

# Network Bandwidth Measurements and Ratio Analysis with the HPC Challenge Benchmark Suite (HPCC)

Rolf Rabenseifner, Sunil R. Tiyyagura, and Matthias Müller  
[rabenseifner@hlrs.de](mailto:rabenseifner@hlrs.de), [sunil@hlrs.de](mailto:sunil@hlrs.de), [mueller@hlrs.de](mailto:mueller@hlrs.de)

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

EuroPVM/MPI'05  
Sorrento, Italy, Sep. 18-21, 2005  
(HPCC data status Sep. 8, 2005)



Balance / HPC Challenge Benchmark

Slide 1

Höchstleistungsrechenzentrum Stuttgart

H L R I S



## Balance Analysis with HPC Challenge Benchmark Data

- How HPCC data can be used to analyze the balance of HPC systems
  - Details on ring based benchmarks
- Resource based ratios
  - Inter-node bandwidth and
  - memory bandwidth
  - versus computational speed
- HPCC footprint
  - Comparing the platforms



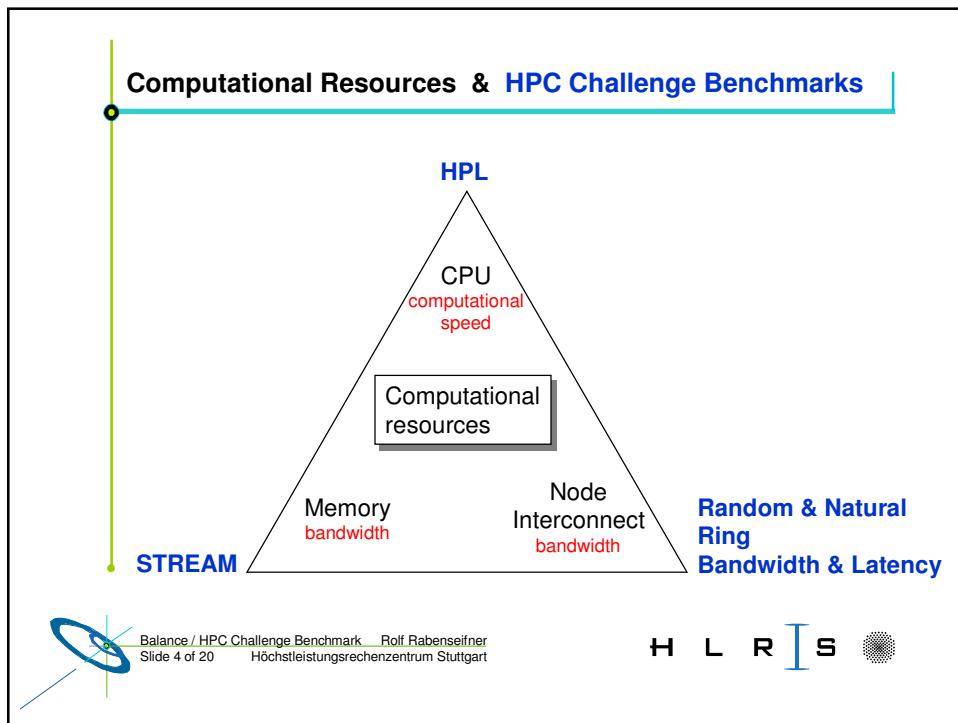
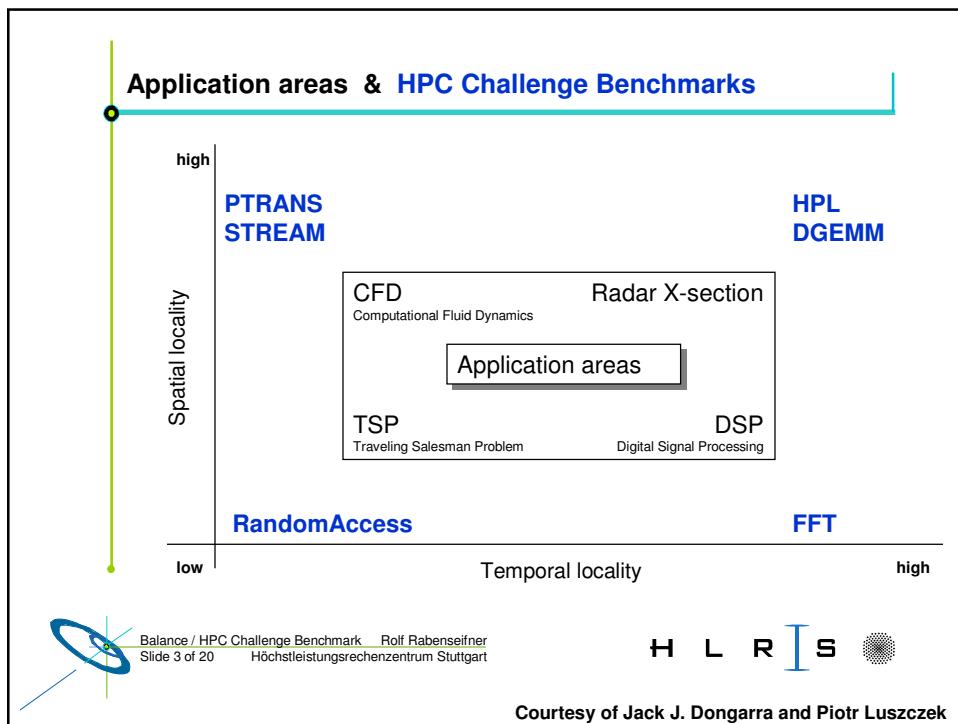
Balance / HPC Challenge Benchmark

