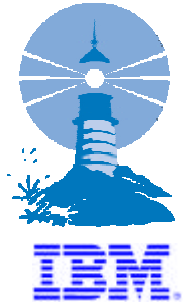# BlueGene/L Supercomputer

George Chiu
IBM Research
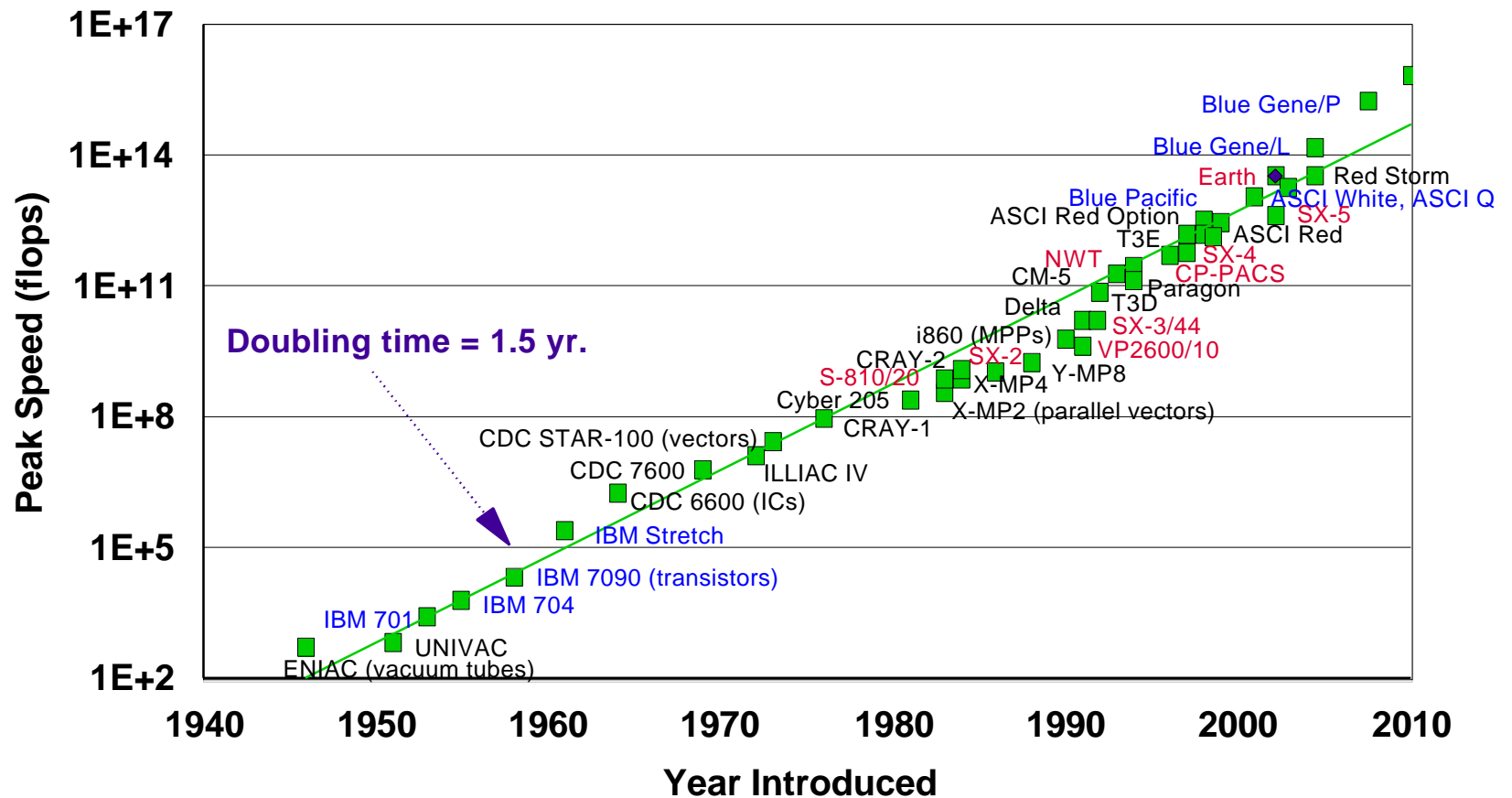
# The Blue Gene Project

- In December 1999, IBM Research announced a 5 year, $100M US, effort to build a petaflop supercomputer to attack problems such as protein folding.
- The Blue Gene project has two primary goals:
  - ƒ Advance the state of the art in computer design and software for extremely large scale systems.
  - ƒ Advance the state of the art of biomolecular simulation.

