

Optimization of Collective Reduction Operations

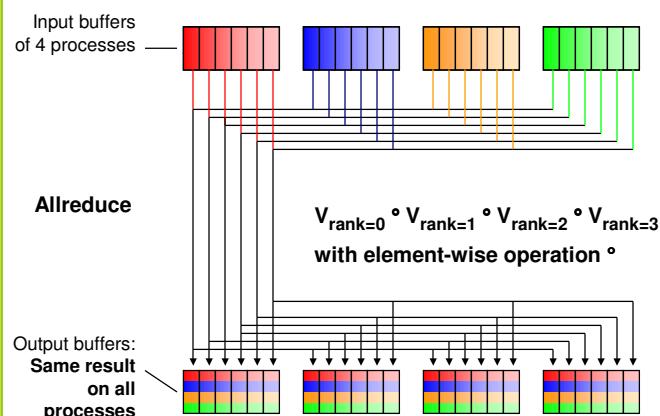
Rolf Rabenseifner
rabenseifner@hlrs.de

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

MPI_Allreduce & MPI_Reduce Optim.
Höchstleistungsrechenzentrum Stuttgart

H L R I S

What do we want to do



MPI_Allreduce & MPI_Reduce Optim.
Rolf Rabenseifner
Höchstleistungsrechenzentrum Stuttgart

H L R I S

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

H L R I 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

H L R I S

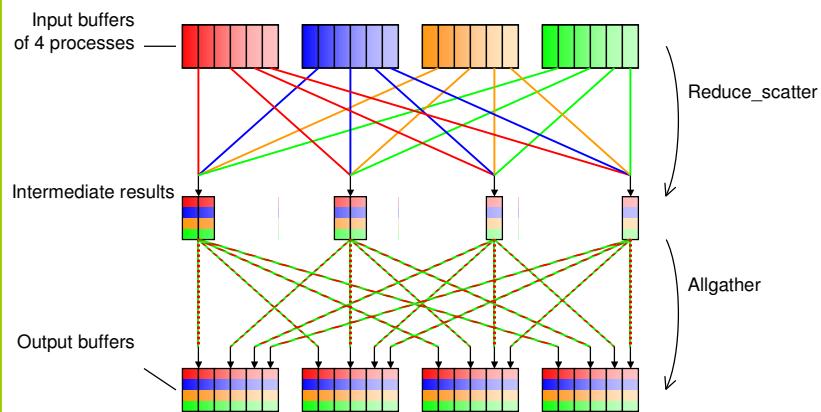
Rabenseifner's Algo., Nov. 1997

- 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 / 19 Höchstleistungsrechenzentrum Stuttgart

H L R I S

Reduce_scatter and Allgather

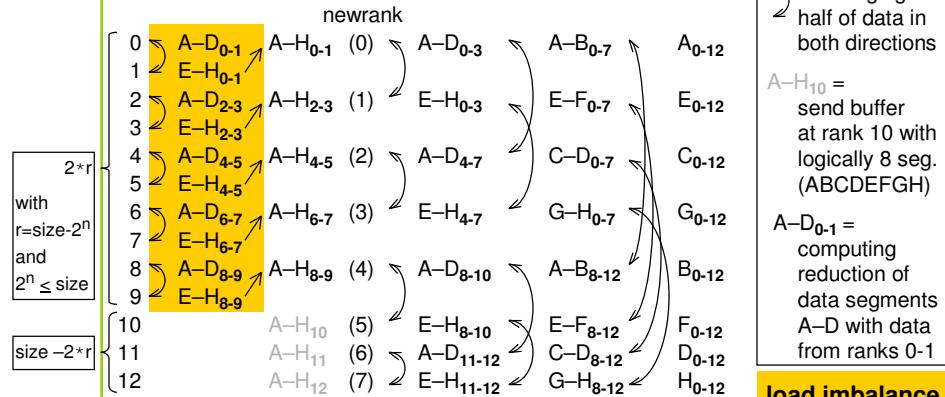


MPI_Allreduce & MPI_Reduce Optim. Rolf Rabenseifner
Slide 6 / 19 Höchstleistungsrechenzentrum Stuttgart