Slide 2 of 20

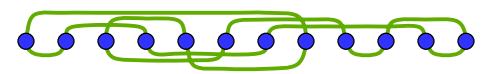
Rolf Rabenseifner  
Höchstleistungsrechenzentrum Stuttgart

H L R I S





## Random & natural ring bandwidth & latency

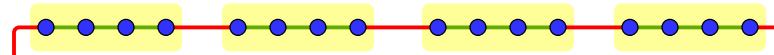
- Parallel communication pattern on all MPI processes (●)
  - Natural ring
  - Random ring
- Bandwidth per process
  - Accumulated message size / wall-clock time / number of processes
  - On each connection messages in both directions
  - With `2xMPI_Sendrecv` and `MPI non-blocking` → best result is used
  - Message size = 2,000,000 bytes
- Latency
  - Same patterns, message size = 8 bytes
  - Wall-clock time / (number of sendrecv per process)

Balance / HPC Challenge Benchmark Rolf Rabenseifner  
Slide 5 of 20 Hochleistungsrechenzentrum Stuttgart

H L R I S

## Inter-node bandwidth on clusters of SMP nodes

- Natural Ring:
  - Only one incoming and one outgoing message per node at the same time
  - Accumulated bandwidth := bandwidth per process × #processes
  - Reflects the communication bandwidth in the **first** dimension of a Cartesian domain decomposition
  - Does **not** reflect the accumulated inter-node bandwidth (nor the bi-section bandwidth) of an HPC system

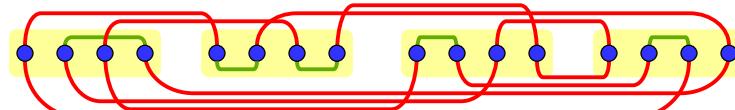


Balance / HPC Challenge Benchmark Rolf Rabenseifner  
Slide 6 of 20 Hochleistungsrechenzentrum Stuttgart

H L R I S

## Inter-node bandwidth on clusters of SMP nodes

- Random Ring
  - Reflects the other dimension of a Cartesian domain decomposition and
  - Communication patterns in unstructured grids
  - Some connections are inside of the nodes
  - Most connections are inter-node
  - Depends on #nodes and #MPI processes per node

