

# Effective Communication and File-I/O Bandwidth Benchmarks

## b\_eff and b\_eff\_io

Rolf Rabenseifner

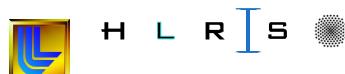
High-Performance Computing-Center Stuttgart (HLRS), University of Stuttgart,  
[rabenseifner@hlrs.de](mailto:rabenseifner@hlrs.de) [www.hlrs.de/people/rabenseifner](http://www.hlrs.de/people/rabenseifner)

Alice E. Koniges

Lawrence Livermore National Laboratory,  
[koniges@llnl.gov](mailto:koniges@llnl.gov) [www.rzg.mpg.de/~ack](http://www.rzg.mpg.de/~ack)

European PVM / MPI Users' Group Meeting 2001  
Santorini, Greece, April 18-20, 2002 (postponed from Sep. 23-26, 2001)

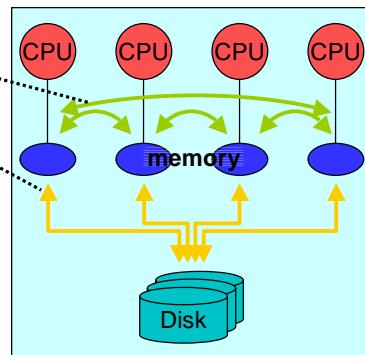
Communication & I/O Benchmarks  
Slide 1 Hochleistungsrechenzentrum Stuttgart



## Effective Communication & I/O Bandwidth Benchmarks

### Goals

- Parallel Communication Benchmark
- Parallel File-I/O Benchmark
  - each process is involved!
- Detailed insight
  - bandwidth experiments of several
    - I/O or communication patterns
    - chunk or message sizes
- One characteristic value
  - based on experiments above
  - averaging
- Appropriate execution time for rapid benchmarking



Communication & I/O Benchmarks  
Slide 2 of 23 Rolf Rabenseifner  
Hochleistungsrechenzentrum Stuttgart



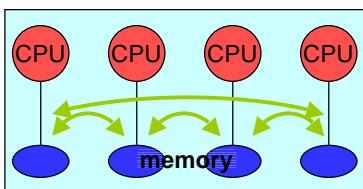
# **b\_eff**

the  
**effective communication bandwidth**  
benchmark

Communication & I/O Benchmarks Rolf Rabenseifner  
Slide 3 of 23 Höchstleistungsrechenzentrum Stuttgart



## **Definition of the Effective Communication Bandwidth Benchmark: b\_eff**



- www.hlrs.de/mpi/b\_eff/
- Authors: Karl Solchenbach, Hans-Joachim Plum, and Gero Ritzenhoefer (Pallas), Rolf Rabenseifner (HLRS)

- 6 ring patterns
- 30 random patterns
- 13 additional patterns
- 21 message sizes
- 3 communication methods
- 3 times repeated
- automatically controlled measurement-loop length, i.e., time driven approach
- 5 - 20 msec / experiment → benchmark completes in **a few minutes**

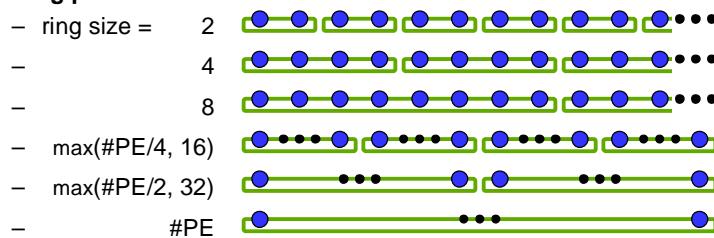
$$(6+30+13) \times 21 \times 3 \times 3 = \text{9261 experiments}$$

Communication & I/O Benchmarks Rolf Rabenseifner  
Slide 4 of 23 Höchstleistungsrechenzentrum Stuttgart



## Definition of $b_{eff}$ — communication patterns and sizes

- 6 ring patterns



- 30 random patterns 

- 21 message sizes

- 1, 2, 4, 8, 16, 32, 64, 128, 256, 512 byte, 1kB, 2kB, (12 sizes)
- 9 logarithmic equidistant sizes: 4kB, ...,  $L_{max}$  = memory per PE / 128