H L R I S

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

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



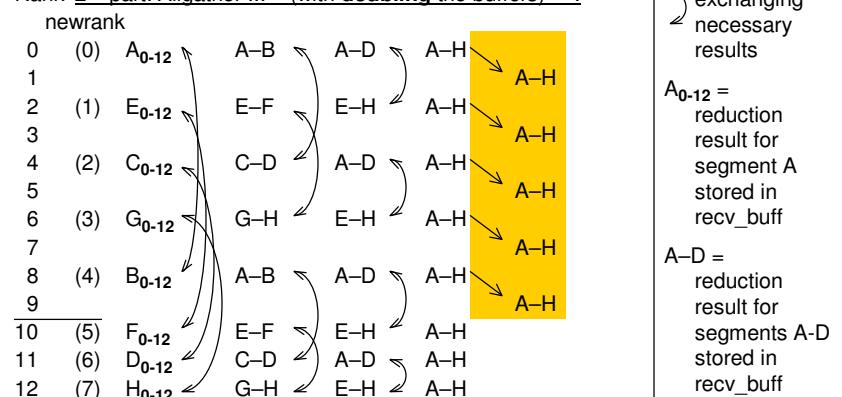
Always computing: { [(0+1)+(2+3)] + [(4+5)+(6+7)] } + { [(8+9)+(10)] + [(11)+(12)] }

MPI_Allreduce & MPI_Reduce Optim. Rolf Rabenseifner
Slide 7 / 19 Höchstleistungsrechenzentrum Stuttgart

H L R I S

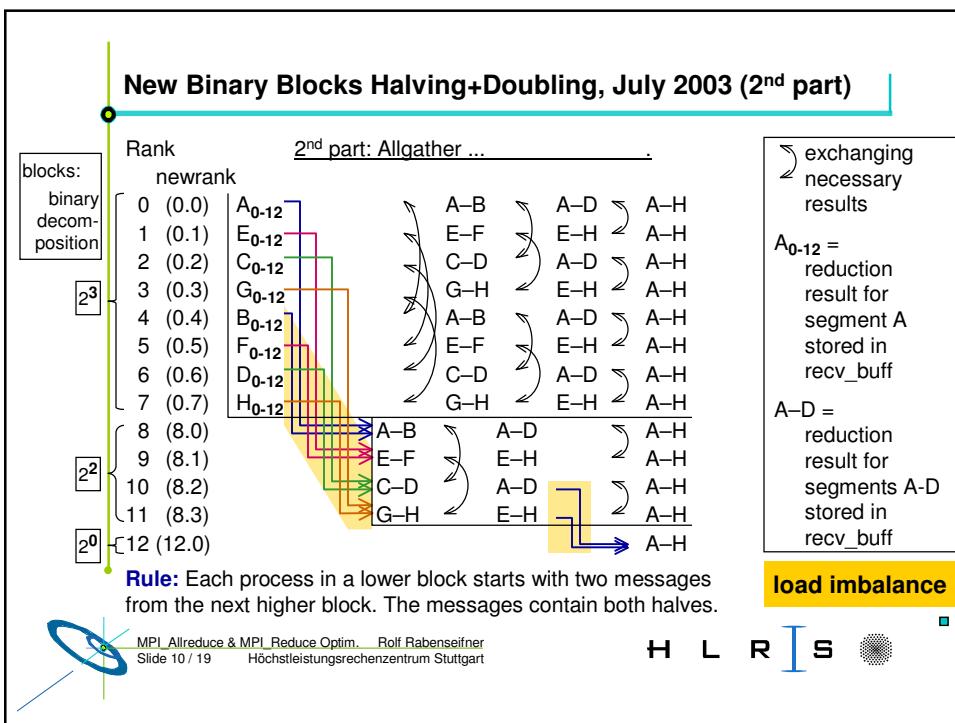
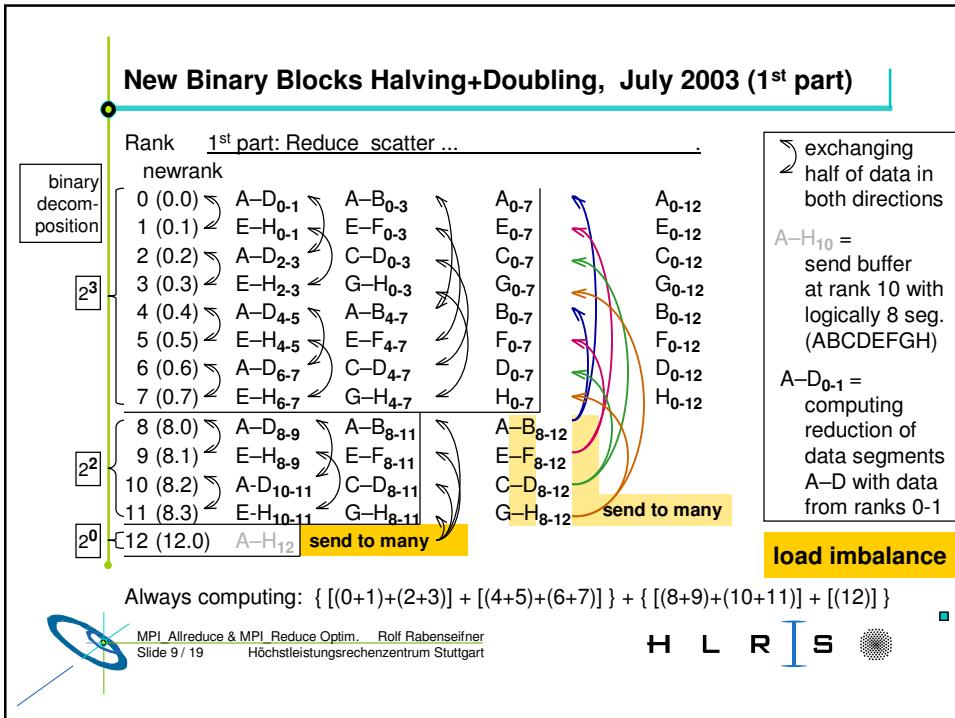
Scheme with Rabenseifner's Algo., Nov. 1997 (2nd part)

