



Penguins Can Fly - The Altix Servers

Steve Caruso
HPC Systems Engineer
scc@sgi.com

Linux Users' Group of Davis

April 19, 2004

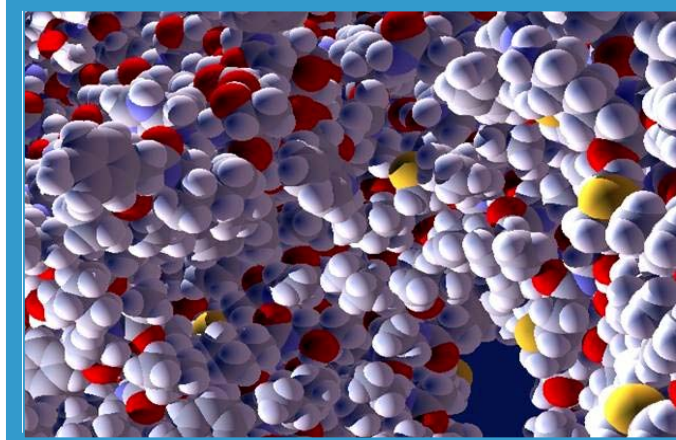


- **Intro**
- **Altix System Architecture**
- **Shared vs Distributed Memory (Clusters)**
- **Altix Linux[®] Environment**
- **Roadmap**

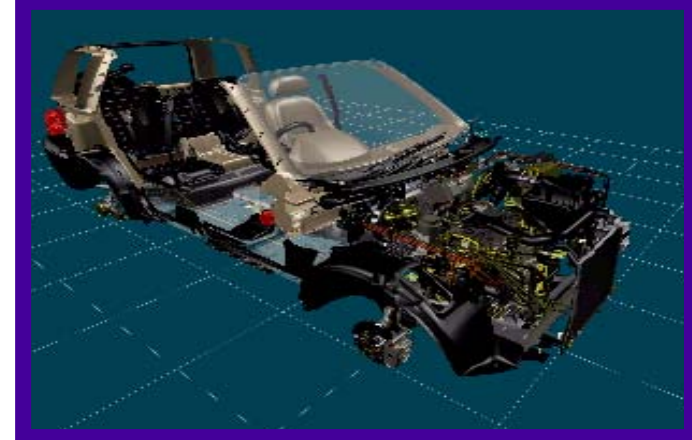
Silicon Graphics' Target Markets



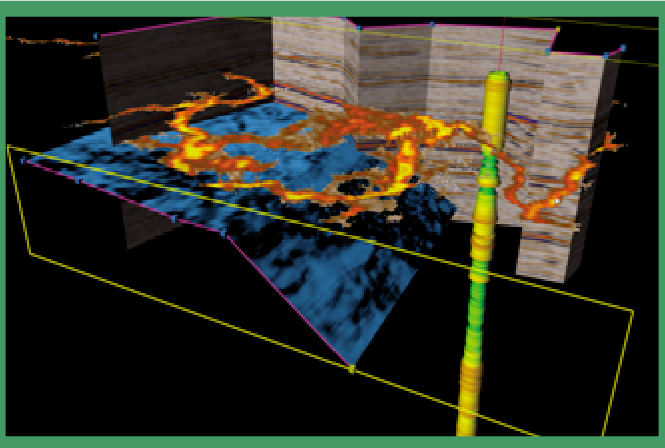
Defense



Science



Manufacturing



Energy

SGI delivers high-performance computing, storage, and visualization solutions that address scientific, engineering, and creative challenges.



Media

Strategic Focus Areas



High Performance Computing



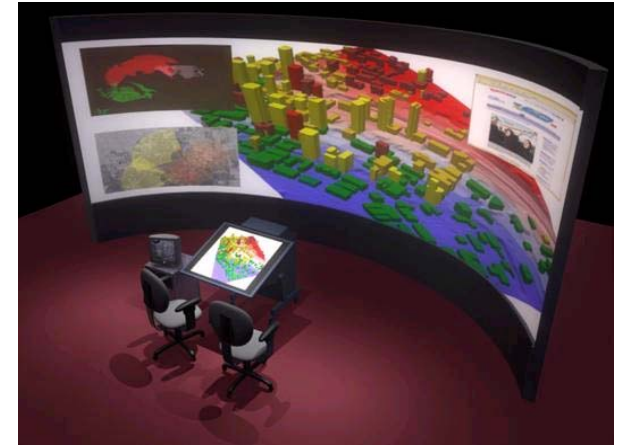
- SGI® Altix™ family
- SGI® Origin® family
- NUMALink™
- Scalability
- Performance
- Multi-paradigm computing

Storage



- SGI® TP9000 storage arrays
- CXFS™
- DMF/TMF
- SAN and NAS servers
- Data lifecycle management
- Heterogeneous file sharing

Advanced Visualization



- Onyx®
- Tezro™
- SGI Reality Center®
- InfiniteReality®
- InfinitePerformance™
- Visual Area Networking

Choice in Deployment with NUMAflex™ HPC Solutions



NUMAflex Global Shared-Memory Architecture

Balanced, scalable performance

Operating environment optimized for HPC

Low-latency memory access

Easily deployable

MIPS® and IRIX®



**SGI® Origin®
3000**

**SGI® Origin®
300**

Itanium® 2 and Linux®



SGI® Altix™ 3000



**SGI® Altix™
350**

SGI® Altix™ : Scaling Linux® to New Altitudes



Intel® Itanium® 2 processors

- + Linux operating system
- + SGI NUMAflex Shared Memory Architecture

The World's Most Scalable Linux
Supercomputer

Altix Momentum

- Altix 3000 introduced January, 2003 & Altix 350 in January, 2004
- First Linux server to scale a single kernel to scale beyond 32 processors
- 256-processor single system image supported configuration
- 512-processor single system image demonstrated
- Delivering a 32P system with 4TB shared memory
- Over 10,000 processors shipped worldwide to 100+ customers
- 11 Systems in the latest Top500 listing
- World record benchmarks including STREAMS performance in excess of 1TB/sec; Fastest Linux IO performance: 7 GB/sec
- Growing applications momentum
- SUSE LINUX distribution agreement
- Graphics solutions in development
- InfiniteStorage data management solution stack complete
- Delivering the benefits of shared memory, scalability, and open source software to the HPC community

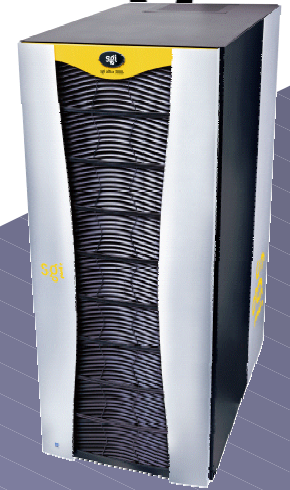


SGI® Altix™ Family

sgi®

SGI® Altix™ 3700

- 4-256P SSI per system
- Up to 8TB of shared memory
- Expandable to 1000+P supercluster
- Infinitely scalable using commercial interconnects



SGI® Altix™ 3300

- 4-12P SSI per system
- Upgradeable to SGI Altix 3700



SGI® Altix™ 350

- 1-16 P SSI per system
- 2-192 GB max memory per system
- Intel® Itanium® 2 and LV Intel® Itanium® 2 processor options



Unified architecture
Scales from 2P to 1000s of processors
Excellent price and performance at every level

Departmental HPC & technical
database applications

Entry-level HPC server

World's most complex,
demanding HPC systems.

Representative Altix Customers



Nationally Funded Supercomputing/GRID Centers

SARA – Dutch National Center	416P Altix
University of Manchester	256P Altix
University of Queensland/QUAKES	225P Altix
University of Cambridge	152P Altix

US Federal R&D Agencies

NASA	512P Altix	CFD, Climate
Oak Ridge	256P Altix	Biology, Environment
PNNL	128P Altix	Chemistry, Biology
NRL	128P Altix	CFD, Climate, Chem
NCI	64P Altix	Bioinformatics/Proteomics

State Funded Supercomputing Centers

Minnesota Supercomputing Center	52P Altix	Chemistry, Proteomics
Ohio Supercomputing Center	32P Altix	Life, Physical Sciences
North Carolina Supercomputing Center	32P Altix	Life, Physical Sciences

Industry

Total – France Oil & Gas	256P Altix	O&G Exploration
Marathon Oil	64P Altix	O&G Exploration
GM	64P Altix	CAE Codes
Honda Americas	32P Altix	CAE Codes
GE Power	32P Altix	CAE Codes
Boeing	8P Altix	CFD Codes

Representative Altix Customers



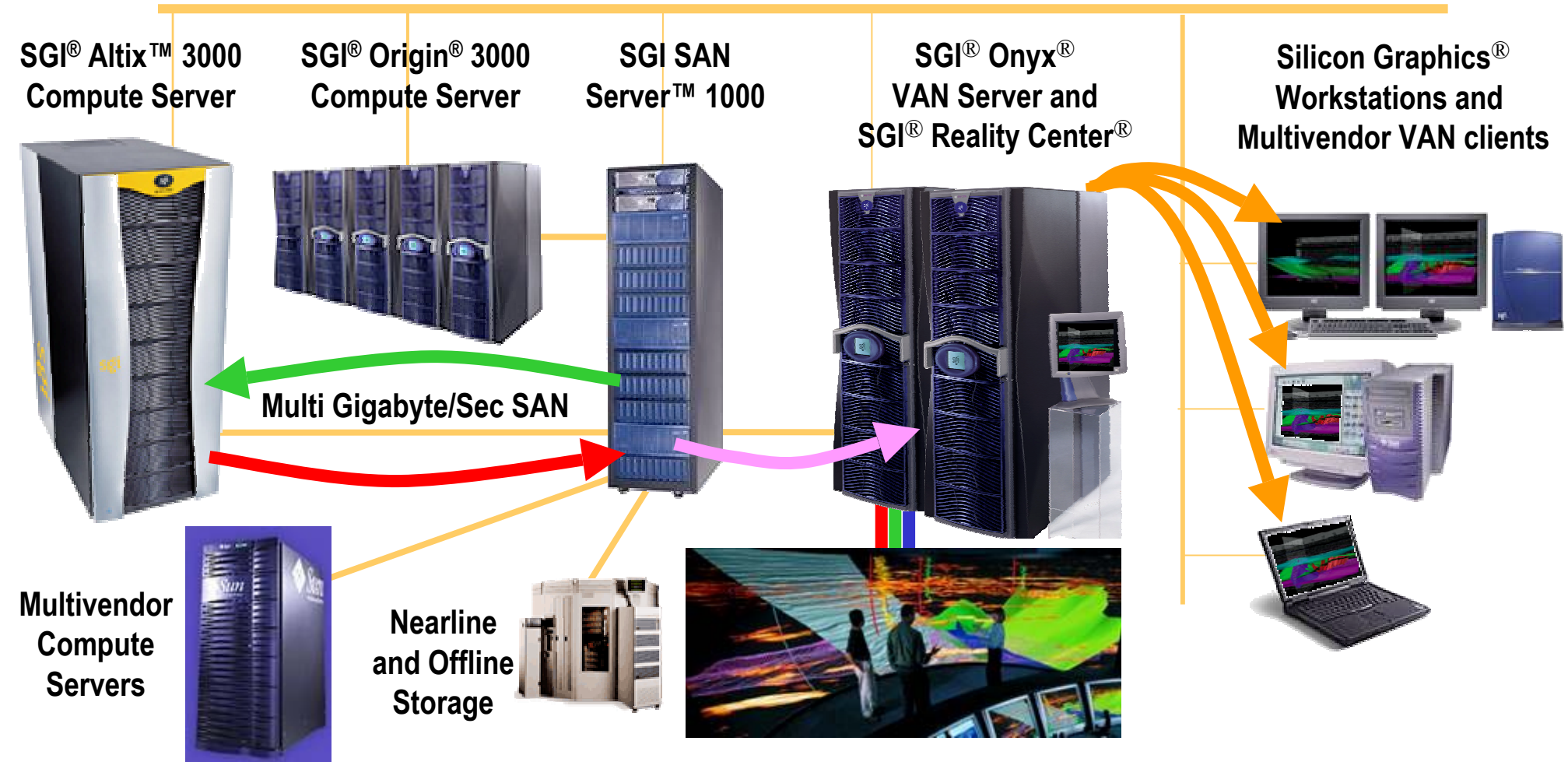
Universities

Washington University	128P Altix	Bio, Chem, Astrophysics
Denmark Technical University	128P Altix	Bio Sequence Analysis
University of Tokyo	108P Altix	Earthquake, Biochem
University of Queensland	64P Altix	Earth, Chem, Life Sci.
Weizmann Institute of Science	44P Altix	Physical Chemistry
Wichita State University	32P Altix	CAE, CFD, Chemistry
University of Nevada – Reno (Desert Research)	32P Altix	Weather & Climate
University of Hawaii – IPRC	32P Altix	Weather & Climate
University of Washington	24P Altix	Astrophysics
U Wisconsin, Madison	24P Altix	NMR Research, Chem
Cornell Weill School of Medicine	16P Altix	Computational BioMed
Georgia Tech	16P Altix	Materials Research
University of Florida	16P Altix	CFD, Aerospace Eng.
Georgia Tech	12P Altix	Atmospheric Science
North Dakota State University	12P Altix	Materials Sciences
Harvard University	12P Altix	Earth Sciences
Memorial Sloan Kettering Cancer Center	12P Altix	Cancer Research
Virginia Tech	12P Altix	Computer Science
University of Southern California	8P Altix	Rendering
Massachusetts Institute of Technology (M.I.T.)	8P Altix	Optics Research
Stanford University	8P Altix	Earthquake, Physics
Wesleyan University	4P Altix	Chemistry

Altix in Today's SGI Environment



Integration of VAN, SAN, and Altix technology allows users to run their simulations, manage their data, and visualize their results faster and more effectively than ever.



Altix Platform Intro

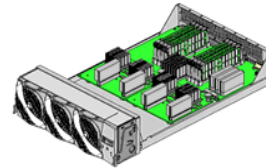
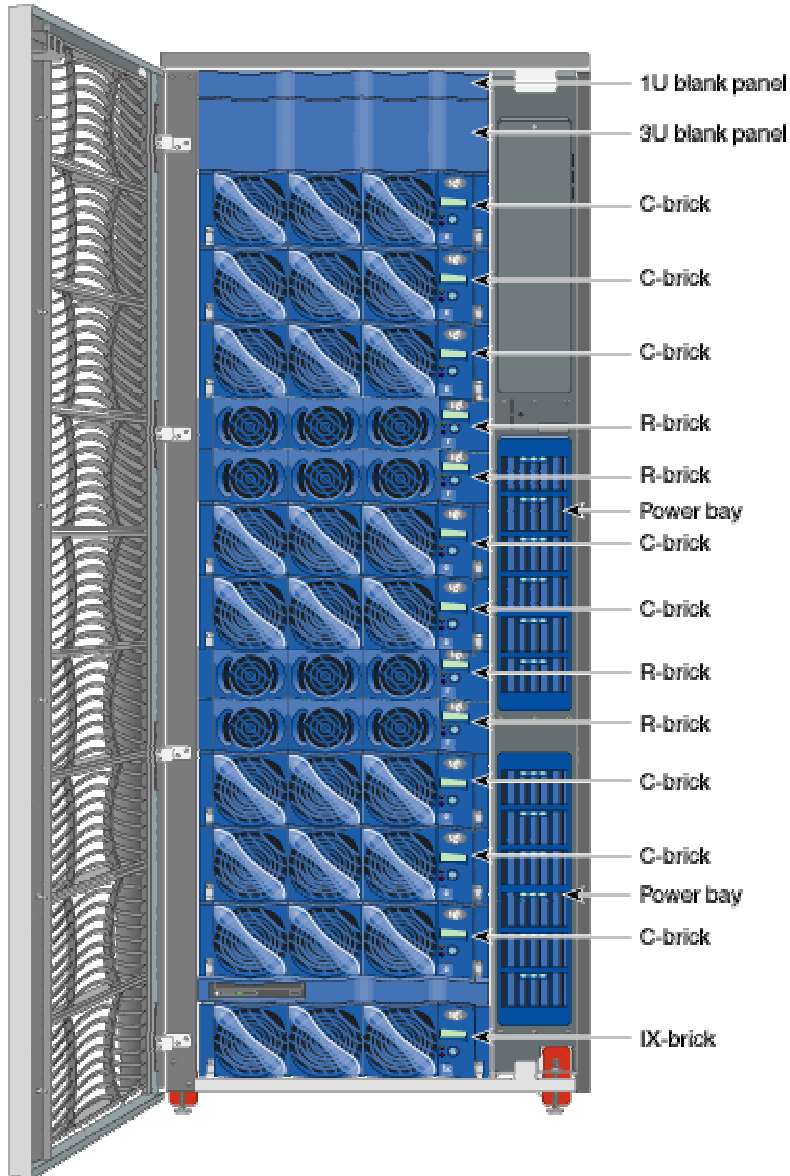
Altix System Architecture