Communication & I/O Benchmarks Rolf Rabenseifner  
Slide 5 of 23 Höchstleistungsrechenzentrum Stuttgart



## Definition of $b_{eff}$ — averaging

One characteristic accumulated communication bandwidth number  
 $\quad :=$  average bandwidth on several communication patterns  
 $\quad \quad$  average on different message sizes  
 $\quad \quad$  maximum over different MPI programming methods

$$b_{eff} = \text{logavg}(\text{logavg}_{\text{ringpat}}(\text{avg}_L(\text{max}_{\text{method}}(\text{max}_{\text{rep}}(b_{\text{pat},L,\text{method},\text{rep}})))), \text{logavg}_{\text{randompat}}(\text{avg}_L(\text{max}_{\text{method}}(\text{max}_{\text{rep}}(b_{\text{pat},L,\text{method},\text{rep}})))))$$

with

- $b_{\text{pat},L,\text{method},\text{rep}}$  = accumulated bandwidth of each experiment  
 $\quad \quad \quad$  over all processes
- methods: MPI\_Sendrecv, MPI\_Alltoallv, and nonblocking Irecv&Isend&Waitall
- pat & L: patterns and message sizes, see previous slide
- rep: repetition number = 1..3
- avg: arithmetic mean
- logavg: geometric mean

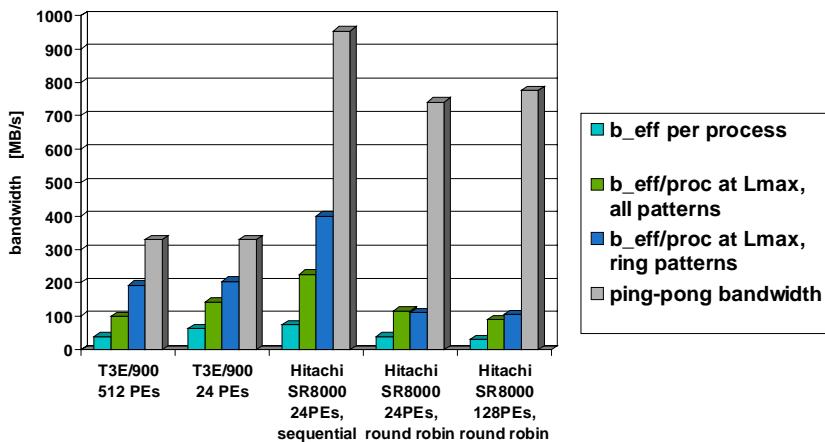
Communication & I/O Benchmarks Rolf Rabenseifner  
Slide 6 of 23 Höchstleistungsrechenzentrum Stuttgart



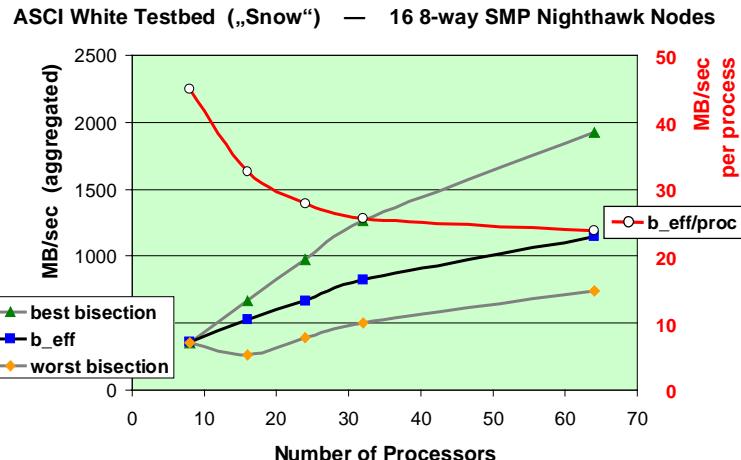
## Features of Effective Bandwidth benchmark

- Based on MPI, source code is available
- Measures total architecture, not only point-to-point
- Checks performance of architecture and not the quality of the MPI implementation
- Suited for MPP-architectures and clusters
- Runs on any number of processors
- Results are easy to understand
- Generates a single number  $b_{eff}$  (like LINPACK  $R_{max}$ )
- Aggregate bandwidth of the total system

