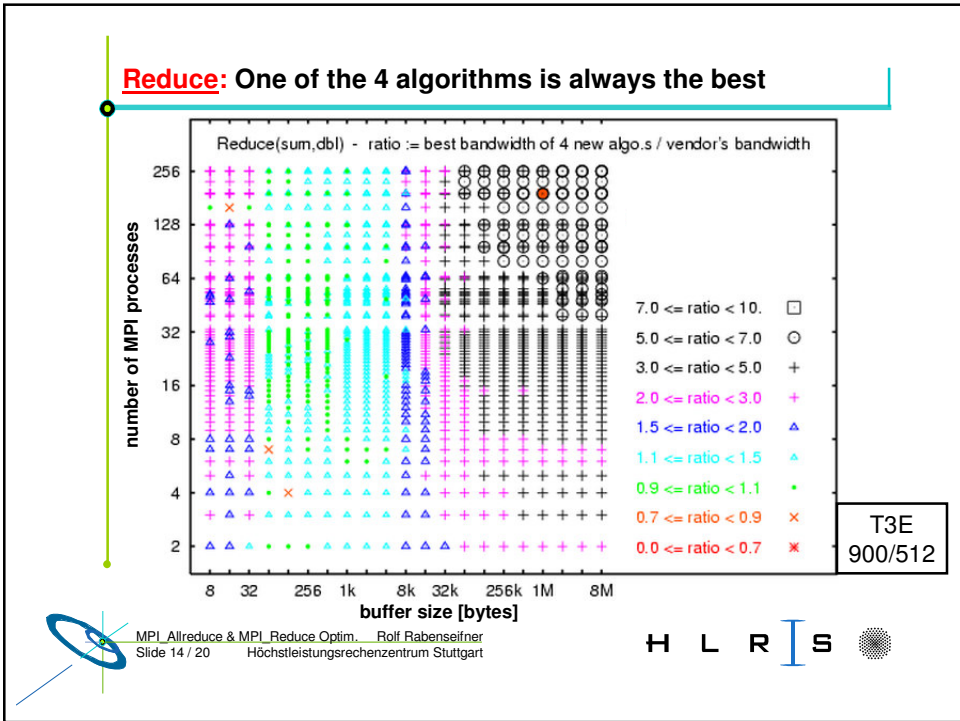
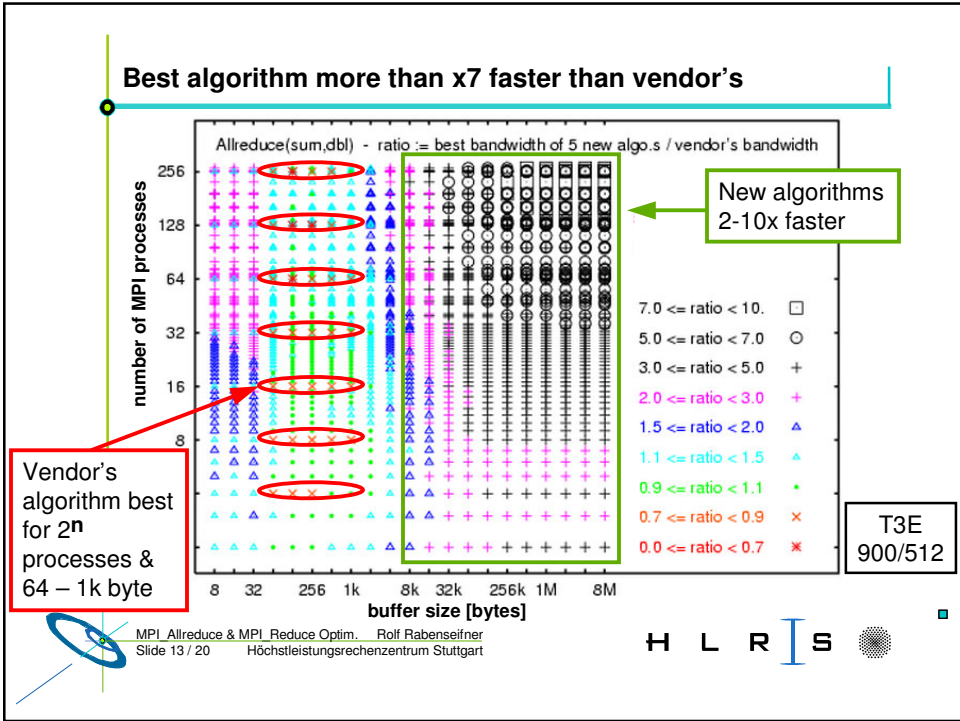
Comparison: Fastest Protocol on T3E 900/512



