



THE SUPERCOMPUTER COMPANY

Performance Lessons from the Cray XT3

Jeff Larkin

Cray, Inc.

larkin@cray.com

2/13/2006

The 7th LCI International Conference on Clusters

EXPLORE SIMULATE CREATE



Overview

- Designing a HPC system
- Cray XT3: About the Parts
 - AMD Opteron Processor
 - Cray SeaStar Network
 - Unicos/lc Kernel
- Benchmark Results
- Application Results

Designing a System

- Architectural decisions must be made very early
 - Processor
 - Network
 - OS
- Early decisions are often difficult/expensive to change
- Months or years may pass between initial design decisions and final results

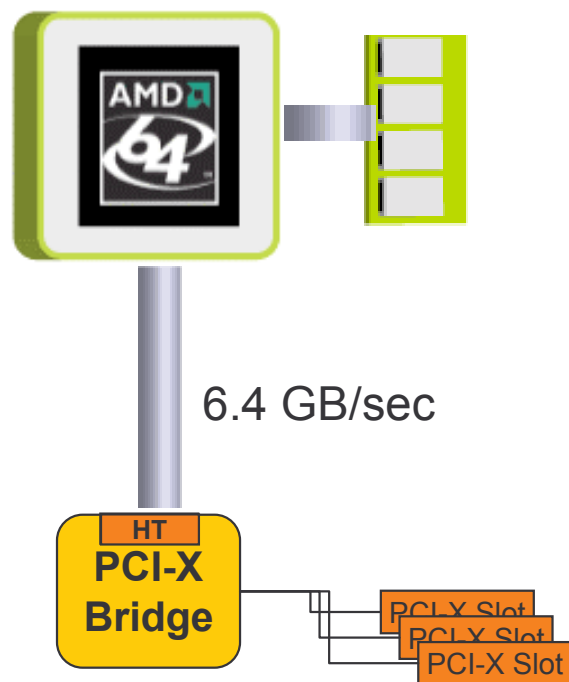
AMD Opteron: The Basics

- 64-bit x86 Architecture
- 128k Fully-associative L1 Cache
- 1MB 4-way Associative L2 Cache
- Integrated Memory Controller
 - No Northbridge
 - Low Memory Latency
- HyperTransport
 - High-bandwidth from the processor
 - Open Standard

AMD Opteron: The Basics

- 64-bit x86 Architecture
- 128k Fully-associative L1 Cache
- 1MB 4-way Associative L2 Cache
- Integrated Memory Controller
 - No Northbridge
 - Low Memory Latency
- HyperTransport
 - High-bandwidth from the processor
 - Open Standard

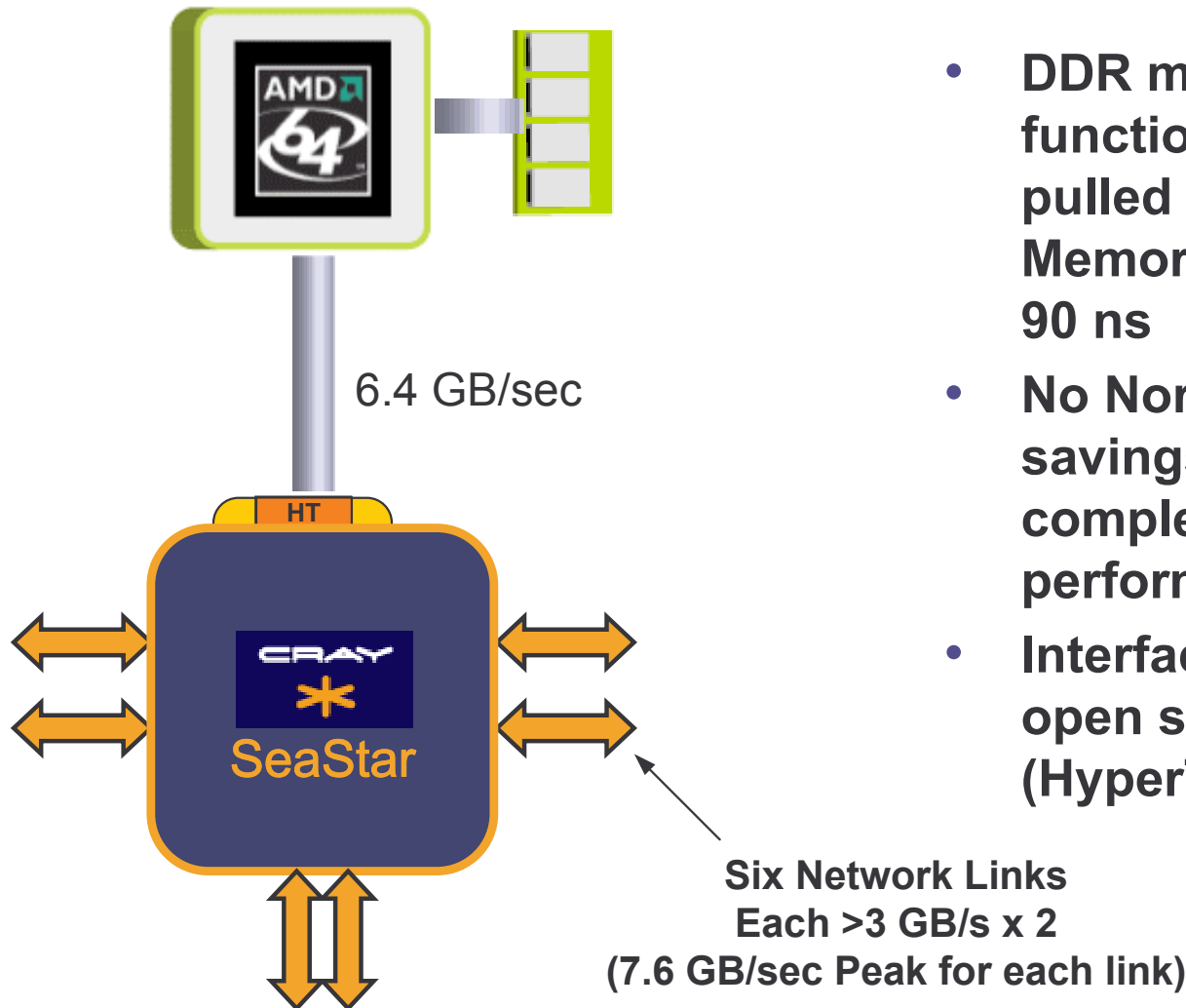
AMD Opteron: Generic System



- DDR memory controller and function of Northbridge is pulled onto the Opteron die. Memory latency reduced to 60-90 ns
- No Northbridge chip results in savings in heat, power, complexity and an increase in performance
- Interface off the chip is an open standard (HyperTransport)

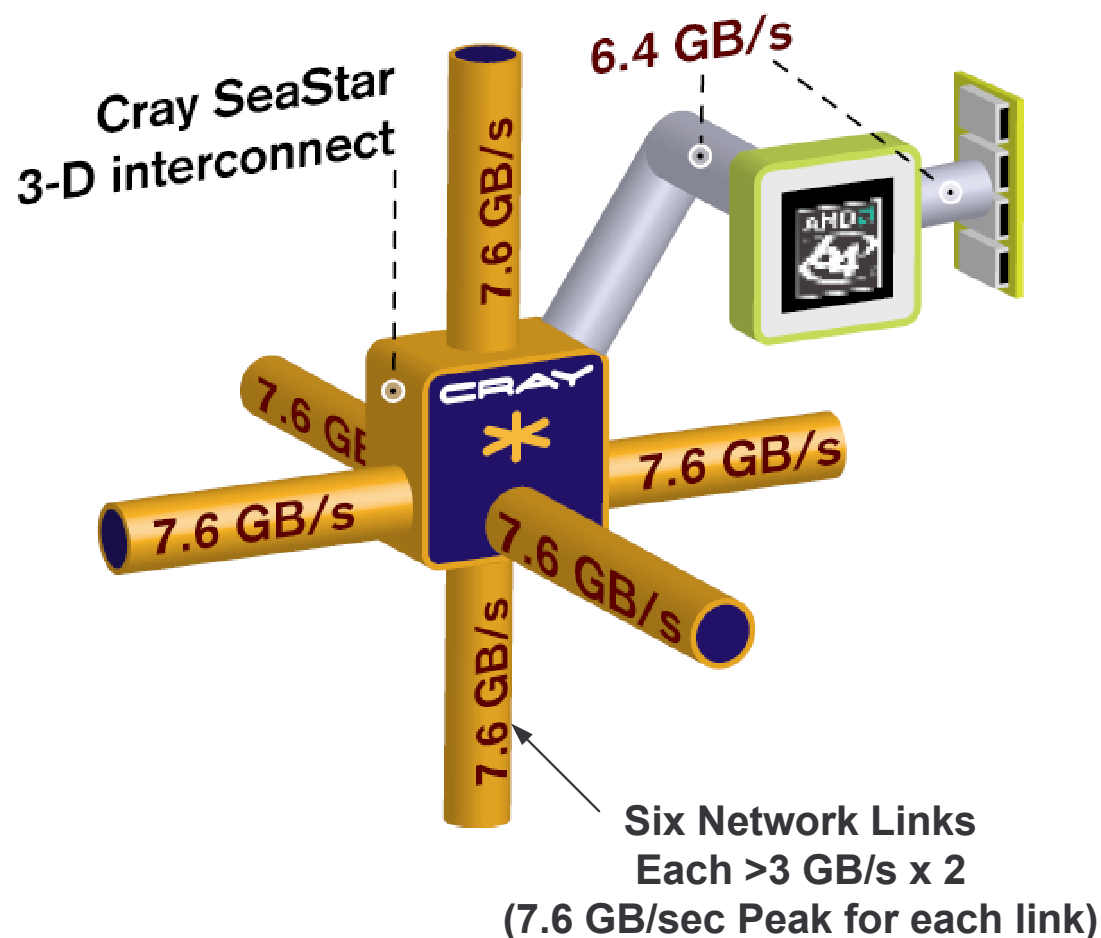
AMD Opteron: Generic System

CRAY XT3 PE



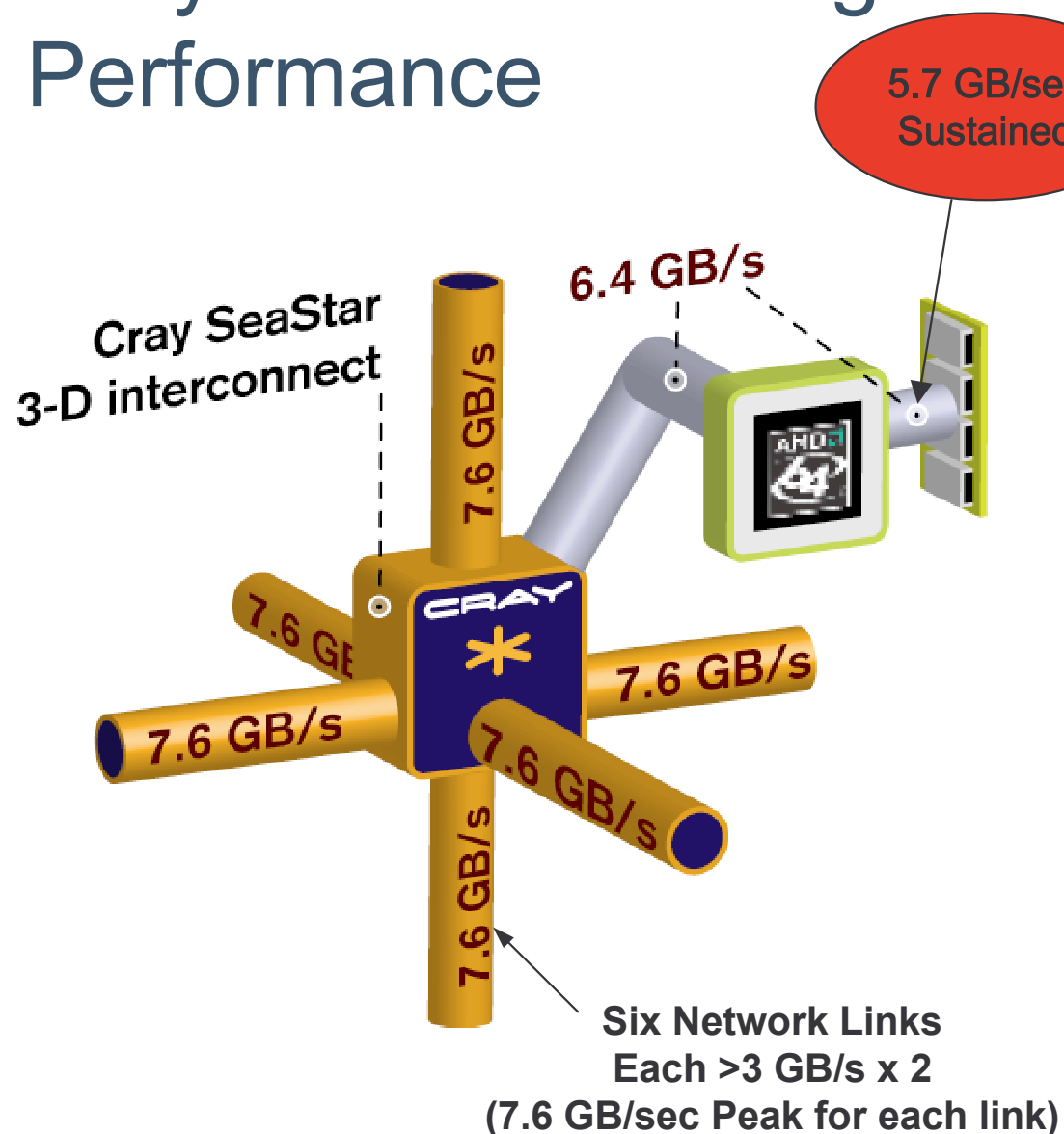
- DDR memory controller and function of Northbridge is pulled onto the Opteron die. Memory latency reduced to 60-90 ns
- No Northbridge chip results in savings in heat, power, complexity and an increase in performance
- Interface off the chip is an open standard (HyperTransport)

Cray XT3 Processing Element: Measured Performance



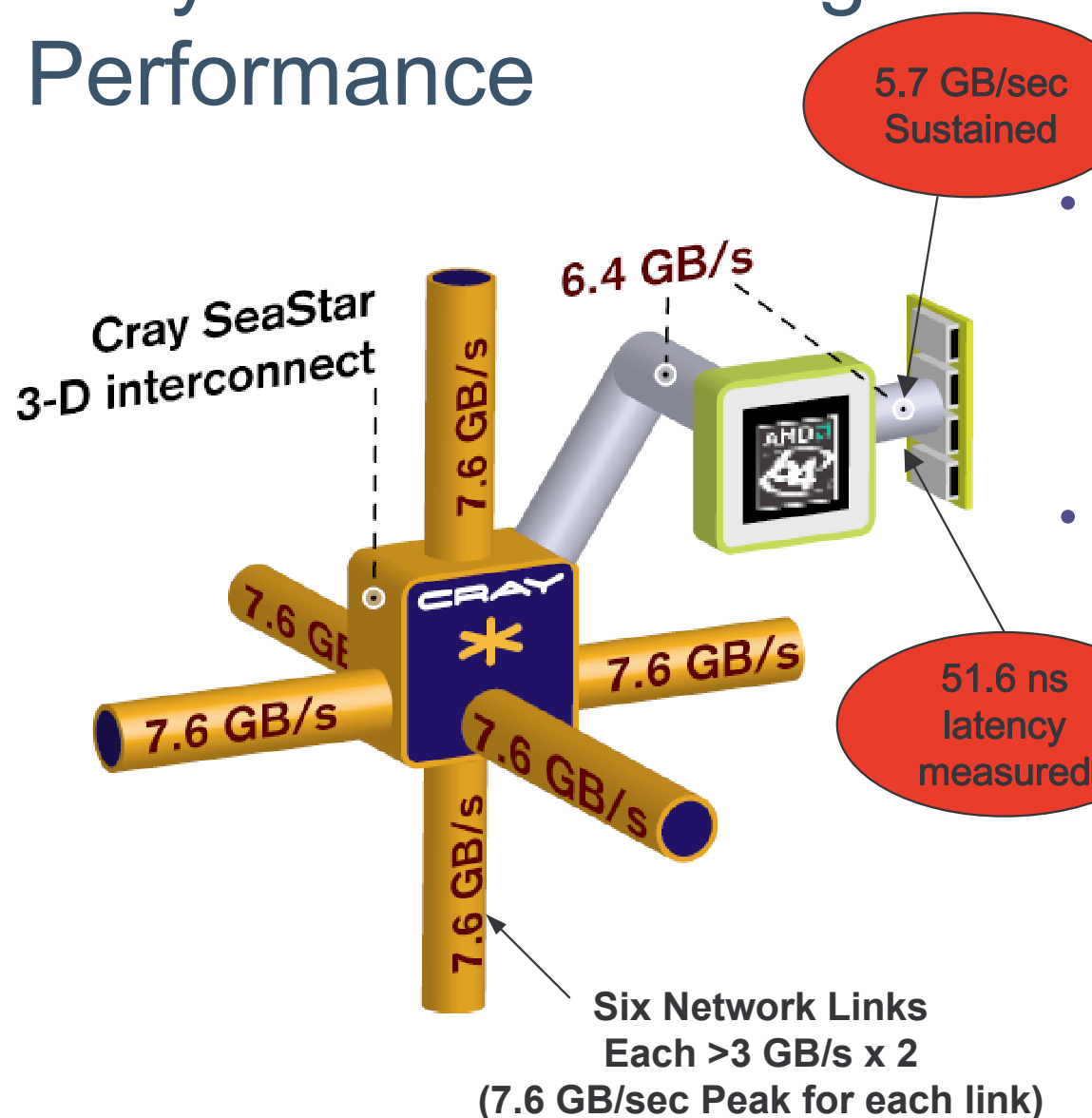
- DDR memory controller and function of Northbridge is pulled onto the Opteron die. Memory latency reduced to <60 ns
- No Northbridge chip results in savings in heat, power, complexity and an increase in performance
- Interface off the chip is an open standard (HyperTransport)