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



MPI_Allreduce & MPI_Reduce Optim. Rolf Rabenseifner
Slide 8 / 19 Höchstleistungsrechenzentrum Stuttgart

H L R I S



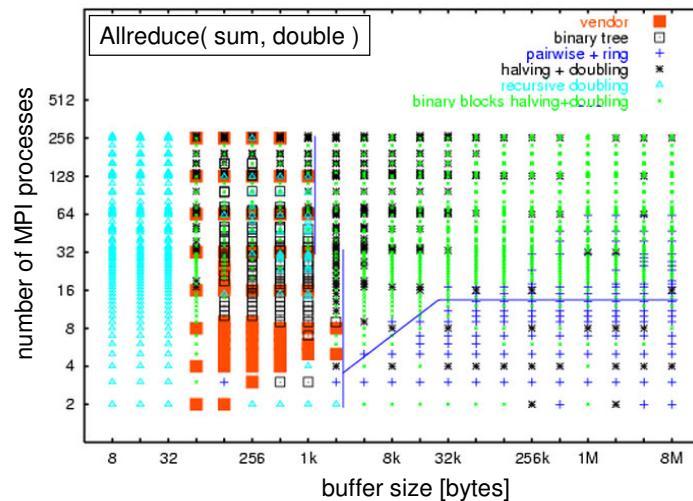
Compared Protocols

- **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 devided into (#proc.) pieces of same size
 - optimal load balance but high latency
 - $O(2x \#processes) + O(2x \text{ vector size})$
 - **Halving & Doubling**
 - $O(2x \lceil \log(\#processes) \rceil) + O(4x \text{ vector size})$
 - **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 / 19 Höchstleistungsrechenzentrum Stuttgart

