

# Optimization of Collective Reduction Operations

Rolf Rabenseifner  
rabenseifner@hls.de

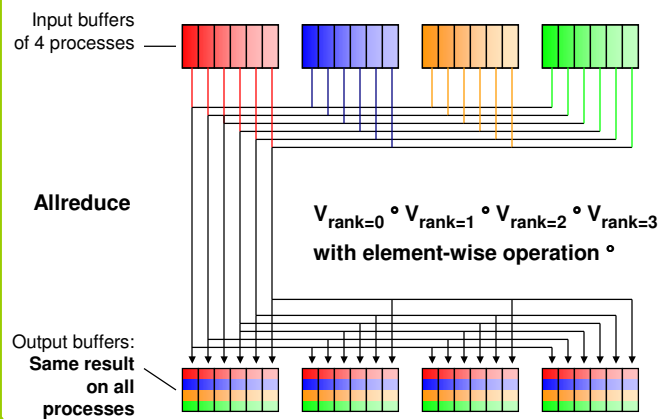
University of Stuttgart  
High-Performance Computing-Center Stuttgart (HLRS)  
www.hls.de



MPI\_Allreduce & MPI\_Reduce Optim.  
Slide 1  
Hochleistungsrechenzentrum Stuttgart



## What do we want to do



MPI\_Allreduce & MPI\_Reduce Optim. Rolf Rabenseifner  
Slide 2 / 19  
Hochleistungsrechenzentrum Stuttgart



## 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 / 19 Höchstleistungsrechenzentrum Stuttgart

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



MPI\_Allreduce & MPI\_Reduce Optim. Rolf Rabenseifner  
Slide 4 / 19 Höchstleistungsrechenzentrum Stuttgart

H L R I S

## Rabenseifner's Algo., Nov. 1997

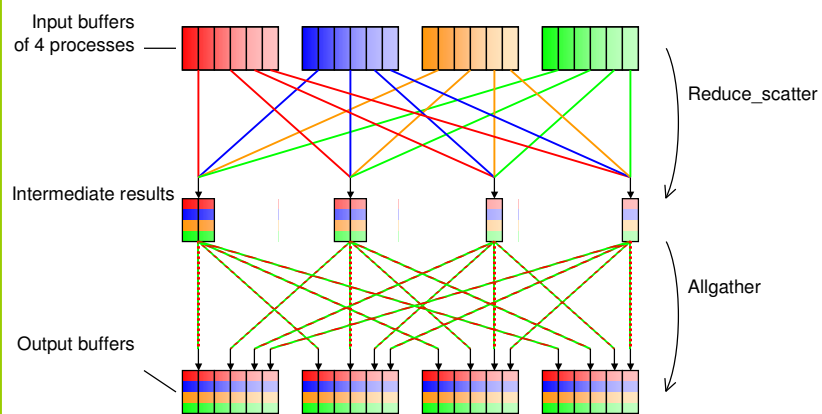
- Standard algorithm used in mpich1:
  - MPI\_Reduce = binomial tree
  - MPI\_Allreduce = binomial tree + MPI\_Bcast
  - Binomial tree is inefficient
    - logarithmic behavior but in each iteration, **half of the processes gets inactive** → **bad load balancing**
- Better algorithms (butterfly-algorithms):
  - 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

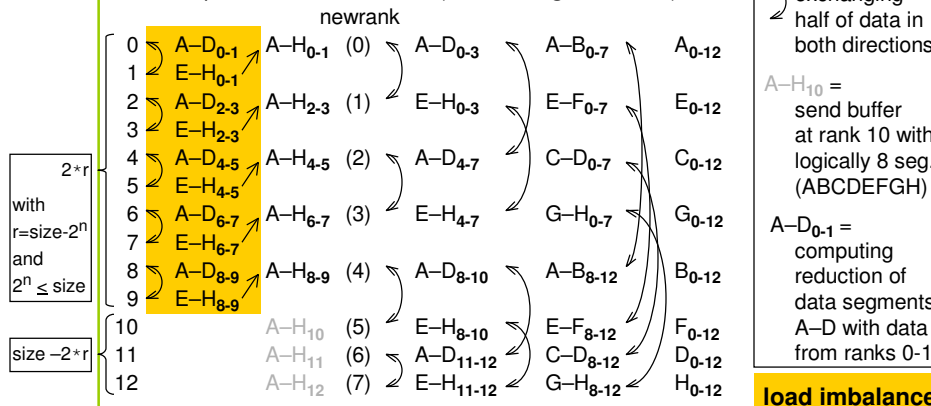


MPI\_Allreduce & MPI\_Reduce Optim. Rolf Rabenseifner  
Slide 6 / 19 Höchstleistungsrechenzentrum Stuttgart

