

More Efficient Reduction Algorithms for Non-power-of-two Number of Processors in Message-Passing Parallel Systems

Rolf Rabenseifner (HLRS),

rabenseifner@hlrs.de

Jesper Larsson Träff (NEC)

traff@ccrl-nece.de

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

C&C Research Laboratories
NEC Europe Ltd.
St. Augustin (www.ccrl-nece.de)

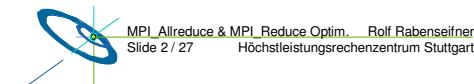
EuroPVM/MPI 2004, Budapest, Sep. 19-22, 2004



H L R I S

Methods

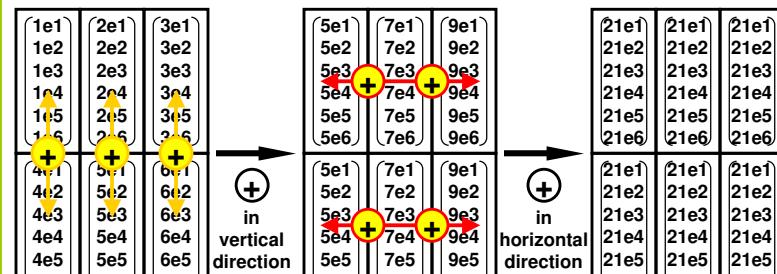
- Factorizing #processes = $2^n \cdot q_1 \cdot q_2 \cdot q_3 \cdots q_m$ ($n \geq 0, m \geq 0, \text{ odd } q_i$)
 - Processes logically in a hypercube with dim=n+m
- Latency-optimization → with full buffer exchange
- Bandwidth-optimization → with buffer halving (or splitting) [**butterfly alg.**]
- Mixing both methods for medium vector sizes
- Factorization implemented with recursive distance doubling ($\times 2$) or $\times q_i$
- Factor $q=2$: optimal butterfly algorithm
- **Odd factors q_i :** several **new** algorithms
 - 2 elimination algorithms:
 - only with 3-to-2 elimination steps
 - with 3-to-2 and 2-to-1 elimination steps
 - Ring algorithm
- Buffer handling
- Further methods



H L R I S

Processes ordered as hypercube

- Reduction separately in each direction

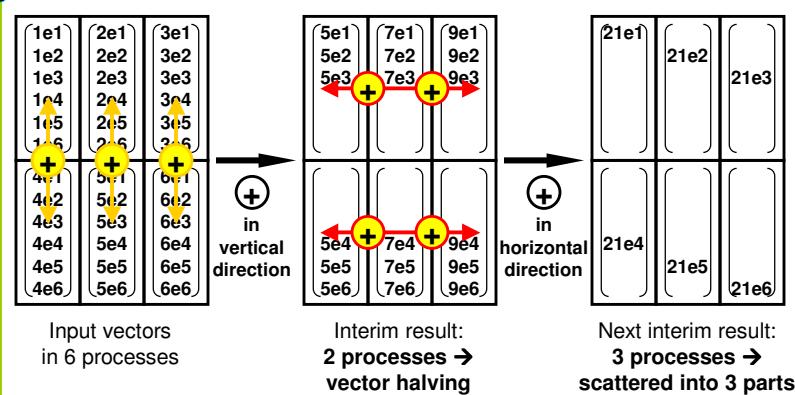


- Pro:** Minimal number of messages – Goal: $O(\log \#proc)$
- Con:** Full size vector transfer in each direction

MPI_Allreduce & MPI_Reduce Optim. Rolf Rabenseifner
Slide 3 / 27 Höchstleistungsrechenzentrum Stuttgart