- In November 2001, an R & D partnership with Lawrence Livermore National Laboratory was announced.
- In November 2002, a contract with Lawrence Livermore National Laboratory was signed to build two systems: ASCI Purple and BlueGene/L.

- In October 2003, a 512-way prototype achieves 1,413 GFlop/s on Linpack and permits entry onto TOP500.

- Deliverable: a 180/360 TF/s BlueGene/L system to Lawrence Livermore National Laboratory in 1H2005
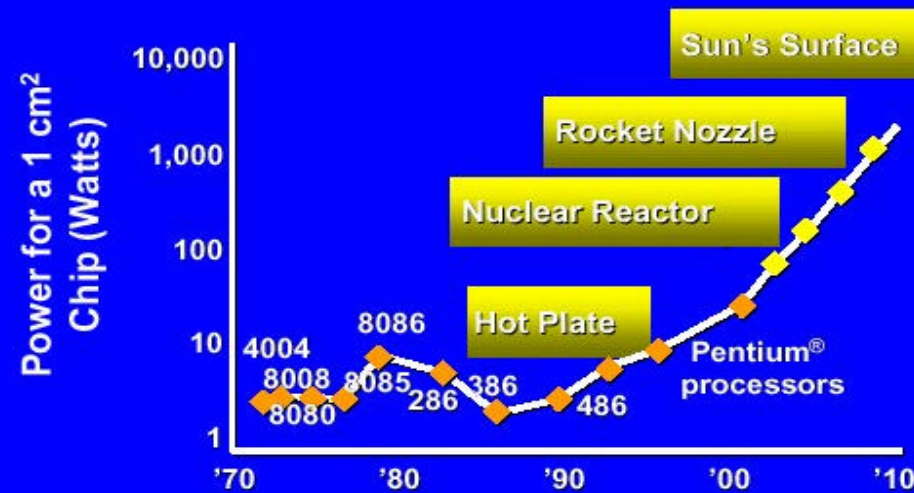
# Supercomputer Peak Performance



**Peak Speed (flops)** vs **Year Introduced**

Doubling time = 1.5 yr.

ENIAC (vacuum tubes)
IBM 701
UNIVAC
IBM 704
IBM 7090 (transistors)
IBM Stretch
CDC 6600 (ICs)
CDC 7600
ILLIAC IV
CDC STAR-100 (vectors)
CRAY-1
Cyber 205
S-810/20
X-MP4
SX-2
CRAY-2
X-MP2 (parallel vectors)
Y-MP8
i860 (MPPs)
VP2600/10
SX-3/44
Delta
T3D
CM-5
Paragon
NWT
CP-PACS
SX-4
T3E
ASCI Red Option
ASCI Red
Blue Pacific
Earth
SX-5
ASCI White, ASCI Q
Blue Gene/L
Red Storm
Blue Gene/P