H L R I S

## Scheme with Rabenseifner's Algo., Nov. 1997 (1<sup>st</sup> part)

Rank 1<sup>st</sup> part: Reduce scatter ... (with **halving** the buffers)



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

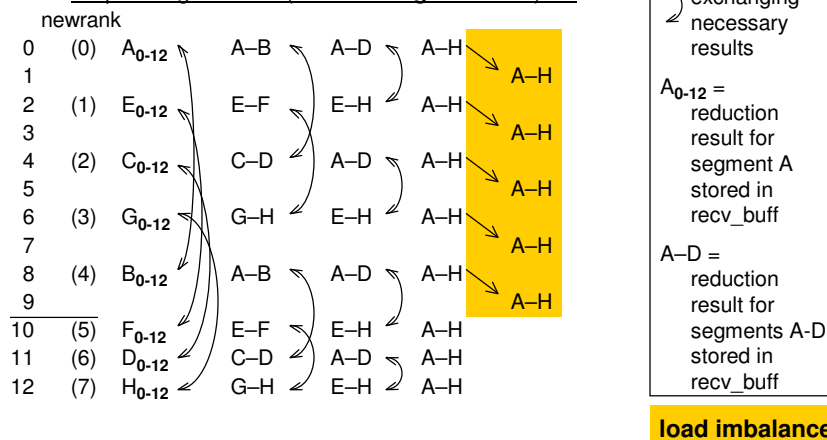


MPI\_Allreduce & MPI\_Reduce Optim. Rolf Rabenseifner  
Slide 7 / 19 Höchstleistungsrechenzentrum Stuttgart

H L R I S

## Scheme with Rabenseifner's Algo., Nov. 1997 (2<sup>nd</sup> part)

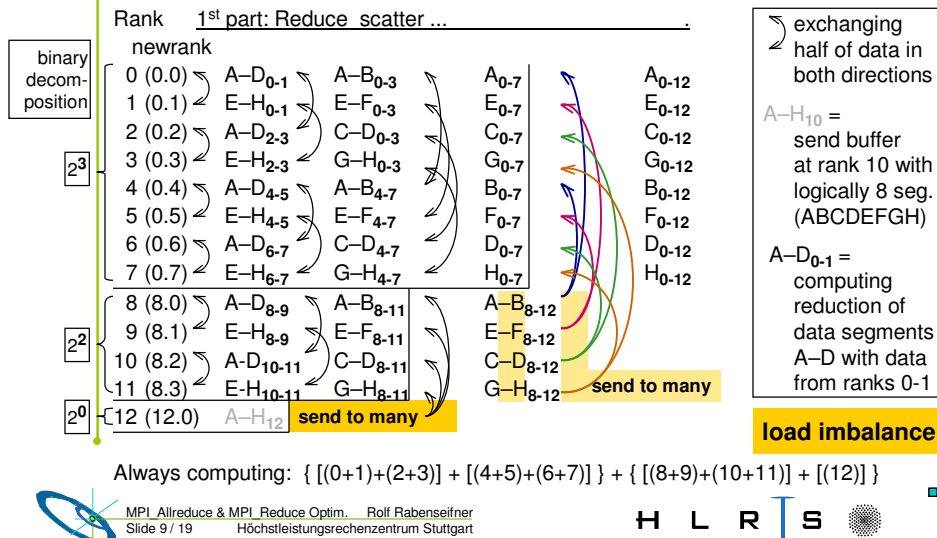
Rank 2<sup>nd</sup> part: Allgather ... (with **doubling** the buffers)