Cray XT3 Processing Element: Measured Performance



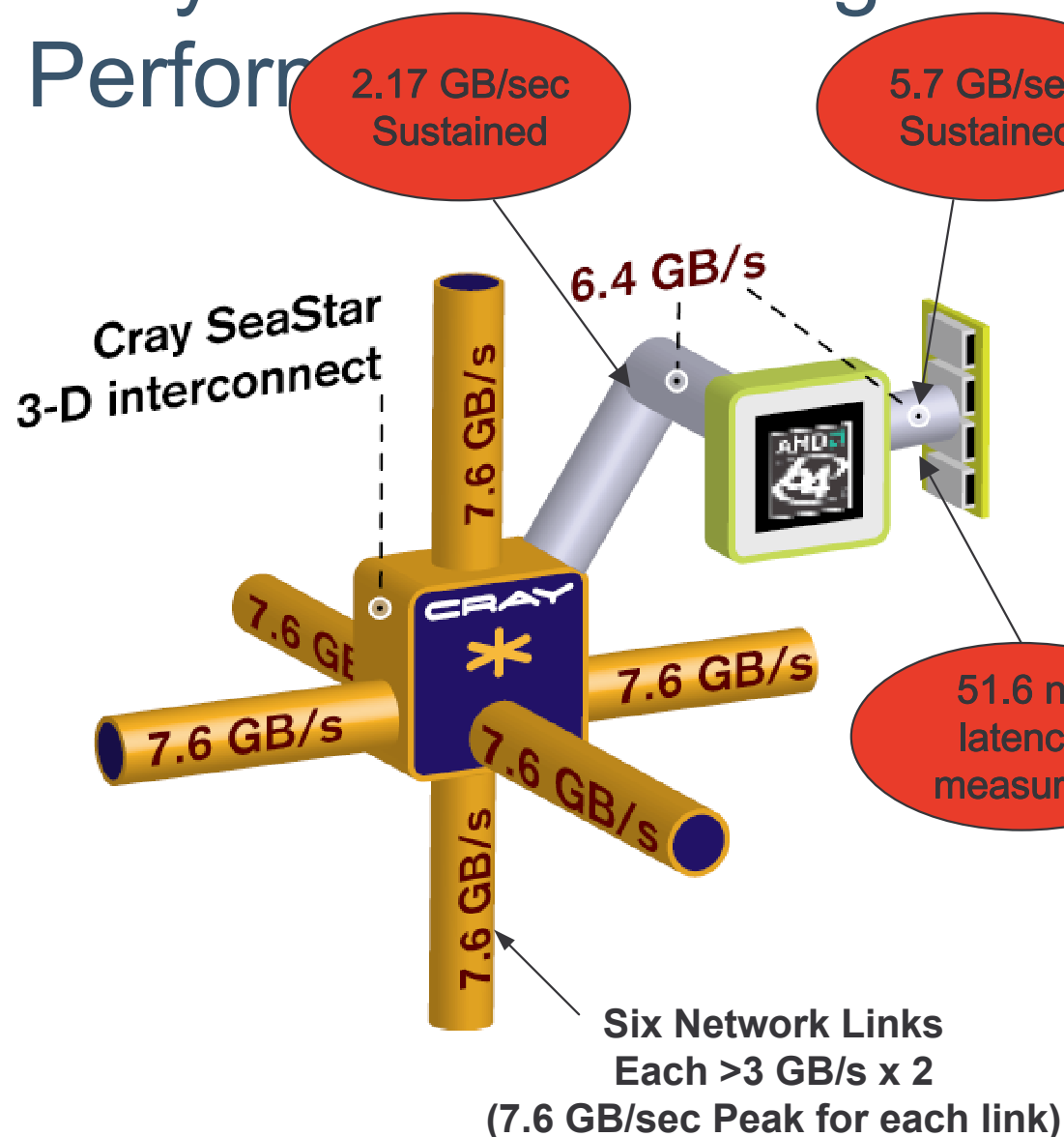
- DDR memory controller and function of Northbridge is pulled onto the Opteron die. Memory latency reduced to <60 ns
- No Northbridge chip results in savings in heat, power, complexity and an increase in performance
- Interface off the chip is an open standard (HyperTransport)

Cray XT3 Processing Element: Measured Performance



- DDR memory controller and function of Northbridge is pulled onto the Opteron die. Memory latency reduced to <60 ns
- No Northbridge chip results in savings in heat, power, complexity and an increase in performance
- Interface off the chip is an open standard (HyperTransport)

Cray XT3 Processing Element: Measured Performance



- DDR memory controller and function of Northbridge is pulled onto the Opteron die. Memory latency reduced to <60 ns
- No Northbridge chip results in savings in heat, power, complexity and an increase in performance
- Interface off the chip is an open standard (HyperTransport)

**2.17 GB/sec
Sustained**

**5.7 GB/sec
Sustained**

**Cray SeaStar
3-D interconnect**


6.4 GB/s

7.6 GB/s

7.6 GB/s

76

/s/

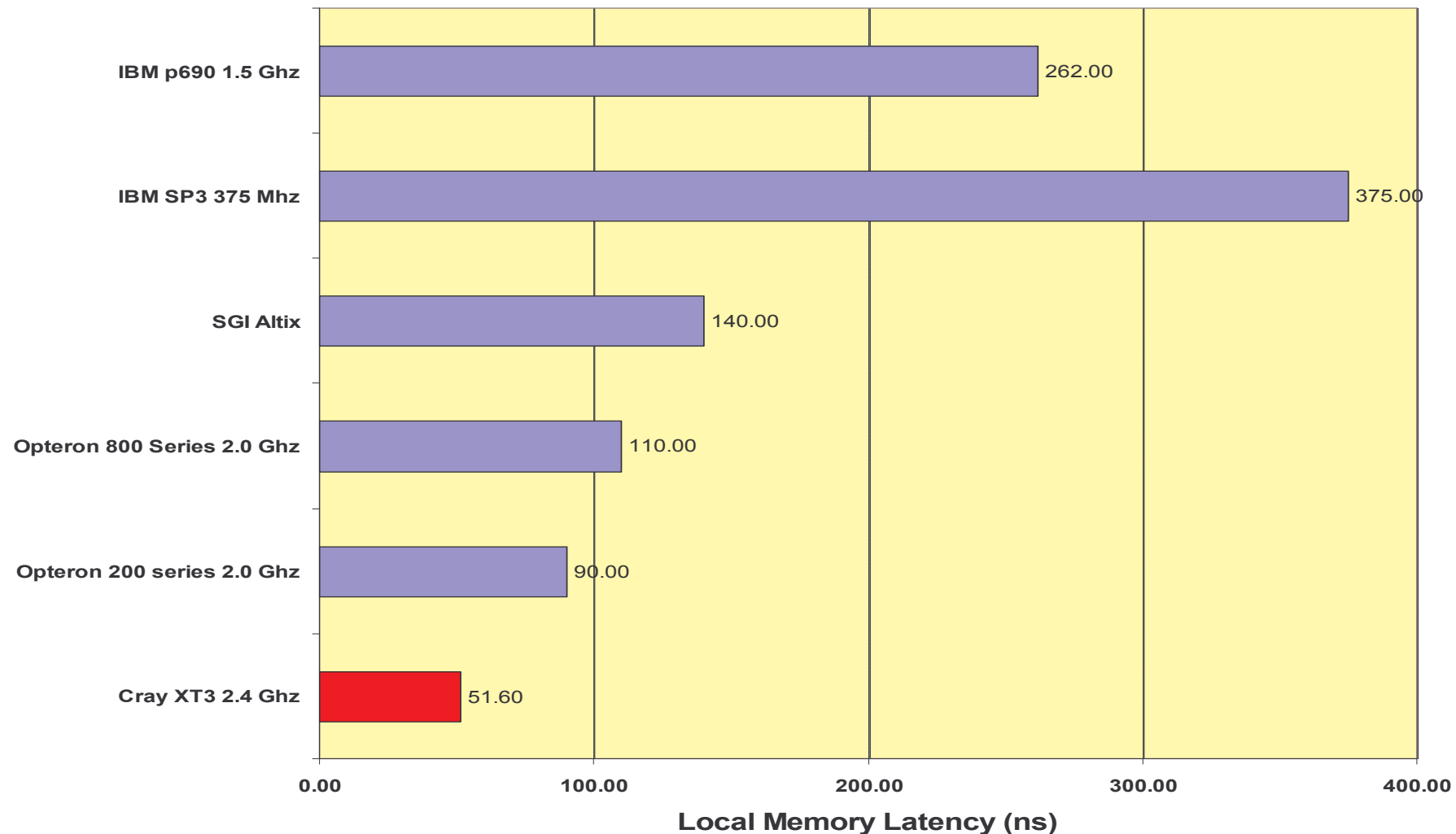


6.5 GB/sec Sustained

Six Network Links
Each >3 GB/s x 2
(7.6 GB/sec Peak for each link)

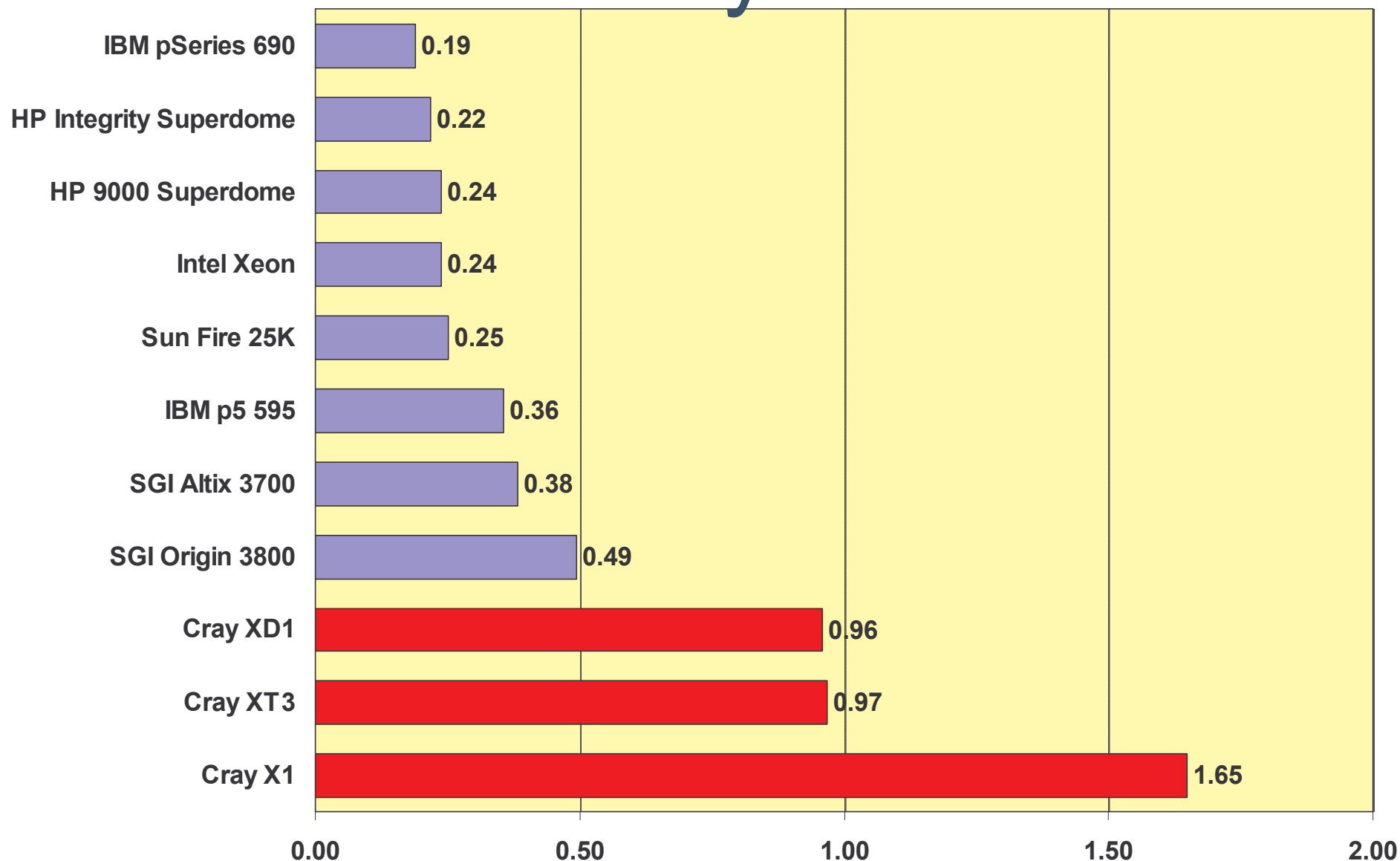
- **DDR memory controller and function of Northbridge is pulled onto the Opteron die. Memory latency reduced to <60 ns**
- **No Northbridge chip results in savings in heat, power, complexity and an increase in performance**
- **HyperTransport interface off the chip is an open standard (HyperTransport)**

Memory Latency



Single Processor architecture yields lowest memory latency

Measured Memory Balance

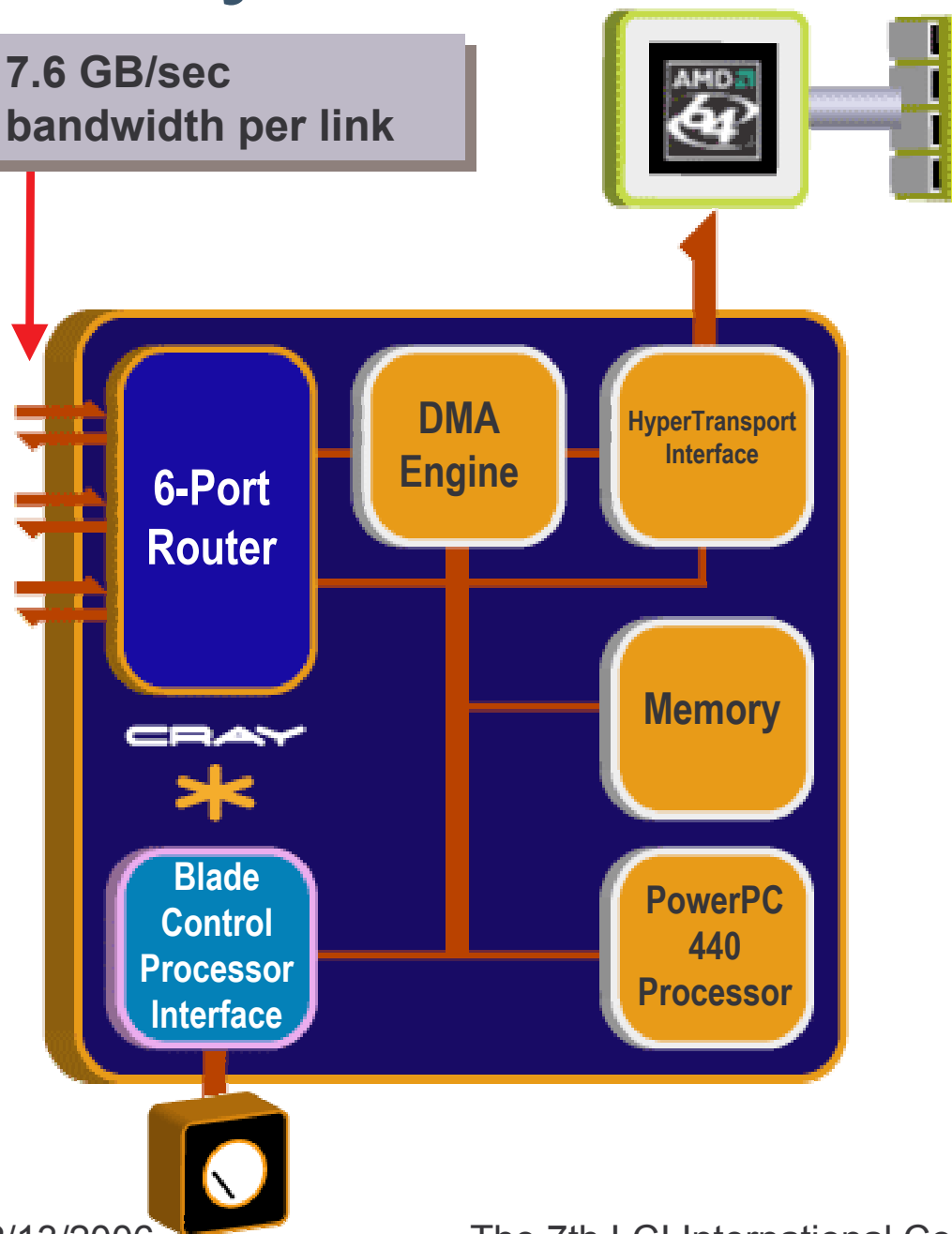


B/F calculated from memory bandwidth
measured via STREAM Triad benchmark

Memory/Computation Balance (B/F)