## $b_{eff}$ per process and ping-pong bandwidth



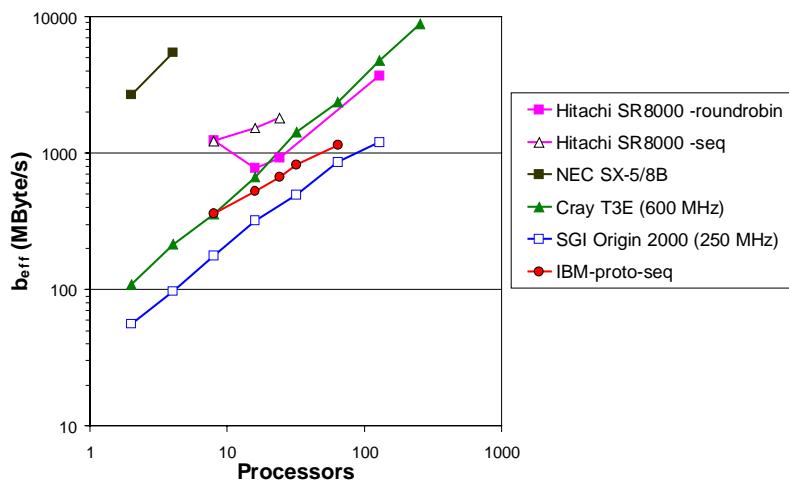
**$B_{eff}$  is monotonic.  $B_{eff}/proc$  is roughly constant indicating scalable balance (see next slide).**



Communication & I/O Benchmarks Rolf Rabenseifner  
Slide 9 of 23 Höchstleistungsrechenzentrum Stuttgart



### **$B_{eff}$ Scaling: current systems**

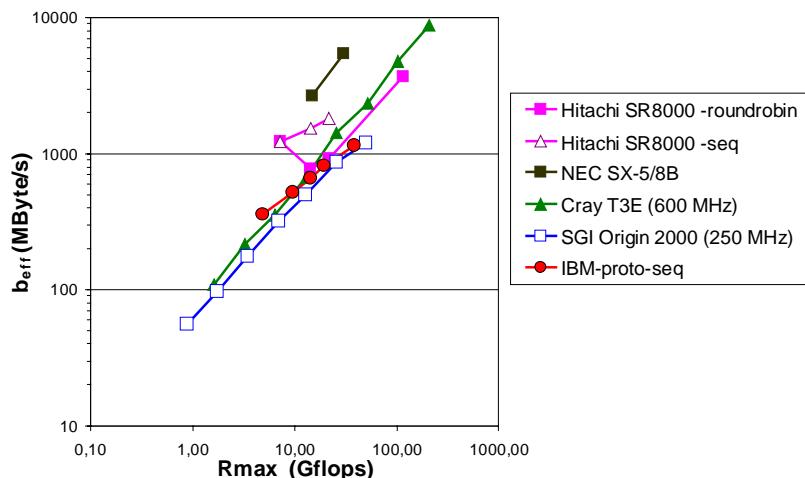


Communication & I/O Benchmarks Rolf Rabenseifner  
Slide 10 of 23 Höchstleistungsrechenzentrum Stuttgart



Effective Communication and File-I/O Bandwidth Benchmarks:  $b_{eff}$  and  $b_{eff\_io}$   
Euro PVM/MPI 2001 – April 18-20, 2002 (postponed from Sep. 23-26, 2001)

### B\_eff Scaling: current systems



Communication & I/O Benchmarks Rolf Rabenseifner

Slide 11 of 23 Höchstleistungsrechenzentrum Stuttgart



**$b_{eff\_io}$**   
the  
effective MPI-I/O bandwidth  
benchmark

Communication & I/O Benchmarks Rolf Rabenseifner

Slide 12 of 23 Höchstleistungsrechenzentrum Stuttgart



Effective Communication and File-I/O Bandwidth Benchmarks:  $b_{eff}$  and  $b_{eff\_io}$   
Euro PVM/MPI 2001 – April 18-20, 2002 (postponed from Sep. 23-26, 2001)