H L R I S

Splitting the vectors

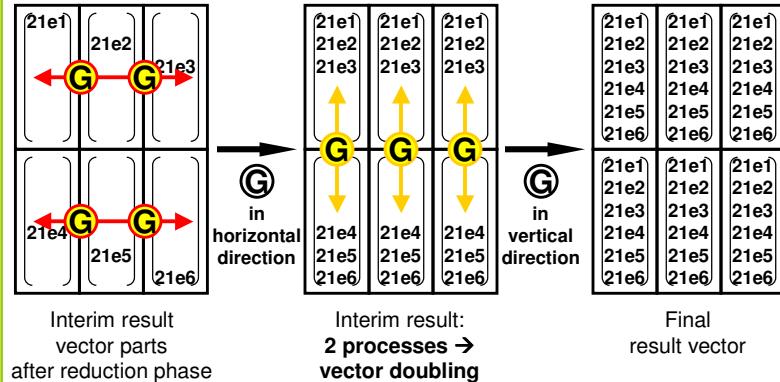


- Reduce_scatter phase
- Additional allgather phase → next slide

MPI_Allreduce & MPI_Reduce Optim. Rolf Rabenseifner
Slide 4 / 27 Höchstleistungsrechenzentrum Stuttgart

H L R I S

In reverse sequence, with MPI_Allgather: Reduction results gathered to all processes



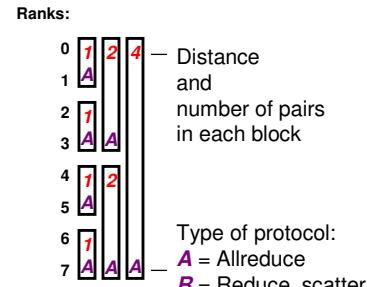
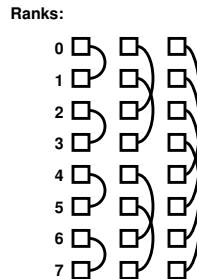
- **Pro:** Sum of all message sizes < 2 * vector size
- **Con:** Doubling the number of messages

MPI Allreduce & MPI Reduce Optim. Rolf Rabenseifner
Slide 5 / 27 Hochleistungsrechenzentrum Stuttgart

H L R I S

Allreduce / Reduce_scatter phase for power-of-two processes

- Communication and reduction protocol with **distance doubling**



- Latency optimized – total reduction result on all processes

Send total buffer receive and reduce it with own total buffer \oplus
 Send total buffer receive and reduce it with own total buffer \oplus

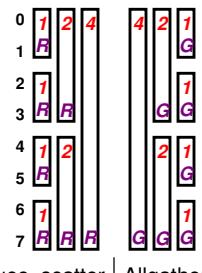
- Bandwidth optimized – on each process: only $1/\#processes$ of reduction result

Send 2nd half of buffer receive and reduce it with own 1st half of own buffer \oplus
 Send 1st half of buffer receive and reduce it with own 2nd half of own buffer \oplus

Allreduce = Reduce_scatter + Allgather (Butterfly)

- Pure Protocol:

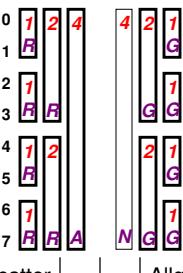
Ranks:



Reduce_scatter | Allgather

- Mixed Protocol:

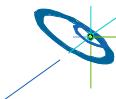
Ranks:



Reduce_scatter | Allgather

Allreduce | No operation

- Allgather protocol:



= Send provisional result receive and store after own provisional result

= Send provisional result receive and store before own provisional result

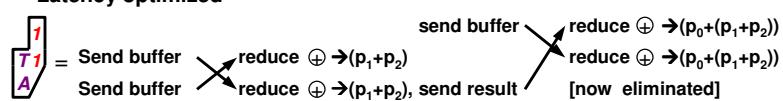
Methods

- Factorizing #processes = $2^n \cdot q_1 \cdot q_2 \cdot q_3 \cdots q_m$ ($n \geq 0, m \geq 0$, odd q_i)
 - Processes logically in a hypercube with dim=n+m
- Full buffer exchange for best latency
- Buffer halving (or splitting) for best bandwidth (butterfly algorithm)
- Mixing both methods for medium vector sizes
- Factorization implemented with recursive distance doubling ($\times 2$) or $\times q_i$
- Factor $q=2$: optimal butterfly algorithm
- Handling odd factors q_i :** several new algorithms
 - 2 elimination algorithms with doubling&halving:**
 - only with 3-to-2 elimination steps
 - with 3-to-2 and 2-1 elimination steps
 - Ring algorithm**
- Buffer handling
- Further methods