H L R I S

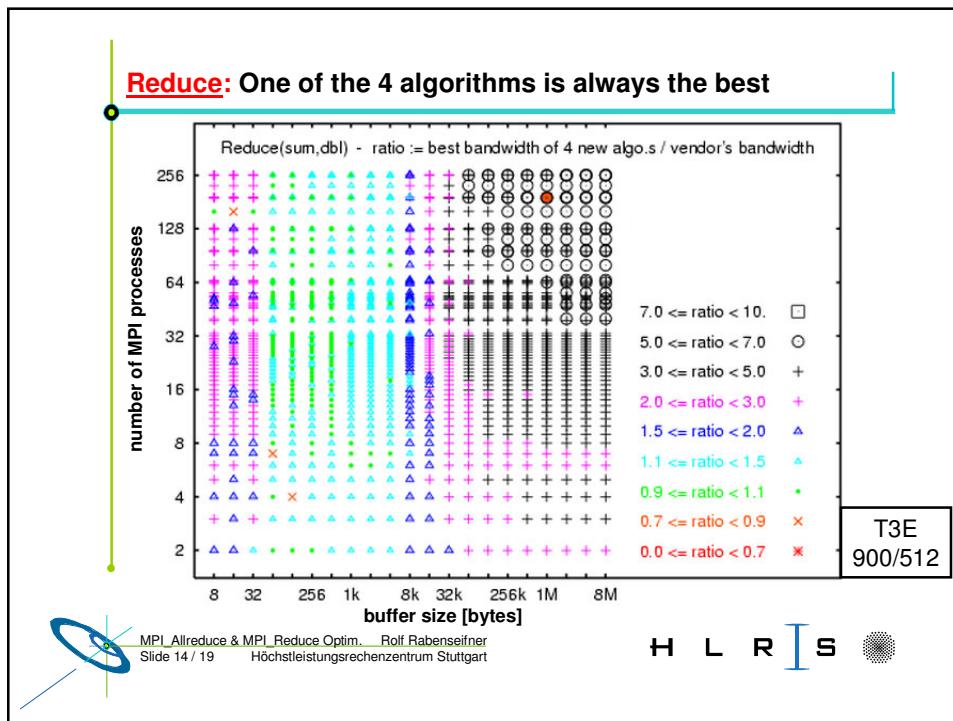
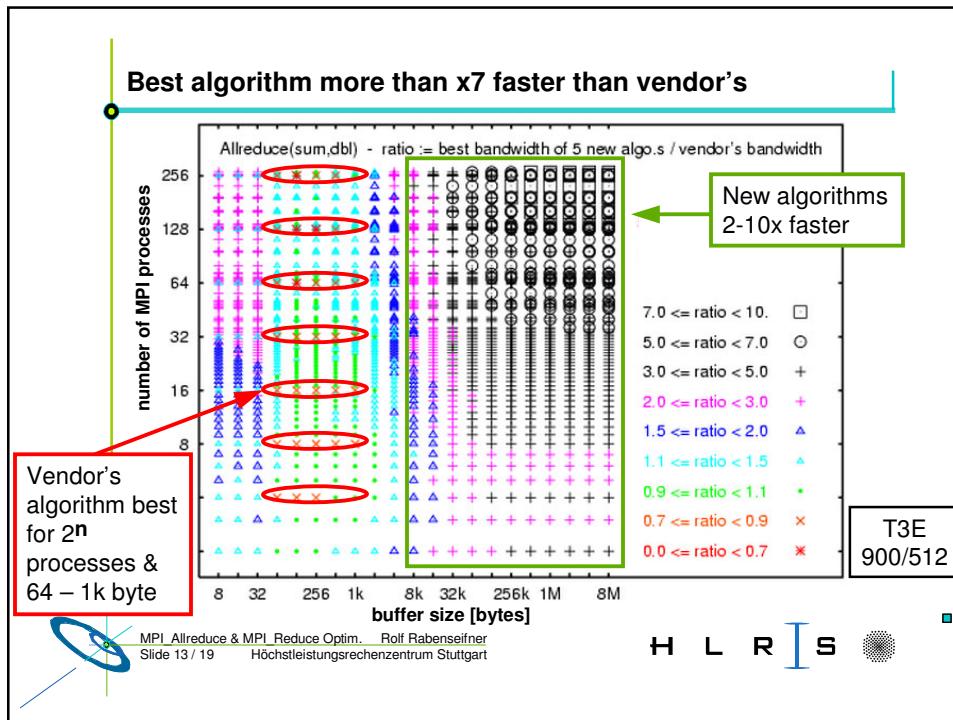
Comparison: Fastest Protocol on T3E 900/512