### Goals for b\_eff\_io

- portable
- characterizing parallel I/O capabilities of a system
- many patterns
- ... that can be optimized
- going beyond the caching capabilities  
--> measuring the real disk I/O

Rule: Balanced HPC systems should be able to write the total memory in 10 minutes to disk

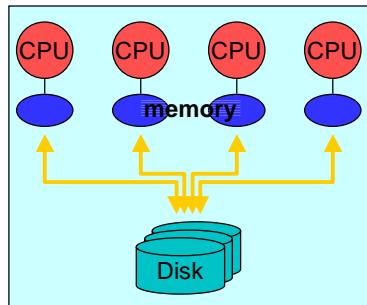
**==> An I/O benchmark should not need hours!  
— 10 minutes may be enough to overrun any cache!**

- time driven approach / automatically controlled repetition factors
- rapid benchmarking (30-60 min.)

### Starting-Points — the I/O Parameter Space

- How to define and measure one characteristic I/O bandwidth value?
  - The I/O parameter space — 20 orthogonal parameters:
    - Application parameters:
      - (a) the size of contiguous chunks in the memory, (b) on disk, (c) ... (f)
    - Usage aspects:
      - (a) how many processes are used
      - (b) how many parallel processors and threads are used for each process.
    - I/O interface:
      - (a) Posix I/O buffered or (b) raw,
      - (c) special filesystem I/O of the vendor filesystem,
      - (d) MPI-I/O.
    - MPI-I/O aspects:
      - (a) access methods, i.e., first writing of a file, rewriting or reading, (b) ...
      - (c) coordination, i.e., collectively or noncollectively, (d) ... (f)
    - Filesystem parameters:
      - (a) which filesystem is used,
      - (b) how many nodes are used as I/O servers, (c) ... (f)
- (full list, see paper)

## Definition of the Effective File-I/O Bandwidth Benchmark: **b\_eff\_io**

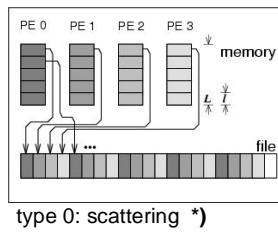


- 5 I/O patterns
- 7 chunk sizes
- 3 accesses  
(initial write, rewrite, read)
- 3 compute partition sizes  
(number of parallel benchmark processes)

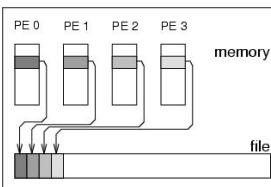
**315 different measurements**

- benchmark completes in **30-60 minutes**
- [www.hlrs.de/mpi/b\\_eff\\_io/](http://www.hlrs.de/mpi/b_eff_io/)

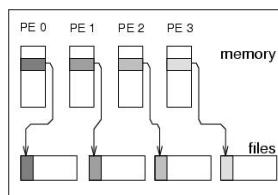
## Definition of b\_eff\_io — the Pattern Types



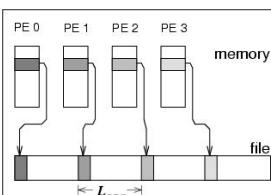
type 0: scattering \*)



type 1: strided coll.



type 2: noncoll., separated



type 3 noncoll. / 4 coll. segmented

Pattern that can be optimized

### Chunk sizes on disk:

- $L_{\max} = \max(2\text{MB}, \text{memory of one node}/128)$  \*)
- **wellformed:**  
1MB, \*,  
32 kB,  
1 kB,
- **non-wellformed:**  
1MB+8B, \*)  
32 kB+8B,  
1kB+8B

\*) double weighted

## Definition of b\_eff\_io

(Release 1.0)

$b_{eff\_io}$  := Maximum over all usage and filesystem parameters } manually  
Average on write, rewrite, read      } automatically,  
Average on five access pattern types      } in time  
Average on several chunk size values\*)      }  $T=30$  min.  
of measured bandwidth