Shared vs Distributed Memory (Clusters)

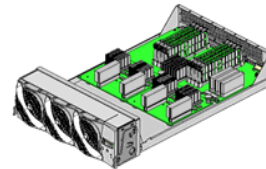
Linux[®] Environment

Roadmap

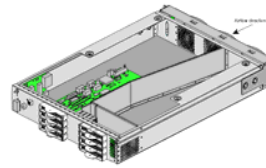
SGI® Altix™ 3000 Hardware Overview



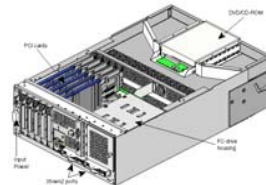
Itanium2™ C-brick
CPU and memory



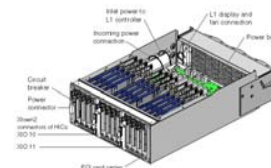
M-brick
Memory



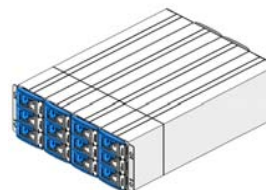
R-brick
Router interconnect



IX-brick
Base I/O module



PX-brick
PCI-X expansion



D-brick2
Disk expansion

Altix™ 3700 Full Rack Configurations



16–512P configurations

- 1.3 GHz/3MB Itanium® 2
- 1.5 GHz/6MB Itanium 2

Memory configuration

- Commodity PC2100 DDR ECC Memory
- 256GB @64P, 2TB @512P with 512MB DIMMS
- 512GB @64P, 4TB @512P with 1GB DIMMS

Dual plane NUMALink™ interconnect

- “Fat tree” topology

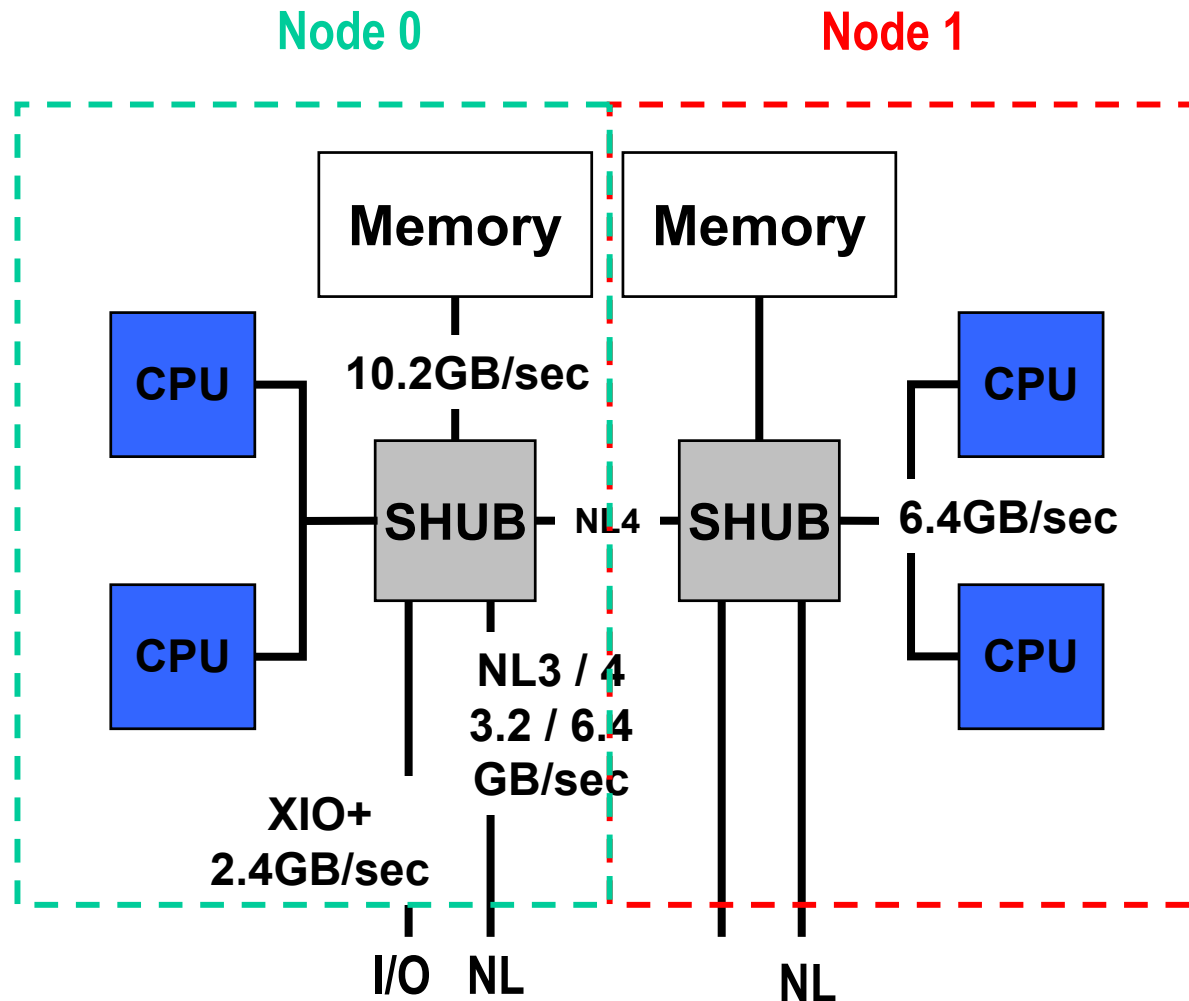
I/O configuration

- $IX + PX \leq 8$ per 64P SSI
- Up to 94 PCIX slots across 47 buses
+ 1 PCI slot

M-brick support

- Up to 15 M-bricks in a 4P configuration
- $M+C \leq 16$
- 128GB per processor max

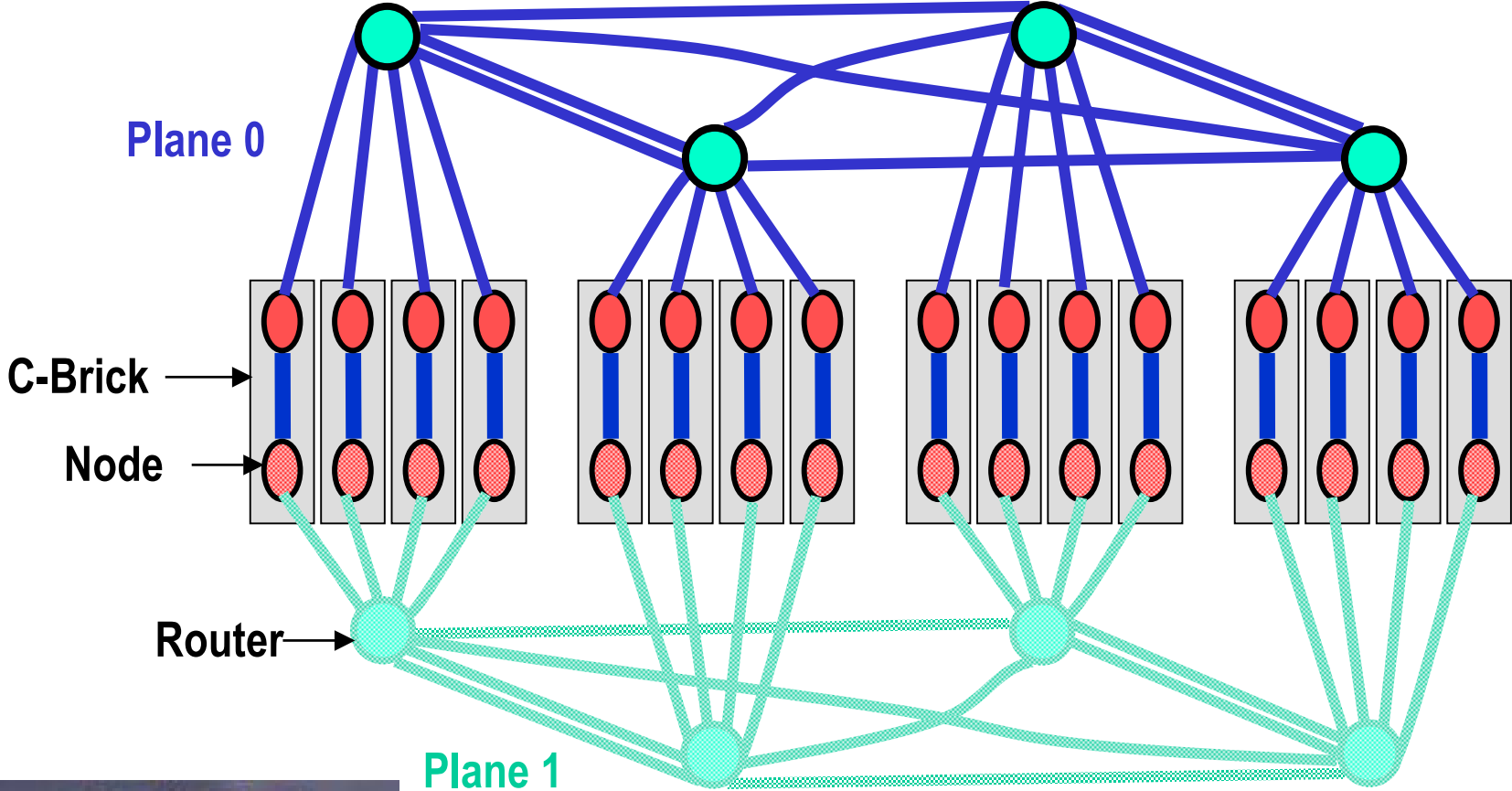
Itanium 2 Based C-brick



NUMALink Interconnect: Dual-plane, fat tree topology

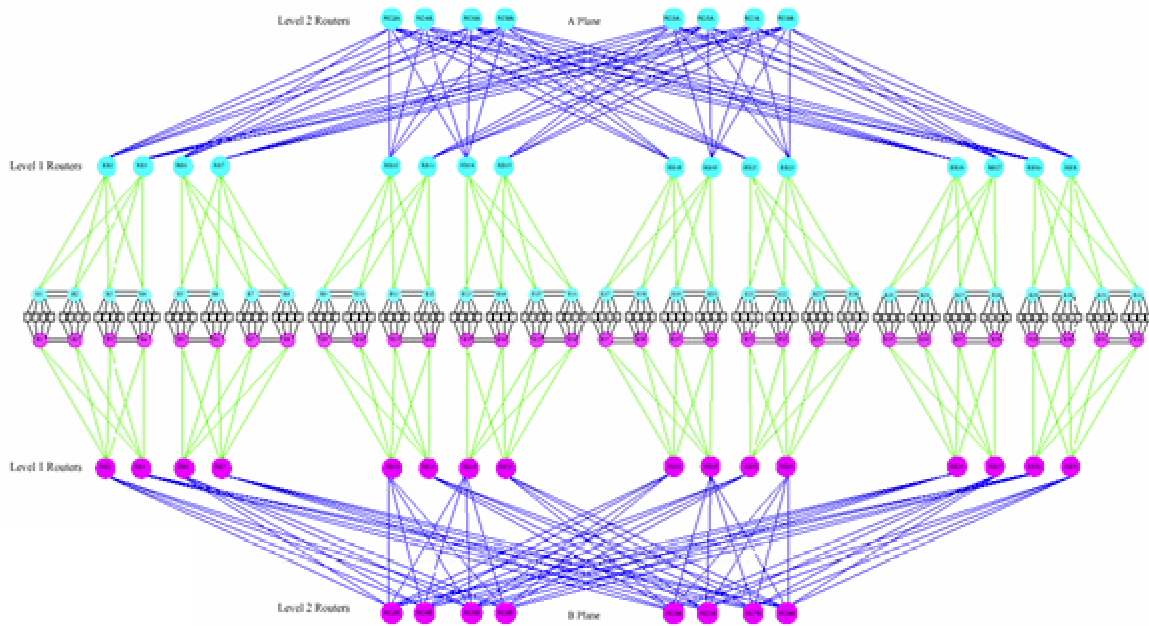


NUMALink Interconnect Fabric



64 PROCESSOR TOPOLOGY

512P Configuration Detail



L2 Controller	L2 Controller	L2 Controller	L2 Controller	L2 Controller	L2 Controller	L2 Controller	L2 Controller
IX	RB23A RB7A	RB22A RB6A	RC5A RC3A	RC9A RC7A	RB30A RB14A	RB31A RB15A	IX
C	C	C	C	C	C	C	C
C	C	C	C	C	C	C	C
C	R22A R22B	R24A R24B	R6A R6B	R8A R8B	R14A R14B	R16A R16B	R30A R30B
C	C	C	C	C	C	C	C
C	C	C	C	C	C	C	C
C	R21A R21B	R23A R23B	R5A R5B	R7A R7B	R13A R13B	R15A R15B	R29A R29B
C	C	C	C	C	C	C	C
C	C	C	C	C	C	C	C
C	C	C	C	C	C	C	C
IX	RB11B RB23B	RB10B RB22B	RC5B RC3B	RC9B RC7B	RB14B RB30B	RB15B RB31B	IX
011	012	013	014	015	016	017	018

L2 Controller	L2 Controller	L2 Controller	L2 Controller	L2 Controller	L2 Controller	L2 Controller	L2 Controller
IX	RB18A RB2A	RB19A RB3A	RC4A RC2A	RC8A RC6A	RB27A RB11A	RB26A RB10A	IX
C	C	C	C	C	C	C	C
C	C	C	C	C	C	C	C
C	R18A R18B	R20A R20B	R2A R2B	R4A R4B	R10A R10B	R12A R12B	R26A R26B
C	C	C	C	C	C	C	C
C	C	C	C	C	C	C	C
C	R17A R17B	R19A R19B	R1A R1B	R3A R3B	R9A R9B	R11A R11B	R25A R25B
C	C	C	C	C	C	C	C
C	C	C	C	C	C	C	C
IX	RB20 RB18B	RB3B RB19B	RC4B RC2B	RC8B RC6B	RB27B RB11B	RB26B RB10B	IX
001	002	003	004	005	006	007	008

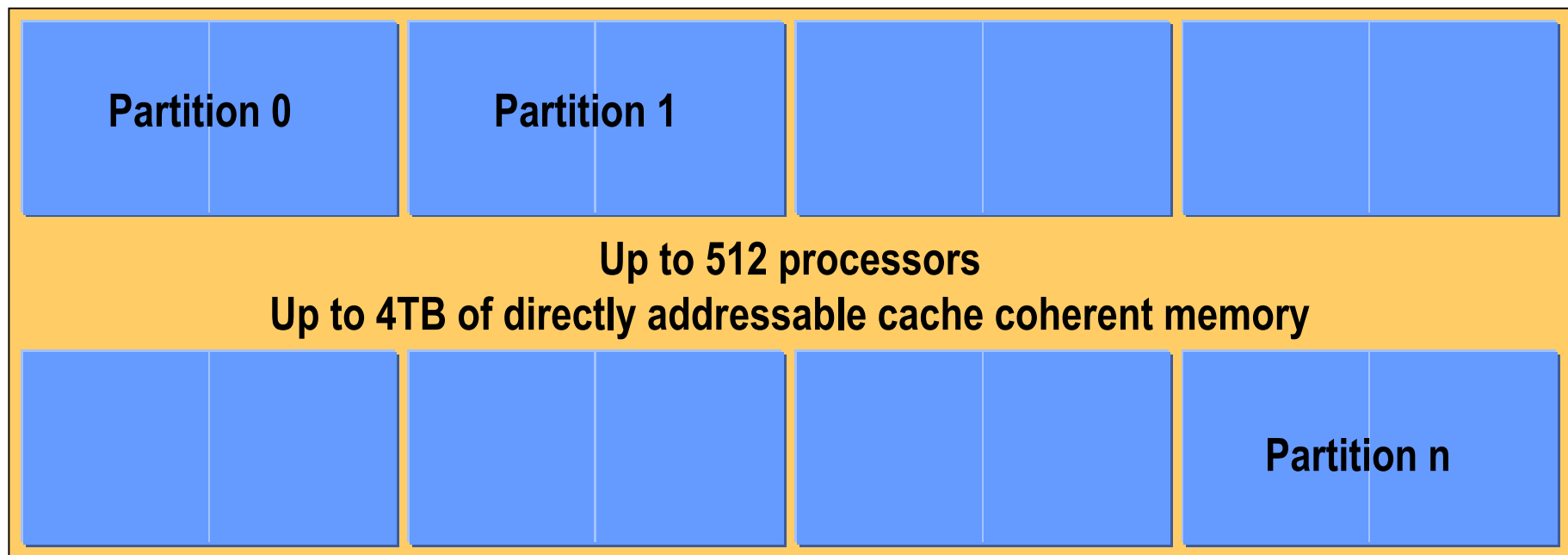
Altix System Nomenclature



Single NUMAlinked system is software-partitioned into multiple sub-systems or partitions each running its own copy of the OS