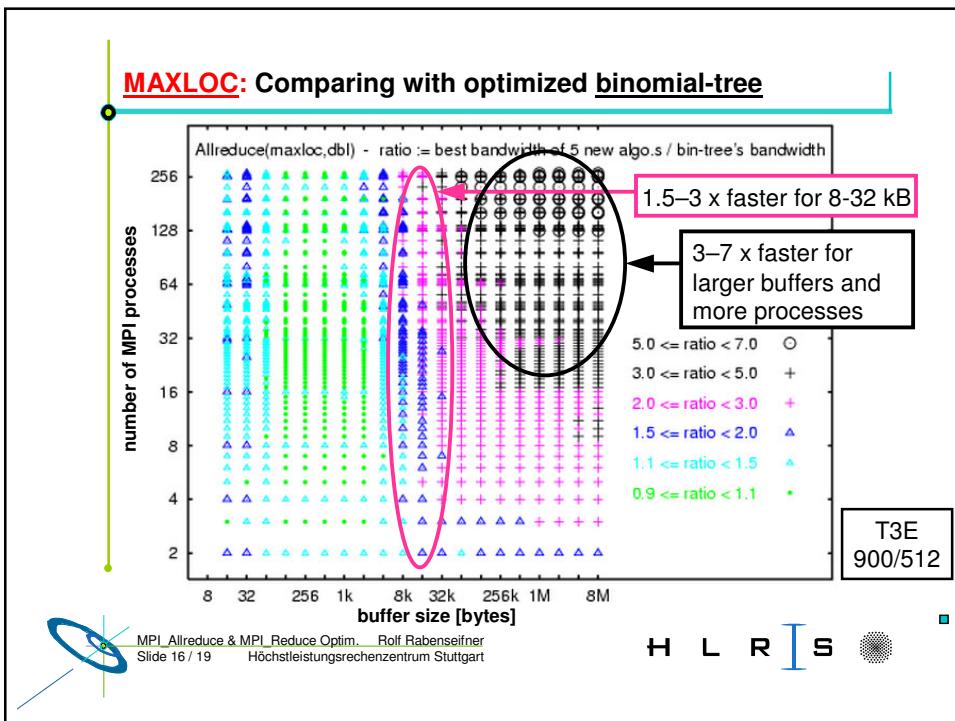
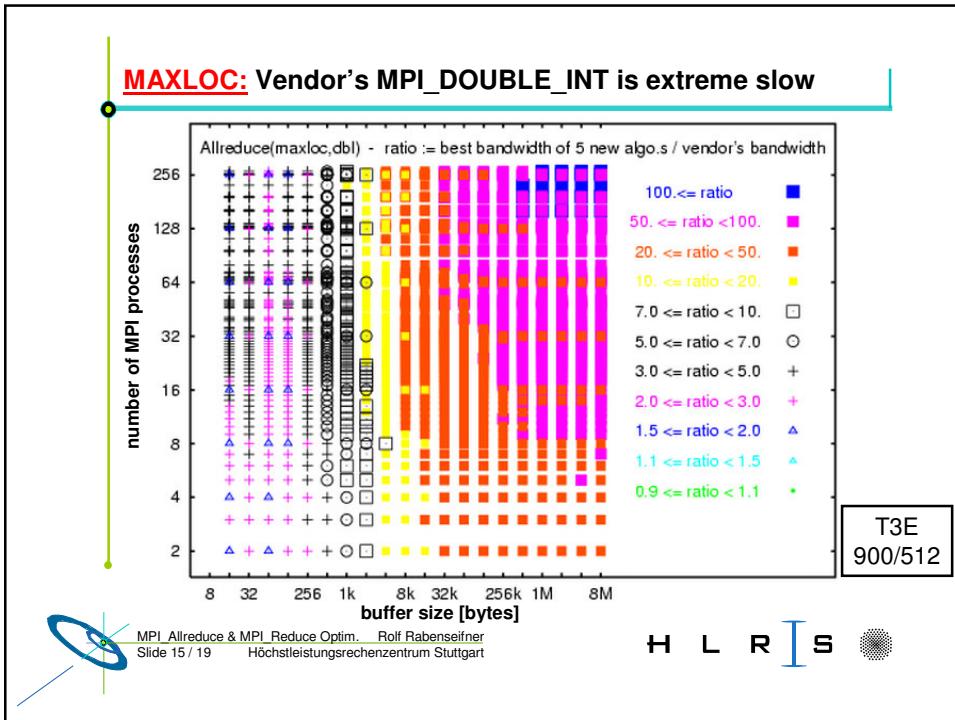