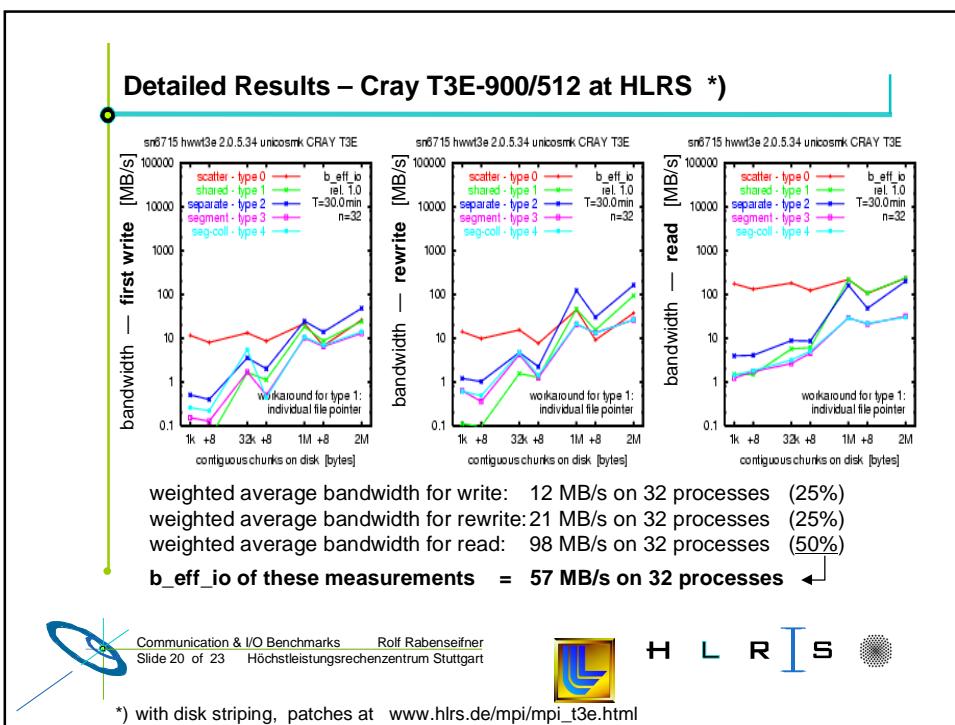
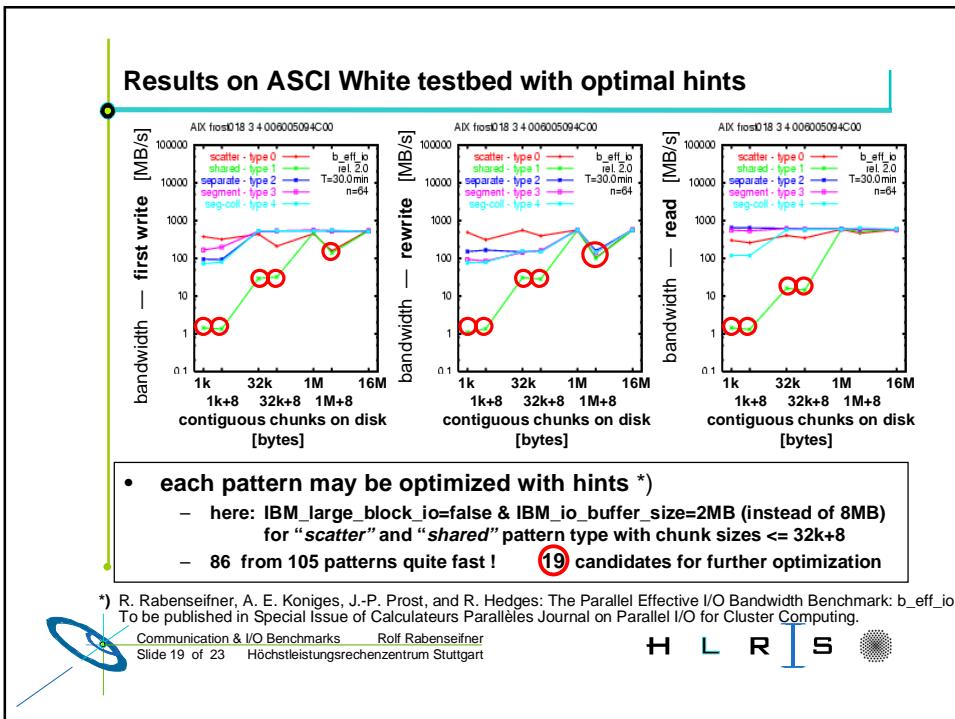
\*) defines the size of contiguous chunks written to disk and the contiguous chunk in memory written by each MPI call