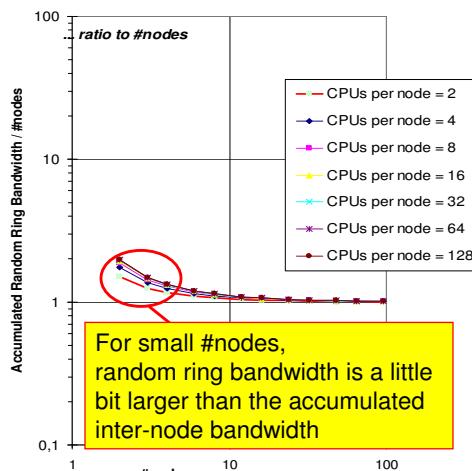
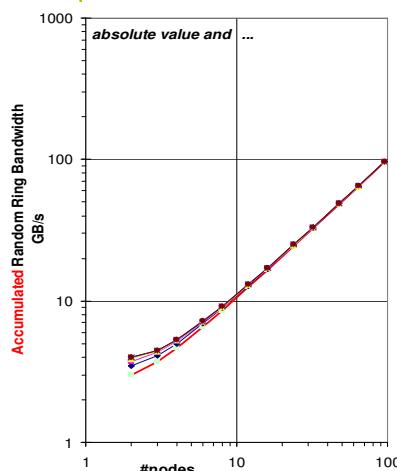


- Accumulated bandwidth  
:= bandwidth per process  $\times$  #processes
- $\approx$  accumulated inter-node bandwidth  $\times (1 - 1 / \#nodes)^{-1}$

Balance / HPC Challenge Benchmark Rolf Rabenseifner  
Slide 7 of 20 Hochleistungsrechenzentrum Stuttgart

H L R I S

## On an ideally switched cluster ...



Balance / HPC Challenge Benchmark Rolf Rabenseifner  
Slide 8 of 20 Hochleistungsrechenzentrum Stuttgart

H L R I S

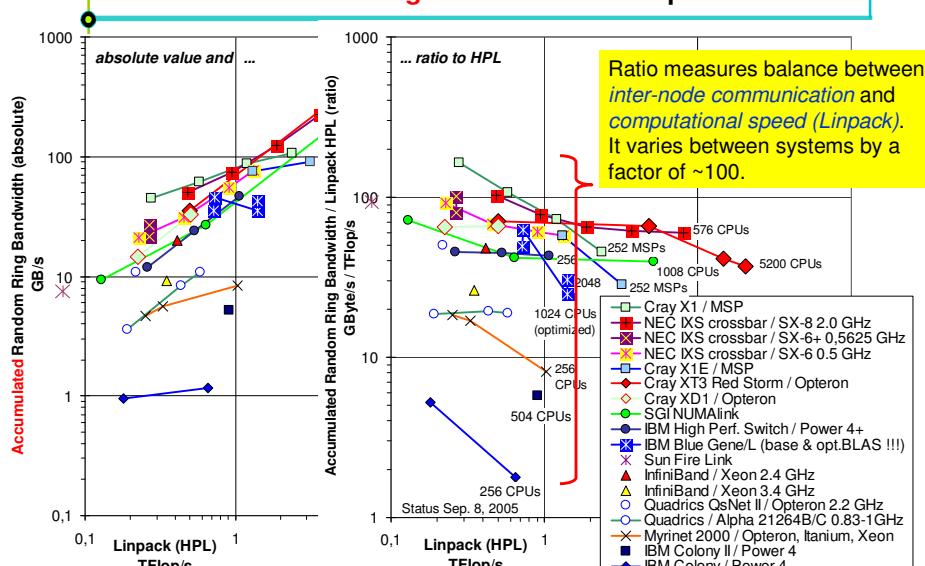
## Balance Analysis with HPC Challenge Benchmark Data