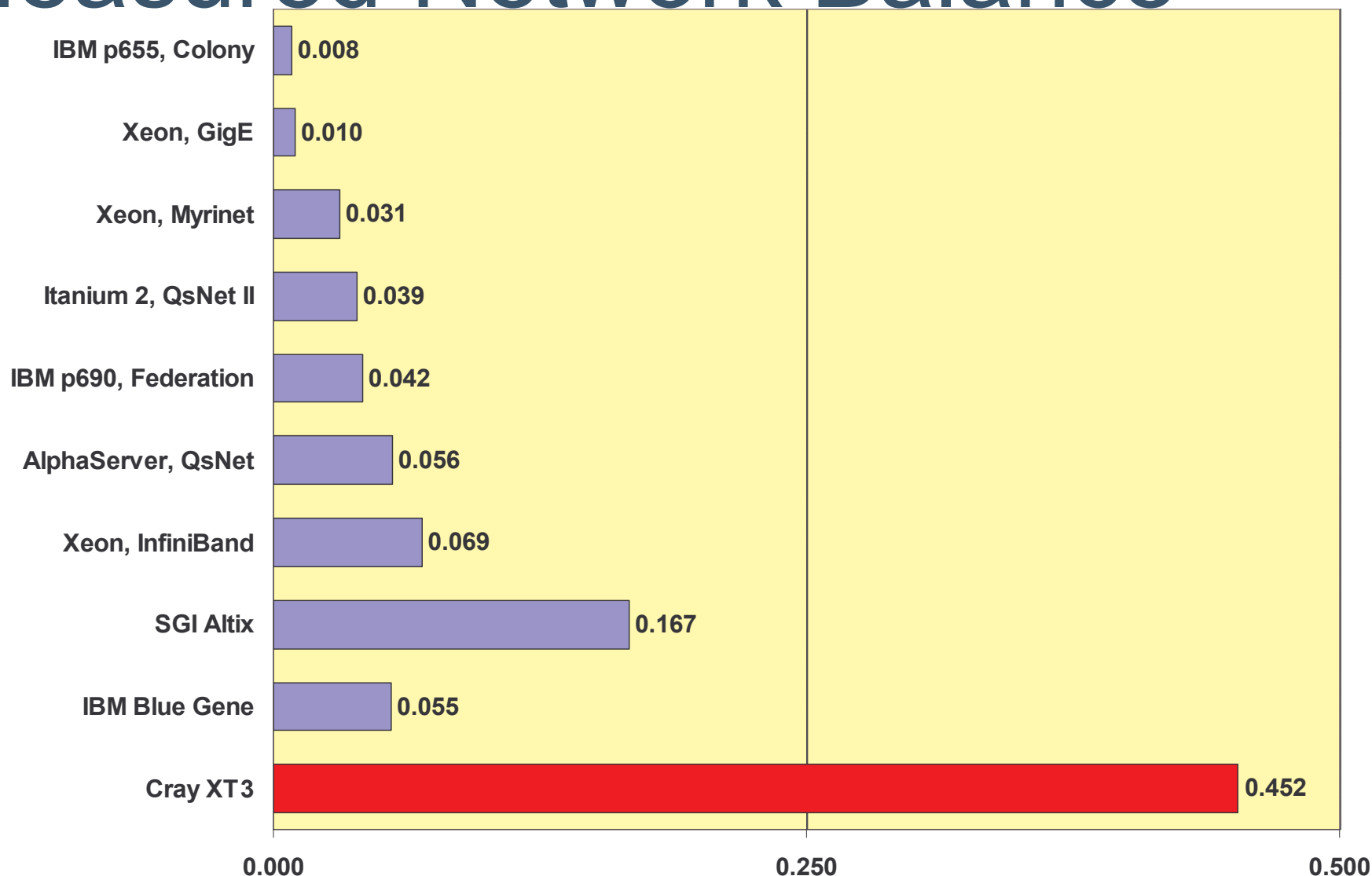
Cray SeaStar Internals

7.6 GB/sec
bandwidth per link



- Each Processor is directly connected to a dedicated SeaStar
- Each SeaStar contains a 6-Port router *and* communications engine
- Provides serial connection to the Cray RAS and Management System

Measured Network Balance



Network bandwidth is the maximum
bidirectional data exchange rate
between two nodes using MPI

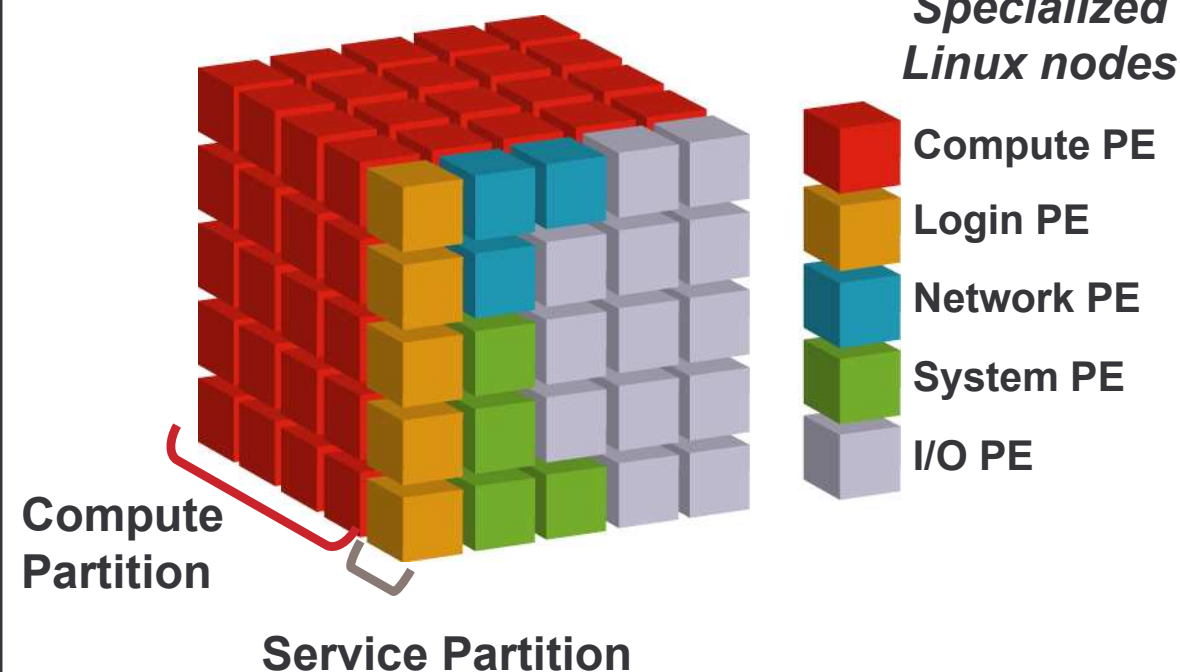
Communication/Computation Balance (B/F)

What is OS Jitter?

- In order to provide certain services, the OS must wake-up periodically
 - Network and I/O Requests
 - Daemon Processes
 - System Calls and Threads
- These wake up calls interrupt user time
- As more processors are added, more interruptions occur across the machine
- This white noise on the system is known as “OS Jitter” and is often the limiting factor for system scaling.

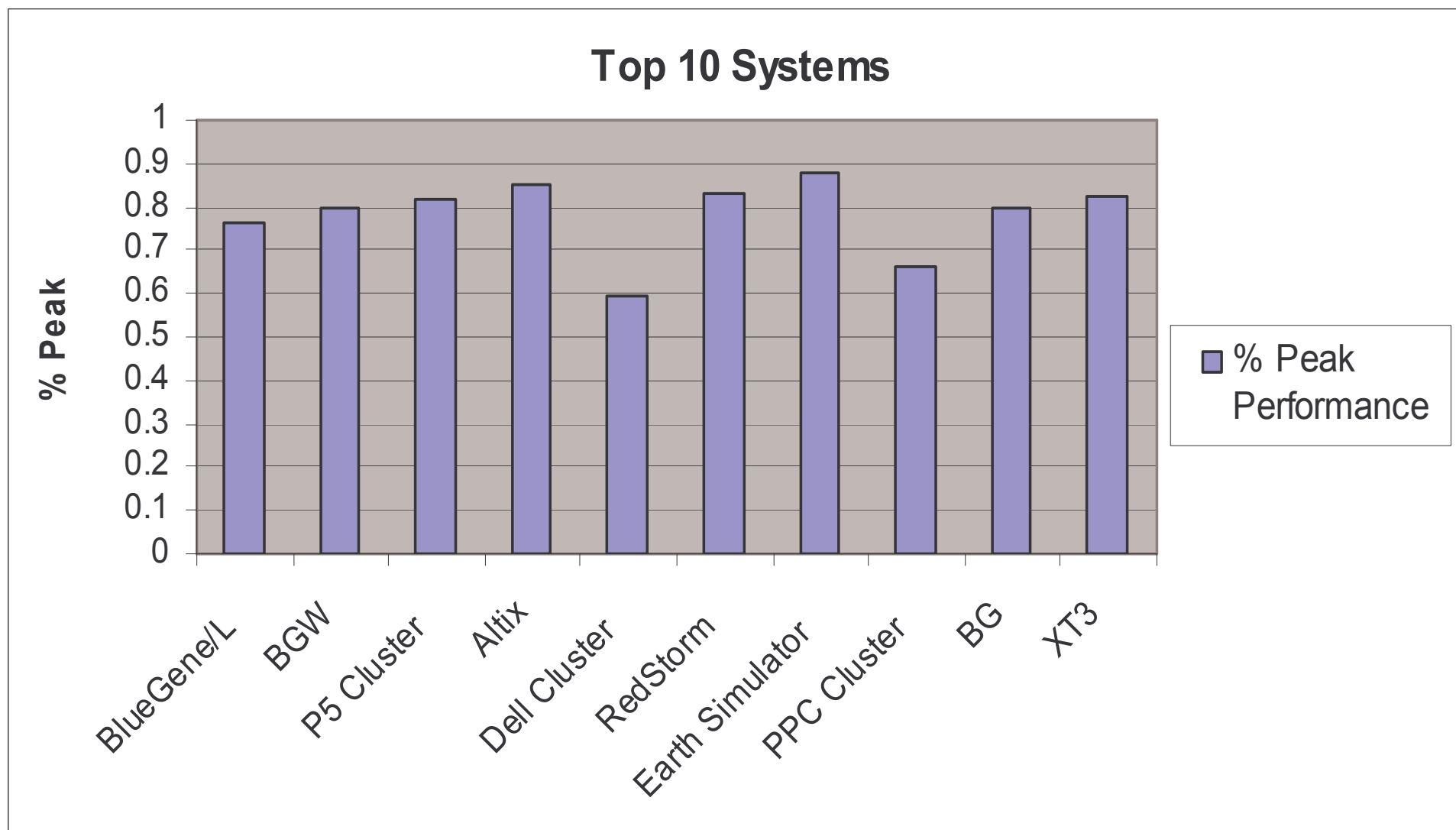


Scalable Software Architecture: UNICOS/Ic



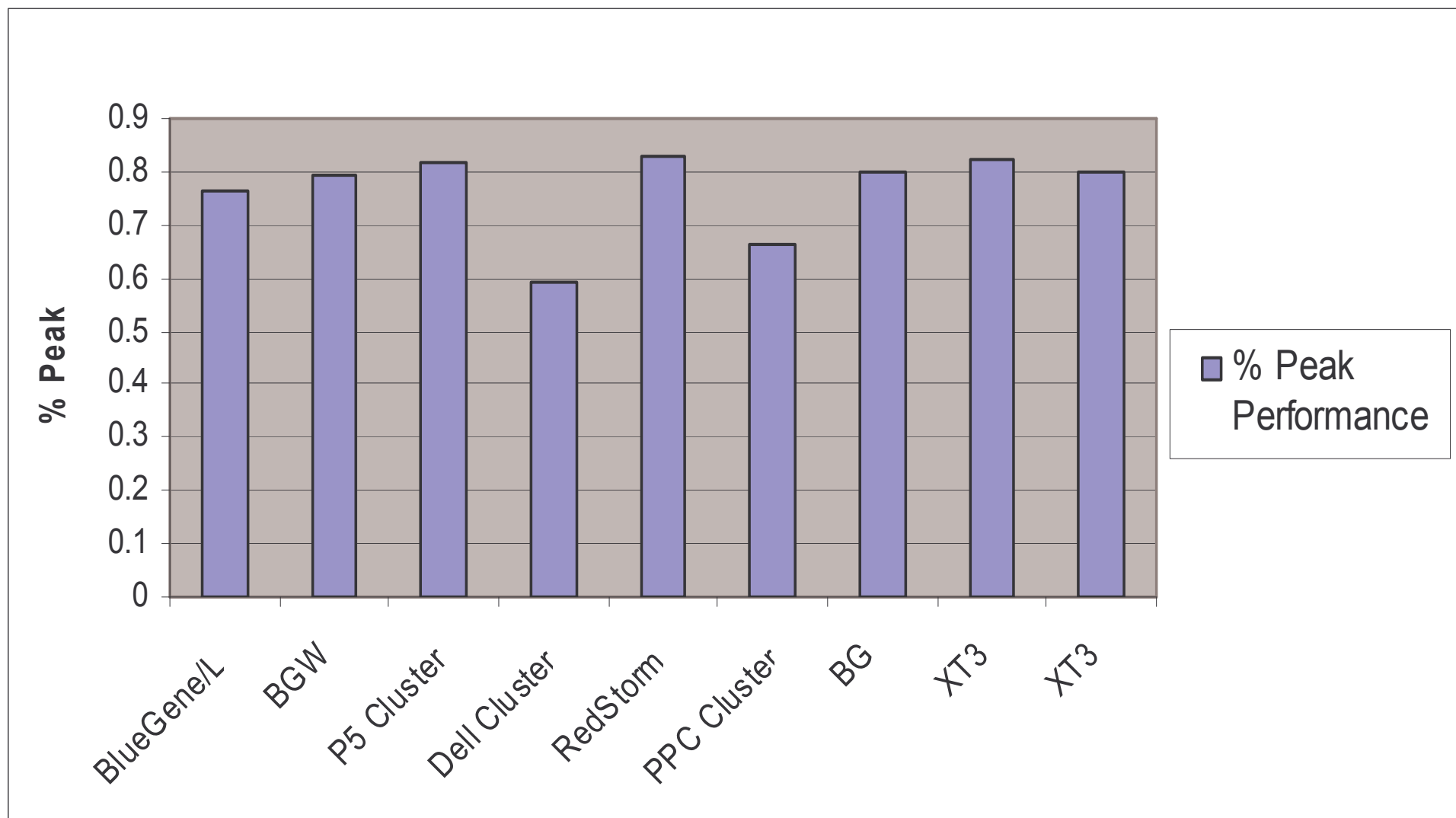
- Full featured Linux on Service PEs, Microkernel on Compute PEs.
- Service PEs specialize by function
- Contiguous memory layout used on compute processors to streamline communications
- Software Architecture eliminates OS "Jitter"
 - 100 ms interrupt times
 - Will be synchronized if required
 - OS heartbeat checked once per second.
- Software Architecture enables reproducible run times

Top500 Linpack Scaling



Data from 11/2005 Top500 list (top500.org)

Top500 Linpack Scaling



Data from 11/2005 Top500 list (top500.org)