Total system memory can be shared and globally adressed from any partition

SGI Marketing calls the whole thing a supercluster, the architecture NUMAflex



Special Engineering Project 512P SSI Altix



- NASA Ames and SGI have a 20-year history of systems collaboration
- 1024 processor Origin 3000 currently installed at NASA
- 512 processor Altix 3700 with 1TB of memory installed in late-October
- Outstanding customer applications performance
- World Record STREAM Triad benchmark result, first system to break 1,000 GB/sec
- 2.4 TB/sec LINPACK NxN Rmax result would rank #26 on current top 500
- 80% efficient on LINPACK NxN
- Lowest processor count of any system in the top 25

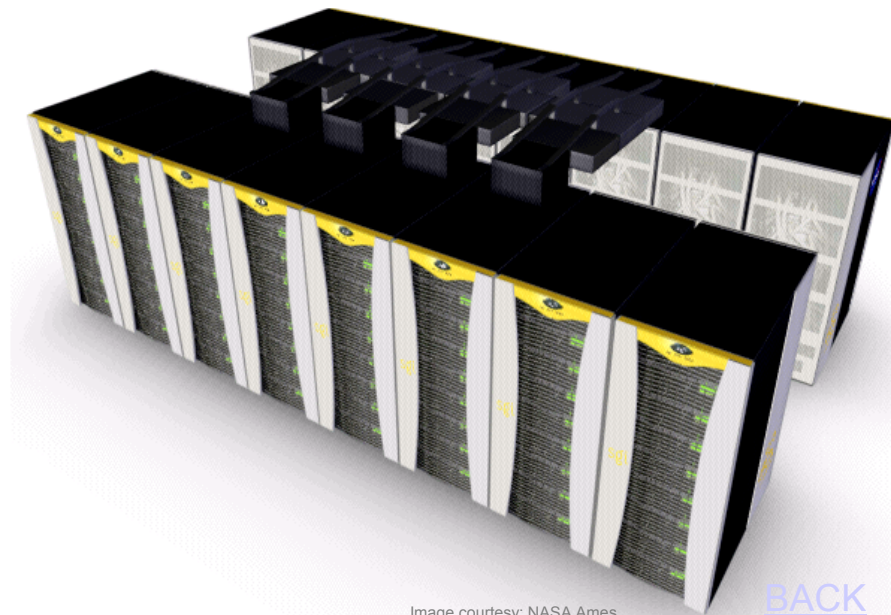


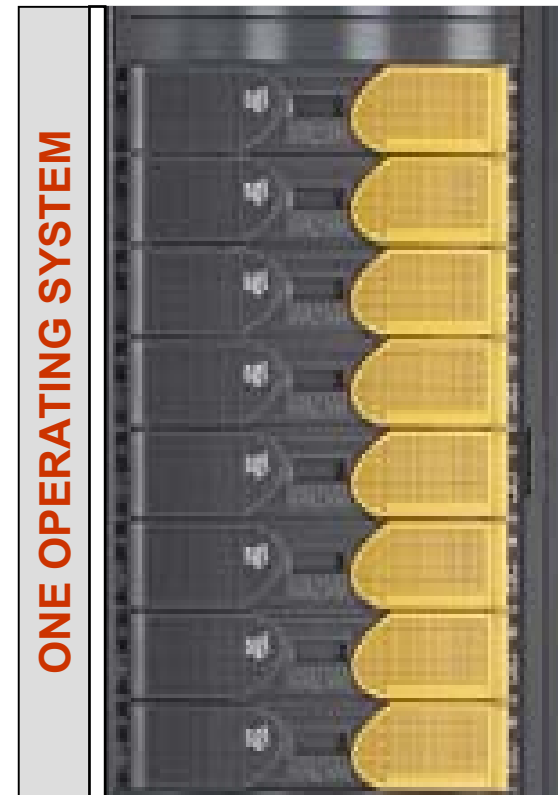
Image courtesy: NASA Ames

[BACK](#)

The SGI® Altix™ 350



- Good performance & scalability from 1 to 16 processors
- Good price/performance mid-range server
- Ideal for departmental application servers, technical database, throughput clusters
- Incrementally scale I/O, processors, memory
- One Linux instance to manage
- NUMALink 4 ring topology
- Cluster to 1000s of processors using industry-standard interconnects



1 - 16P / 2-192GB

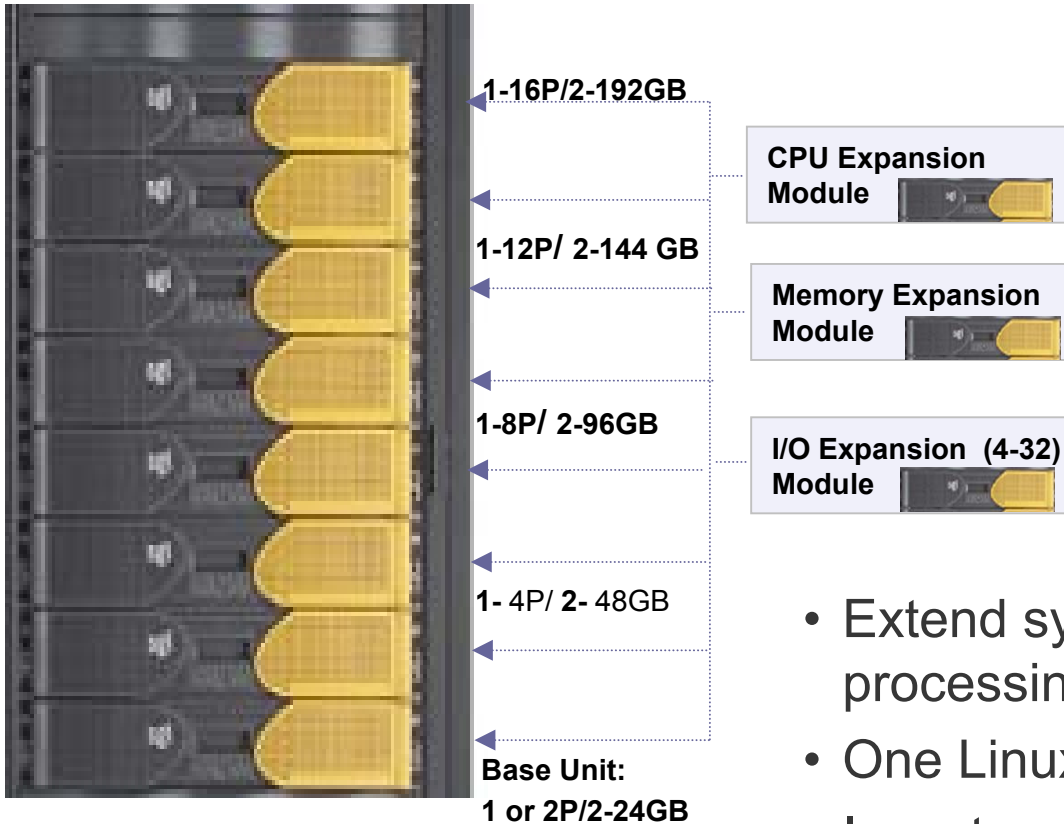
1- 12P / 2-144 GB

1 - 8P / 2-96GB

Base Unit:
1 or 2P / 2-24GB

System features both Standard and new Low Voltage Intel® Itanium® 2 processors

Altix 350 Expansion Options



- Extend system with I/O, memory, and/or processing power as required
- One Linux[®] instance to manage
- Investment protection & leverage current assets
- Allocate budget and resources to ongoing needs

Altix™ 350 Front & Rear

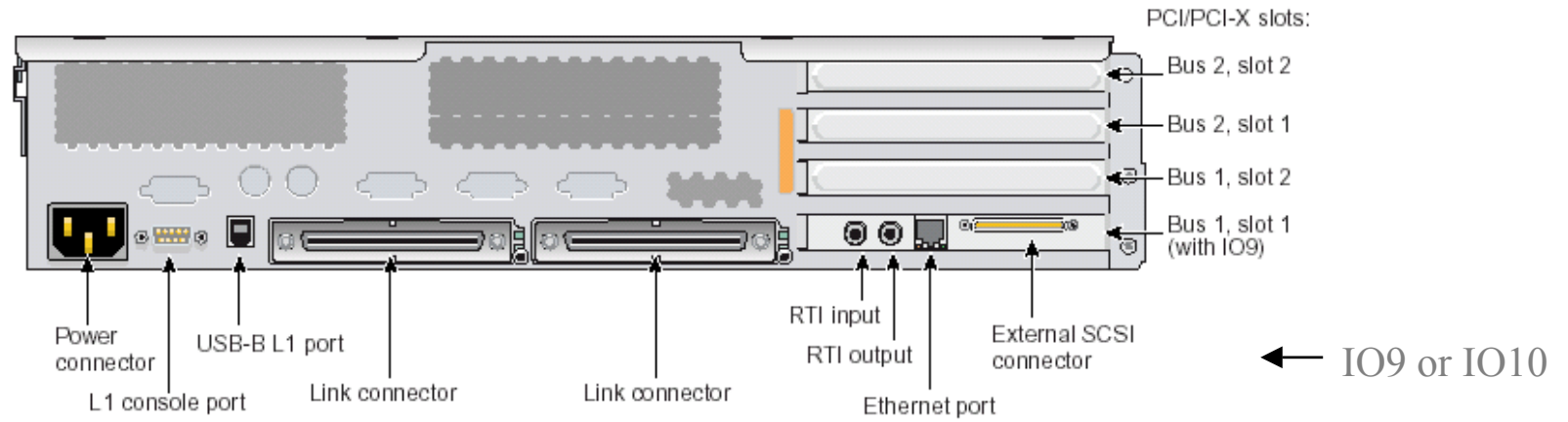
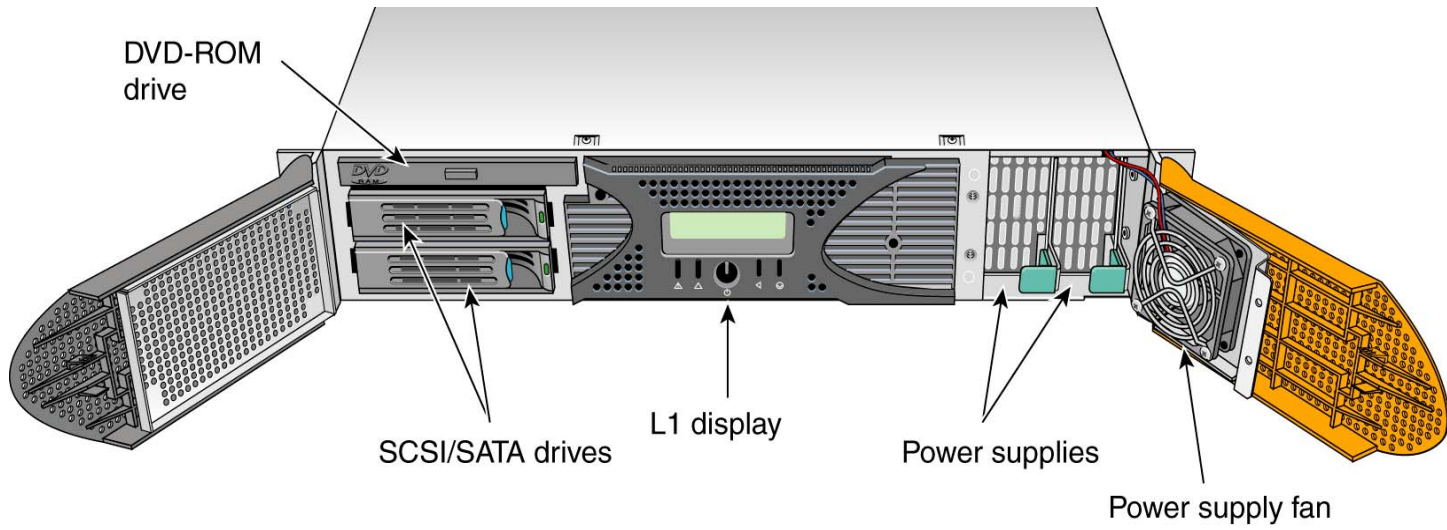


Figure 5-4 Base Compute Module (Rear View)

External Storage Options



HBA interfaces

- 1Gb Fibre Channel, 100MB/sec peak bandwidth
- 2Gb Fibre Channel, 200MB/sec peak bandwidth
- Ultra 160 SCSI, 160MB/sec peak bandwidth
- Gigabit Ethernet copper and optical

JBOD

- SGI® TP900 (Ultra160 SCSI)

RAID

- 1Gb SGI® TP9100 (1Gb Fibre Channel)
- 2Gb SGI TP9100 (2Gb Fibre Channel)
- SGI® TP9400 (2Gb Fibre Channel)
- SGI® TP9500 (2Gb Fibre Channel)

Data servers

- SGI® File Server 830 and SGI® File Server 850 (Gigabit Ethernet)
- SGI SAN Server™ 1000 (1Gb Fibre Channel)

Tape and libraries

- STK L20, L40, L80, L700
- STK 9840, 9940, LTO
- ADIC® Scalar® 100, Scalar® 1000, Scalar® 10000ADIC® AIT

SGI InfiniteStorage Overview



- **Software**

- XFS
- CXFS
- High availability
- Data migration facility
- Backup & Restore

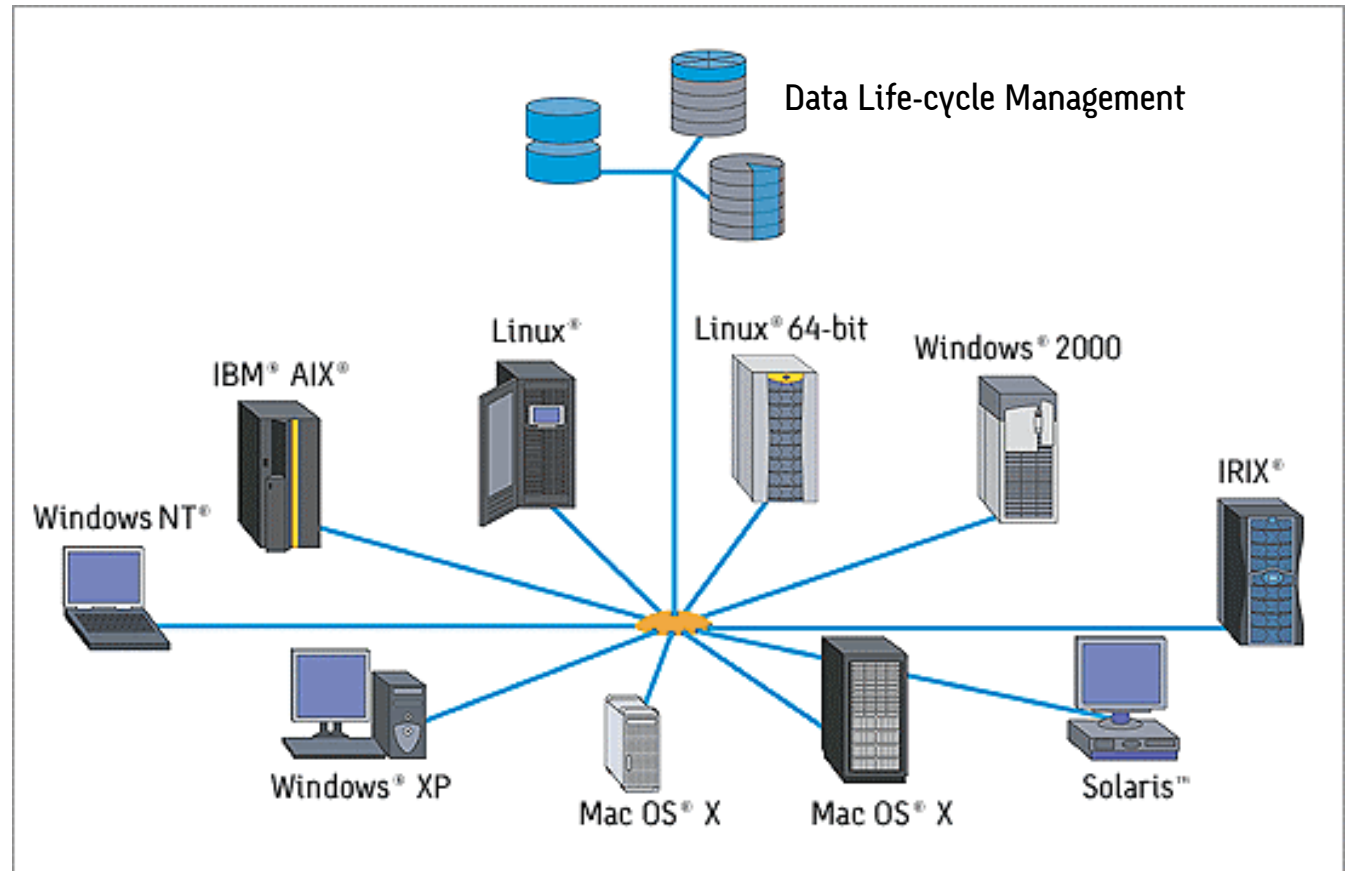
- **Solutions**

- NAS
- SAN

- **Disk Arrays**

- **SAN Infrastructure**

- **Tape Libraries**



Altix Platform Intro

Altix System Architecture

Shared vs Distributed Memory (Clusters)