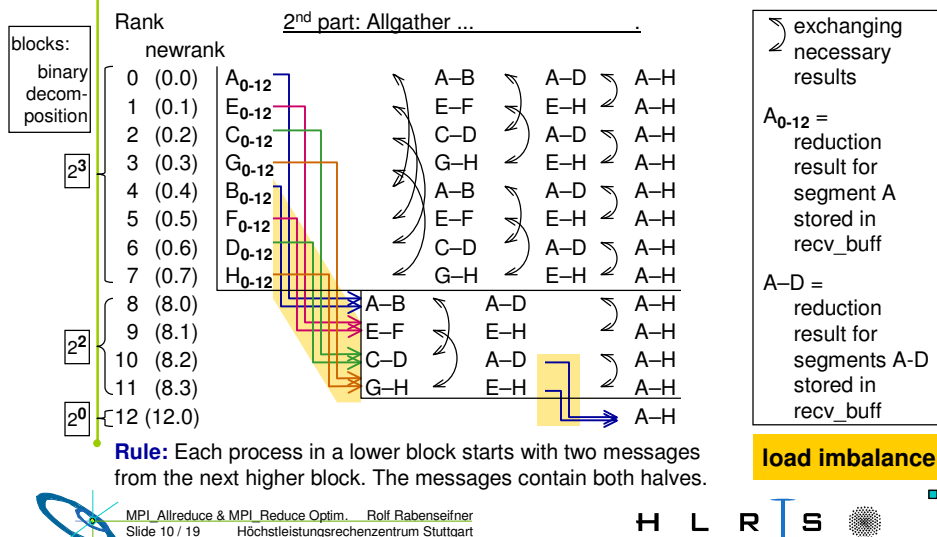
MPI\_Allreduce & MPI\_Reduce Optim. Rolf Rabenseifner  
Slide 8 / 19 Höchstleistungsrechenzentrum Stuttgart

H L R I S

## New Binary Blocks Halving+Doubling, July 2003 (1<sup>st</sup> part)



## New Binary Blocks Halving+Doubling, July 2003 (2<sup>nd</sup> part)



## Compared Protocols

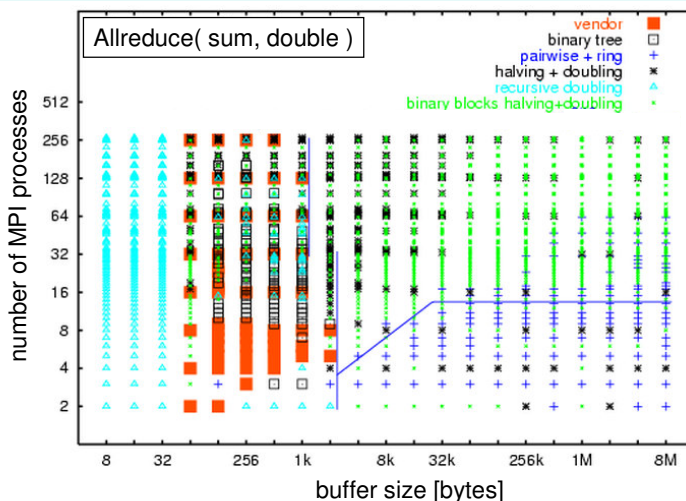
- **Vendor** (MPI\_Allreduce and MPI\_Reduce of the used MPI library)
- **Binomial tree** + Bcast (i.e., without latency optimization)
- **Recursive doubling** with full buffers (i.e., with latency optimization)
- *Reduce\_scatter + Allgather (or Gather)*
  - **Pairwise & Ring**
    - input buffer is 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 / 19 Höchstleistungsrechenzentrum Stuttgart