# Microprocessor Power Density Growth



## Power Extrapolation

Sun's Surface

Rocket Nozzle

Nuclear Reactor

Hot Plate

Pentium® processors

4004
8008
8080
8085
8086
286
386
486

Power for a 1 cm² Chip (Watts)
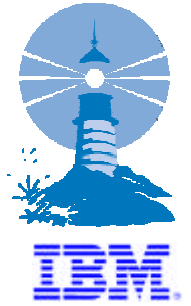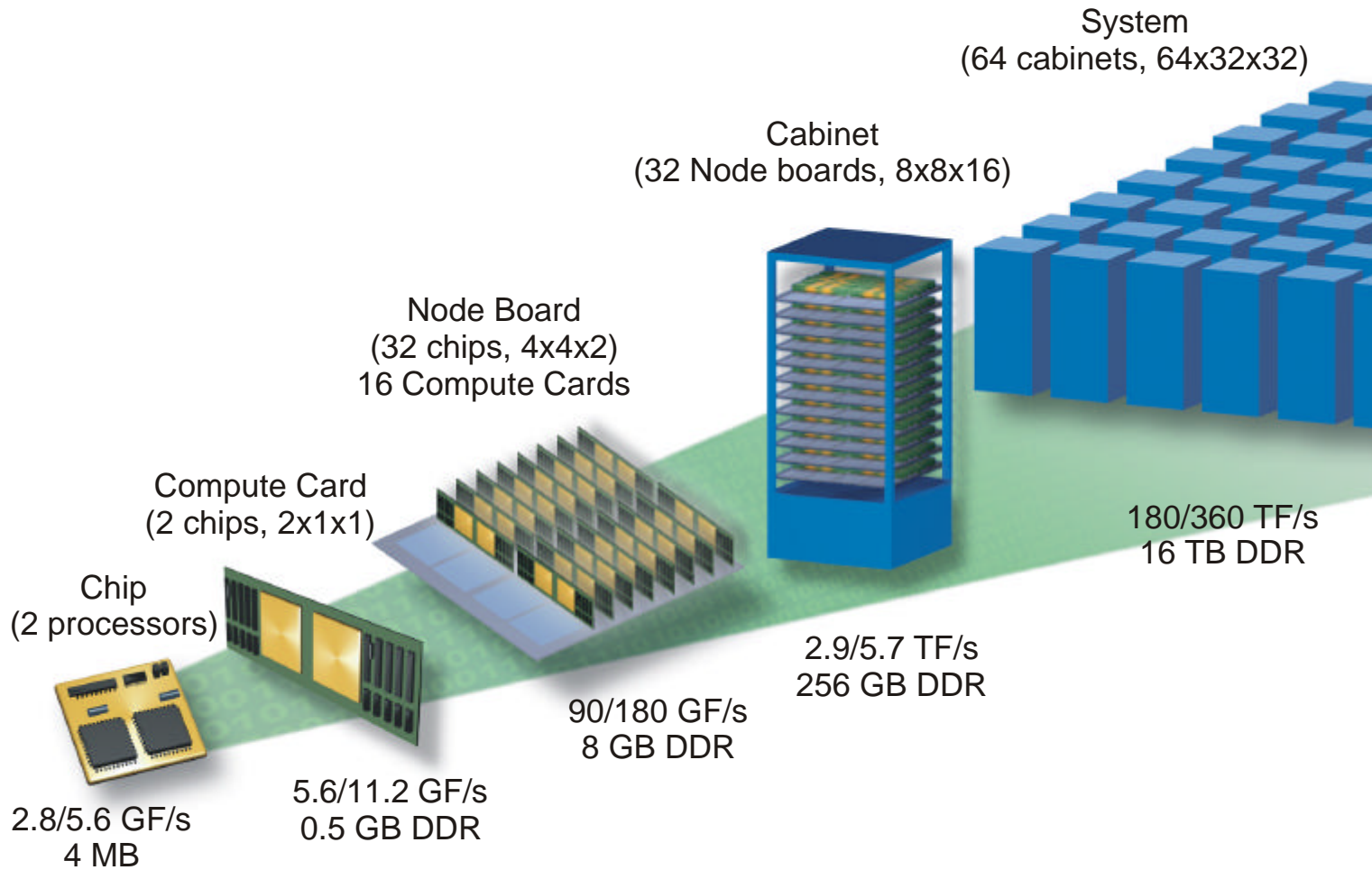
Pat Gelsinger's Slide from ISSCC 2001

*If nothing is done, power will get out of hand and Moore's Law cannot continue. Future CPUs will not be feasible.*
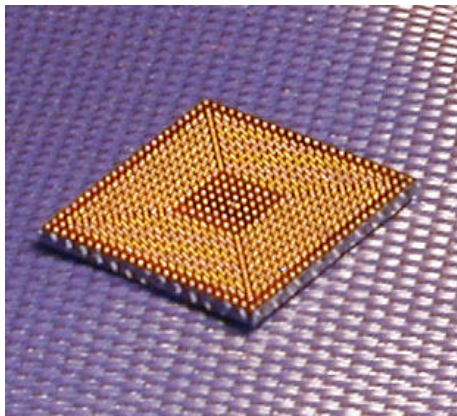
# What is the Blue Gene/L Project?