XT3 Benchmark Results

2/13/2006

The 7th LCI International Conference on Clusters

EXPLORE SIMULATE CREATE

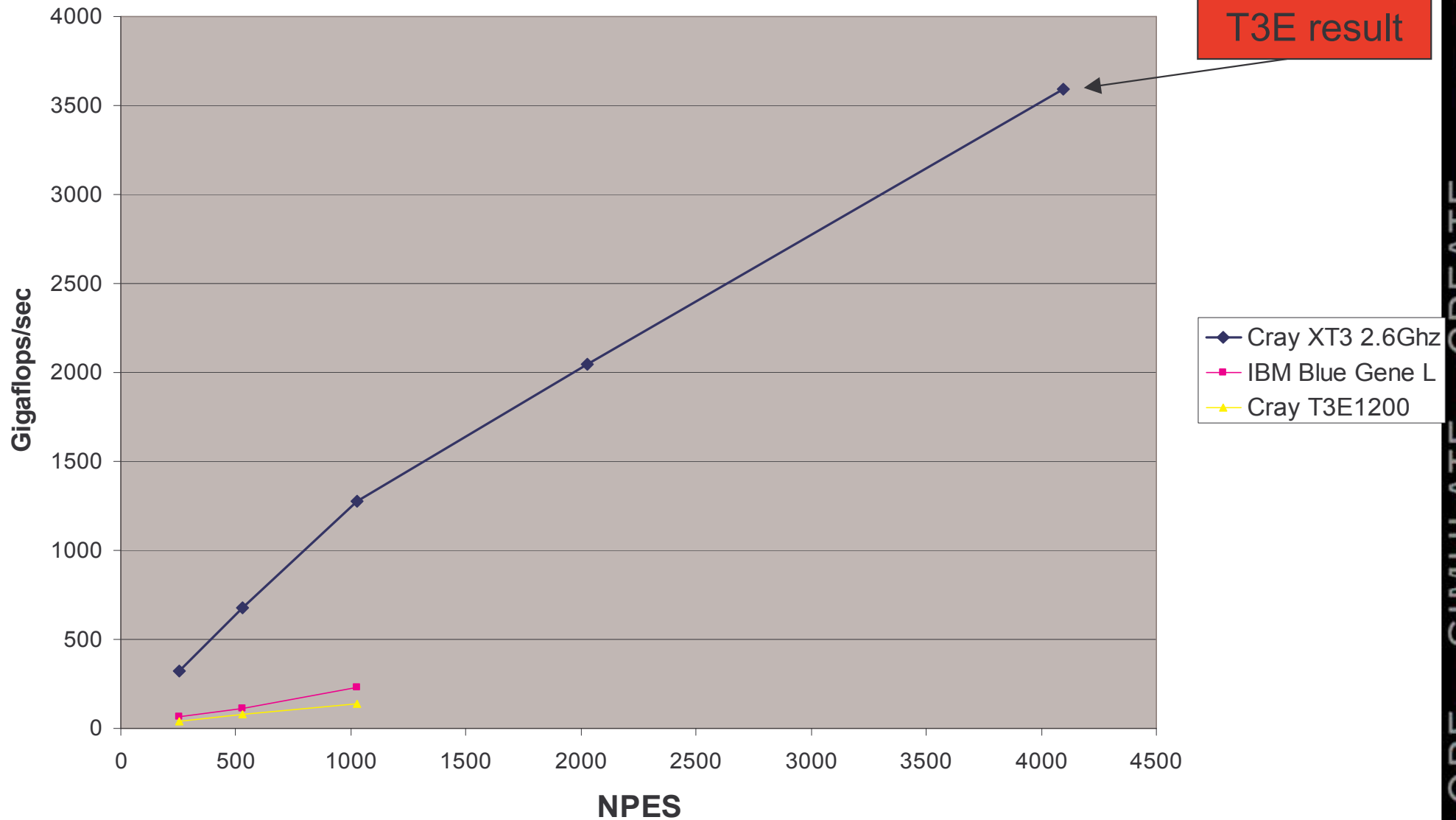


NAS Parallel Benchmarks

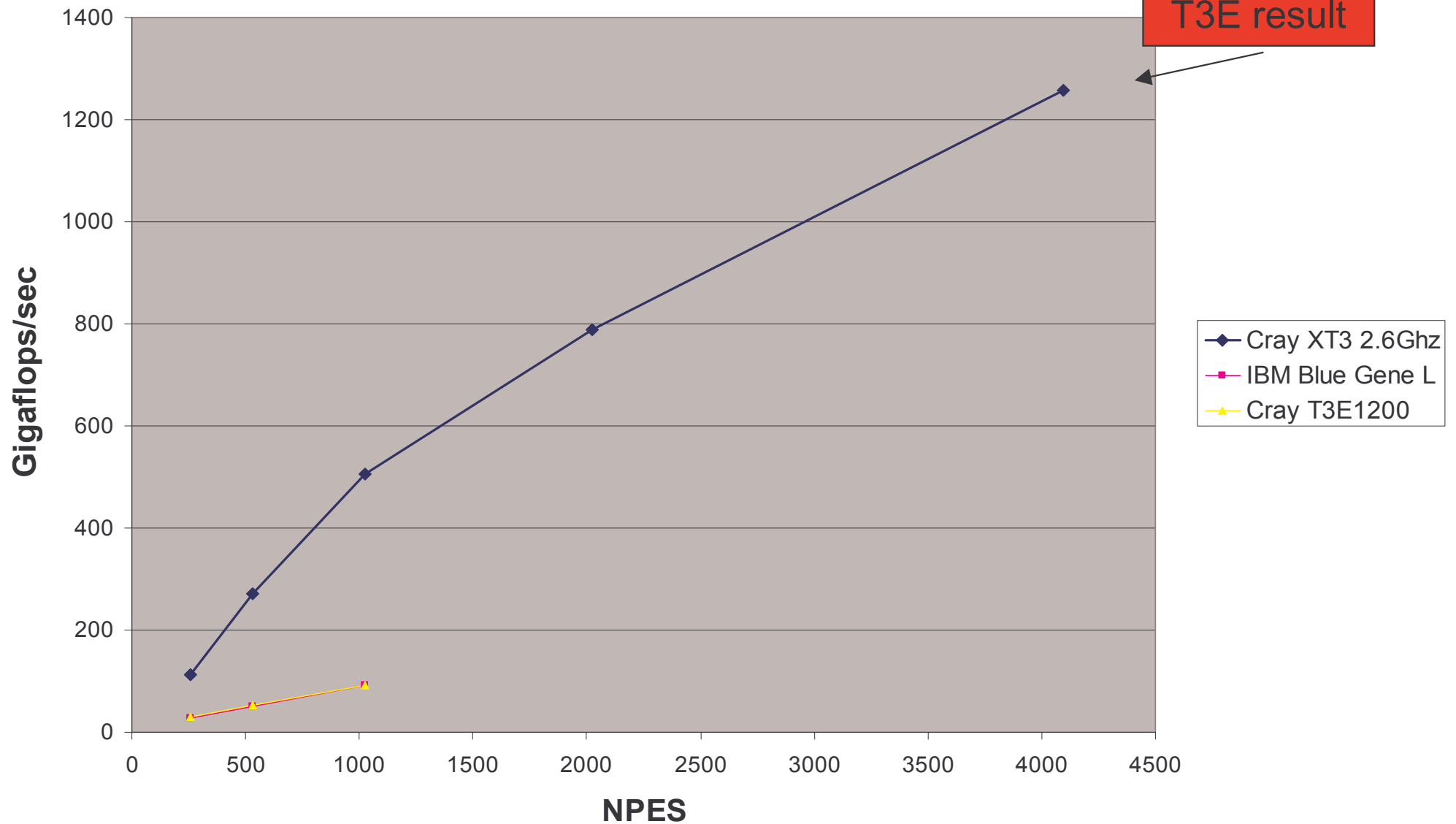


- The first time we looked at these was in 1991 on a Cray Y-MP
- The last time we looked at these was in 1999 on a 1024 Processor Cray T3E 1200
 - We tuned codes for E-registers, shmем, etc.
- We recently ran these on a 4096 PE Cray XT3
 - Codes were the ASIS MPI version (3.2) (one exception)
 - -O3 -fastsse
- IBM Blue Gene Results from Argonne
 - http://www-unix.mcs.anl.gov/~kaushik/bgl/npb_results.htm
 - Run in co-processor mode
 - -O3 -qarch=440d -qtune=440 -qbgl -qmaxmem=64000