Methods for non-power-of-two number of processes

- 3-to-2-process (**Triple**) and 2-to-1-process (**Double**) reduction eliminates processes in the allreduce/reduce_scatter epoch, i.e., those processes are not used in further reduction steps
- 3-to-2-process reduction (**Triple**):
 - Allreduce/reduce_scatter phase:
 - Latency optimized



- Bandwidth optimized



MPI_Allreduce & MPI_Reduce Optim. Rolf Rabenseifner
Slide 9 / 27 Höchstleistungsrechenzentrum Stuttgart

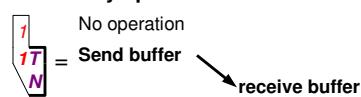
H L R I S

This work was done together with Jesper Larsson Träff, NEC.

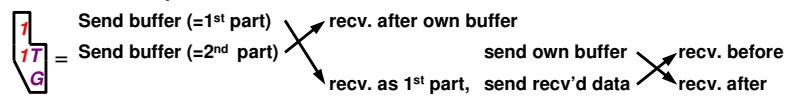
Methods for non-power-of-two number of processes

- 3-to-2-process reduction (**Triple**):
 - Allgather phase:

- Latency optimized



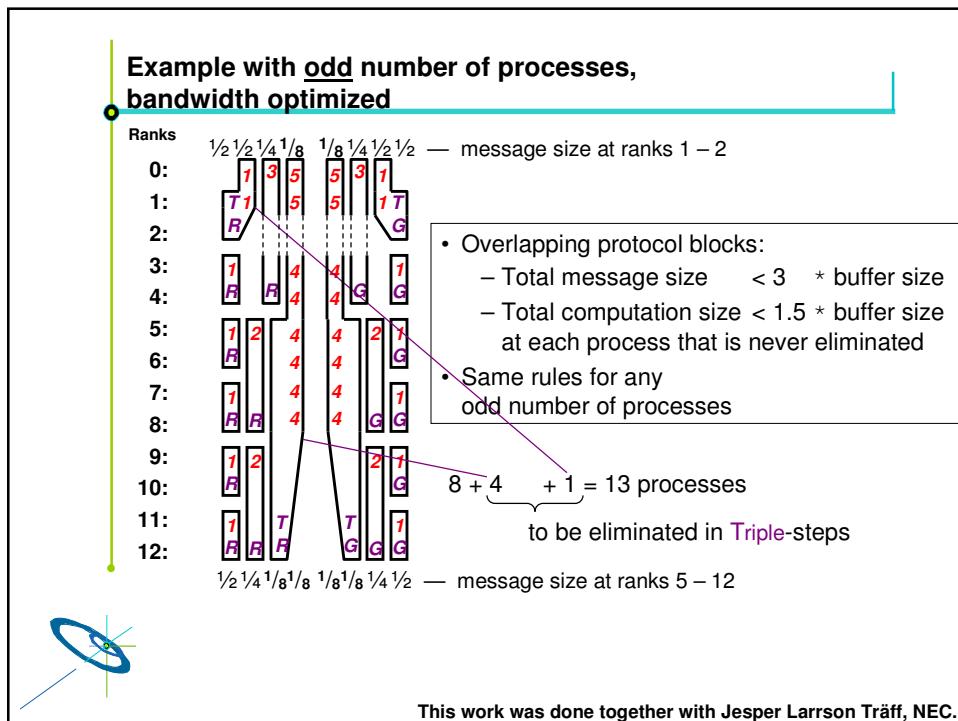
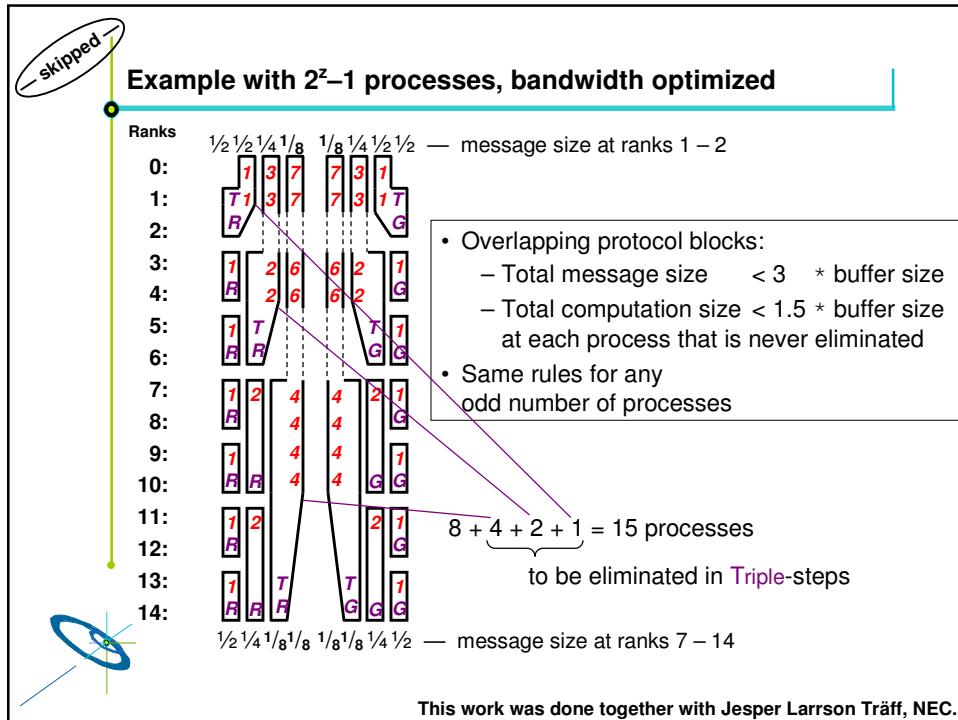
- Bandwidth optimized