- Partnership between IBM, ASCI-Trilab and Universities to develop and build a 180-360 TFlops/s computer.
- Focus is on numerically intensive scientific problems.
- A 64k node highly integrated supercomputer based on system-on-a-chip technology. Two ASICs were designed:
  - BlueGene/L Compute (BLC)
  - BlueGene/L Link (BLL)
- LLNL contributions
  - Validation and optimization of architecture based on real applications.
  - Researchers are accustomed to "new architectures" and will work hard to adapt to constraints. Assists us in the investigation of the reach of this machine
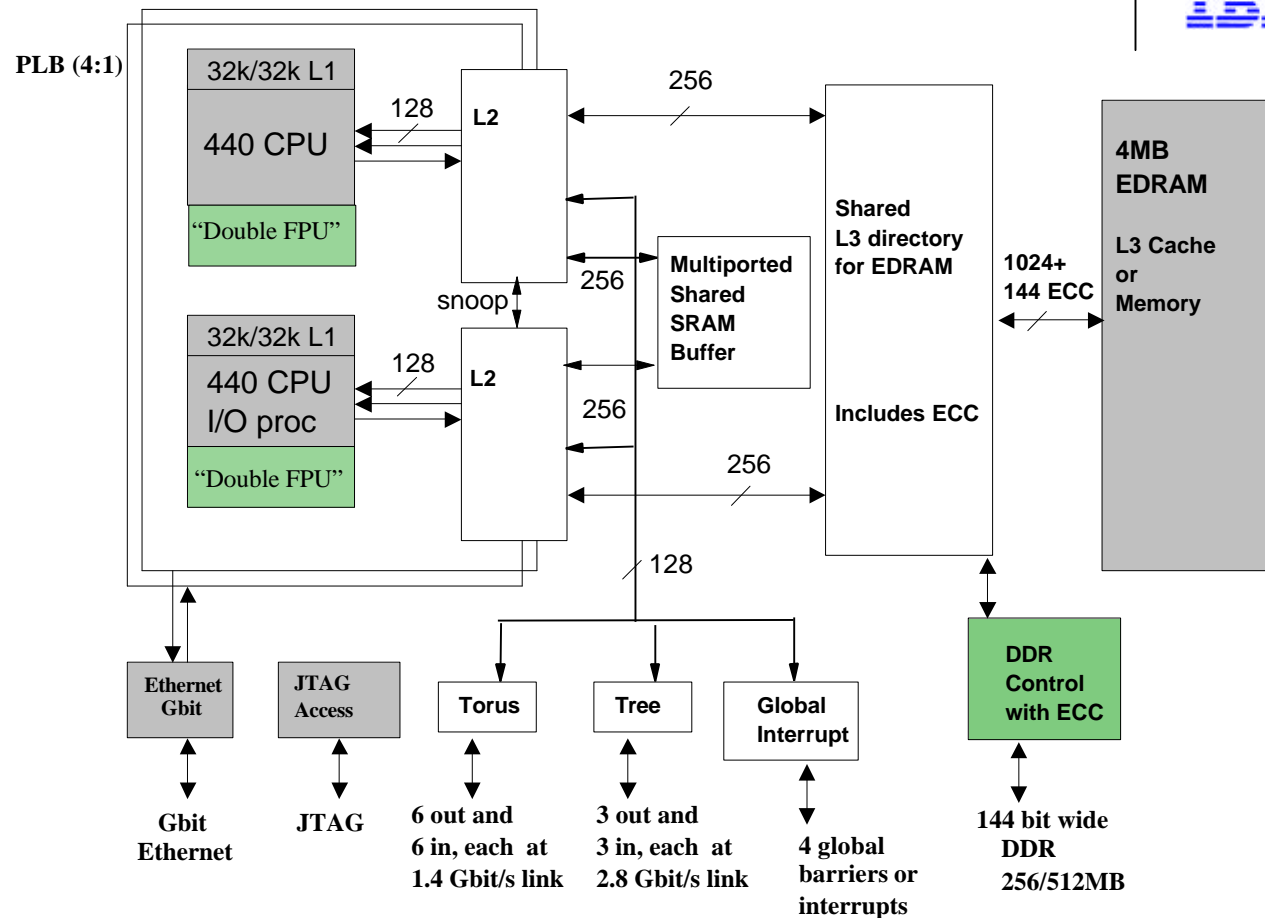- Grand challenge science requirements stress I/O, memory (bandwidth, size and latency), and processing power.
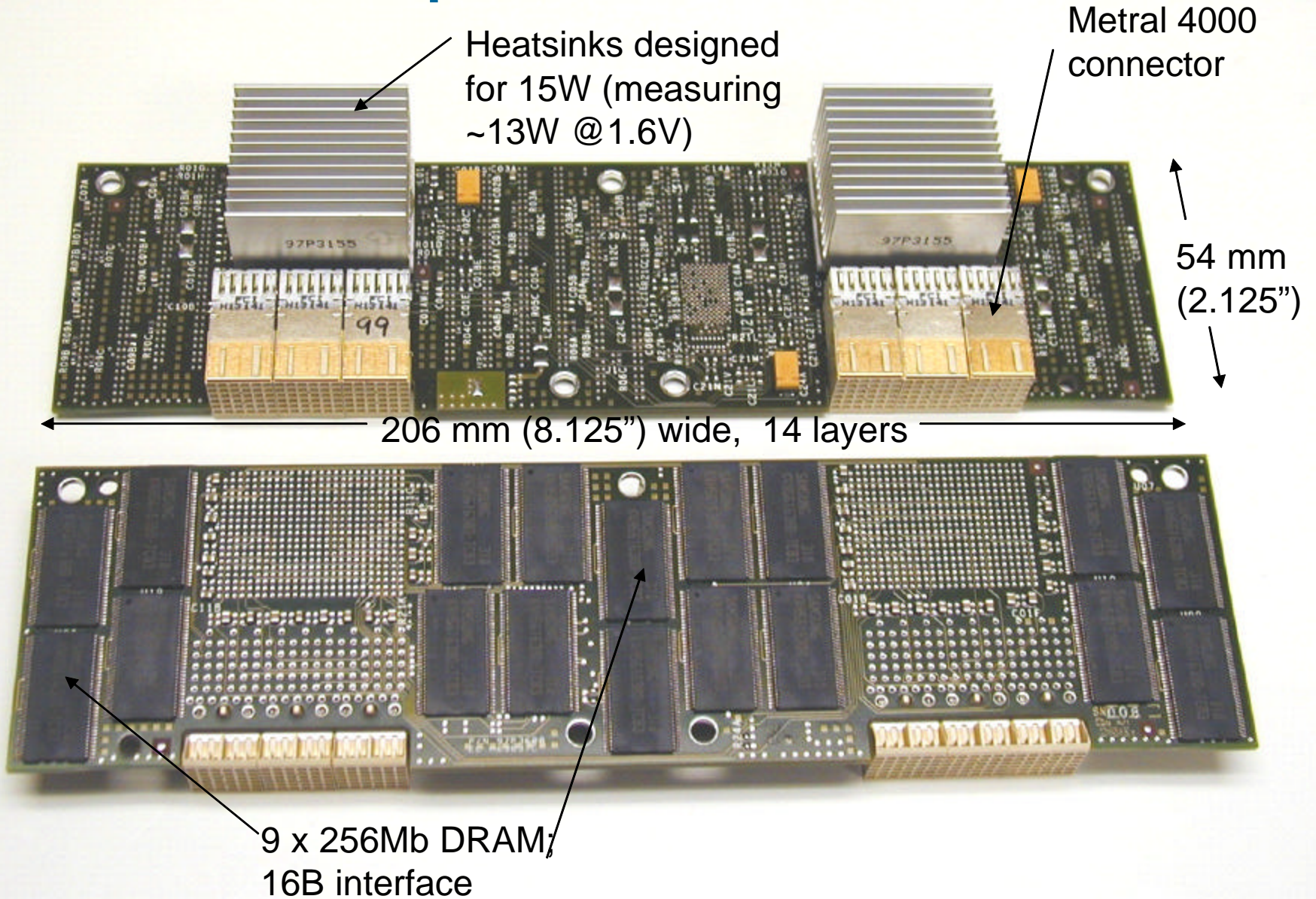
# BlueGene/L