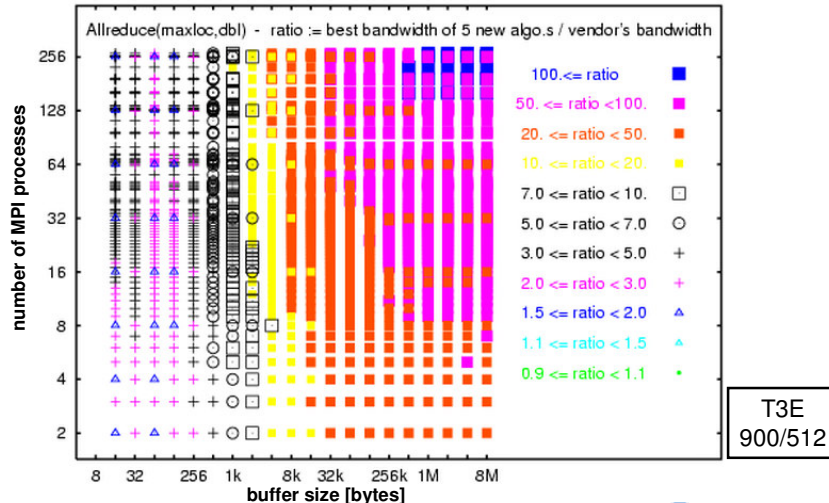
MPI_Allreduce & MPI_Reduce Optim. Rolf Rabenseifner
Slide 12 / 20 Höchstleistungsrechenzentrum Stuttgart

HLRS

Benchmarks on T3E 900/512, sum of doubles, bandwidth := buffersize / wallclock time



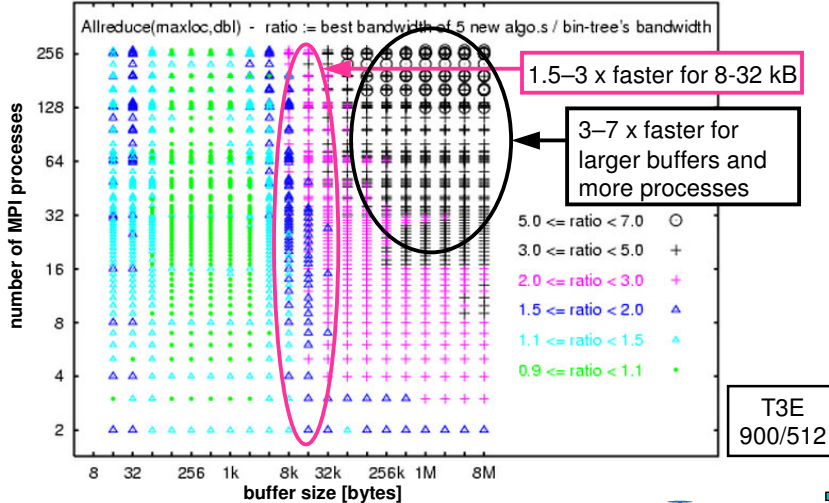
MAXLOC: Vendor's MPI_DOUBLE_INT is extreme slow



MPI_Allreduce & MPI_Reduce Optim. Rolf Rabenseifner
Slide 15 / 20 Höchstleistungsrechenzentrum Stuttgart