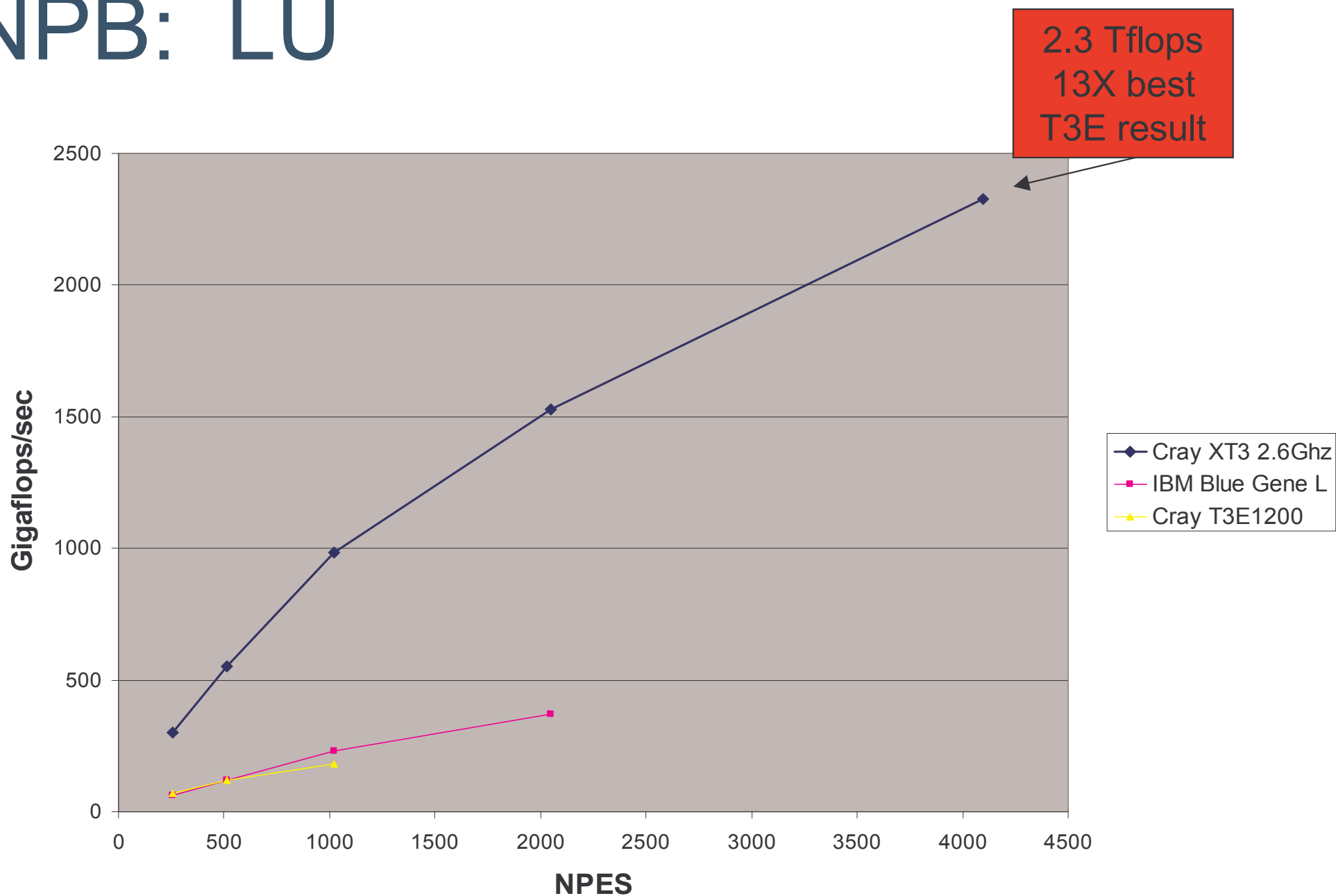
NPB: BT



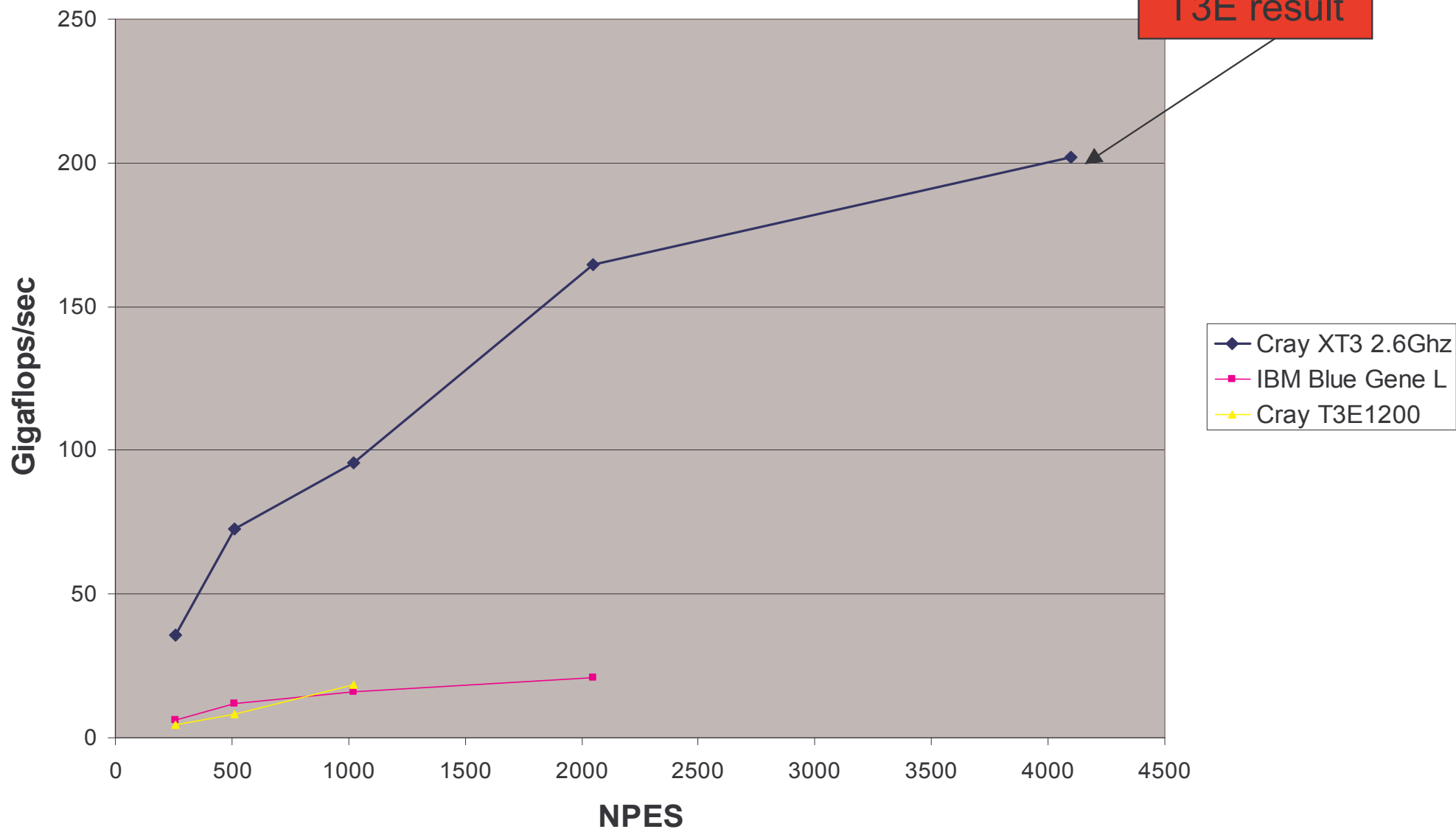
NPB: SP



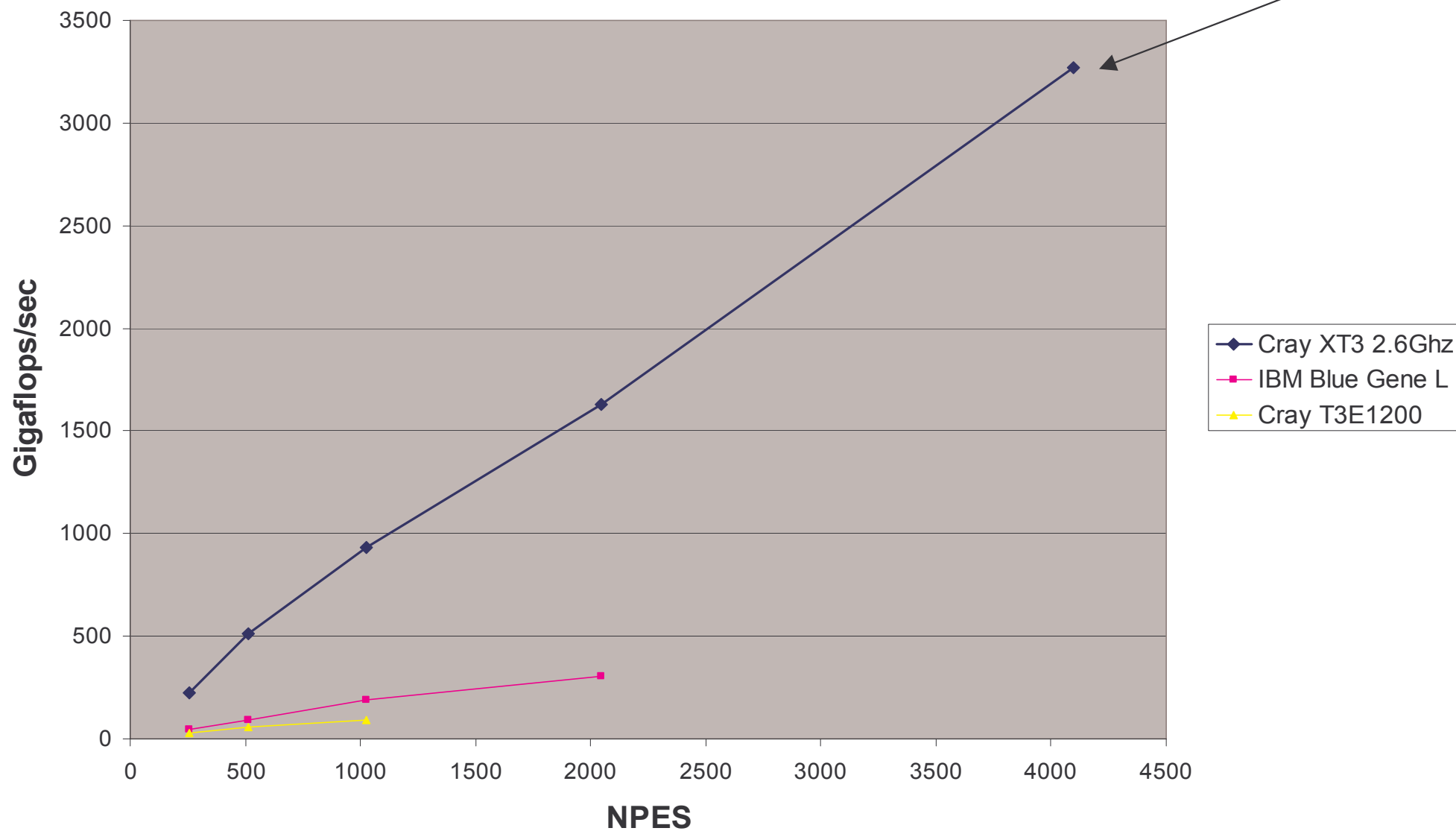
NPB: LU



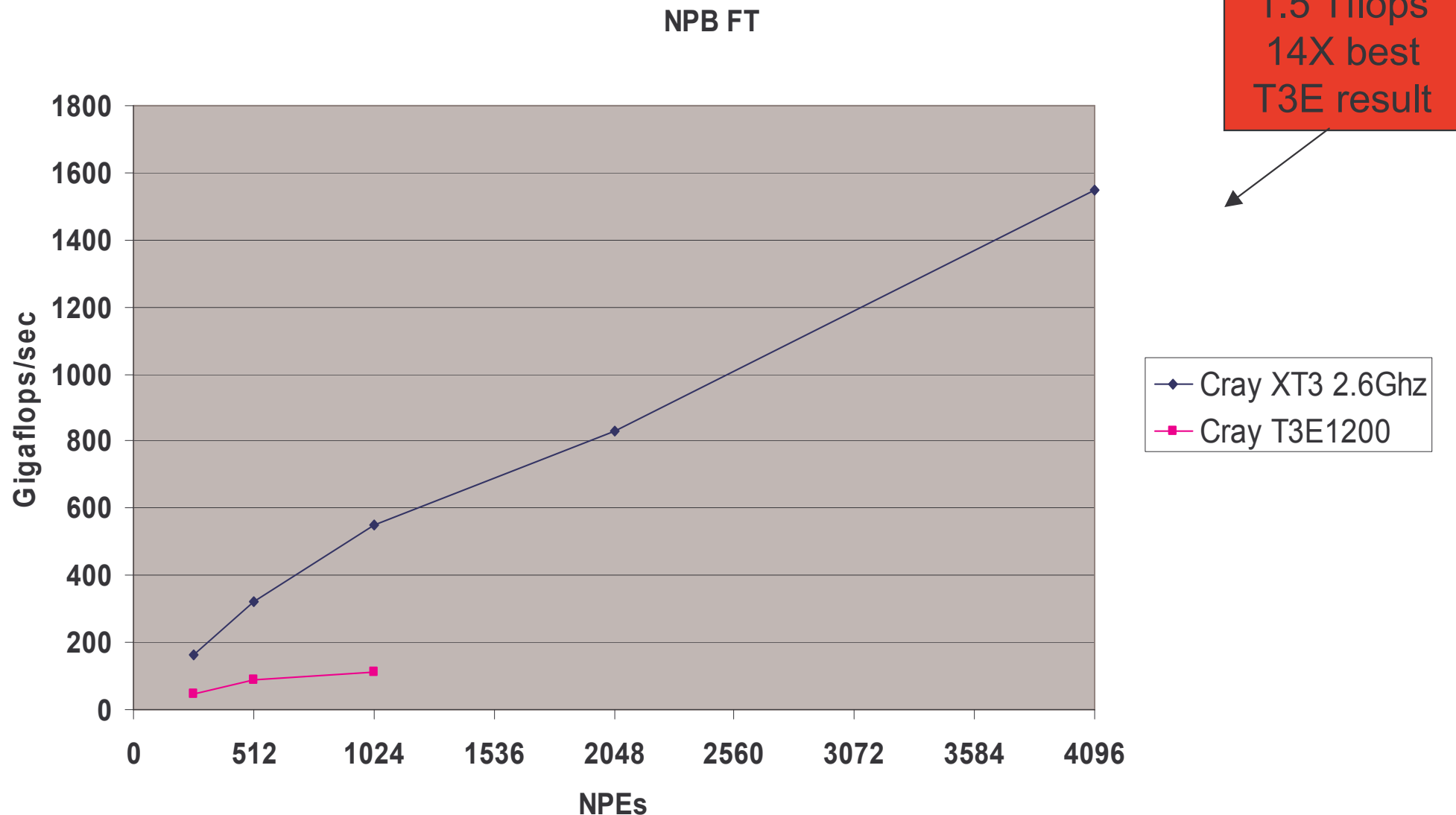
NPB: CG



NPB: MG



NPB: FT



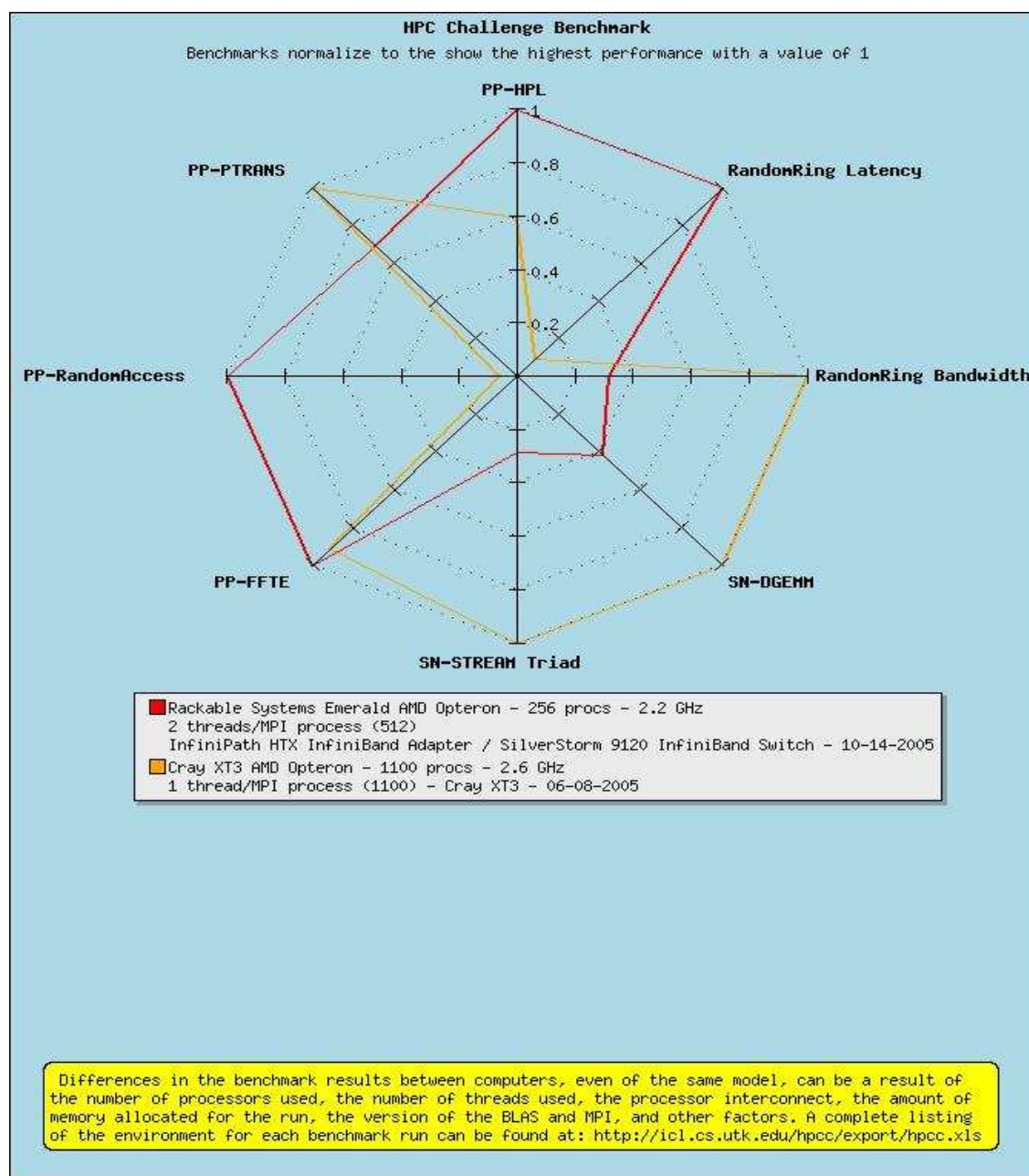
HPC Challenge

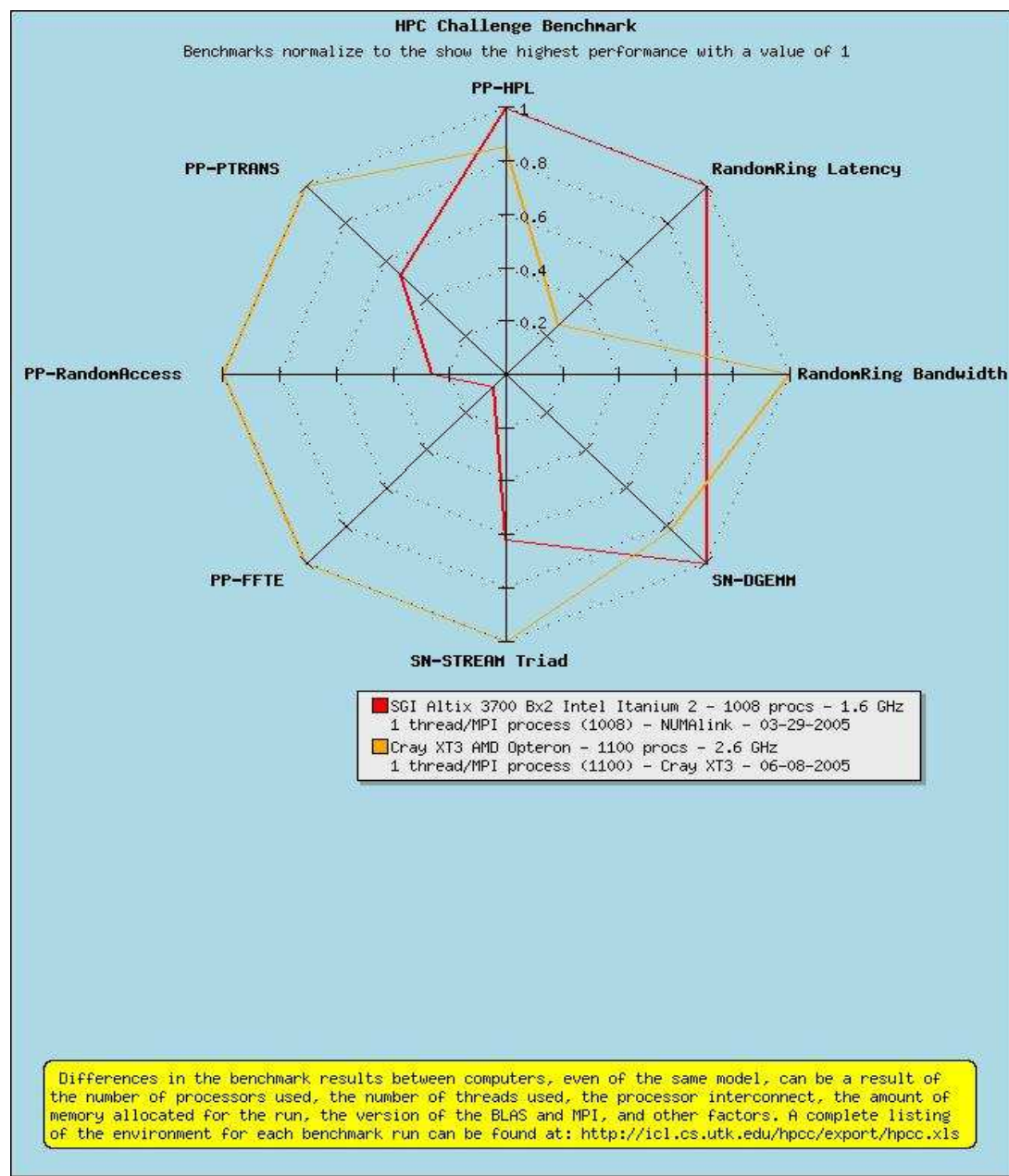
2/13/2006

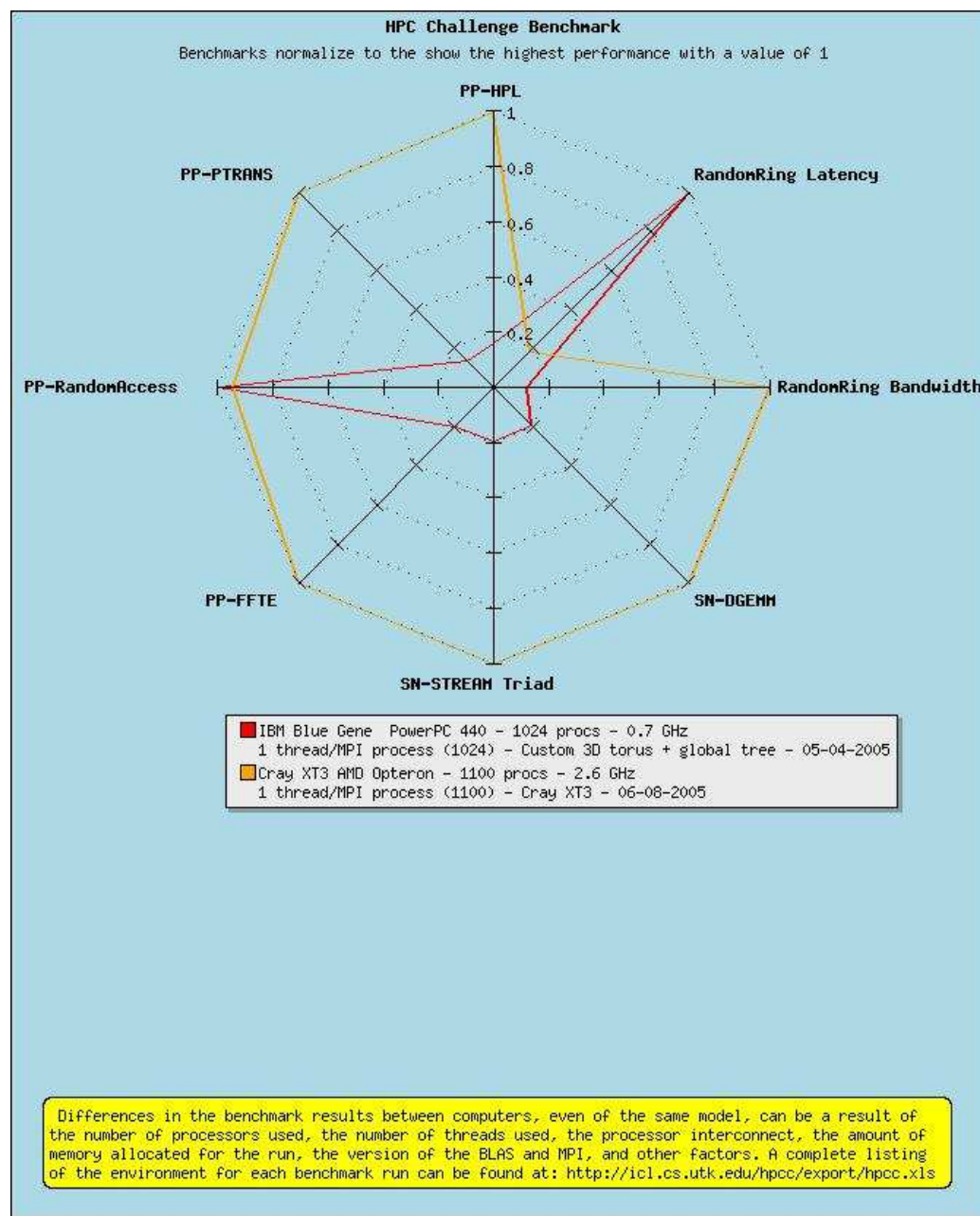
The 7th LCI International Conference on Clusters

EXPLORE SIMULATE CREATE



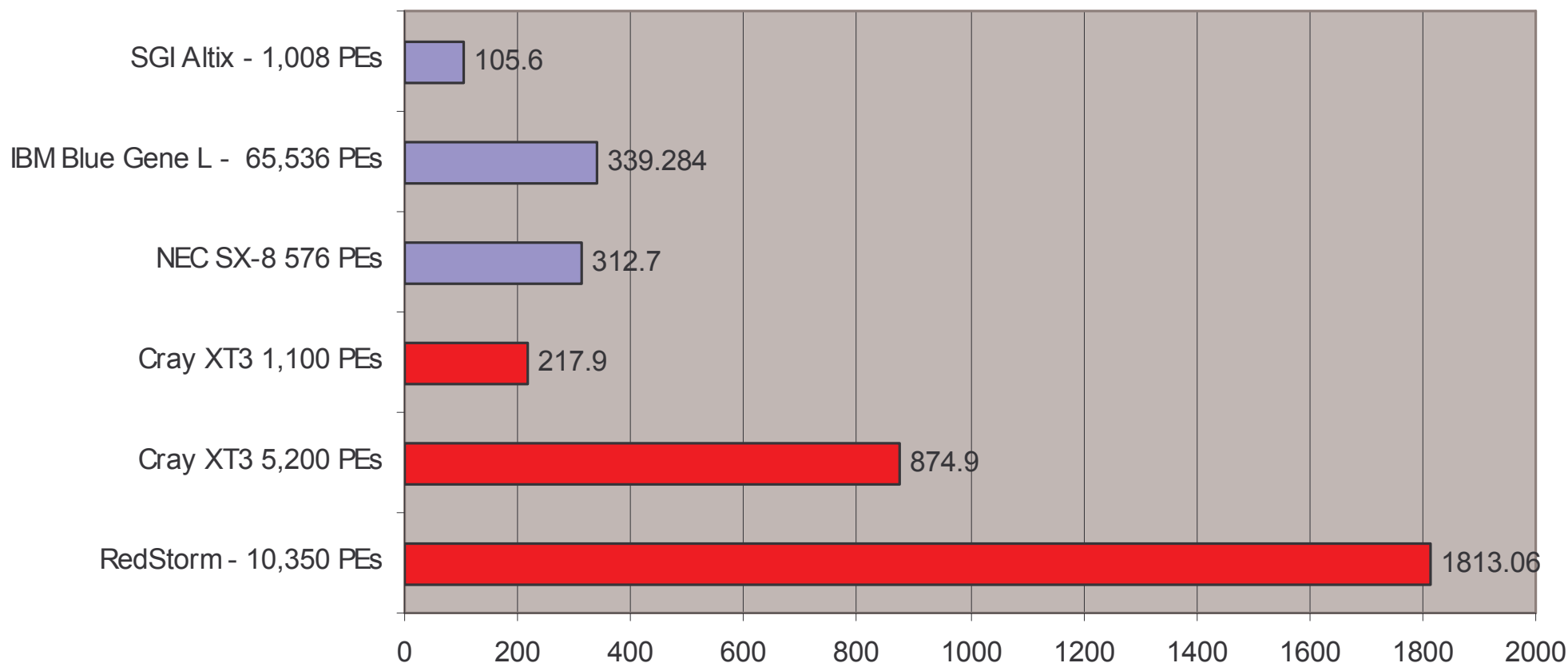






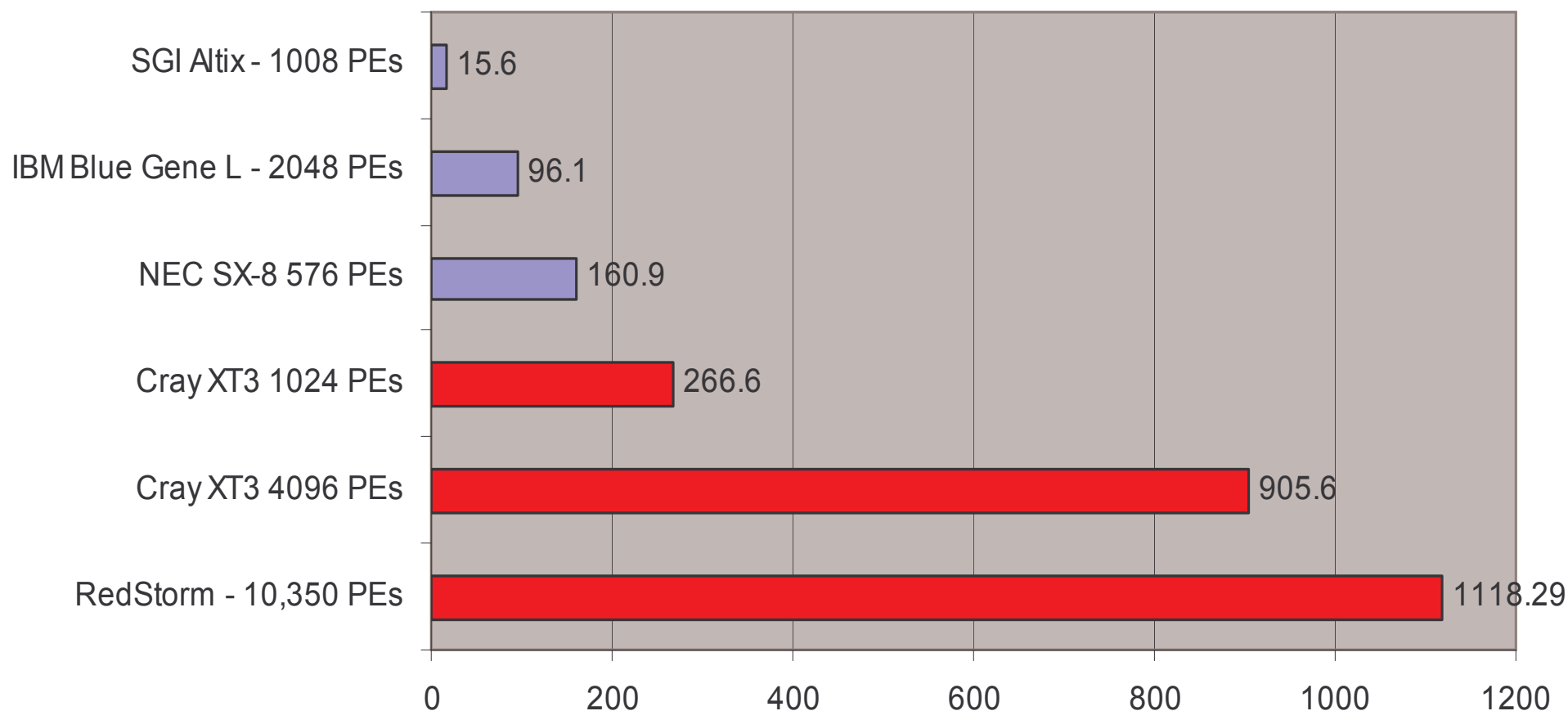
Ptrans (Global Bandwidth)