- Balance can be expressed as a set of ratios
  - e.g., accumulated memory bandwidth / accumulated Tflop/s rate
- Basis
  - Linpack (HPL) → Computational Speed
  - Random Ring Bandwidth → Inter-node communication
  - Parallel STREAM Copy or Triad → Memory bandwidth
- Be careful:
  - Some data are presented for the **total system**
  - Some per **MPI process** (HPL processes)
  - i.e., balance calculation always
    - with accumulated data on total system, or
    - with divided data to one MPI process

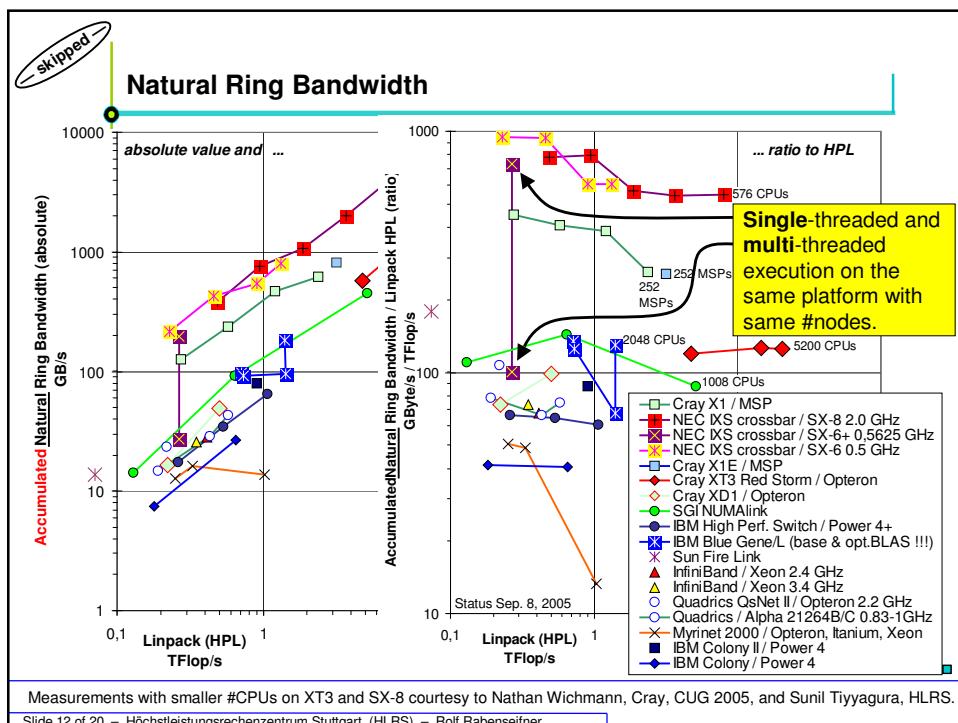
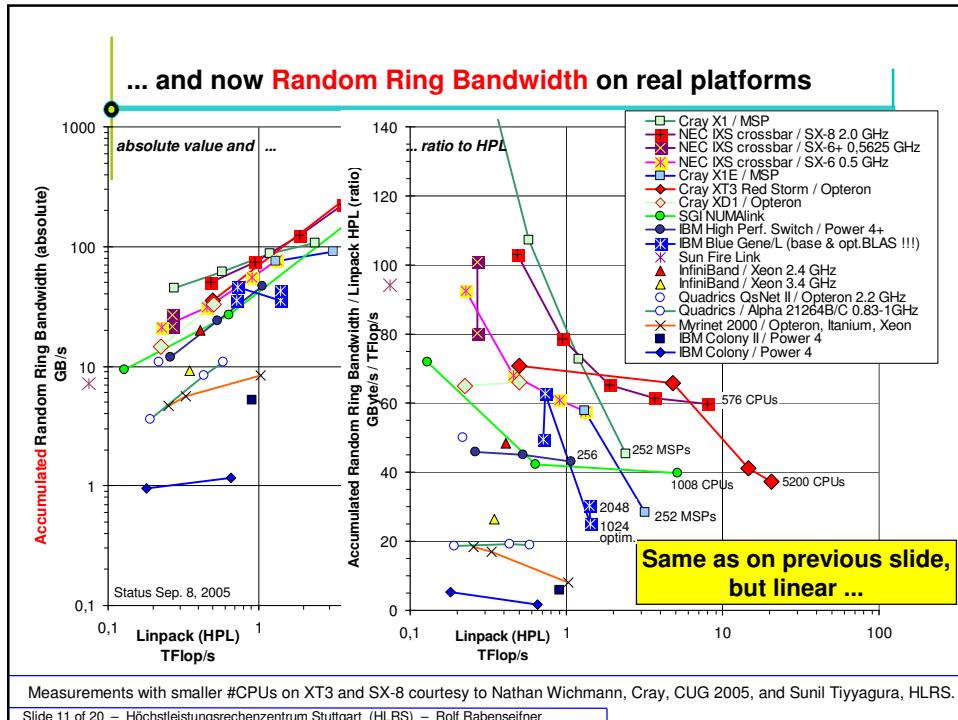
Balance / HPC Challenge Benchmark Rolf Rabenseifner  
Slide 9 of 20 Hochleistungsrechenzentrum Stuttgart