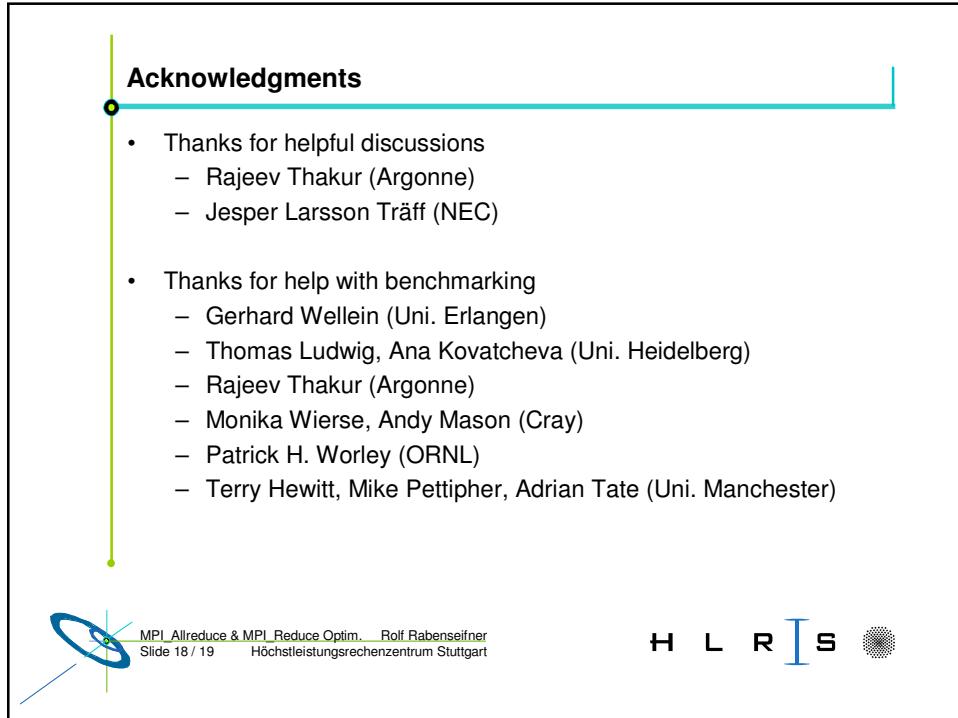
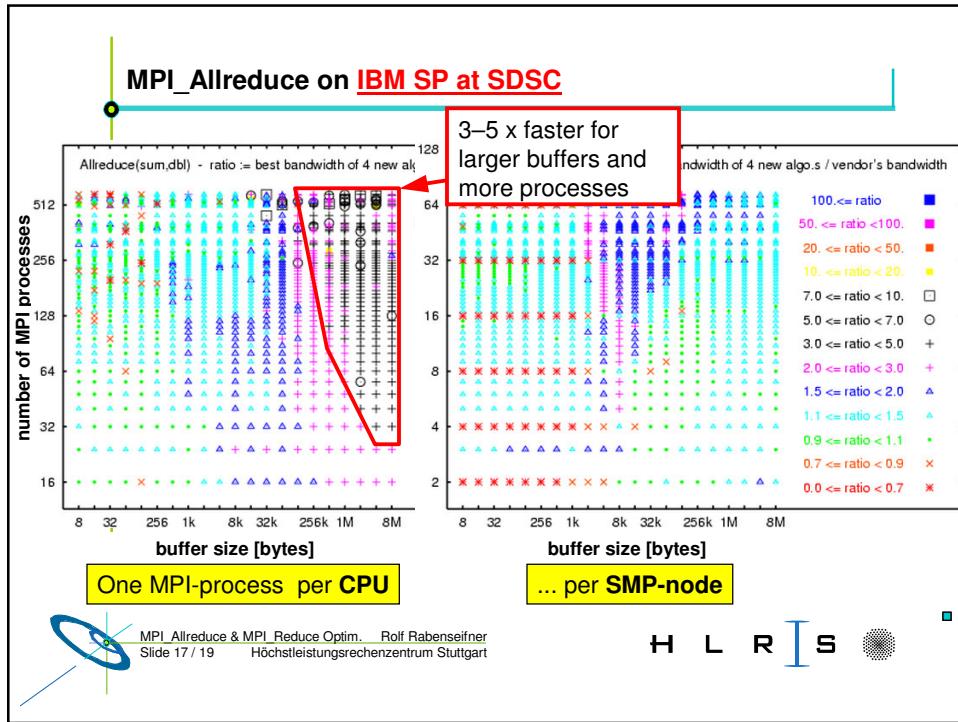
MPI_Allreduce & MPI_Reduce Optim. Rolf Rabenseifner
Slide 12 / 19 Höchstleistungsrechenzentrum Stuttgart

H L R I S

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







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