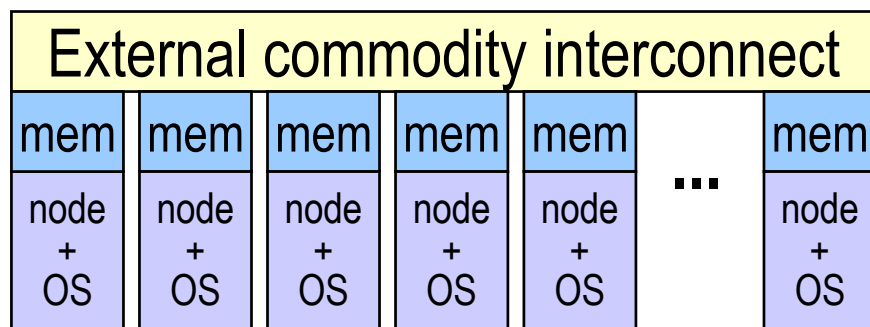
Linux[®] Environment

Roadmap

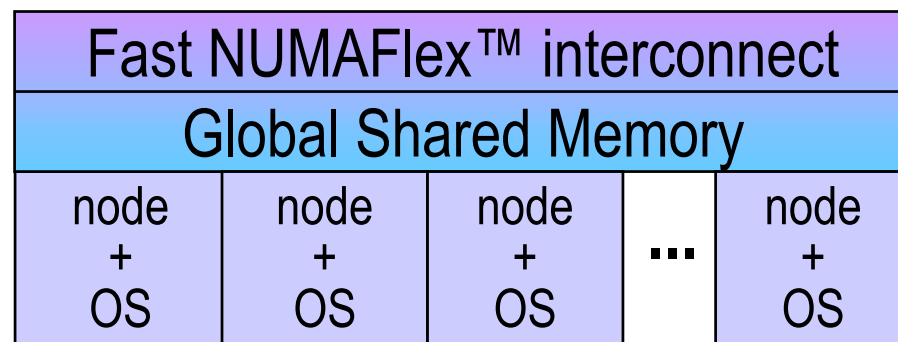
Benefits of Global Shared Memory



Traditional Clusters



SGI® Altix™ 3000



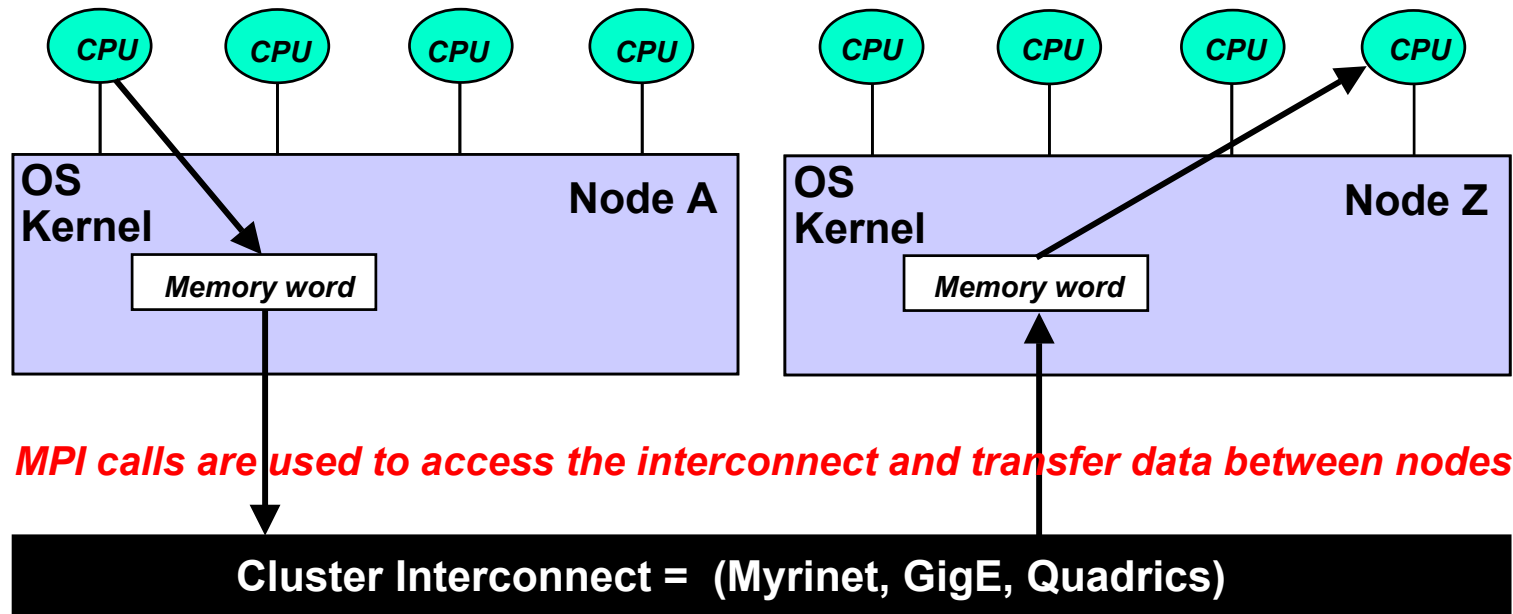
What is global shared memory?

- All nodes operate on one large shared memory space, instead of each node having its own small memory space

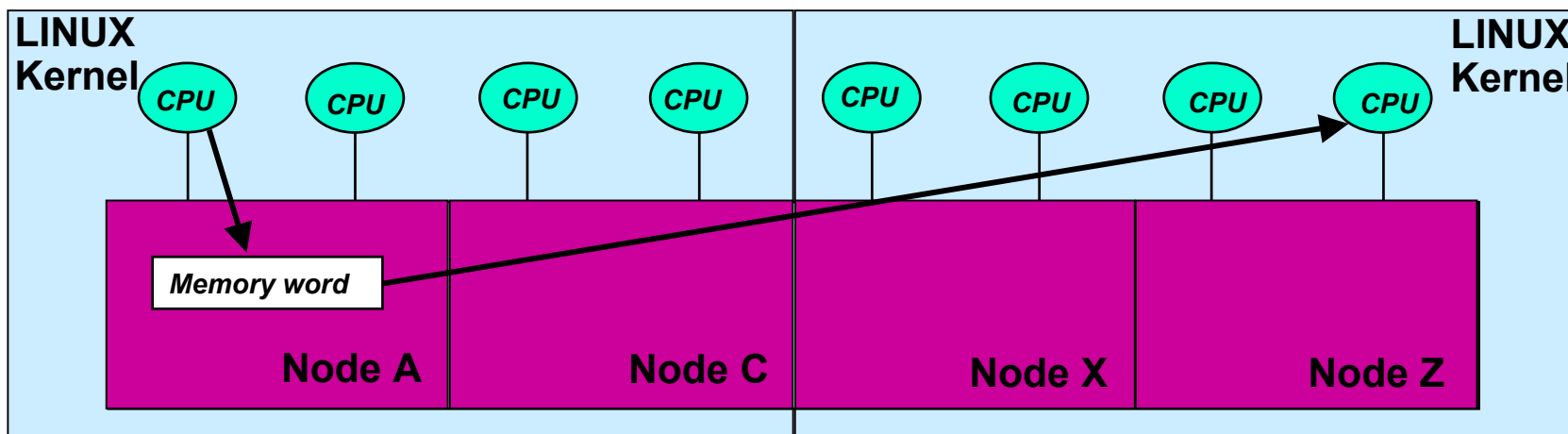
Global shared memory is high performance

- All nodes can access one large memory space efficiently, so complex communication and data passing between nodes isn't needed
- Big data sets fit entirely in memory; less disk I/O is needed
- Global shared memory allows **application performance and scalability**

GSM vs Distributed Memory Cluster



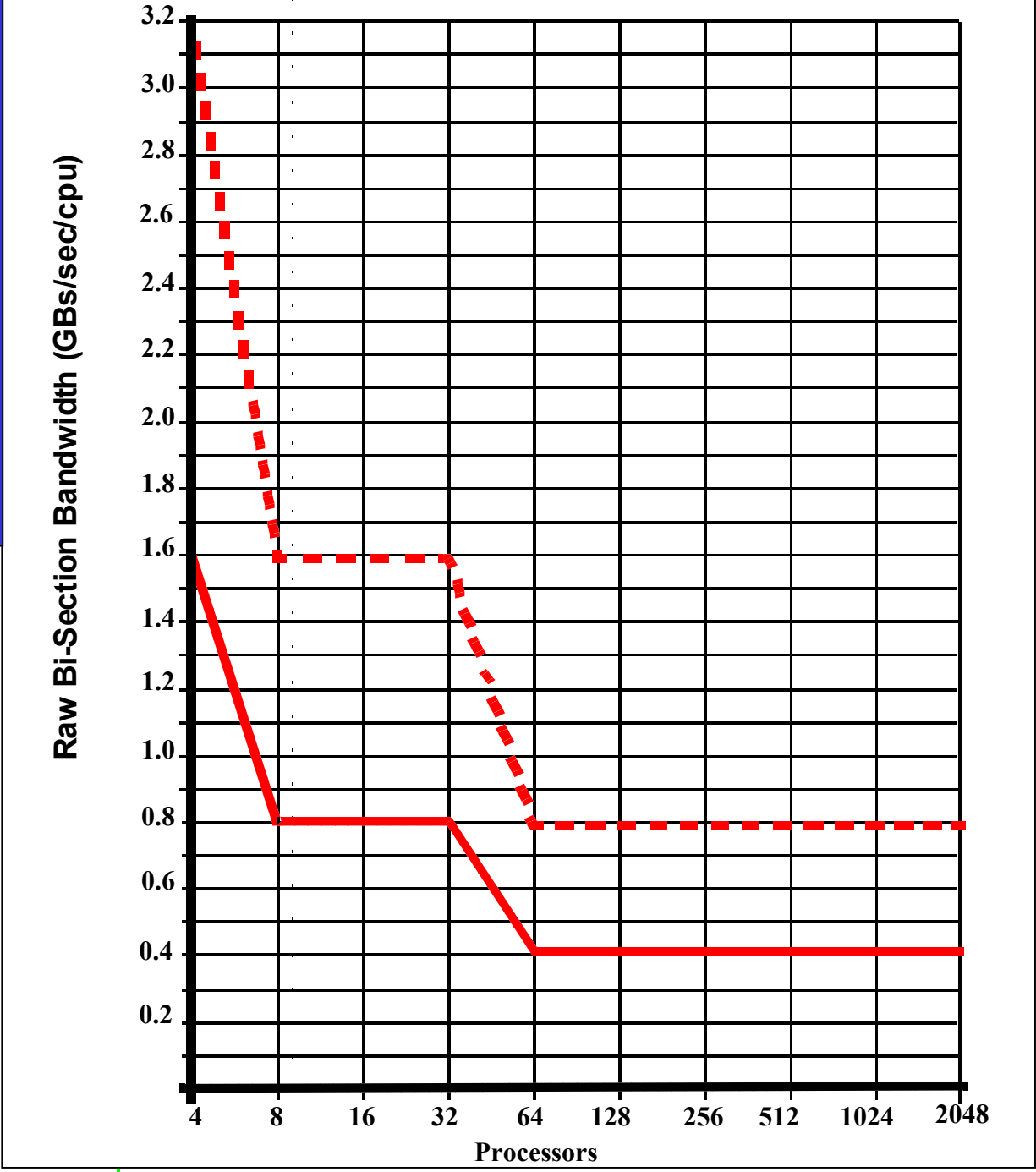
Altix : Memory is global and truly shared and data transfer is simply a memory store then load



Interconnect Topology

Bi-Section Bandwidth Profiles GBs/sec/cpu

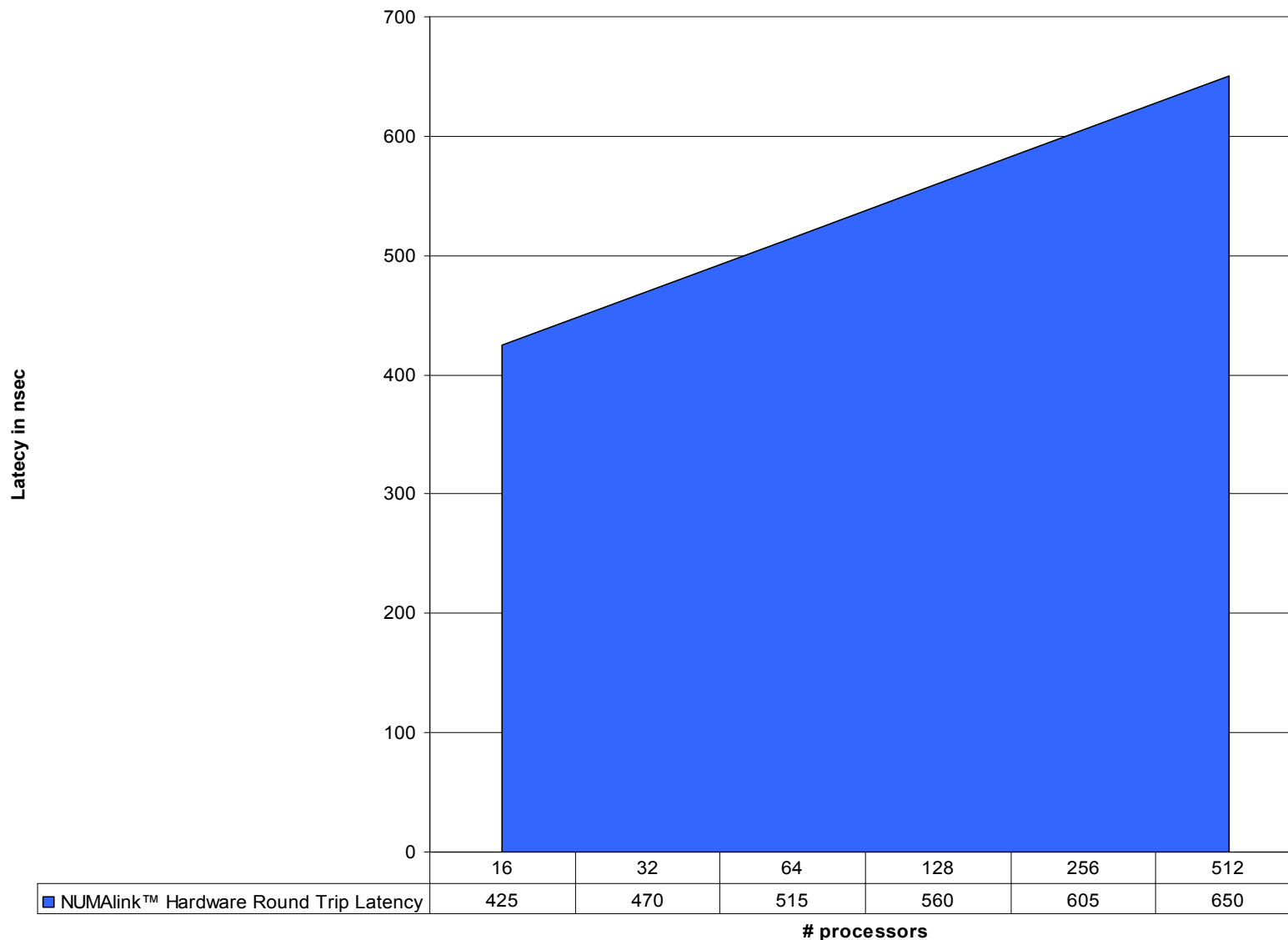
- Dual Plane - NL3 router - 8 port router bricks
- - -** Dual Plane - NL4 router - 8 port router bricks



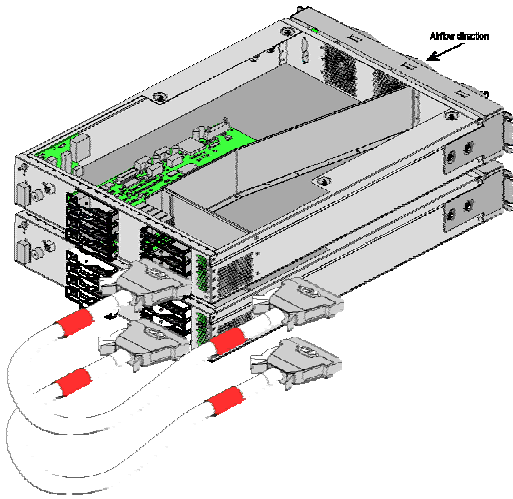
NUMAlink™ 3 Hardware Latency Theoretical Numbers



NUMAlink™ Hardware Round Trip Latency



NUMAlink™ Infrastructure



Communication over NUMAlink3

- 6.4 GB/sec aggregate bandwidth
 - 3.2 GB/sec per plane
 - 1.6 GB/sec per link per direction
- Low hardware memory latency
 - 140–515ns for 64P system
 - MPI partition to partition latency is 135ns
 - Maximum h/w latency for a 512P system from one MPI partition to the farthest partition is 650ns