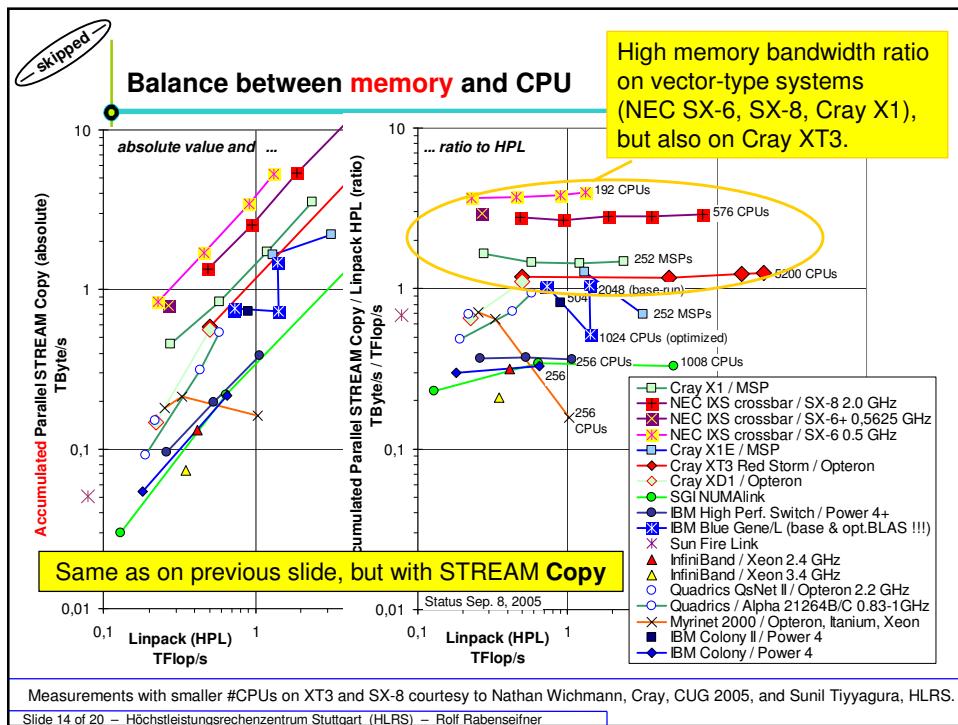
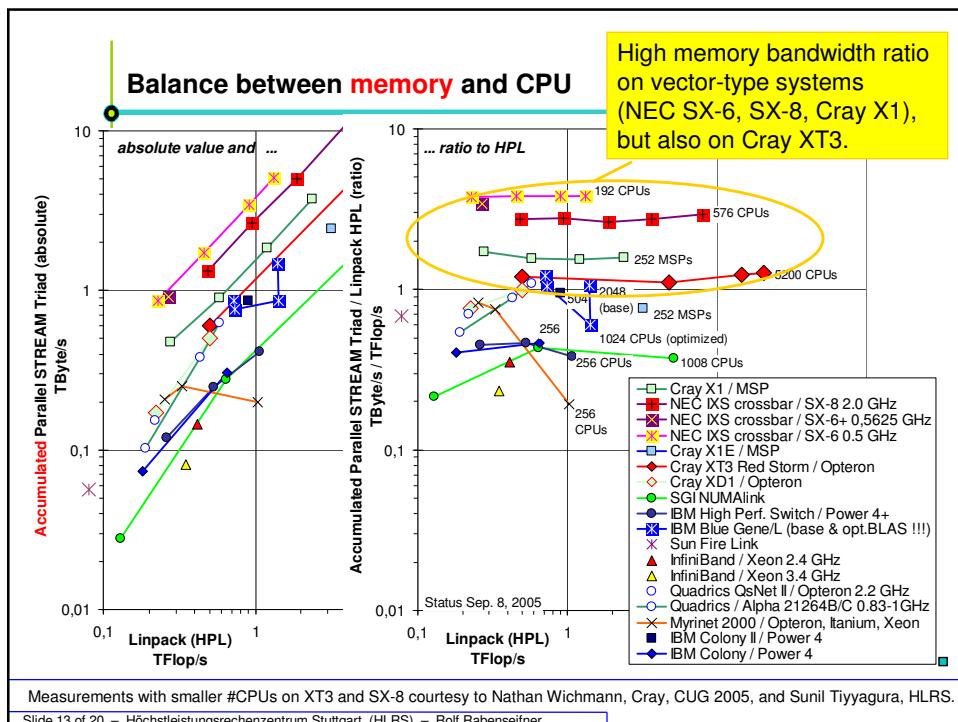
H L R I S