H L R I S

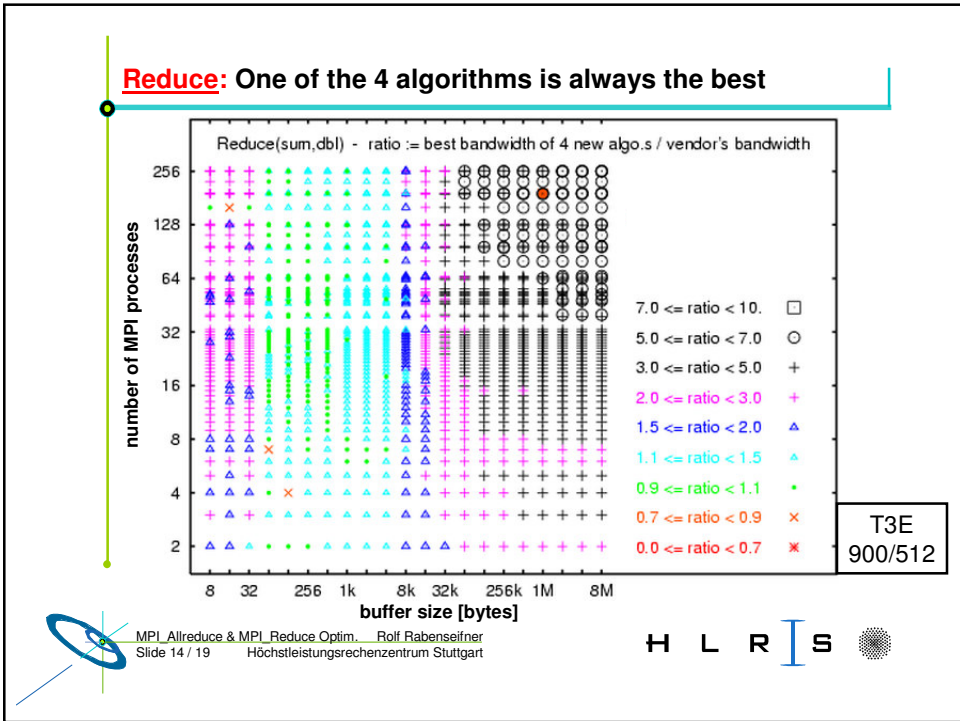
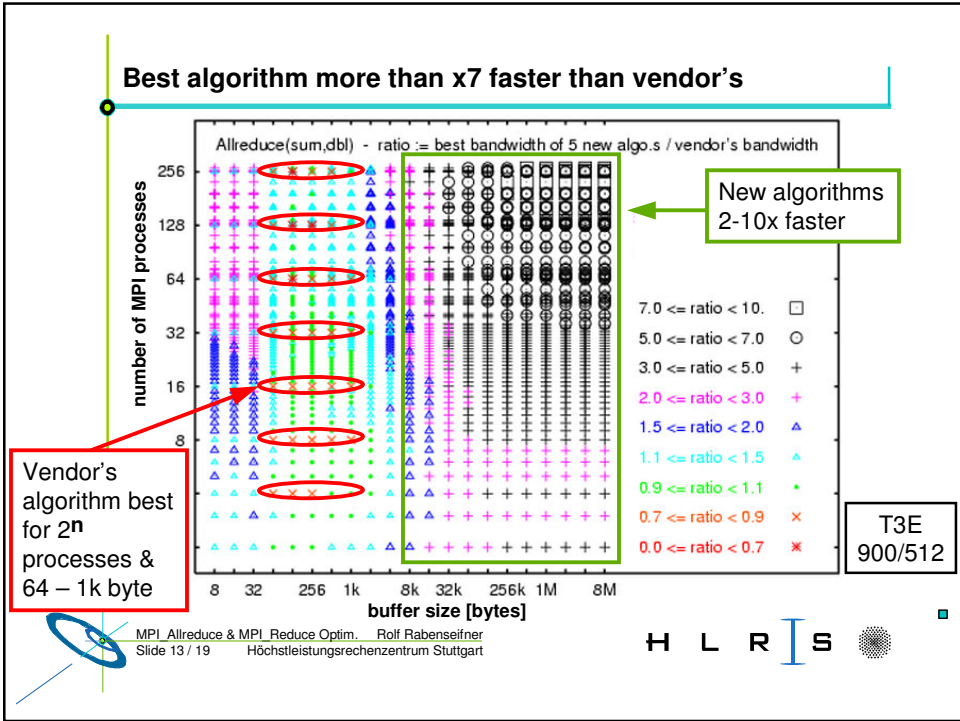
## Comparison: Fastest Protocol on T3E 900/512



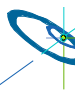
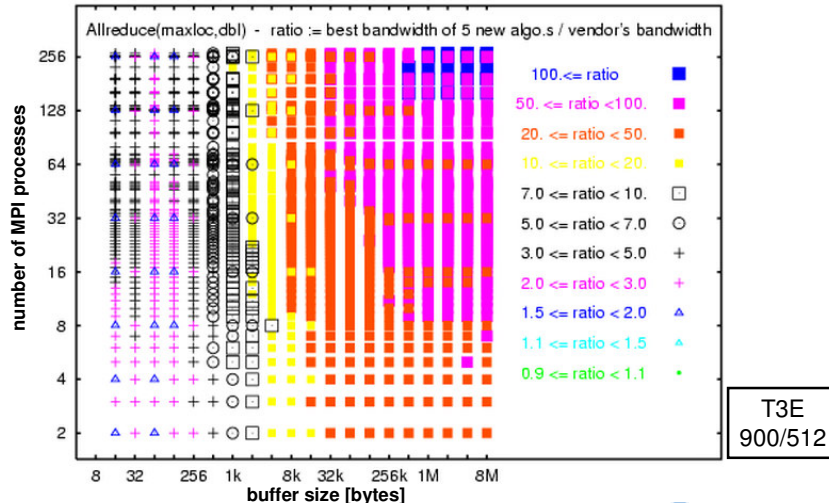
MPI\_Allreduce & MPI\_Reduce Optim. Rolf Rabenseifner  
Slide 12 / 19 Höchstleistungsrechenzentrum Stuttgart