TCP/IP over NUMAlink for high-speed interpartition communication

Interconnect Comparison	
Interconnect	Aggregate Bandwidth (GB/sec)
NUMAlink 3 Dual Plane	6.4
10GigE	1.25
IBM "Federation Switch"	1.0*
Quadrics	0.680
Myrinet	0.500
1000Bt	0.125

*based on current information

Altix MPI Latencies (off node)



- **MPI-1 (two sided send/receive)**

1.8 microseconds

- **MPI-2 (one sided get/put)**

0.6 microseconds

Numalink Compared to Other Cluster Interconnects



MPI Send/Recv Performance

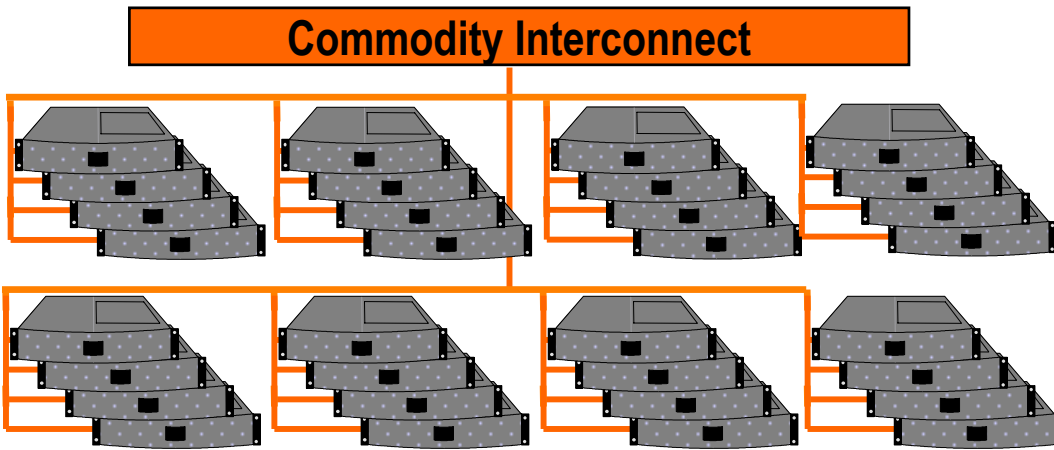
Product	Platform	Bandwidth (1MB xfer)	Latency (4KB)
Numalink 4	Altix	1.50 GB/s	<2 microsec
Numalink 3	Origin3000, 400Mhz	.28	4 microsec
Dolphin	Itanium	.32	14 microsec
Myrinet 2000	PIII, 1GHz, PCI64C	.30	30 microsec

Benefits of Shared Memory: Accelerating Time to Solution



Traditional Clusters

Commodity Interconnect



No support for shared-memory programming models

Large data sets require disk swapping

Overprovisioning of hardware, memory, and software

Load balancing requires communication between nodes

Support for all major parallel programming models

Large data sets fit into shared memory

Larger nodes mean lower total costs

Efficient load balancing; no need to move data

SGI® Altix™ 3000 Family

Fast SGI® NUMAflex™ Interconnect

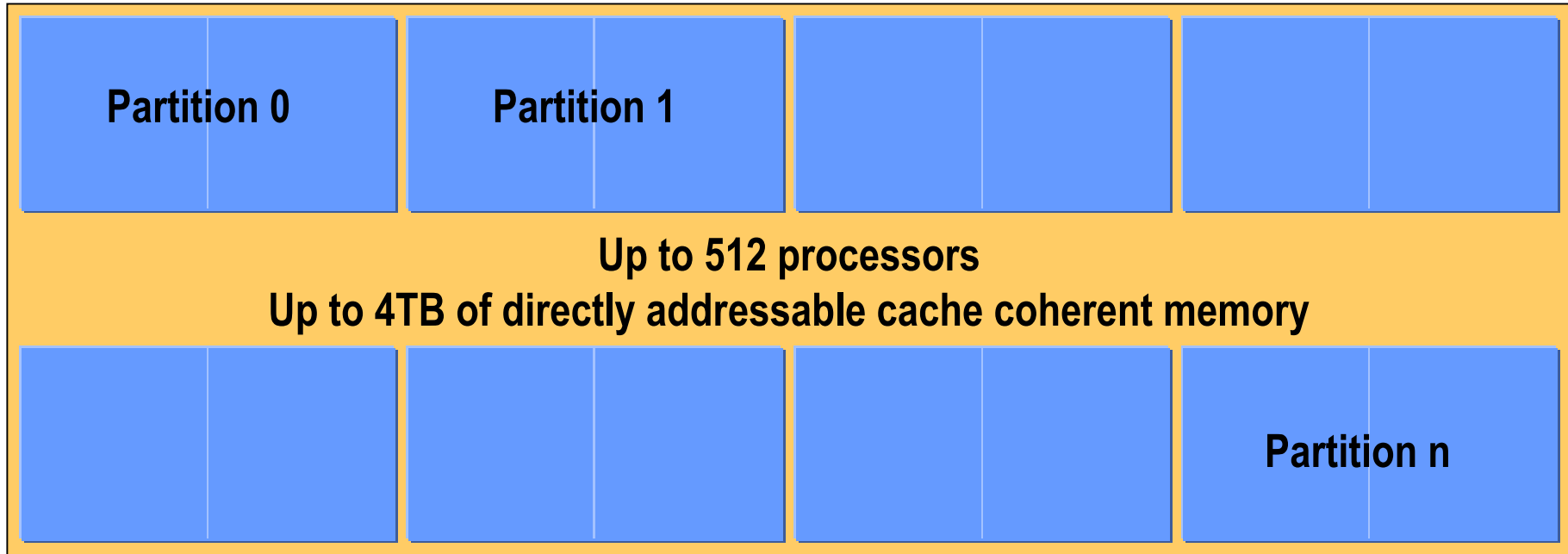
Global Shared Memory



CXFS™ Shared Filesystem

Ultrafast SAN Storage

Benefits of Shared Memory: Multiple Programming Models



Programming models within a partition

- OpenMP™
- MPI
- SHERM
- MPI and SHMEM can utilize put/get

Programming models across partitions

- Hybrid
- MPI
- SHERM
- MPI and SHMEM can utilize put/get
- Global pointers

Benefits of Shared Memory: World-Record Results



Performance, Efficiency, Price/Performance

World-record memory
bandwidth

STREAM Triad

Unsurpassed Linux® scalability
on **real-world applications**

Fastest Linux I/O
performance

7GB/sec

World-record 16, 32, and 64P
compute performance

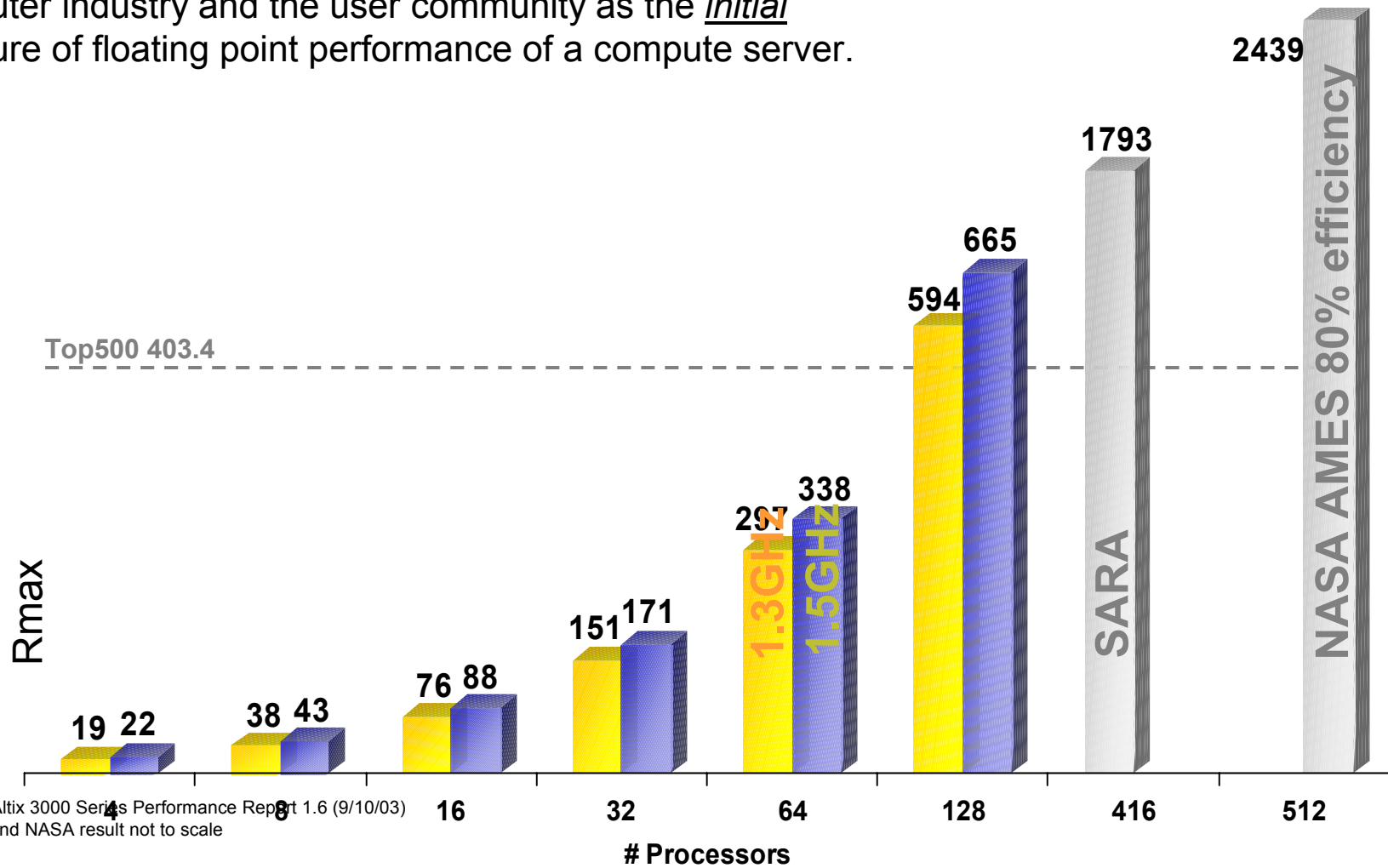
SPEC® fp_rate base 2000
SpecOMPm2001
Linpack NxN

SGI Altix 3000 LINPACK NxN Benchmark



Floating Point Performance Benchmark

The Linpack Benchmark is widely accepted in both the computer industry and the user community as the *initial* measure of floating point performance of a compute server.



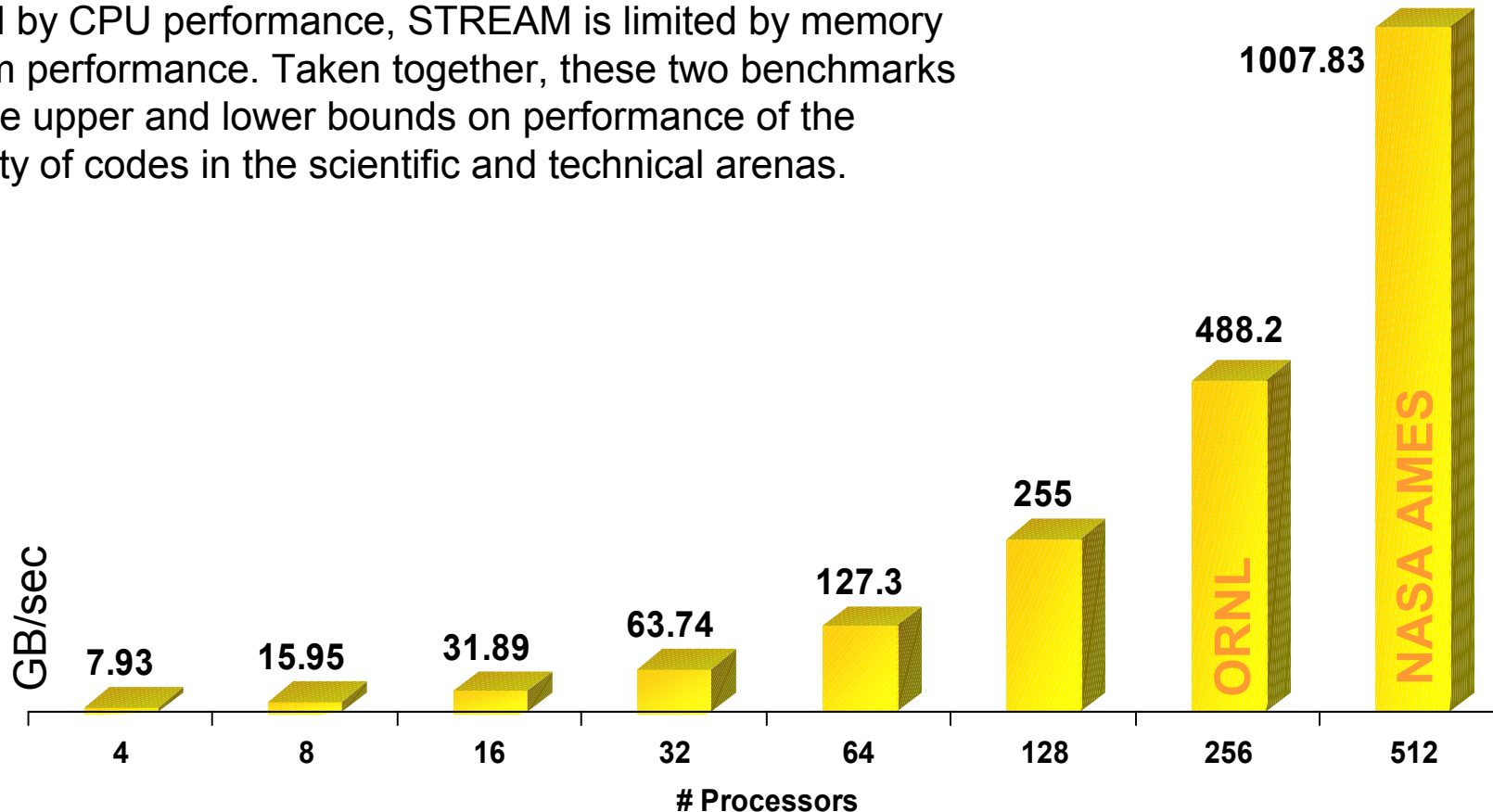
Source: SGI Altix 3000 Series Performance Report 1.6 (9/10/03)
Note: SARA and NASA result not to scale

SGI Altix 3000 STREAM Benchmark



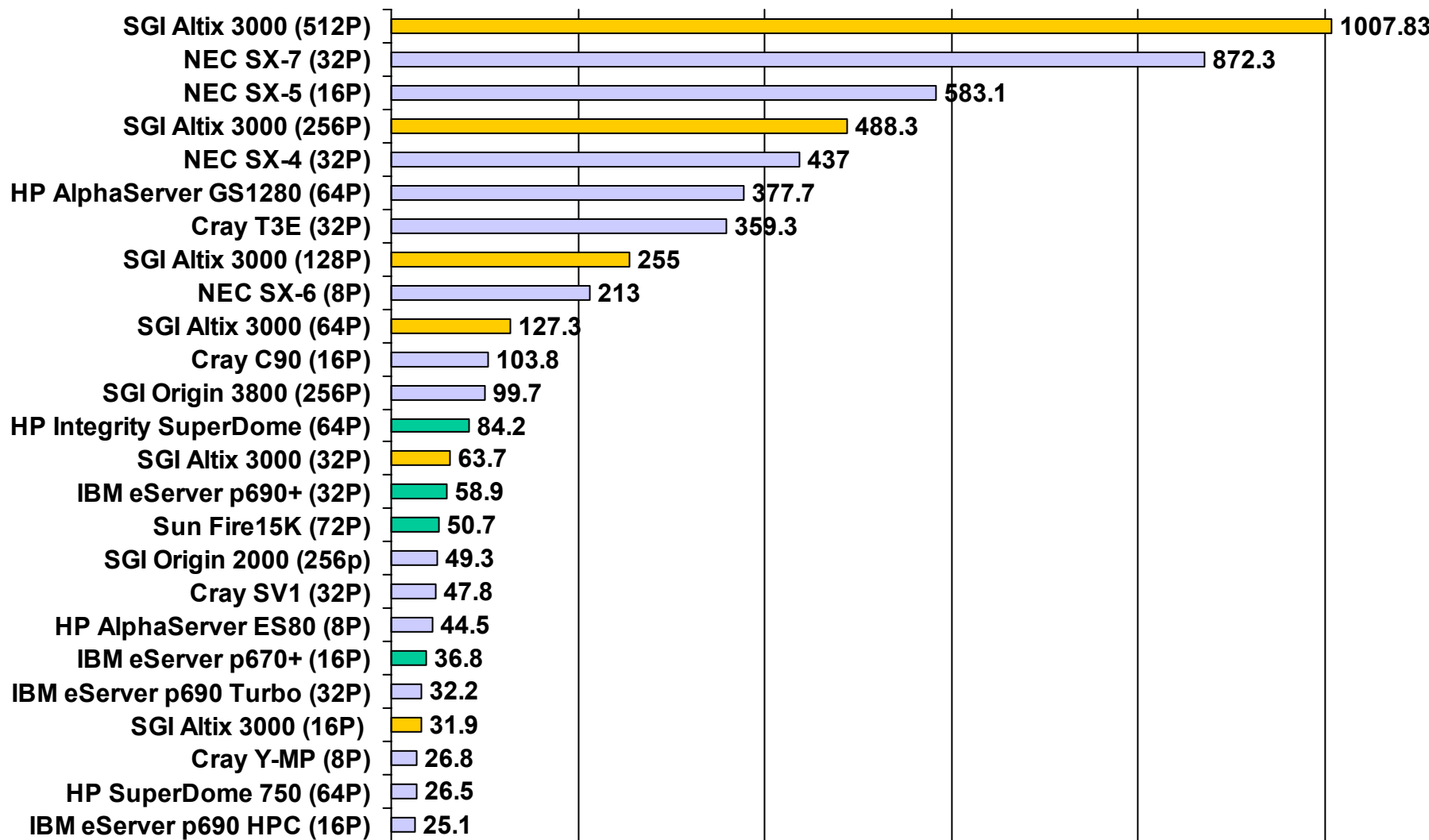
Memory Bandwidth Benchmark

STREAM is especially useful as a counterpoint to the LINPACK benchmark. While LINPACK is almost always limited by CPU performance, STREAM is limited by memory system performance. Taken together, these two benchmarks provide upper and lower bounds on performance of the majority of codes in the scientific and technical arenas.