GP Trans (GB/sec)



Global FFT (Gflops/sec)

Global FFT (Gflops/sec)



XT3 Application Results

2/13/2006

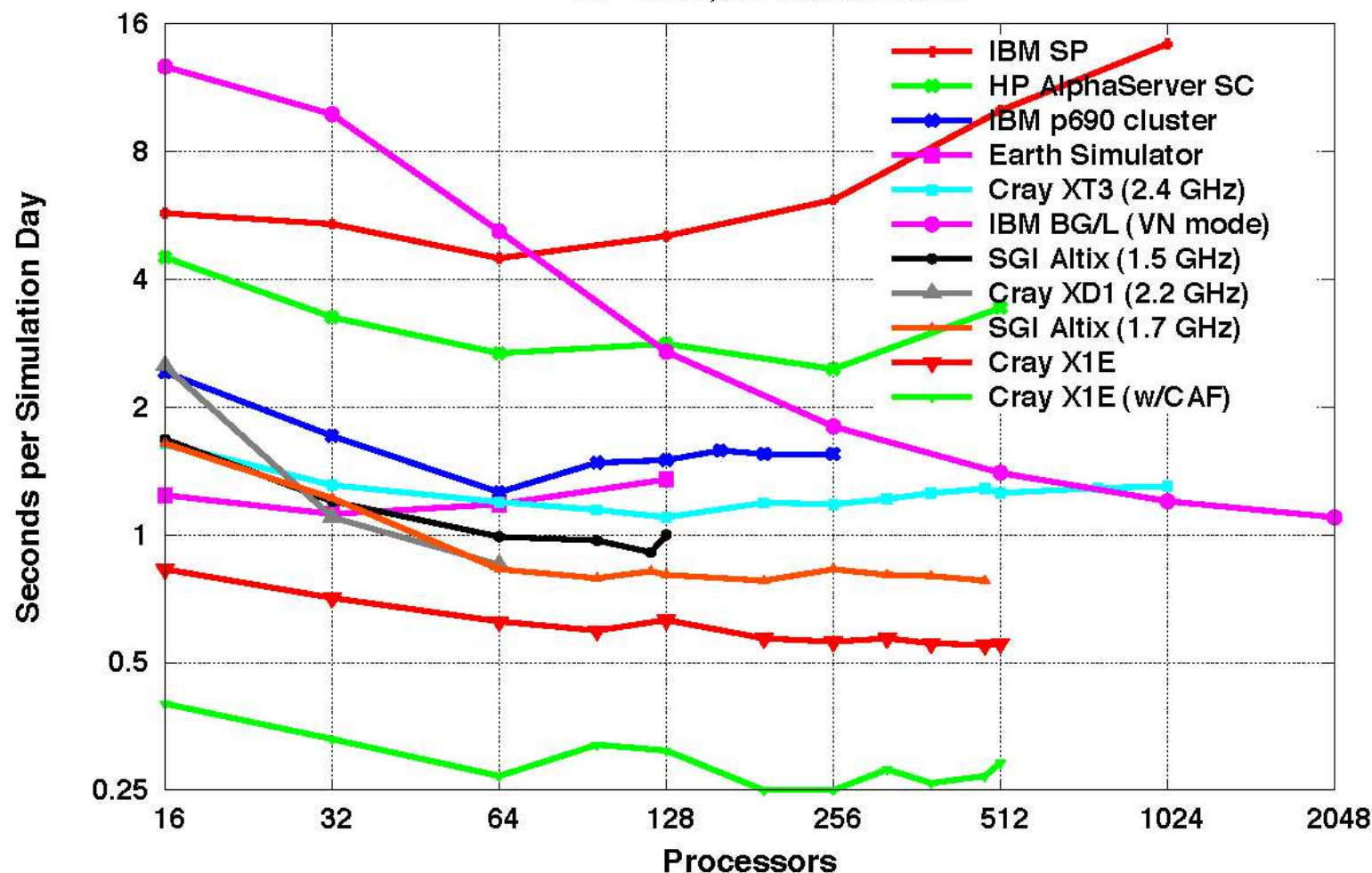
The 7th LCI International Conference on Clusters

EXPLORE SIMULATE CREATE



POP Barotropic

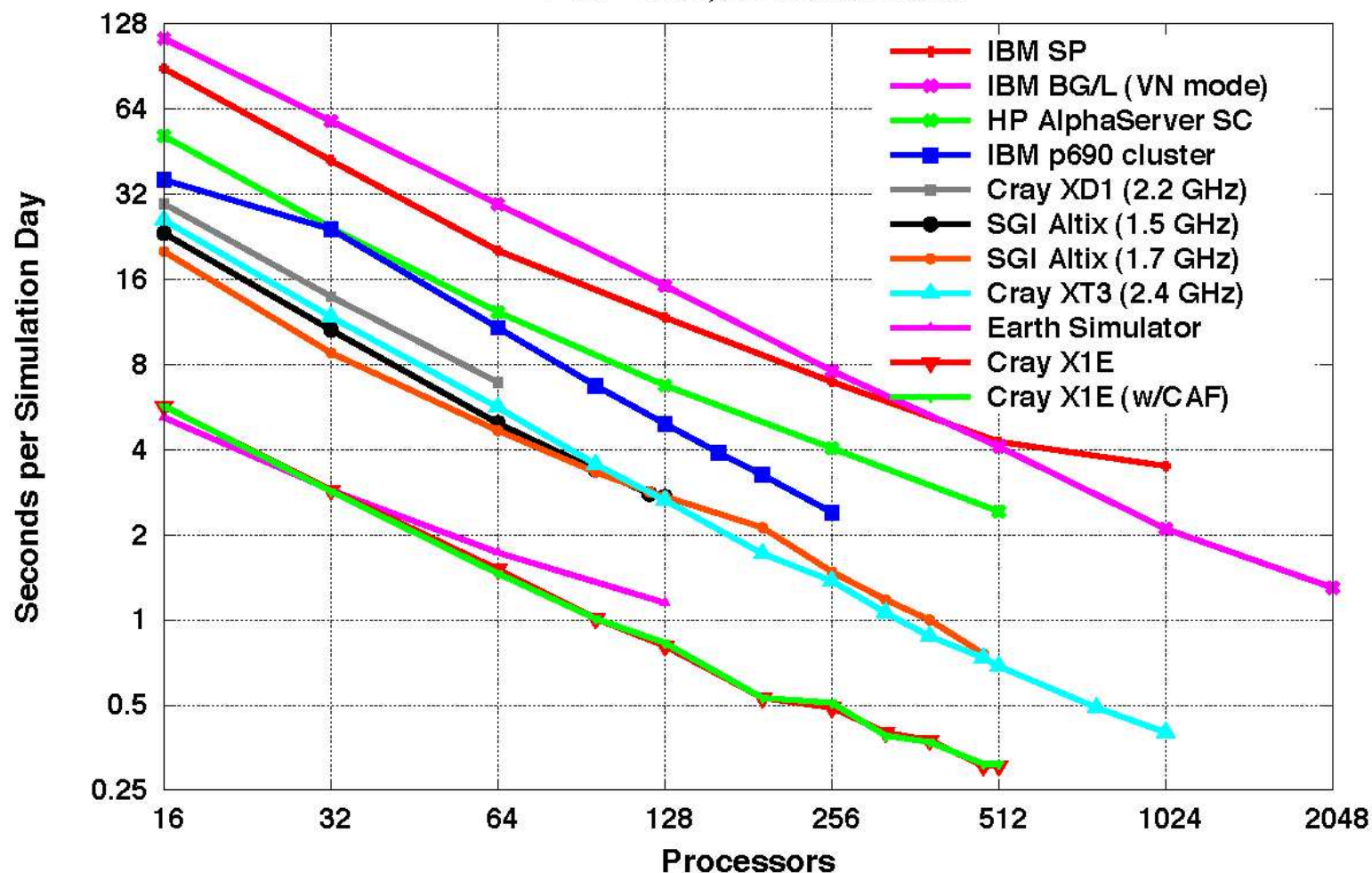
POP Barotropic Timings
POP 1.4.3, x1 benchmark