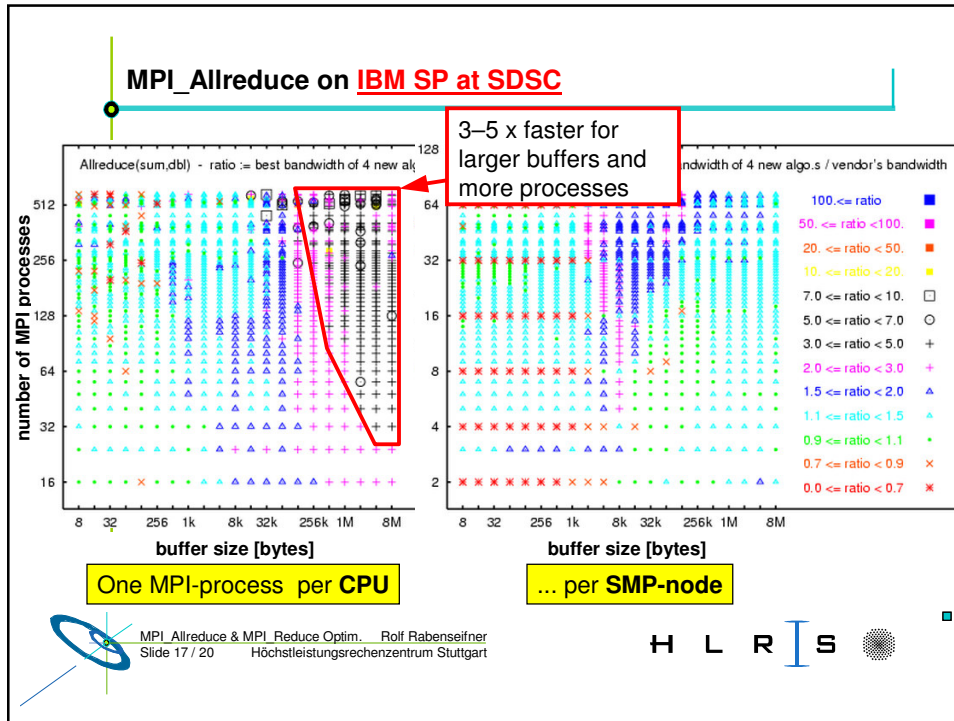
HLRS

MAXLOC: Comparing with optimized binomial-tree



MPI_Allreduce & MPI_Reduce Optim. Rolf Rabenseifner
Slide 16 / 20 Höchstleistungsrechenzentrum Stuttgart

HLRS



- ### Real Benefit
- depends on real usage-pattern
 - # calls
 - # bytes (bandwidth vs. latency)
 - # processes (power-of-two?)
 - operation (MPI_SUM, MPI_MAX, MPI_MAXLOC, user-defined, ...)
 - Allreduce / Reduce
 - Cray X1
 - code is developed on T3E
 - code is tested on CRAY Opteron clusters, Linux clusters, IA32/IA64, IBM
 - code has a problem on Cray X1, need more time for debugging
 - first result on X1:
 - MAXLOC problem with Cray-MPI: **same as on T3E**
- MPI_Allreduce & MPI_Reduce Optim. Rolf Rabenseifner
Slide 18 / 20 Höchstleistungsrechenzentrum Stuttgart
- H L R I S

Acknowledgments

- Thanks for helpful discussions
 - Rajeev Thakur (Argonne)
 - Jesper Larsson Träff (NEC)
- Thanks for help with benchmarking
 - Gerhard Wellein (Uni. Erlangen)
 - Thomas Ludwig, Ana Kovatcheva (Uni. Heidelberg)
 - Rajeev Thakur (Argonne)
 - Monika Wierse, Andy Mason (Cray)
 - Patrick H. Worley (ORNL)
 - Terry Hewitt, Mike Pettipher, Adrian Tate (Uni. Manchester)



MPI_Allreduce & MPI_Reduce Optim. Rolf Rabenseifner
Slide 19 / 20 Höchstleistungsrechenzentrum Stuttgart

HLRS

Conclusion & Future Work

- Latency & Bandwidth optimization of MPI_Allreduce and MPI_Reduce is
 - possible
 - important
 - the '97 algorithm is now part of mpich
- Future work:
Integrated algorithm under construction
 - smooth optimization for any vector size
 - nearly optimal for any # processes
 - again significantly better bandwidth for non-power-of-two



MPI_Allreduce & MPI_Reduce Optim. Rolf Rabenseifner
Slide 20 / 20 Höchstleistungsrechenzentrum Stuttgart

HLRS