Source: SGI Altix 3000 Series Performance Report 1.6 (9/10/03),
STREAM website (10/13/03)

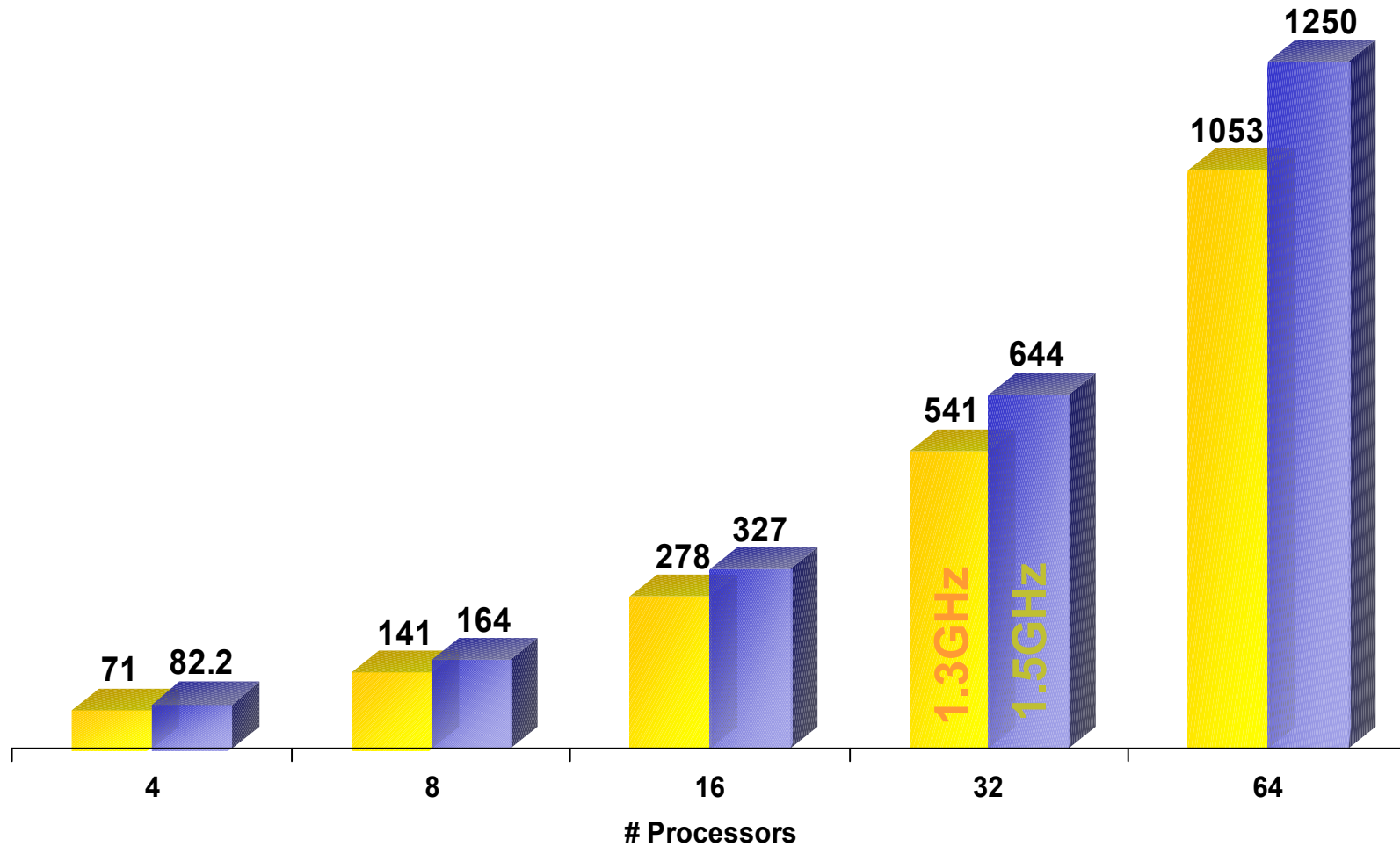
STREAM Benchmark Top 20 Results versus SGI Altix 3000



Source: SGI Altix 3000 Series Performance Report 1.6 (9/10/03), Customer provided data, & STREAM website (11/13/03)

[BACK](#)

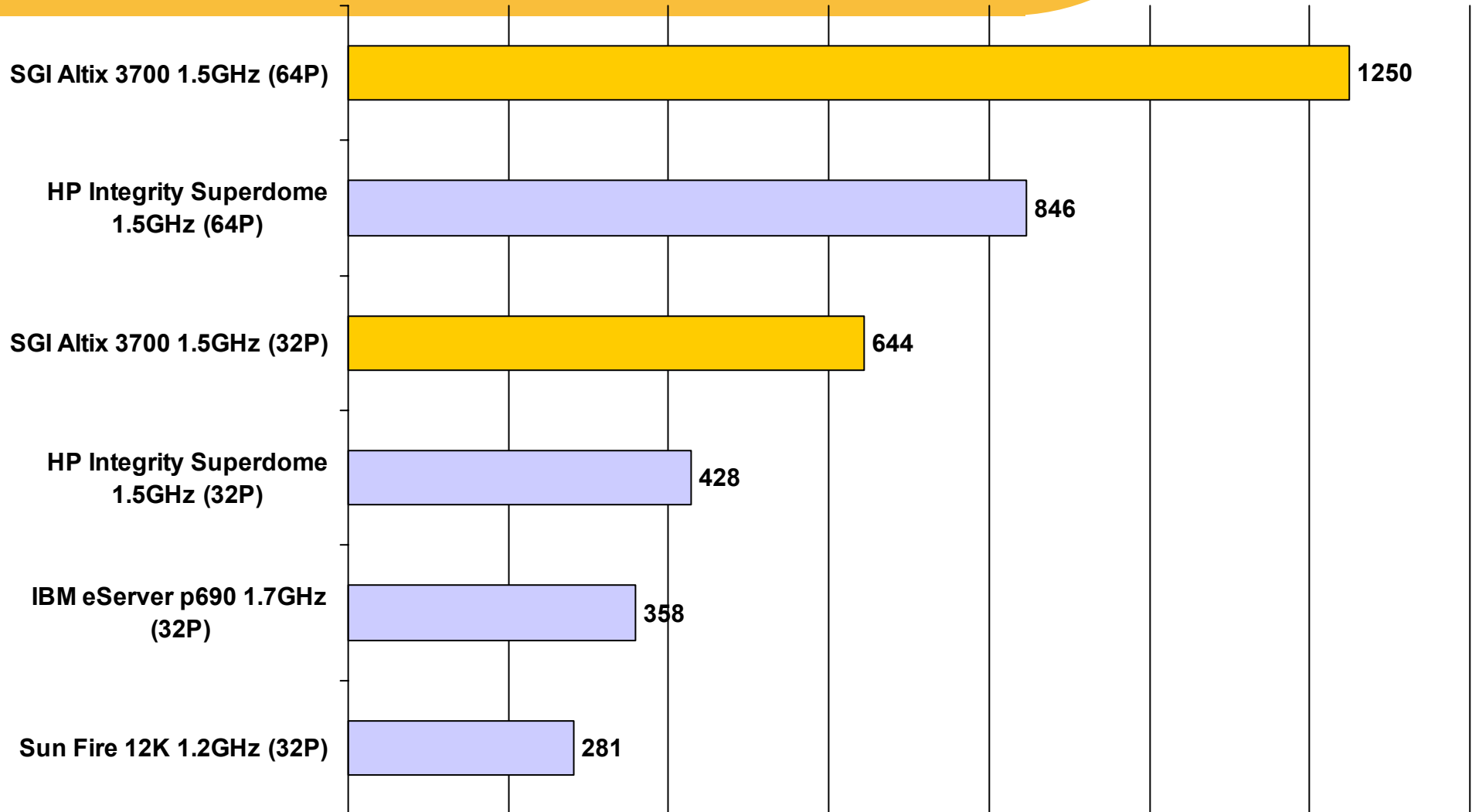
SGI Altix 3000 SPECfp_rate_base2000



Source: SGI Altix 3000 Series Performance Report 1.6 (9/10/03)

[BACK](#)

Competitive SPECfp_rate_base2000

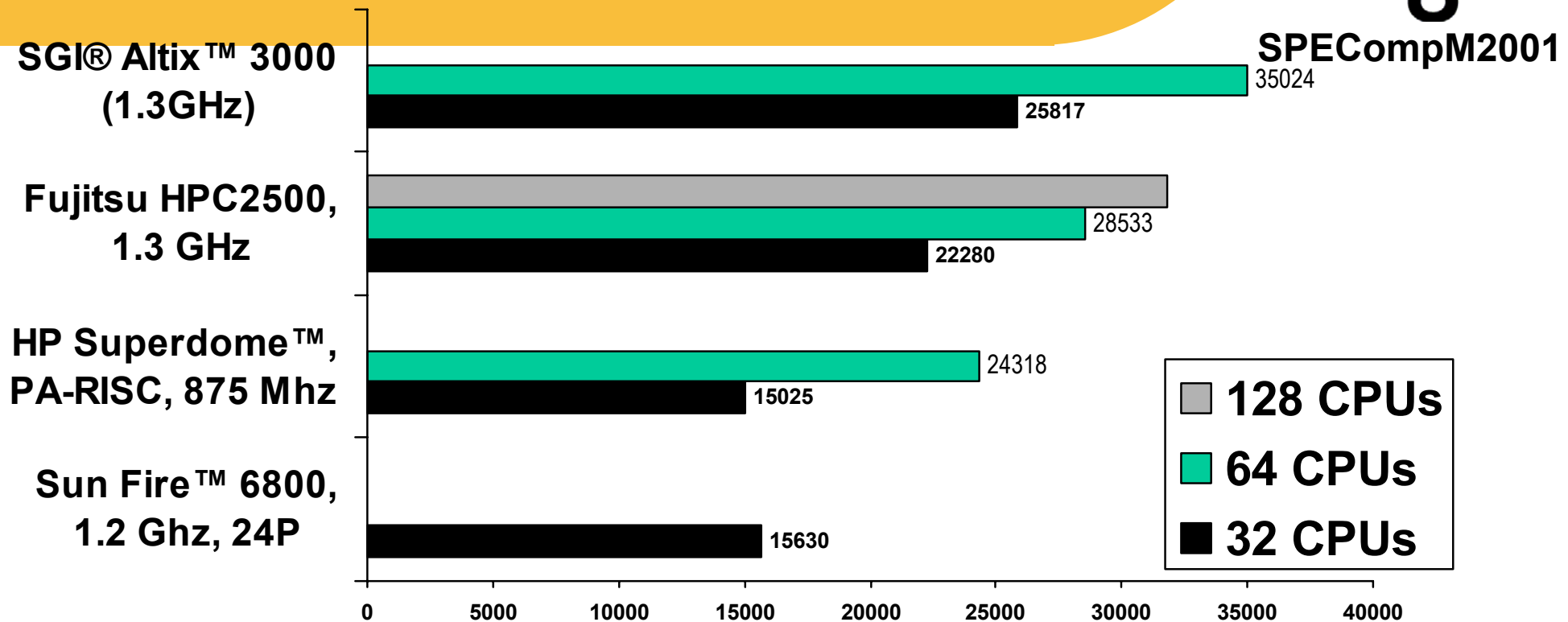


Source: SPEC website (11/9/03)

[BACK](#)

World-Record Parallel Performance: SPECComp® Results

sgi®



- World-record result for 64 and 32-processor systems
- SGI's 1.5Ghz, 64P result is 11% better than Fujitsu's and 44% better than HP Superdome..
- SPECCompM2001 is defined as the higher of SPECCompMpeak2001 and SPECCompM_base2001

Developments in High Performance Computing

A Preliminary Assessment of the NAS SGI 256/512 CPU SSI Altix (1.5 GHz) Systems

SC'03

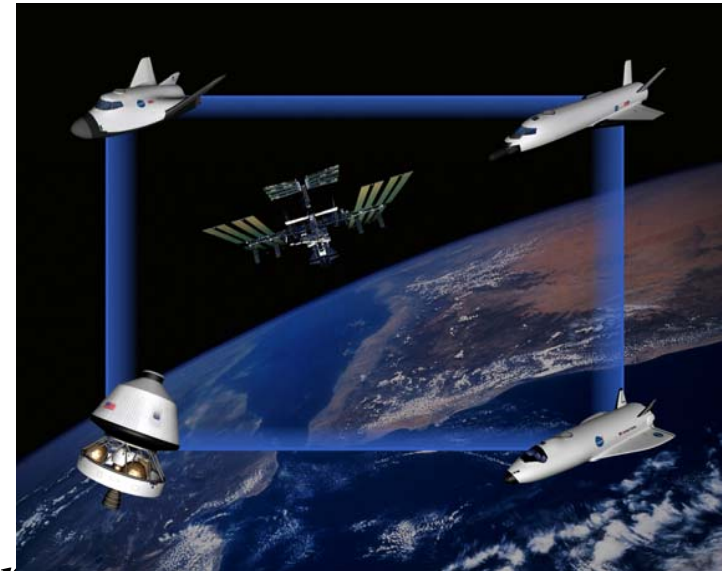
November 17-20, 2003

Jim Taft

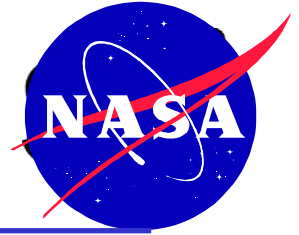
NASA Ames Research Center

Jtaft@nas.nasa.gov

QuickTime™ and a
TIFF (Uncompressed) decompressor
are needed to see this picture.