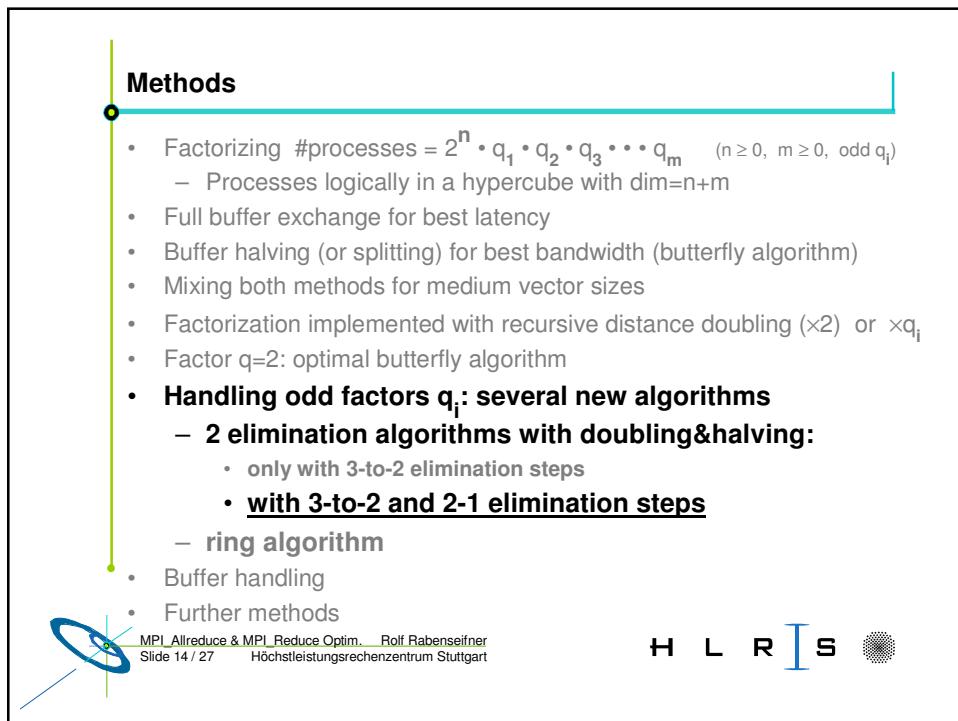
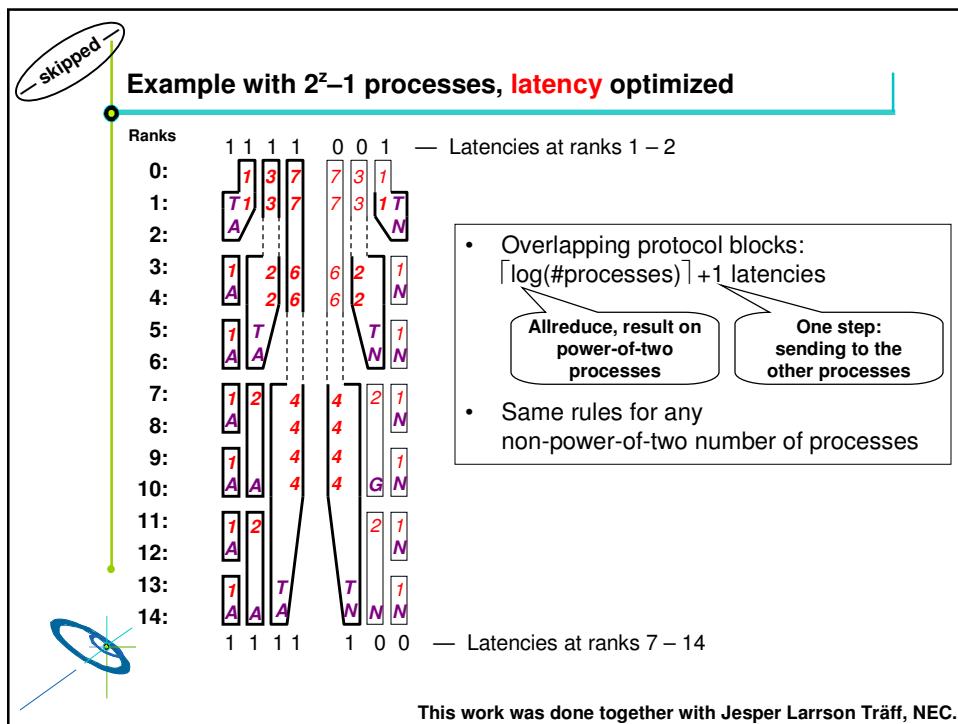