Lower is better

POP Baroclinic

POP Baroclinic Timings
POP 1.4.3, x1 benchmark

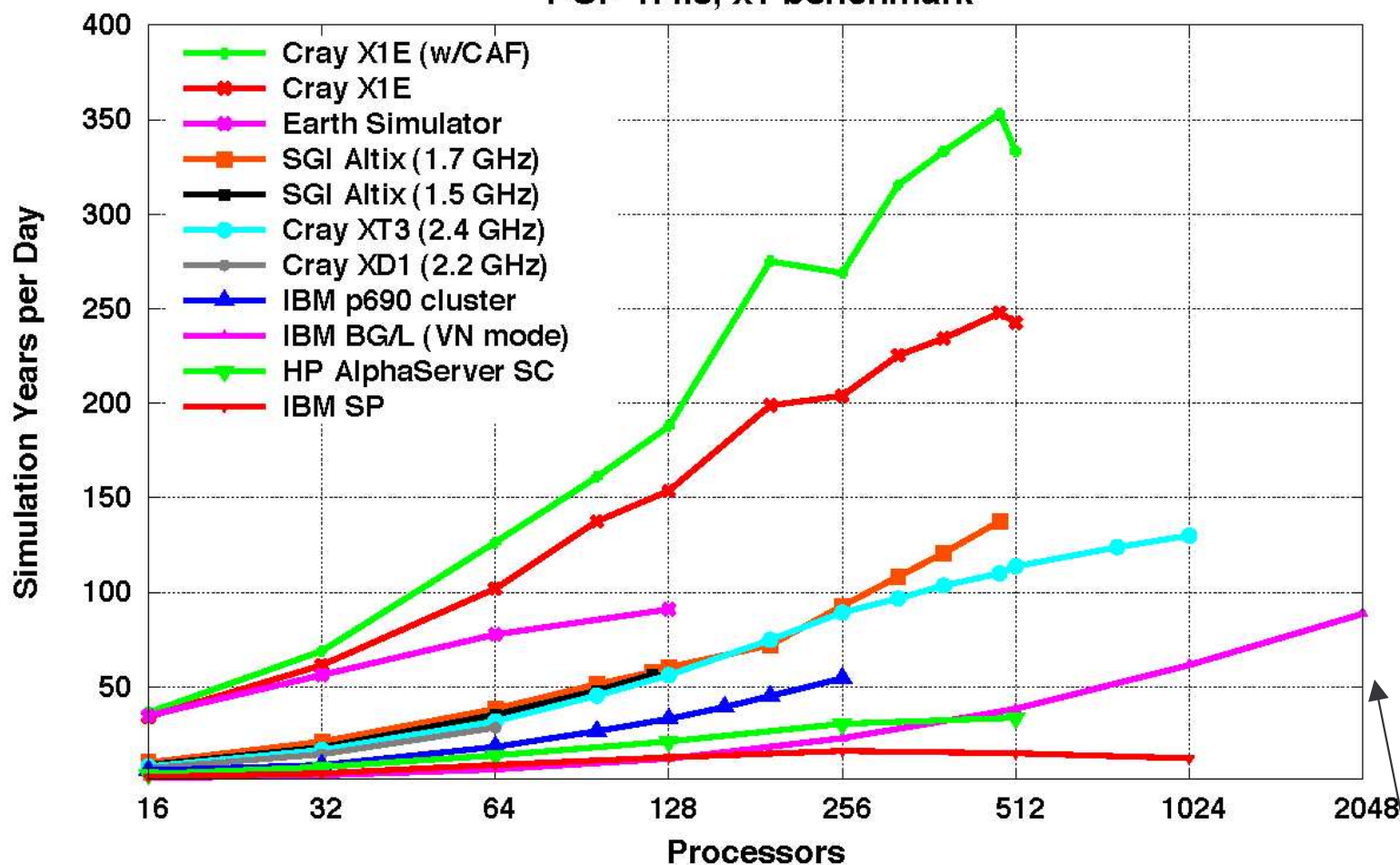


Lower is better

POP Overall Performance

LANL Parallel Ocean Program

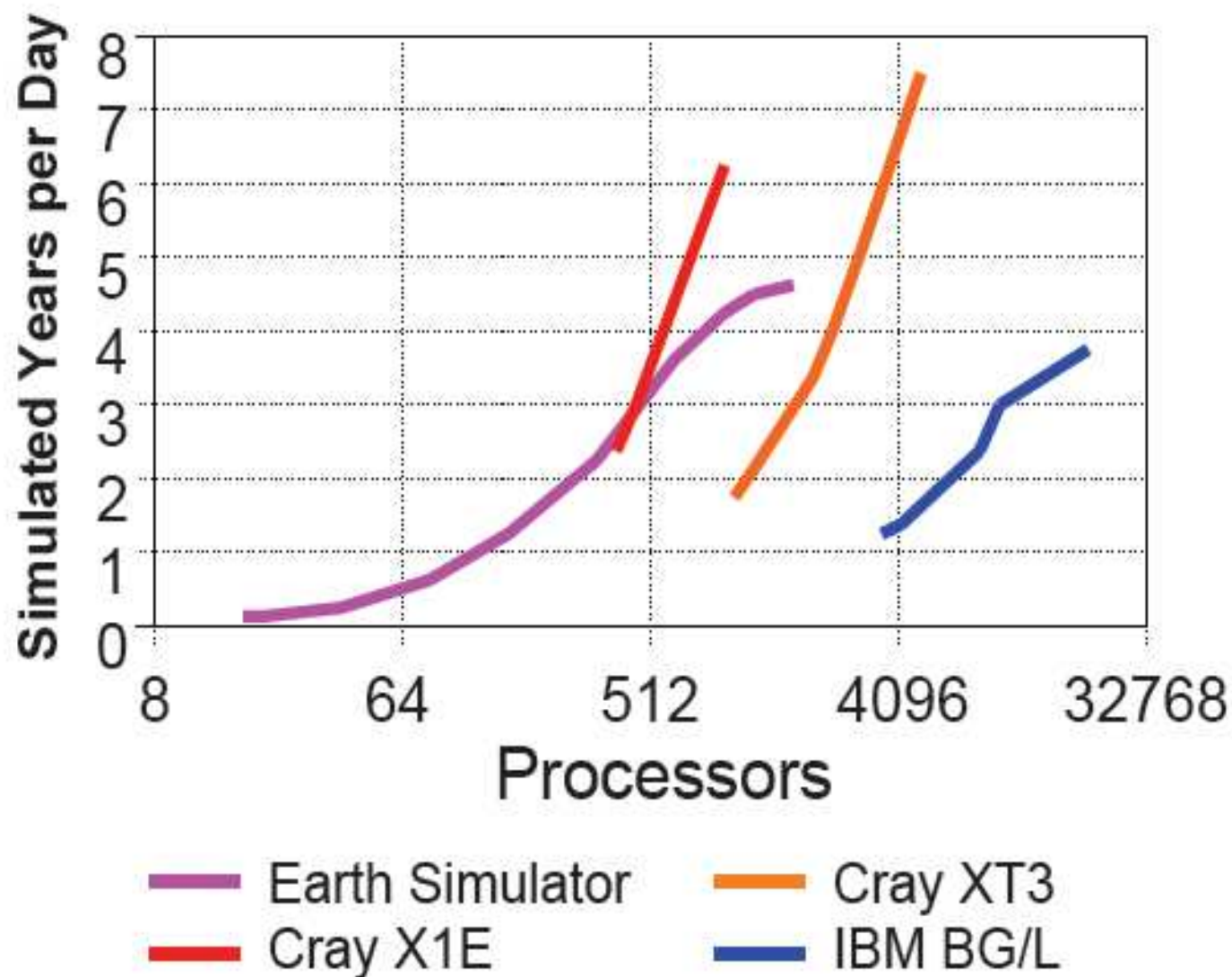
POP 1.4.3, x1 benchmark



Higher is better

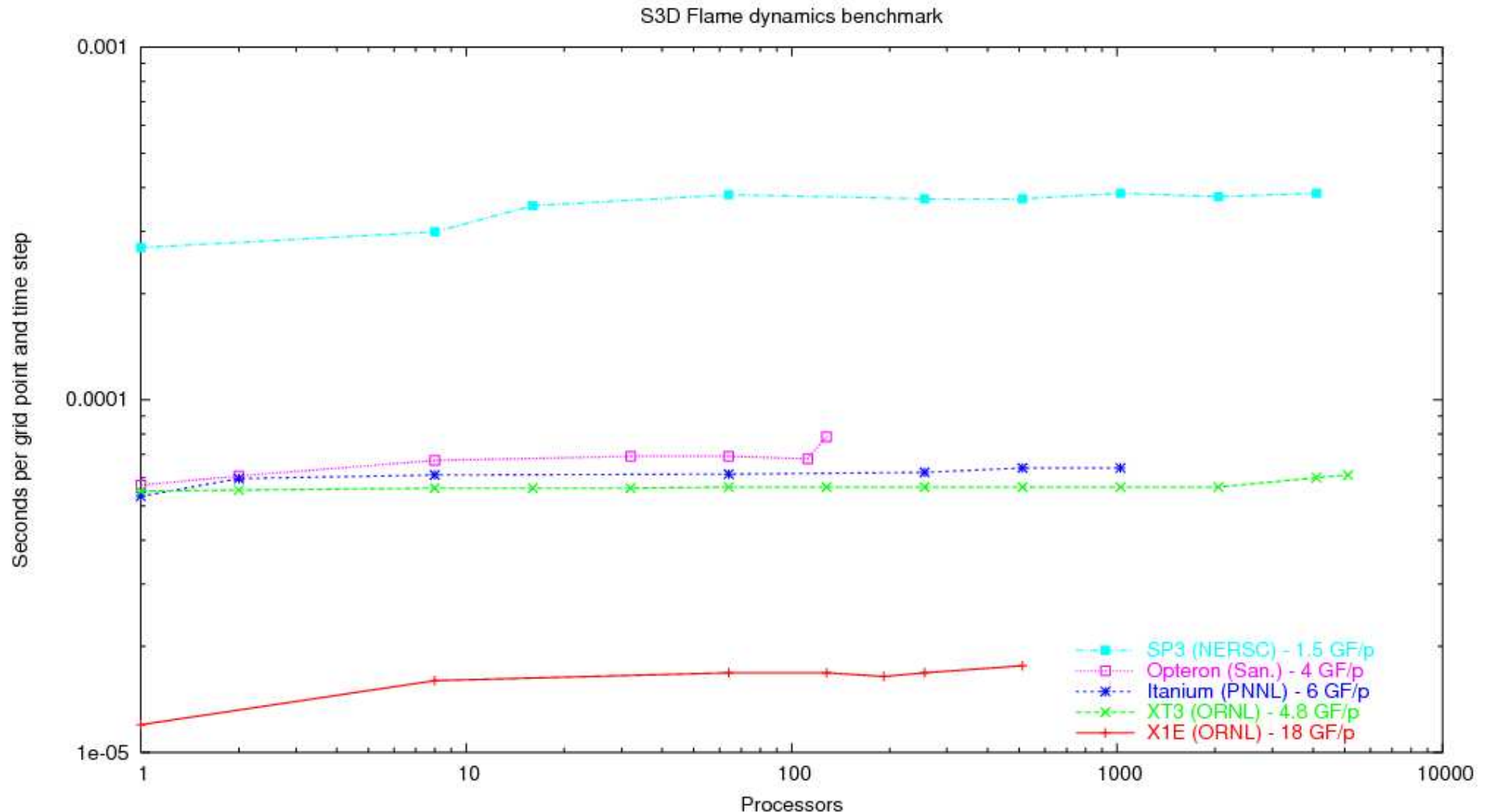
60 surface cells/processor

POP 1.4.3 tenth-degree benchmark



Higher is better

S3D Flame Dynamics

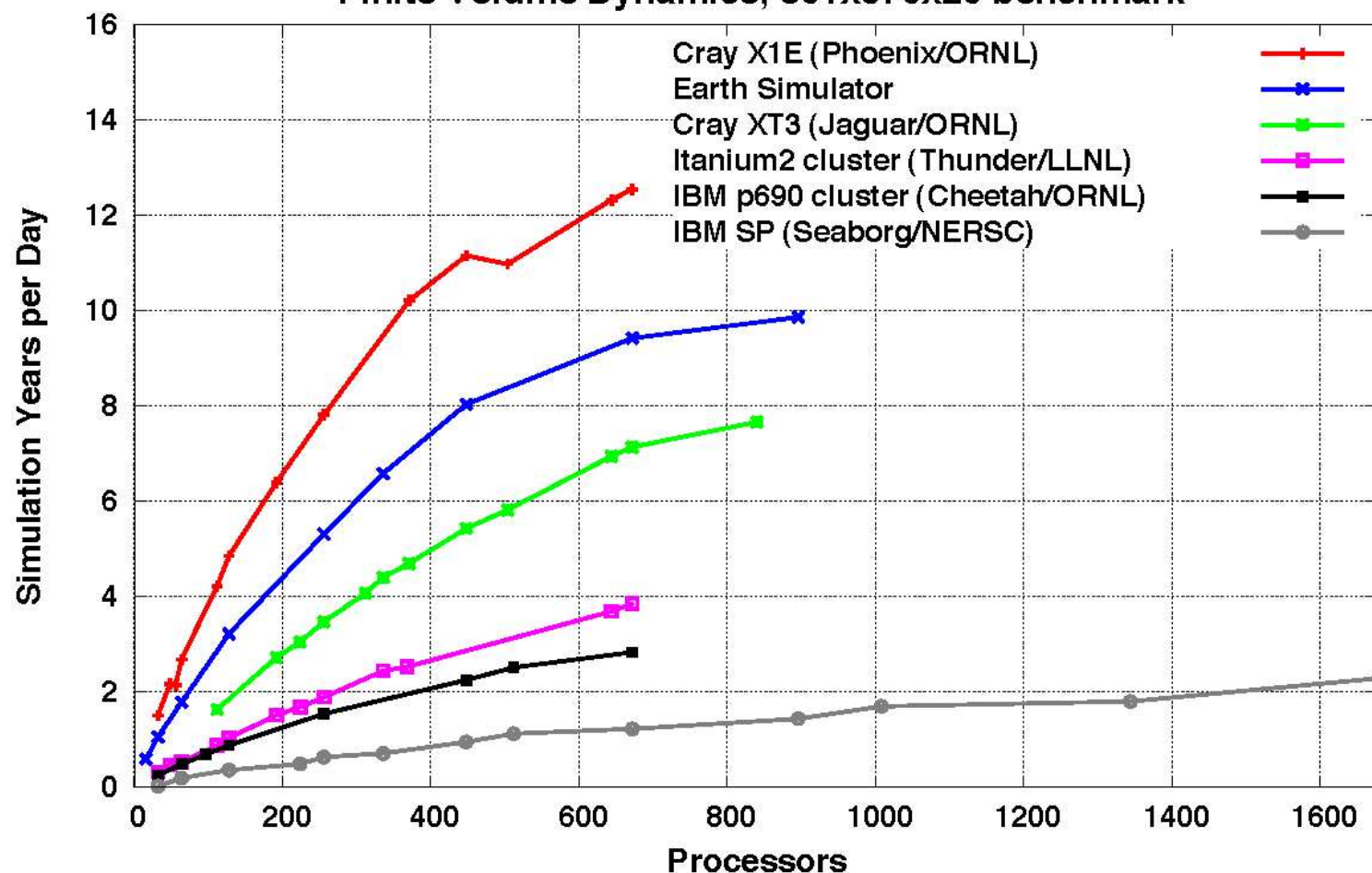


Flat and Lower is better

CAM Atmospheric Model

Performance of the CAM3.1 Atmospheric Model

Finite Volume Dynamics, 361x576x26 benchmark



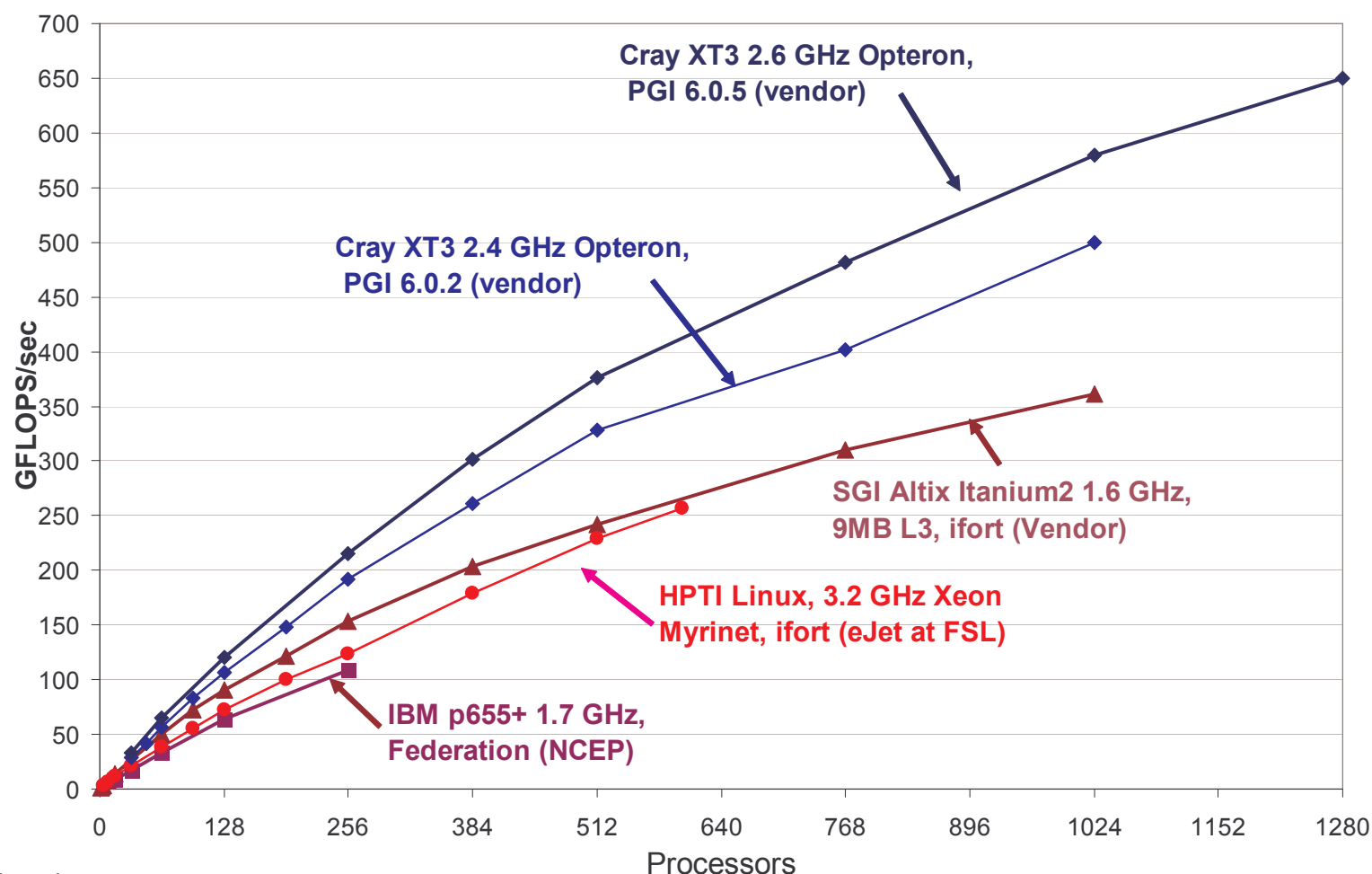
Higher is better

WRF Performance on XT3



THE WEATHER RESEARCH & FORECASTING MODEL

WRF v2 EM Core, 425x300x35, DX=12km, DT=72s

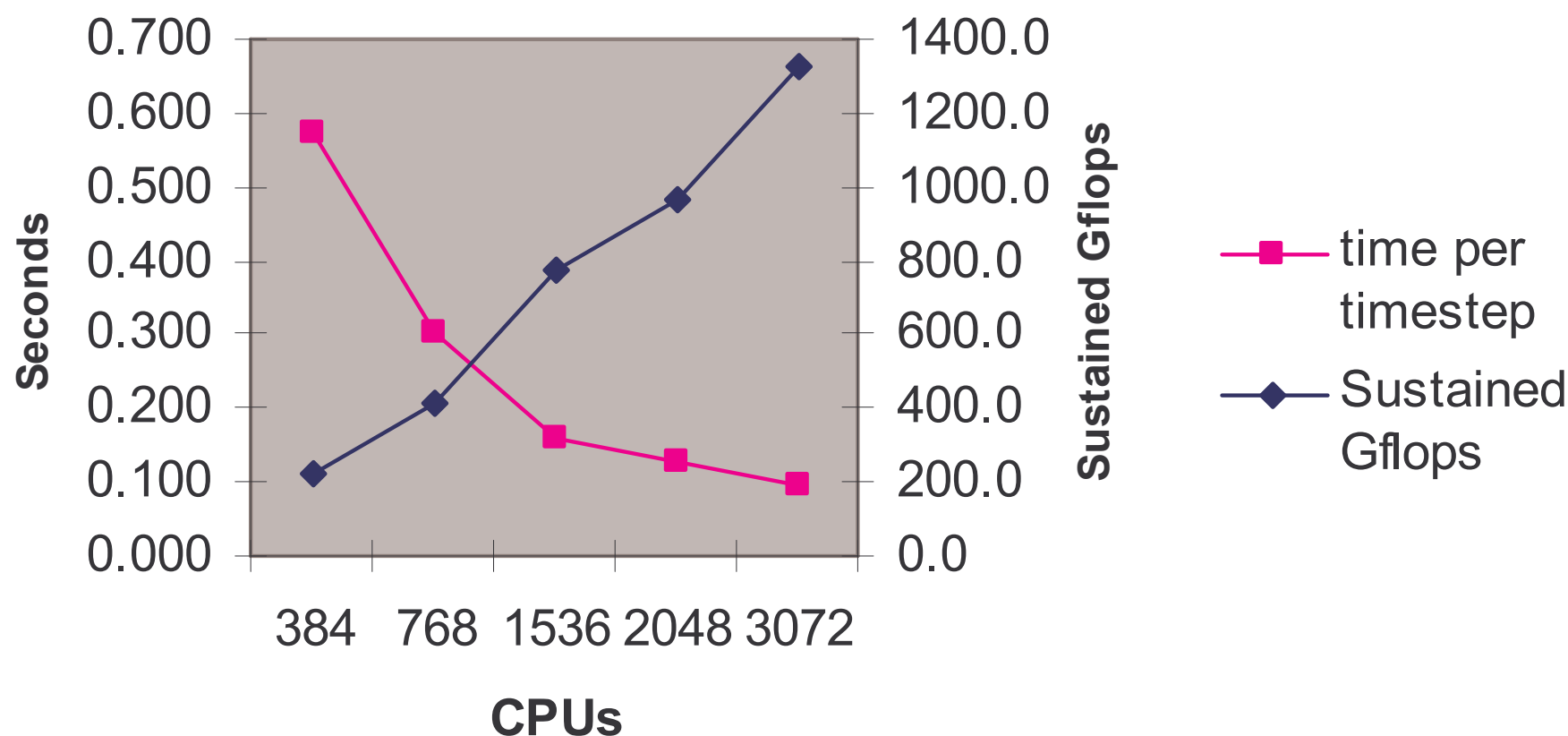


Higher is better