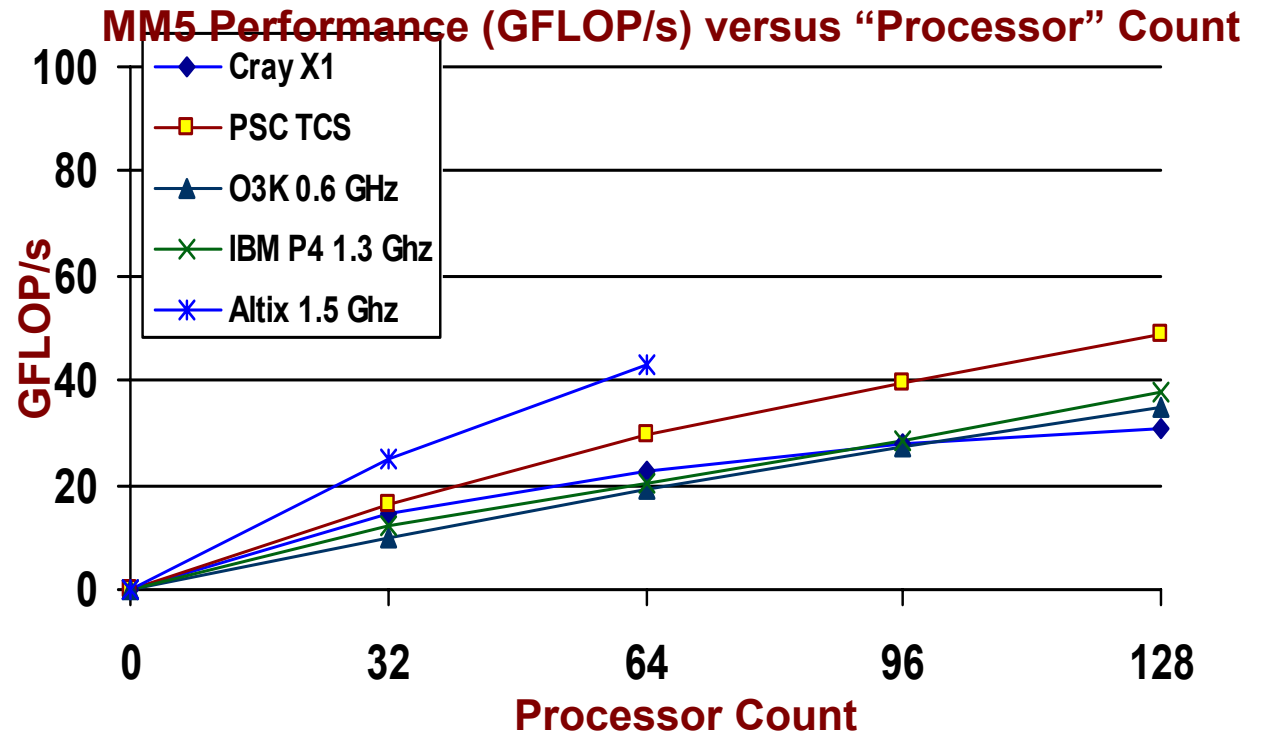


MM5 10Km Performance Results



MM5 is a classic NCAR weather code benchmark. It is known for its excellent scaling on clusters given the right problem size. It is NOT a climate model. It is inappropriate to use MM5 for setting expectations for most climate models, which usually have great difficulty in scaling to large CPU counts on clustered systems without shared memory interconnects

NOTE: The Cray X-1 results have been plotted using SSP count as the “processor” count instead of MSP. This is more of an apples to apples comparison of “processors”. Note X-1 scaling is already falling rapidly. Altix 1.5 GHz is outstanding on this code.

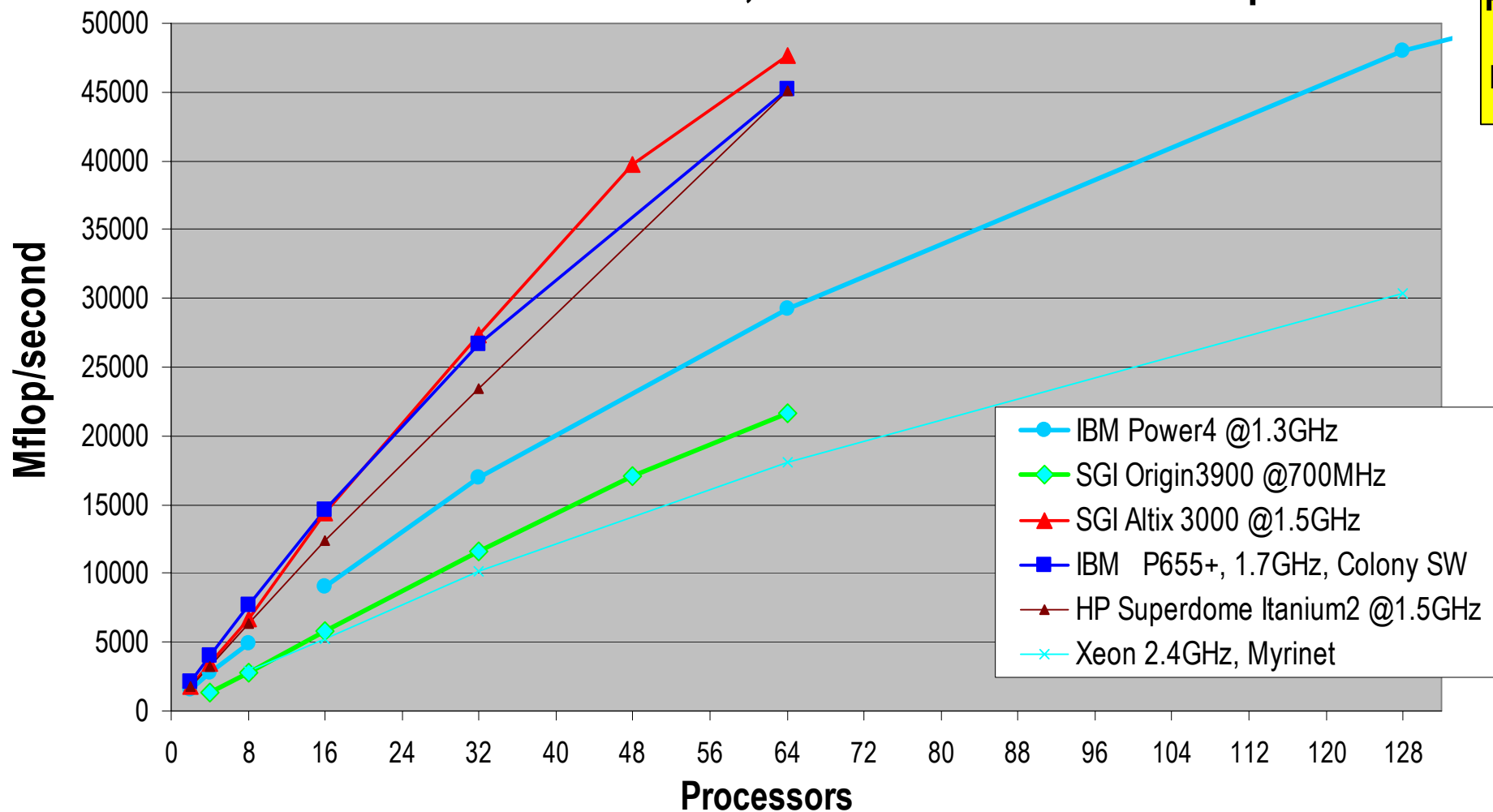


NOTE: Non Altix Data replotted from Paul Muzio charts presented at IDC-Utah

MM5 3.5 Scalability on SGI Altix 3000



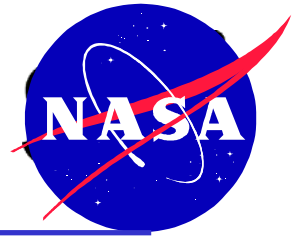
MM5- NCAT t3a benchmark, 3-hour forecast over Europe



Higher
is
Better

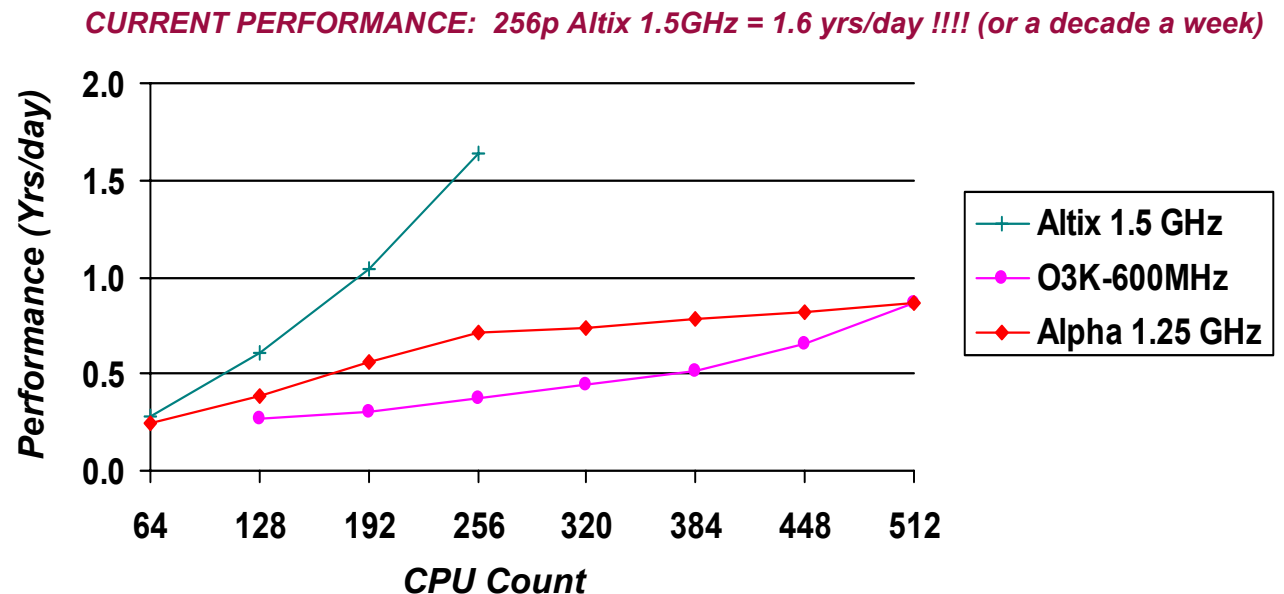
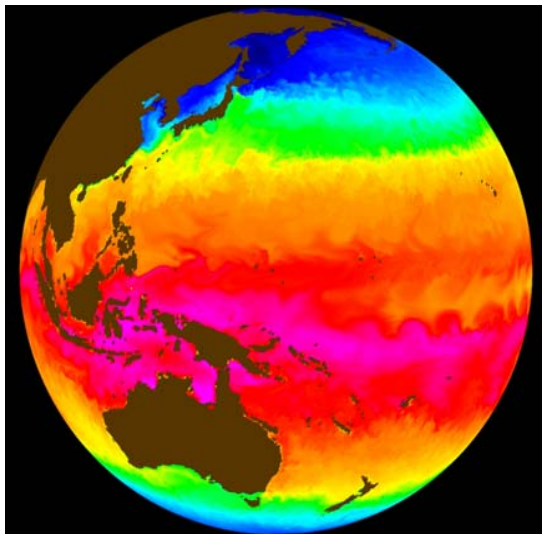
ECCO Code Performance

11/04/03



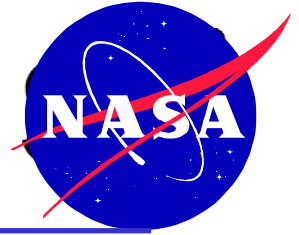
The ECCO code is a well known ocean circulation model, with features that allow it to run in a coupled mode where land, ice, and atmospheric models are run to provide a complete earth system modeling capability. In addition the code can run in a “data assimilation” mode that allows observational data to be used to improve the quality of the various physical sub-models in the calculation. The chart below shows the current performance on the Altix and other platforms for a “1/4 degree” resolution global ocean circulation problem. (in reality, much of the calculation runs at an effective much higher resolution due to grid shrink at the poles).

Note: Virtually no changes to the code have been made across platforms. Only changes needed to make it functional have been done. The preliminary Altix results are very good to date. A number of code modifications have been identified that will significantly improve on this performance number. NOTE: The performance on both Chapman and Altix with full I/O are super-linear. That is, as you add more CPUs you get even faster speedups. The Alpha numbers show a knee at 256 CPUs.



NOTE: Alpha Data re-plotted from Gerhard Theurich charts in NCCS paper

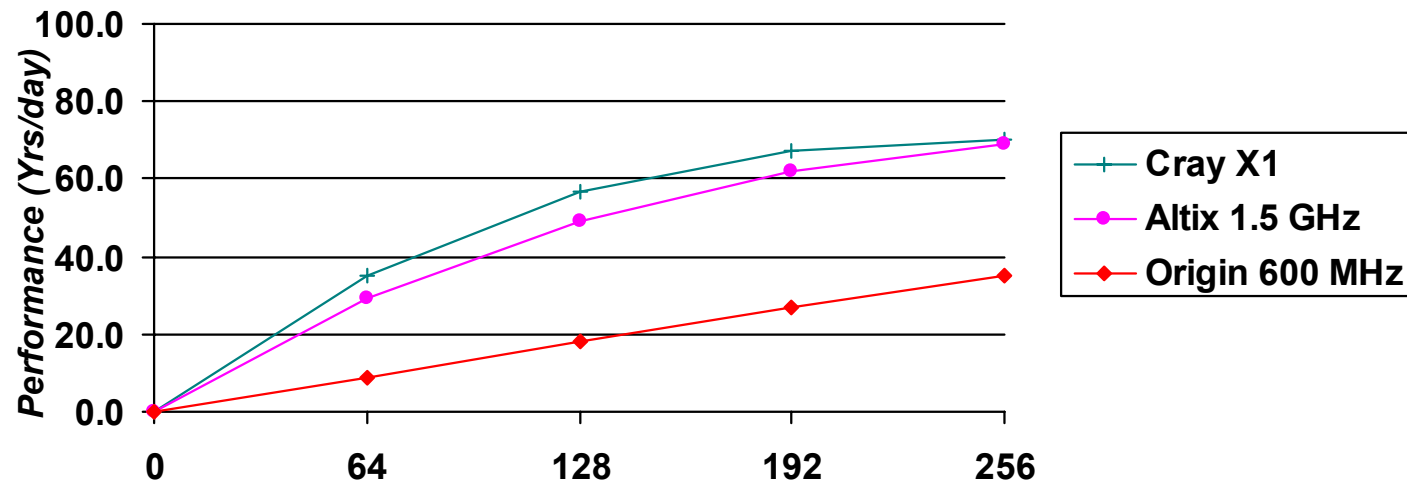
POP 1.4.3 Performance - 1 Degree Global Problem



The POP code is a well known ocean circulation model developed at LANL. It is the ocean model component of the Community Climate Systems Model (CCSM) from NCAR. The chart below shows the current performance on the Altix and other platforms for a “1 degree” resolution global ocean circulation problem.

Note: Virtually no changes to the original code have been for the Altix runs. A total of about 100 lines of code have been modified. Most of the changes are in the boundary routine used in the CG solver. At this point a number of code modifications have been identified that will significantly improve on this performance. In contrast, the vector version has been in development for about 2 years by Japan, and lately Cray.

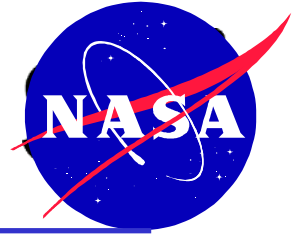
POP 1.4.3 - Performance on 1 Degree “X1” Problem



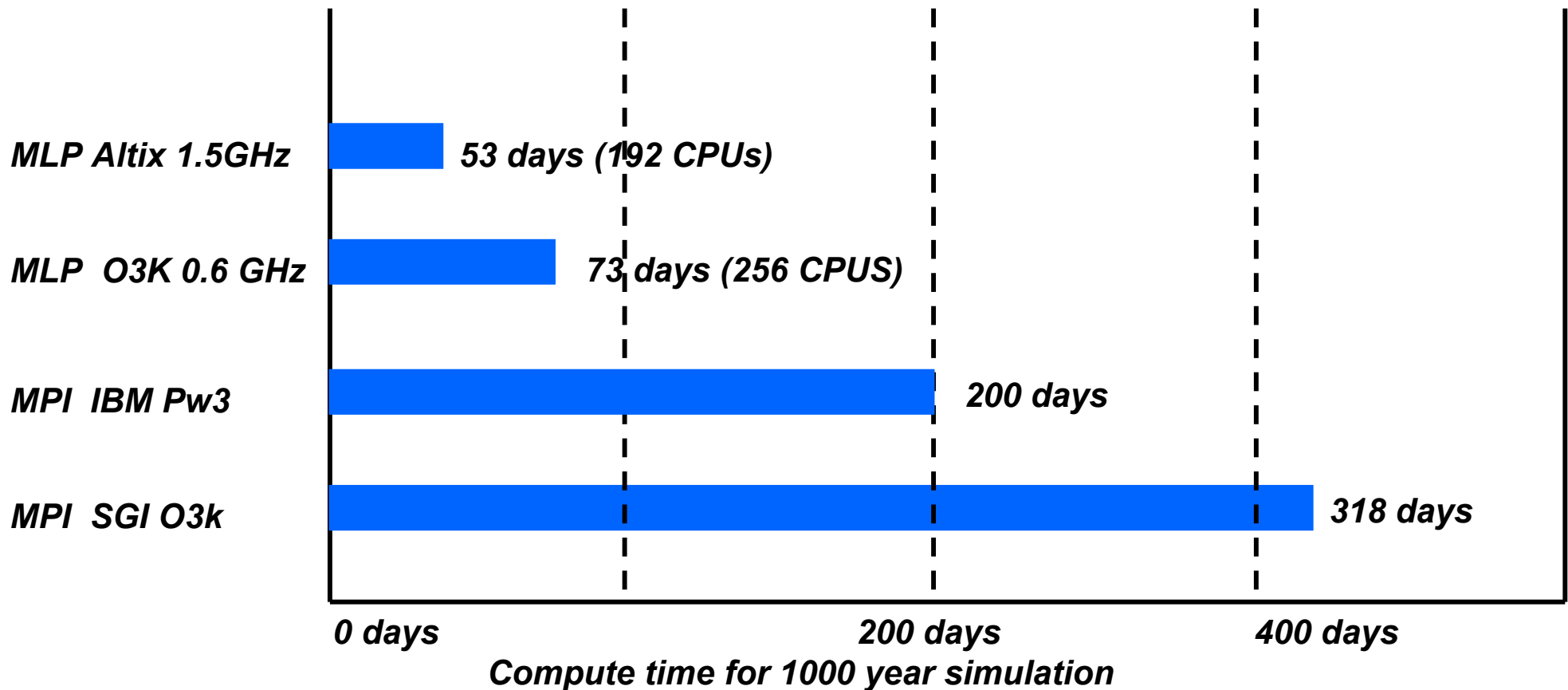
CPU Count (Cray X1 plotted as SSP count)

NOTE: X1 Data re-plotted from Pat Worley charts in X1 Early Performance Evaluation

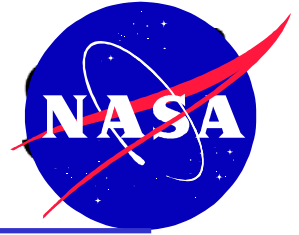
CCSM 2.0 Code Performance - 1000 year simulation



CCSM was used last year by NCAR to conduct a 1000 year global simulation using T42 resolution for the atmosphere and 1 degree resolution for the ocean. The simulation required 200 days of compute time to complete. The Altix code at this point has been partially optimized using MLP for all inter model communications. Some sub-models have been optimized further. About 4 man-months have been devoted to the project.