MPI_Allreduce & MPI_Reduce Optim. Rolf Rabenseifner
Slide 10 / 27 Höchstleistungsrechenzentrum Stuttgart

H L R I S

This work was done together with Jesper Larsson Träff, NEC.





2. Algorithm with Triple and Double Reduction steps

- 2-to-1-process reduction (**Double**):
 - Same protocol for “Latency optimized” and “Bandwidth optimized”
 - Allreduce/reduce_scatter epoch:

 = Send buffer → receive and reduce $\oplus \rightarrow (p_0 + p_1)$
[now eliminated]

- Allgather epoch:

 = Send buffer → receive

Example with 2^{z-1} processes, bandwidth optimized

Ranks $\frac{1}{2} \frac{1}{2} \frac{1}{4} \frac{1}{4} \frac{1}{8} \frac{1}{8} \frac{1}{4} \frac{1}{4} \frac{1}{2} \frac{1}{2}$ — message size at all active ranks

0:	
1:	
2:	
3:	
4:	
5:	
6:	
7:	
8:	
9:	
10:	
11:	
12:	
13:	
14:	

- Non-overlapping protocol blocks:
 - Total message size $< 3 * \text{buffer size}$
 - Total computation size $< 1.5 * \text{buffer size}$ at each process that is never eliminated
- Same costs as with previous algorithm
- Same rules for any odd number of processes

**Example with odd number of processes,
bandwidth optimized**

skipped

Ranks

0:	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{8}$	$\frac{1}{4}$	$\frac{1}{2}$	$\frac{1}{2}$	— message size at all active ranks
1:	1	R	2	4	4	2	1	G	
2:	1	2	R			2	1		
3:	R	2	R			2	1	G	
4:	1	D				D	1		
5:	R					1	G		
6:	1	2	2			2	1		
7:	R	2				2	1	G	
8:	1	D				D	1		
9:	R					1	G		
10:							1		
11:		T	1	R	R	G	G	1	T
12:		R						G	