ECHAM5 T255L60 Performance on XT3

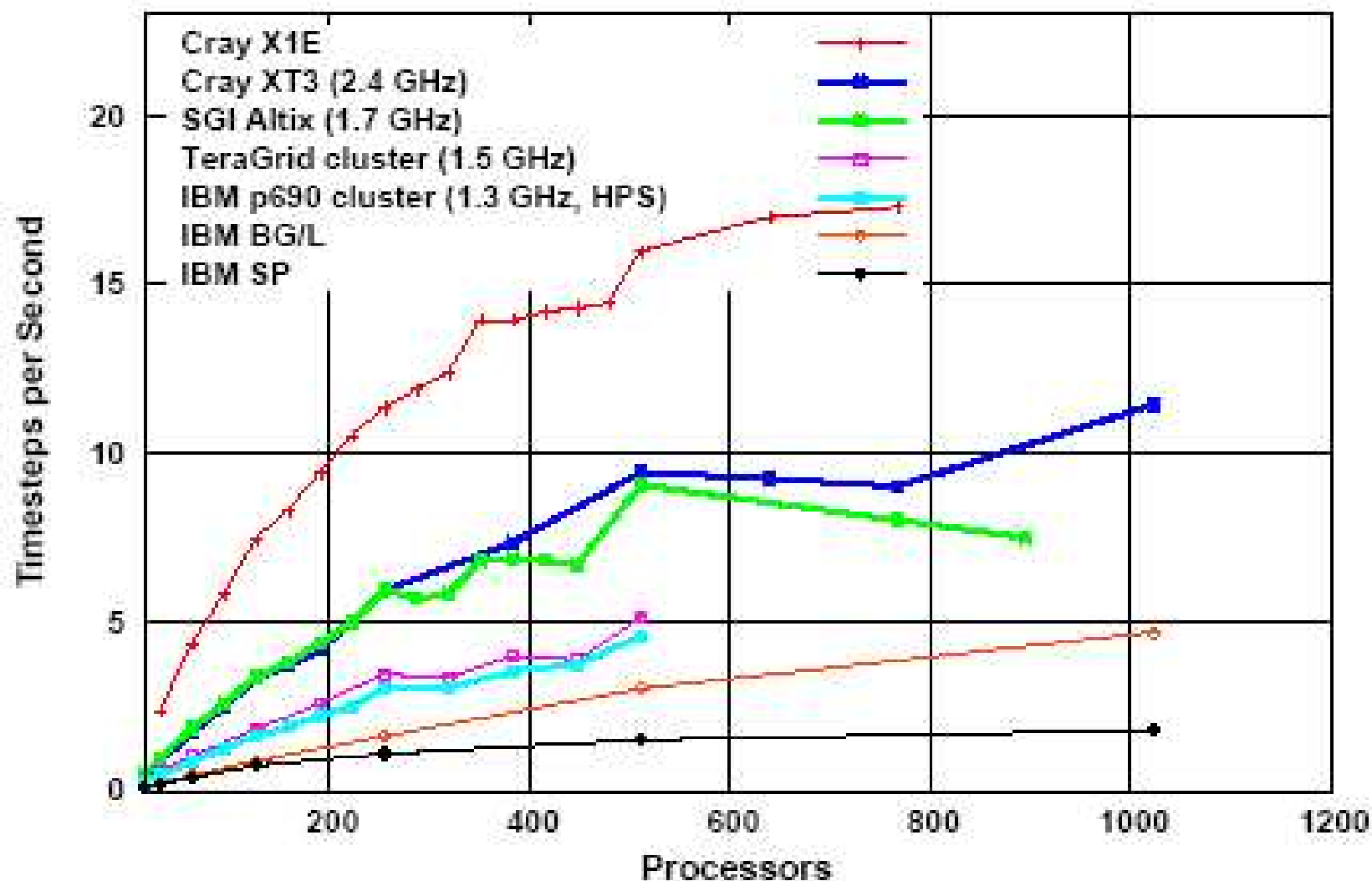


Max-Planck-Institut für Meteorologie
MaxPlanck Institute for Meteorology



GYRO B1

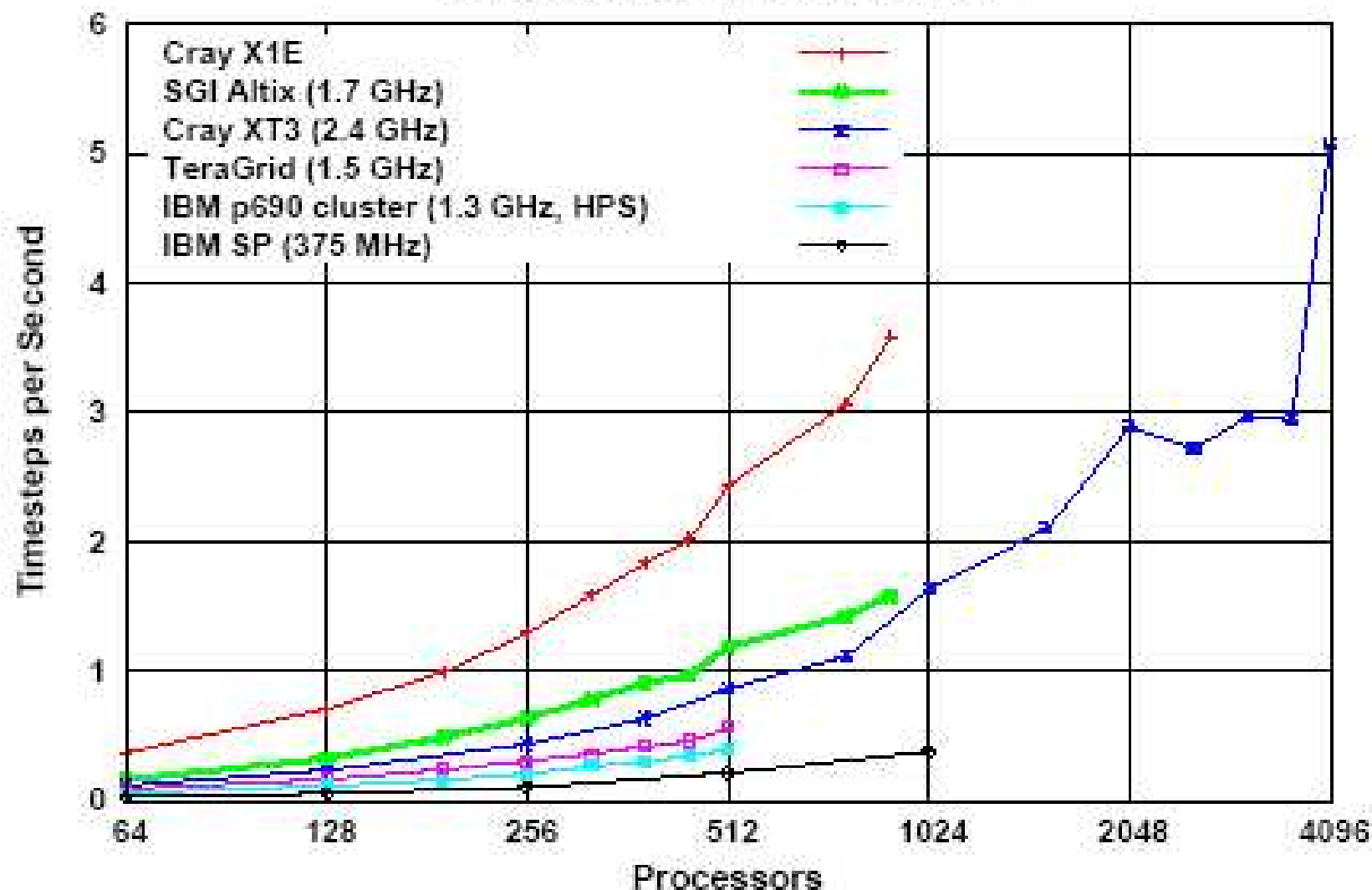
GYRO performance for B1-std



Higher is better

GYRO B3

GYRO performance for B3-gtc



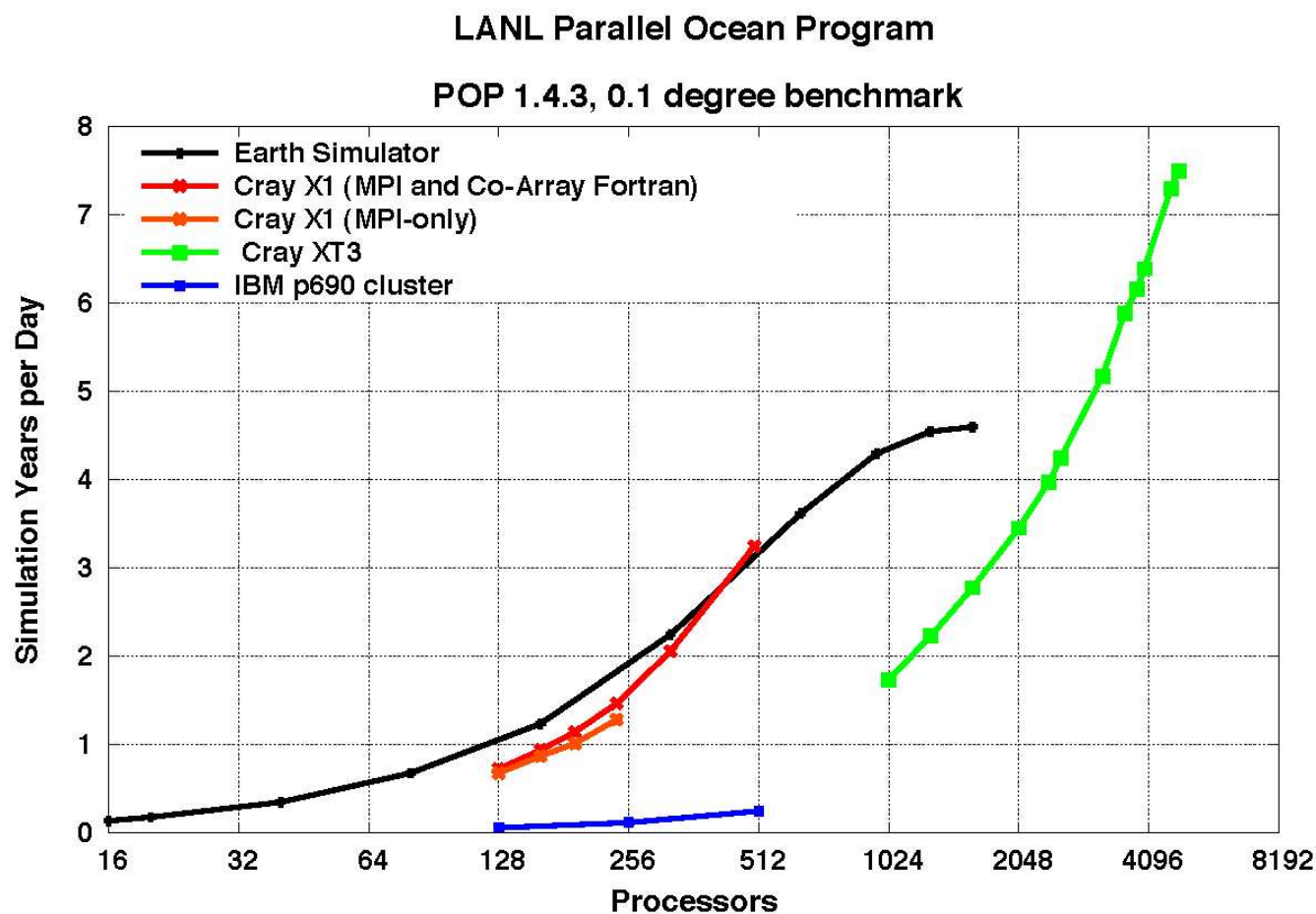
Higher is better

Conclusions

- We believe that we made the right choices
 - AMD Opteron provides excellent memory latency and injection bandwidth
 - SeaStar network gives extremely good network balance
 - Microkernel gives excellent scalability by reducing OS jitter.



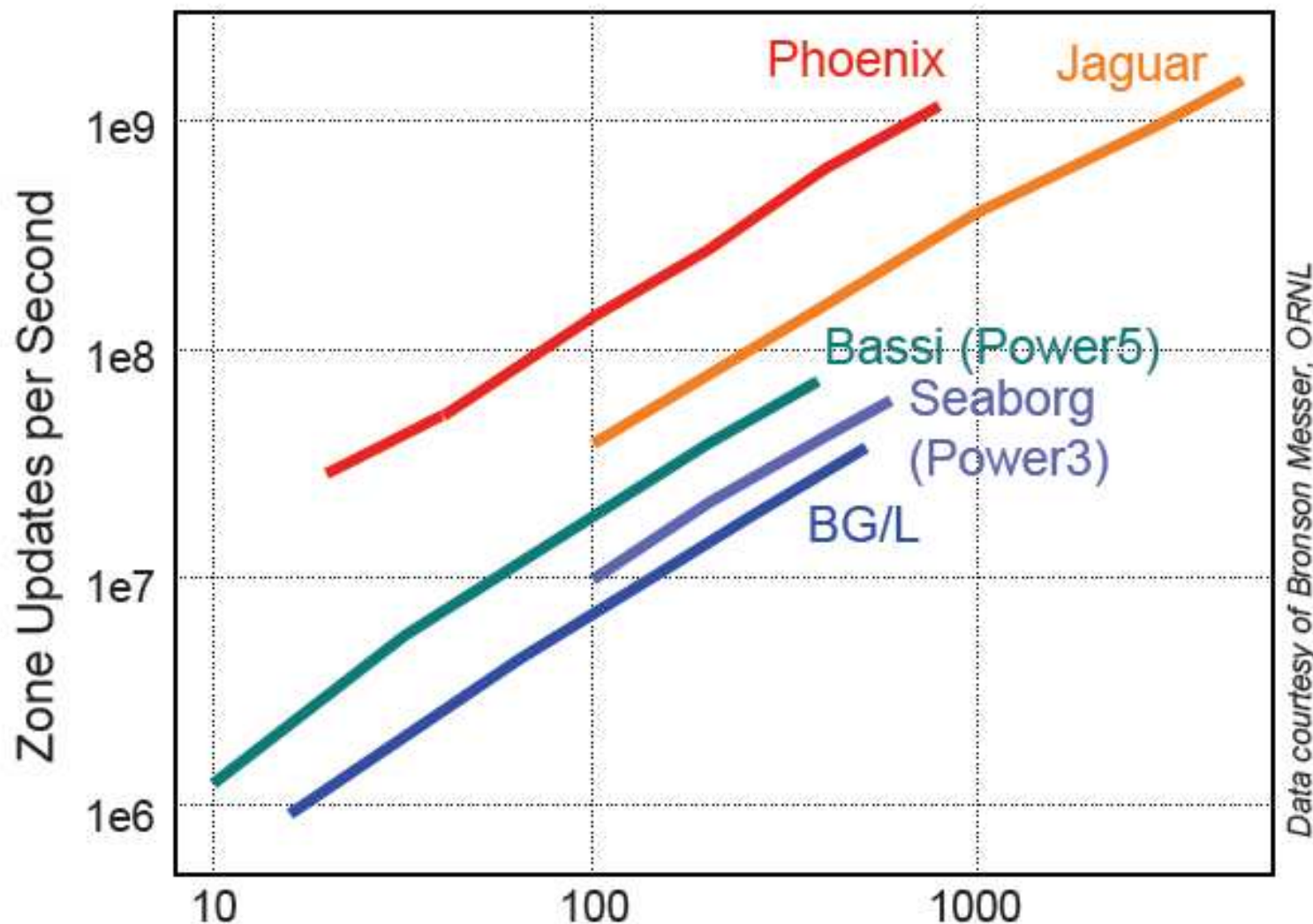
Extra 1



Higher is better

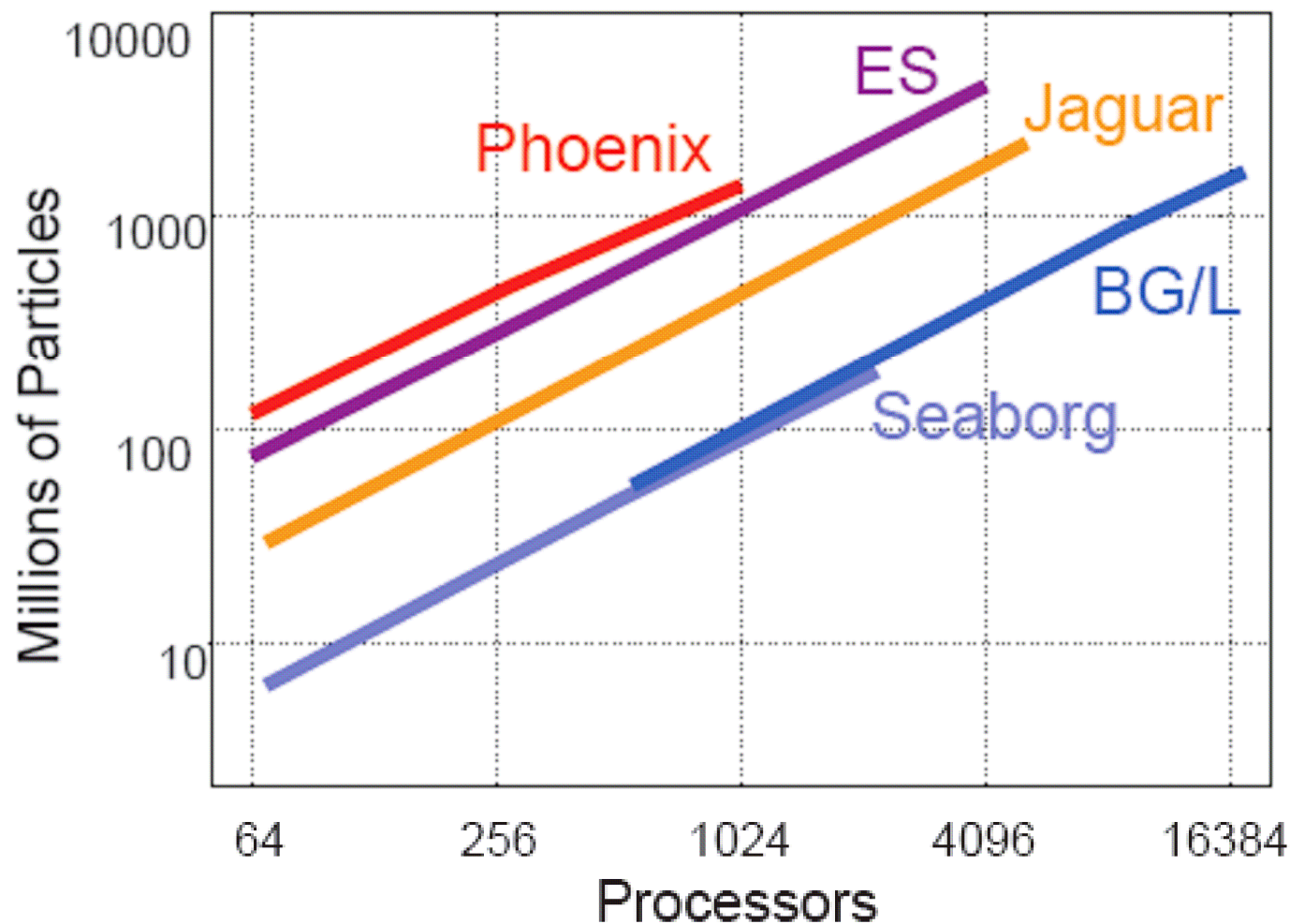
Extra 2

VH1 weak-scaling benchmark



Extra 3

GTC weak-scaling benchmark



Data courtesy of Stephane Ethier, PPPL