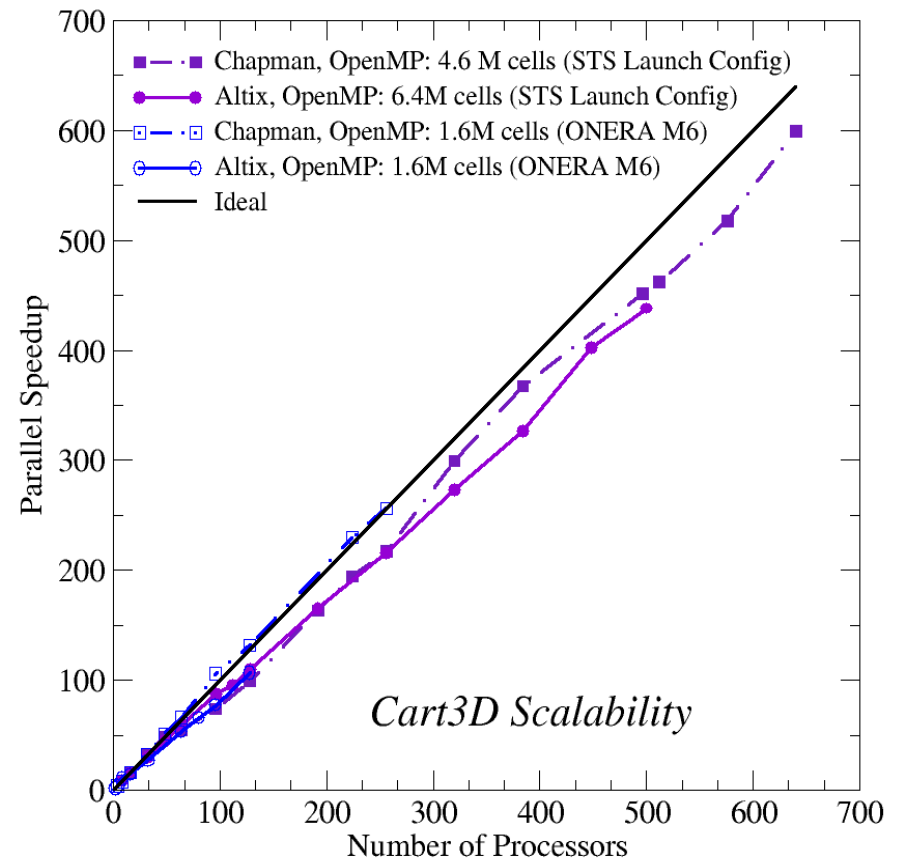


The CART3D Code - OpenMP Test



The CART3D code was the NASA “Software of the Year” winner for 2003. It is routinely used for a number of CFD problems within the agency. It’s most recent application was to assist in the foam impact analysis done for the STS107 accident investigation.

The chart to the right presents the results of executing the OpenMP based CART3D production CFD code on various problems across differing CPU counts on the NAS Altix and O3K systems. As can be seen, the scaling to 500 CPUs on the weeks old Altix 512 CPU system is excellent.



Altix 3000 w/ IPF, GSM and high performance NUMALink interconnect enables:

- low latency, high bandwidth memory access and communications
- world-record performance on std benchmarks
- top performance and scalability on customer applications

Altix Platform Intro

Altix System Architecture

Shared vs Distributed Memory (Clusters)

Linux® Environment

Roadmap

System Software Options for SGI®

Altix™



Differentiating features and functionality

- NUMA tools: cpusets/dplace
- MPT and Array Services
- Partitioning
- XVM, XSCSI
- Performance Co-Pilot™

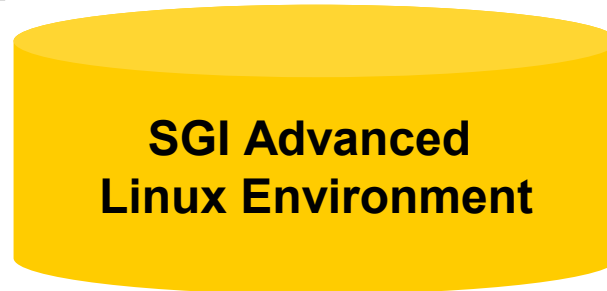


Enabling features and functionality

- Comprehensive System Accounting
- Job Container (PAGG)
- XFS



Base Linux Operating System



HPC-targeted

ISV/Database-targeted

Two OS Options for Altix



Option 1: SGI Advanced Linux Environment -PLUS- SGI ProPack

- SGI Advanced Linux Environment
 - uses Red Hat Enterprise Linux as base (RHAS v2.1)
- SGI ProPack
 - SGI tuned kernel for *world-class and best* performance and scalability
 - HPC libraries and Tools (e.g. MPT, FFIO, SCSL, cpuset, etc.)
 - Both closed-source (SGI proprietary) and open-source components
 - Latest bug fixes (both user and kernel)
- ABI compatibility (user/kernel) with RHEL
- Target markets: HPC focused
- Customer support path: Customer -> SGI

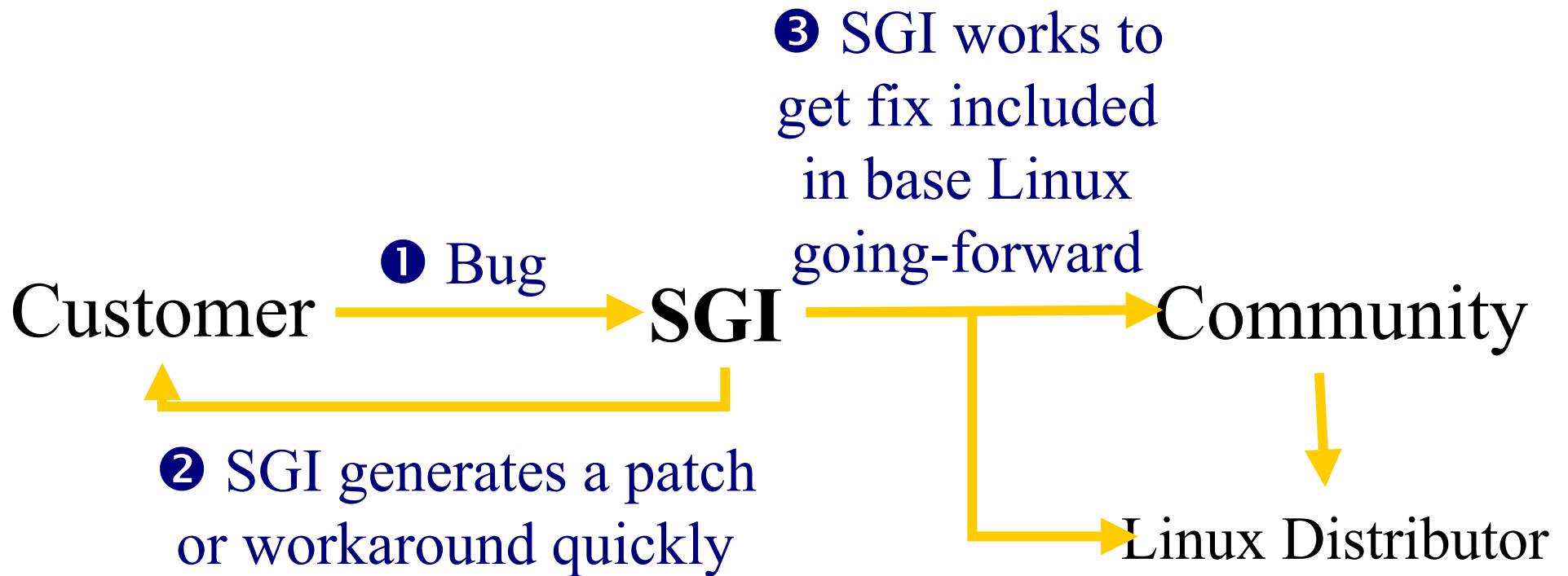
Core Linux

- Platform support
- NUMA support (discontig memory, VM, CPU timers, etc.)
- Big systems (>32P, memory, etc.)
- Error handling, MCAs, etc.
- Partitioning enabling hooks
- Scalability improvements (too many to mention!!)

Other OS features

- Cpumemsets (runon, cpuset, memory placement/dplace, etc.)
- PAGG (process aggregates for jobs containers, job acctg, etc.)
- Pthreads (ngpt), MQ schedulers, kdb, lkcd, etc.

SGI Provides Linux[®] Support Directly for ProPack



Two OS Options for Altix (cont.)



Option 2: SUSE LINUX Enterprise Server 8 for Altix

- SLES8 SP3 on Altix (available now)
- SUSE distribution boots and runs on Altix
 - SUSE *defines* distribution and contents and releases CDs
 - SGI ProPack components are *NOT* provided or supported on SUSE
- Supports smaller Altix config sizes, fewer devices than SGI Propack:
 - 64P SSI max, but not same scaling/performance
 - less IO controller devices, storage, etc
- Target market: ISV focused
 - SUSE certified/ Oracle Unbreakable Linux certified
- Customer support path: Customer -> SGI -> SUSE

Rich HPC/Linux[®] Development Environment



Rapid evolution

- SGI knowledge of compilers
- Intel knowledge of processors

Leverage open source

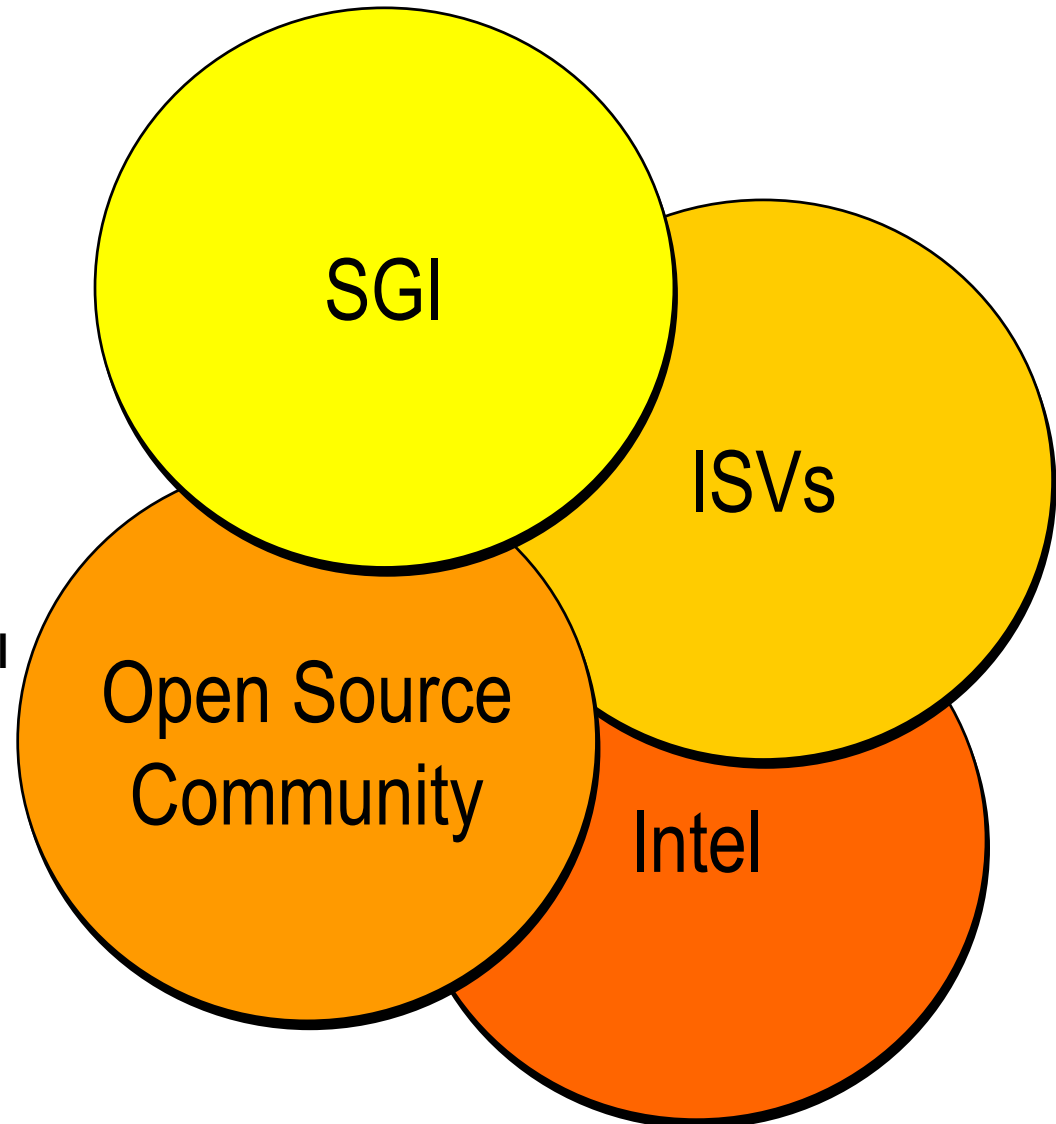
- Many apps available
- We test to verify

Differentiation

- Enhanced ISV app performance from SGI Libraries--MPT and SCSL
- Only on Altix[™]

Engagement with premier tools ISVs

- Etnus (TotalView[®])
- Pallas (Vampir[™])



Intel® C/C++ and Fortran compilers

- **8.0 Compilers**
 - Fortran support for OpenMP 2.0
 - Improved gcc compatibility
 - C99 compatibility (subset)

GNU Fortran and gcc

SGI® MPT (Message Passing Toolkit)

- **MPI and SHMEM parallel programming libraries**
- **Global pointer construct allows jobs to address both local and remote memory regions**
- **Low latency, high bandwidth, NUMA-aware**
- **MPT performance on multipartition systems same as SSI with little performance penalty crossing nodes**

SGI® SCSL

- **Comprehensive science and math functions**
- **Optimized BLAS, FFTs, sparse solvers**
- **Provides performance improvements**

Libraries



SGI® FFIO

- Enables control of specifics of I/O transfers
- Enhances application performance

Intel® MKL 6.0

- Optimized math functions
- LAPACK, BLAS, FFTs and vector transcendental functions

Intel® idb (included with Intel compilers)

- Thread support
- Supports MPI

Gnu gdb (with Fortran extensions)

TotalView from Etnus®

- Excellent C++ and F90 support
- Thread support includes MPI
- Advanced features

Performance and System Analysis



SGI® Performance Co-Pilot™

- System performance analysis
- Visualization of all nodes in a system

pfmon—(open source) provided by SGI

- Command-line binary and system analysis
- Uses IA-64 PMU (Performance Monitoring Unit) to do counting and sampling

SGI® Histx provided at no charge by SGI

- Provides profiling and perfex like tools

Intel® Vtune™

- Remote & native sampling on Linux®
- Multithreaded application and hyper-threaded processor analysis

Other HPC Tools for Analysis and Parallelization



Vampir™ and Vampirtrace™ from Pallas **Performance analysis and visualization**

- For MPI applications
- Graphical analysis of runtime traces
- Indispensable for efficient parallel program development and tuning

Parallel Software Products from ParaWise

- Already ported (formerly CAPTools)
- Generates parallel code from serial code for Fortran
- Computer Aided Parallelization Toolkit

Agenda



Altix Platform Intro

Altix system Architecture

Shared vs Distributed Memory (Clusters)

Linux[®] Environment

Roadmap

Altix Server Roadmap



- **Future Itanium Processor Family (IPF) microprocessors**
- **Maximum system size increase to thousands of processors**
- **Maximum single system image increase**
- **More compact packaging (for increased computational density)**
- **Next-gen Altix architecture and components**

Altix Software Roadmap



- **New SGI hardware support**
(including future IPF microprocessors)
- **Realtime support features (e.g. for viz-sim applications)**
- **2.6 Kernel**
- **512p SSI**
- **thousands of processors via *Supercluster***
- **MPI optimized for Infiniband**
- **RHEL 4.0 base**

sggi[®]