## Output of the b\_eff\_io benchmark program

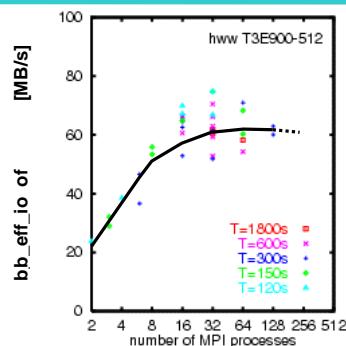
- the  $b_{eff\_io}$  value

```
weighted average bandwidth for write : 422.388 MB/s on 64 processes
weighted average bandwidth for rewrite : 331.833 MB/s on 64 processes
weighted average bandwidth for read : 473.441 MB/s on 64 processes
Total amount of data written/read with each access method: 208558.575 MBytes
= 39.8 percent of the total memory (524288 MBytes)
b_eff_io of these measurements = 451.727 MB/s
on 64 processes with 2048 MByte/PE, scheduled time=30.0 Min,
on AIX frost018 3 4 006005094C00
total memory / b_eff_io = 524288 Mbytes / 451.727 MB/s = 19.3 min.
```

- detailed results
  - as ASCII table
  - one page with 3+5 plots
    - all measurements sorted by access: write / rewrites / reads
    - and same sorted by pattern types: type-0 / type-1 / type-2 / type-3 / type-4

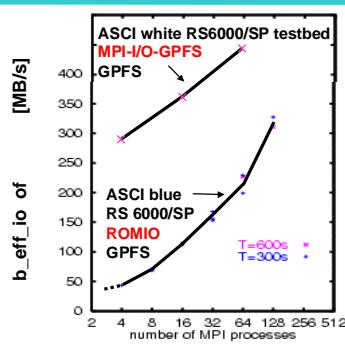


## First Results — Comparing b\_eff\_io (#processes)



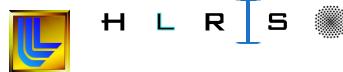
### Cray T3E:

- reached already with 8 processors!
- => optimal for any load:  
many small jobs ... one large job



### IBM SP:

- with optimized MPI-I/O-GPFS
  - high bandwidth for any load
- with ROMIO
  - I/O capability bound to application
  - full bandwidth only with many processes



## Acknowledgements

Thanks to  pallas.

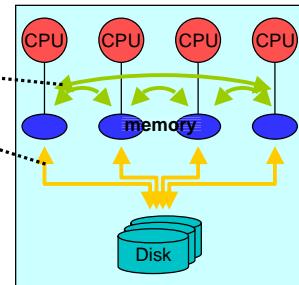
They initiated this project with their bi-section based b\_eff benchmark.

Work by Lawrence Livermore National Laboratory  
is performed under the auspices of the U.S. Department  
of Energy by the University of California  
under Contract W-7405-ENG-48, UCRL-VG-143637.

## Summary

Two parallel benchmarks:

- $b_{eff}$ : Communication
- $b_{eff\_io}$ : File-I/O
- Detailed insight
  - many patterns, chunk sizes, ...
- One characteristic value
  - averaging
- Appropriate execution time for rapid benchmarking
  - time driven approach
  - $b_{eff}$ : 3–5 min.
  - $b_{eff\_io}$ : 30–60 min.
- ASCI white compared with ...



### Further information:

[www.hlrs.de/mpi/b\\_eff](http://www.hlrs.de/mpi/b_eff) & .../[b\\_eff\\_io](http://www.hlrs.de/mpi/b_eff_io)  
[www.hlrs.de/mpi/mpi\\_t3e.html](http://www.hlrs.de/mpi/mpi_t3e.html) & .../[ufs\\_t3e](http://www.hlrs.de/mpi/ufs_t3e)  
[www.hlrs.de/people/rabenseifner/publ/publications.html](http://www.hlrs.de/people/rabenseifner/publ/publications.html)

Communication & I/O Benchmarks Rolf Rabenseifner  
Slide 23 of 23 Höchstleistungsrechenzentrum Stuttgart