- Splitting the number of processes into
 - $L \times (2\text{-process-symmetric-reduction } [R])$
 - $M \times (2 \times 2\text{-process-symmetric-reduction } [R] + 2 \times 2\text{-1-reduction } [D])$
 - $1 \times (3\text{-2-reduction } [T])$
- with $2(L+M+1) = 2^{\lfloor \log \#processes \rfloor}$

MPI_Allreduce & MPI_Reduce Optim. Rolf Rabenseifner
Slide 17 / 27 Höchstleistungsrechenzentrum Stuttgart

H L R I S

Methods

- Factorizing $\#processes = 2^n \cdot q_1 \cdot q_2 \cdot q_3 \cdots q_m$ ($n \geq 0, m \geq 0$, odd q_i)
 - Processes logically in a hypercube with dim=n+m
- Full buffer exchange for best latency
- Buffer halving (or splitting) for best bandwidth (butterfly algorithm)
- Mixing both methods for medium vector sizes
- Factorization implemented with recursive distance doubling ($\times 2$) or $\times q_i$
- Factor $q=2$: optimal butterfly algorithm
- **Handling odd factors q_i : several new algorithms**
 - **2 elimination algorithms with doubling&halving :**
 - only with 3-to-2 elimination steps
 - with 3-to-2 and 2-1 elimination steps
 - **Ring algorithm**
- Buffer handling
- Further methods

MPI_Allreduce & MPI_Reduce Optim. Rolf Rabenseifner
Slide 18 / 27 Höchstleistungsrechenzentrum Stuttgart

H L R I S

Further Optimization of odd #processes with ring-algorithms

- Latency optimization:

- Algorithm:
 - Use Bruck's "Concatenation" = Allgather algorithm (optimal, i.e., $O(\lceil \log p \rceil)$ for any #p)
 - Then reduce all data in each process

- Bandwidth optimization:

- Algorithm:
 - Divide input vector logically into (#processes) parts
 - Alltoall: send j^{th} parts from all input vectors to process j
 - Then local reduction of all j^{th} vector-parts in each process j
 - Then Bruck's Allgather

Further Optimization of odd #processes with ring-algorithms

The inner part for odd #processes can be additionally improved:

- Latency optimization:

- Algorithm:
 - Use Bruck's "Concatenation" = Allgather algorithm
 - Then reduce all data in each process

	Criterion	Ring-algorithm	Elimination-Protocol [p = #processes]
• Latency	$O(\lceil \log(p) \rceil)$ always faster		$O(\lceil \log(p) \rceil + 1)$
• Transfer	$O(p - 1)$ faster if $p=3,5$		$O(\lceil \log(p) \rceil + 1)$ faster if $p \geq 7$
• Operations	$O(p - 1)$ faster if $p=3$		$O(\lceil \log(p) \rceil)$ faster if $p \geq 5$

#processes = p	3	5	7	9	11	13	15
Latencies	$2/3$	$3/4$	$3/4$	$4/5$	$4/5$	$4/5$	$4/5$
Transfer / input_buffer_size	$2/3$	$4/4$	$6/4$	$8/5$	$10/5$	$12/5$	$14/5$
Operations / input_buffer_elements	$2/2$	$4/3$	$6/3$	$8/4$	$10/4$	$12/4$	$14/4$

Further Optimization of odd #processes with ring-algorithms

- Bandwidth optimization:**

- Algorithm:
 - Divide input vector logically into (#processes) parts
 - Alltoall: send j^{th} parts from all input vectors to process j
 - Then local reduction of all j^{th} vector-parts in each process j
 - Then Bruck's Allgather

Criterion	Ring-algorithm	Elimination-protocol	$[p = \#processes]$
Latency	$O(p-1 + \lceil \log(p) \rceil)$	$O(2 \lceil \log(p) \rceil)$	$[p' = 2^{\lfloor \log p \rfloor}]$
Transfer	$O(2(p-1)/p)$	$O(2(1.5-1/p'))$	
Operations	$O((p-1)/p)$	$O(1.5-1/p')$ i.e., $O(\text{half of transfer})$	

better

#processes = p	3	5	7	9	11	13	15
Latencies	4 / 4	7 / 6	9 / 6	12 / 8	14 / 8	16 / 8	18 / 8
Transfer per input_buffer_...	1.33 / 2.00	1.60 / 2.50	1.71 / 2.50	1.78 / 2.75	1.82 / 2.75	1.85 / 2.75	1.87 / 2.75