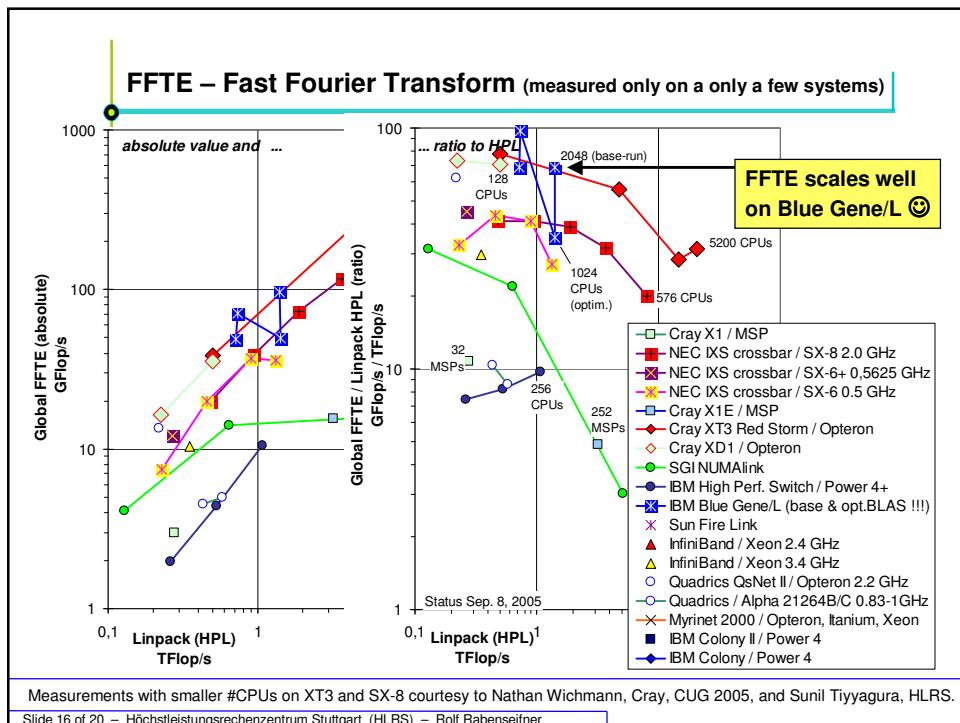
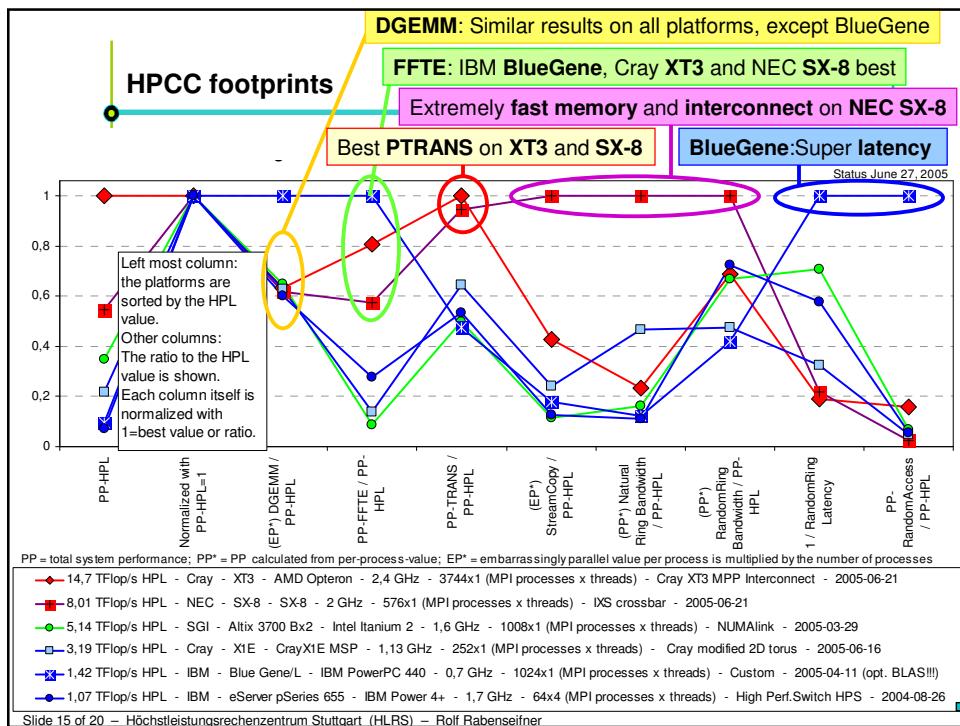
## ... and now Random Ring Bandwidth on real platforms