H L R I S

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



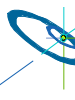
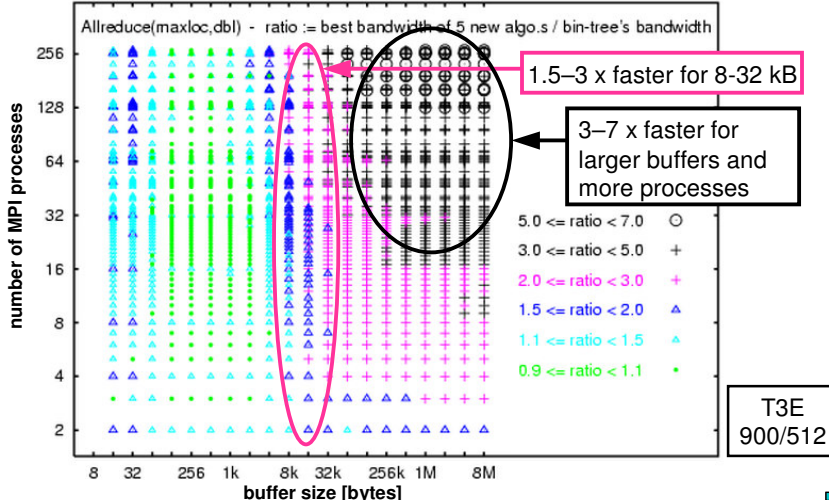
**MAXLOC: Vendor's MPI\_DOUBLE\_INT is extreme slow**



MPI\_Allreduce & MPI\_Reduce Optim. Rolf Rabenseifner  
 Slide 15 / 19 Höchstleistungsrechenzentrum Stuttgart

H L R I S

**MAXLOC: Comparing with optimized binomial-tree**

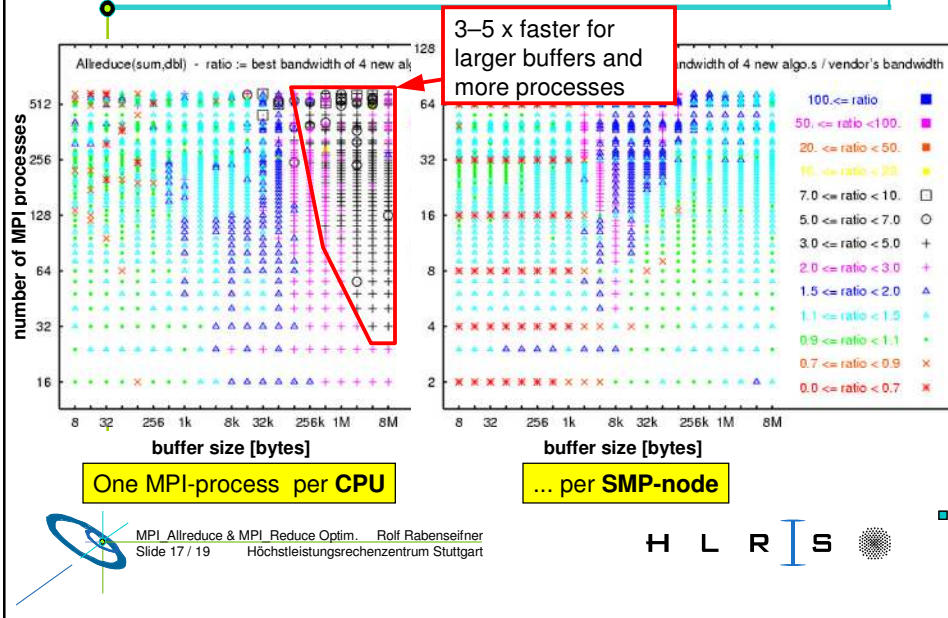


MPI\_Allreduce & MPI\_Reduce Optim. Rolf Rabenseifner  
 Slide 16 / 19 Höchstleistungsrechenzentrum Stuttgart

H L R I S



## MPI\_Allreduce on IBM SP at SDSC



## 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)

## 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 19 / 19 Höchstleistungsrechenzentrum Stuttgart

H L R I S 