MPI_Allreduce & MPI_Reduce Optim. Rolf Rabenseifner
Slide 21 / 27 Höchstleistungsrechenzentrum Stuttgart

H L R I S

Best algorithm for odd factor q_i

Algorithm	Latency	+ Transfer	+ Operations
Lat.-opt. Ring	$O(\lceil \log(p) \rceil)$	$+ s(O(p-1))$	$+ O(p-1)$
Lat.-opt. Elimi.	$O(\lceil \log(p) \rceil + 1)$	$+ s(O(\lceil \log(p) \rceil + 1))$	$+ O(\lceil \log(p) \rceil)$
Bw.-opt. Elimi.	$O(2\lceil \log(p) \rceil)$	$+ s(O(3-2/p'))$	$+ O(1.5-1/p')$
Bw.-opt. Ring	$O(p-1 + \lceil \log(p) \rceil + 1^*)$	$+ s(O(2-2/p))$	$+ O(1-1/p)$

p =number of processes, s =vector size

* some additional messages to prohibit local buffer copy in Bruck's allgather algorithm

Example for $p=125$

Algorithm Latency + Transfer + Operations → Optimization for large p :

Lat.-opt. Ring	$O(7)$	$+ s(O(124))$	$+ O(124)$	→ very small vectors
Lat.-opt. Elimi.	$O(8)$	$+ s(O(8))$	$+ O(7)$	→ small vectors
Bw.-opt. Elimi.	$O(14)$	$+ s(O(3))$	$+ O(1.5)$	→ med.-long vectors
Bw.-opt. Ring	$O(132)^{**}$	$+ s(O(2))$	$+ O(1)$	→ very long vectors

** may be reduced to $O(24)$ with ring(5)*ring(5)*ring(5) instead of ring(125)

MPI_Allreduce & MPI_Reduce Optim. Rolf Rabenseifner
Slide 22 / 27 Höchstleistungsrechenzentrum Stuttgart