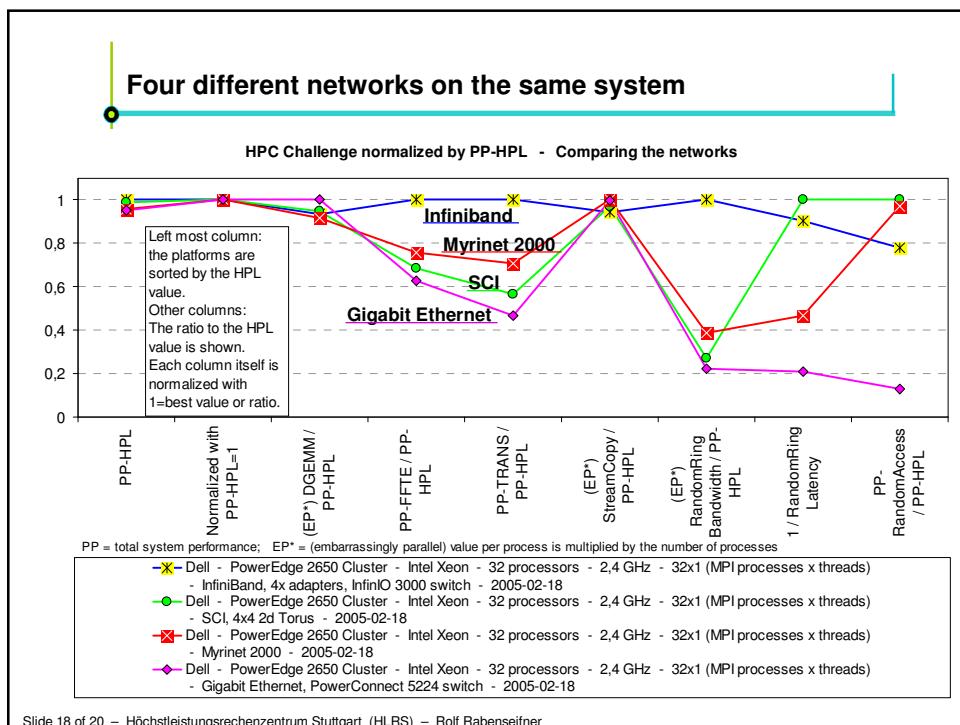
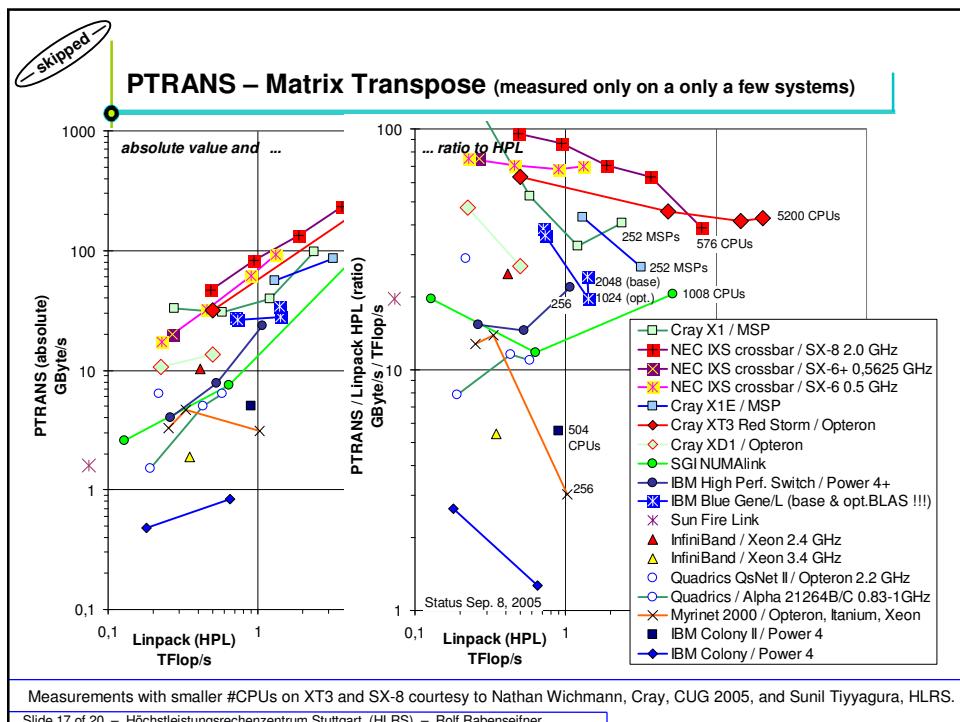


Slide 10 of 20 – Hochleistungsrechenzentrum Stuttgart (HLRS) – Rolf Rabenseifner









## Acknowledgments

- Thanks to
  - all persons and institutions that have uploaded HPCC results.
  - Jack Dongarra and Piotr Luszczek for inviting me into the HPCC development team.
  - Matthias Müller, Sunil Tiyyagura and Holger Berger for benchmarking on the SX-8 and SX-6 and discussions on HPCC.
  - Nathan Wichmann from Cray for Cray XT3 and X1E data.
  - David Koester for his helpful remarks on the HPCC Kiviat diagrams.

## Conclusions

- HPCC is an interesting basis for
  - benchmarking computational resources
  - analyzing the balance of a system
  - scaling with the number of processors
  - with respect to application needs
- HPCC helps to show the strength of
  - vector systems
  - cluster networks

See also HPCC award at SC'06  
<http://icl.cs.utk.edu/hpcc/>  
→ Award