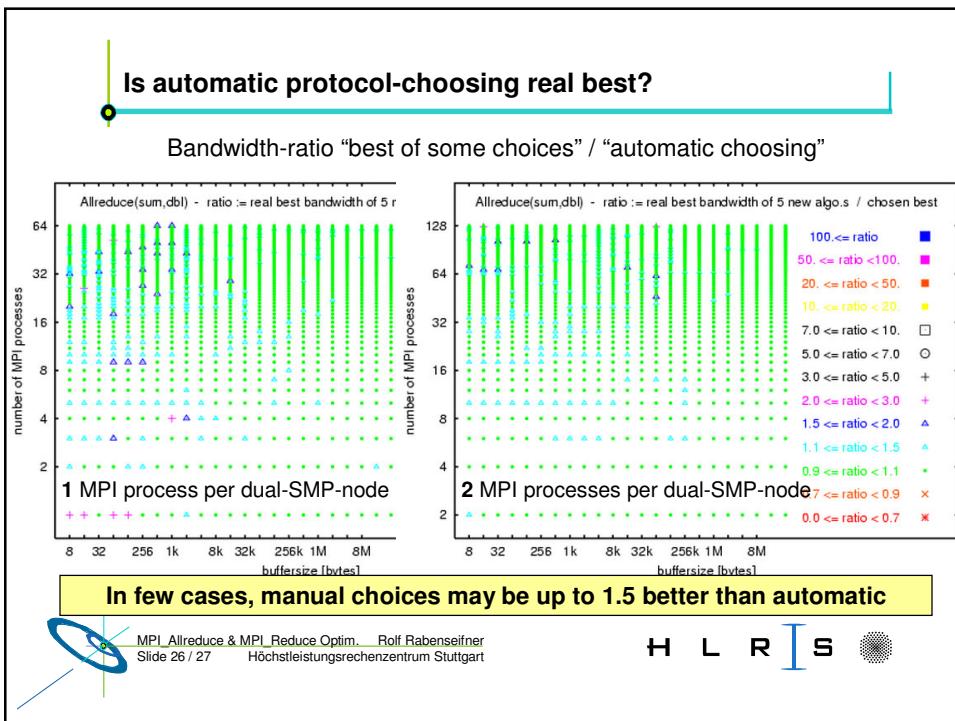
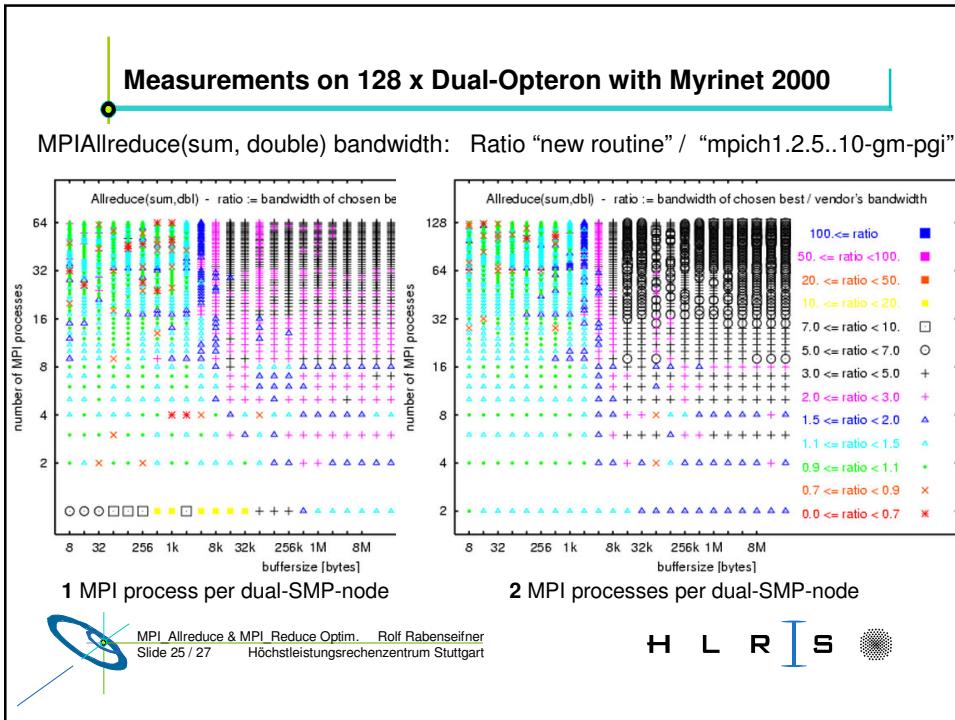
H L R I S

Buffer handling

- If operation is commute
 - No local buffer copying,
 - Except in Bruck's allgather, used in latency optimized ring (i.e., with small vectors)
- If operation is not commute (only with user defined operations)
 - Minimizing buffer copying by receiving data to optimized location

Further optimization methods

- Pipelining:
 - Splitting vectors into chunks
 - Overlaying chunk message transfer with reduction already received chunk
- Segmentation
 - For algorithms that do not use all processors at the same time
 - E.g., MPI_Reduce with binomial tree
 - Executing the operation for several segments of the input vector
 - (Recursive halving and doubling should be faster)
 - Smaller scratch buffers



Summary

- Optimized MPI_Allreduce
 - any number of processes #p
 - especially non-power-of-two #p
 - any vector length
- Optimal latency for $\#p = 2^n$ or $\#p = 2^n \cdot 3 \rightarrow O(\lceil \log_2 p \rceil)$
 - one additional round for other #p $\rightarrow O(\lceil \log_2 p \rceil + 1)$
- Bandwidth optimized: if $\#p = 2^n \cdot q$ then $\rightarrow O(2(1+1/2^{n+1})s)$
 - Instead of previous $O(4s)$ [mpich ≥ 1.2.6] between 2 and 3
- Smooth transition
 - from latency to bandwidth optimized
- Implementations
 - By NEC: (non-disclosed, 3-2-elimination)
 - By HLRS: freely available on request (2-1-elimination and ring)
(e-mail to rabenseifner@hlrs.de)