System
(64 cabinets, 64x32x32)

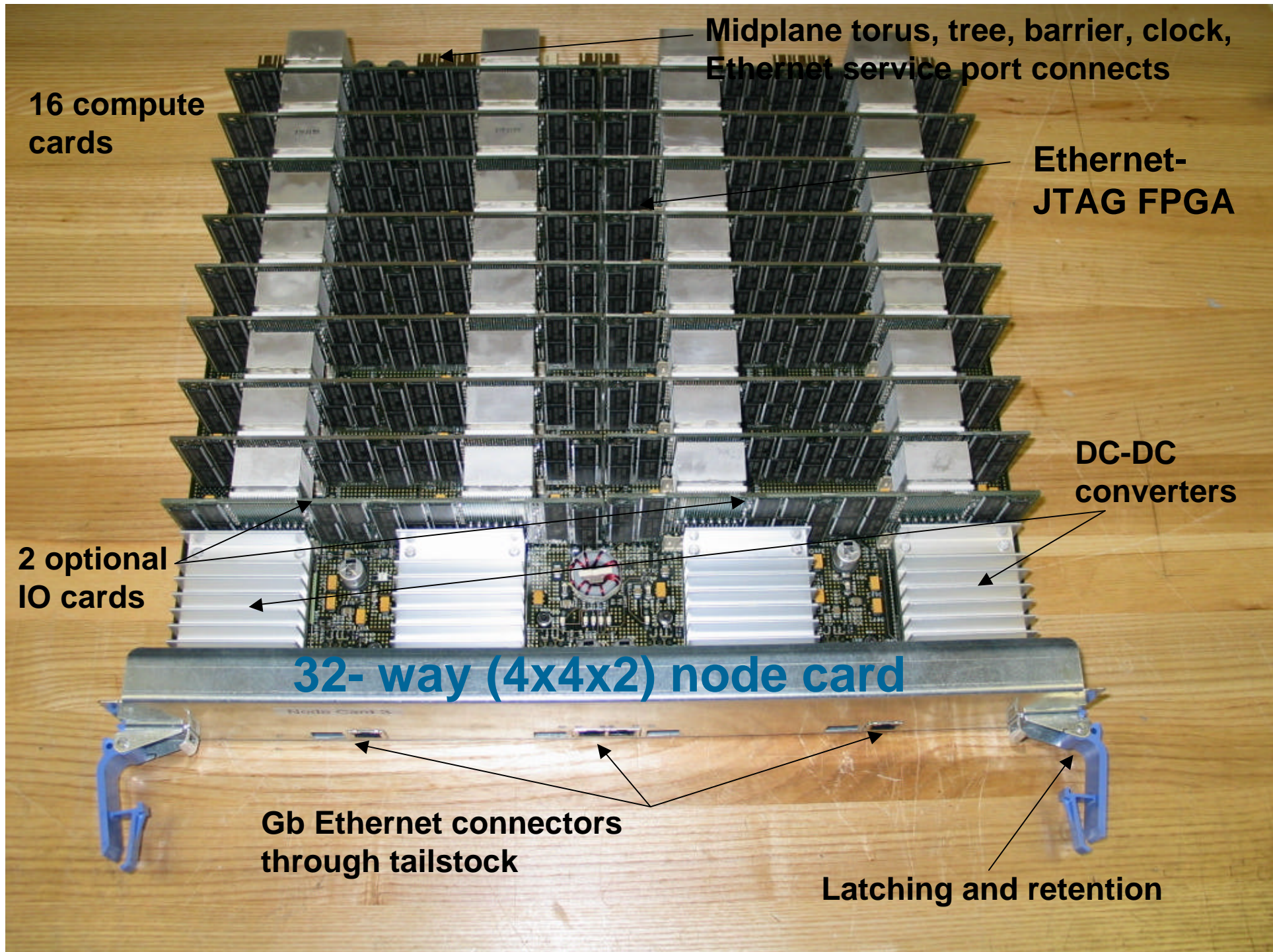Cabinet
(32 Node boards, 8x8x16)

Node Board
(32 chips, 4x4x2)
16 Compute Cards

Compute Card
(2 chips, 2x1x1)

Chip
(2 processors)

180/360 TF/s
16 TB DDR

2.9/5.7 TF/s
256 GB DDR

90/180 GF/s
8 GB DDR

5.6/11.2 GF/s
0.5 GB DDR

2.8/5.6 GF/s
4 MB

# BlueGene/L Compute ASIC

- IBM CU-11, 0.13 µm
- 11 x 11 mm die size
- 25 x 32 mm CBGA
- 474 pins, 328 signal
- 1.5/2.5 Volt

PLB (4:1)

32k/32k L1

440 CPU

"Double FPU"

L2

128

256

snoop

256

Multiported Shared SRAM Buffer

32k/32k L1

440 CPU I/O proc

"Double FPU"

L2

128

256

256

128

Shared L3 directory for EDRAM

Includes ECC

1024+ 144 ECC

4MB EDRAM

L3 Cache or Memory

Ethernet Gbit

JTAG Access

Torus

Tree

Global Interrupt

DDR Control with ECC

Gbit Ethernet

JTAG

6 out and 6 in, each at 1.4 Gbit/s link

3 out and 3 in, each at 2.8 Gbit/s link

4 global barriers or interrupts

144 bit wide DDR 256/512MB

# Dual Node Compute Card



Heatsinks designed for 15W (measuring ~13W @1.6V)

Metral 4000 connector

54 mm (2.125")

206 mm (8.125") wide, 14 layers

9 x 256Mb DRAM; 16B interface

Midplane torus, tree, barrier, clock, Ethernet service port connects

16 compute cards

Ethernet-JTAG FPGA

DC-DC converters

2 optional IO cards

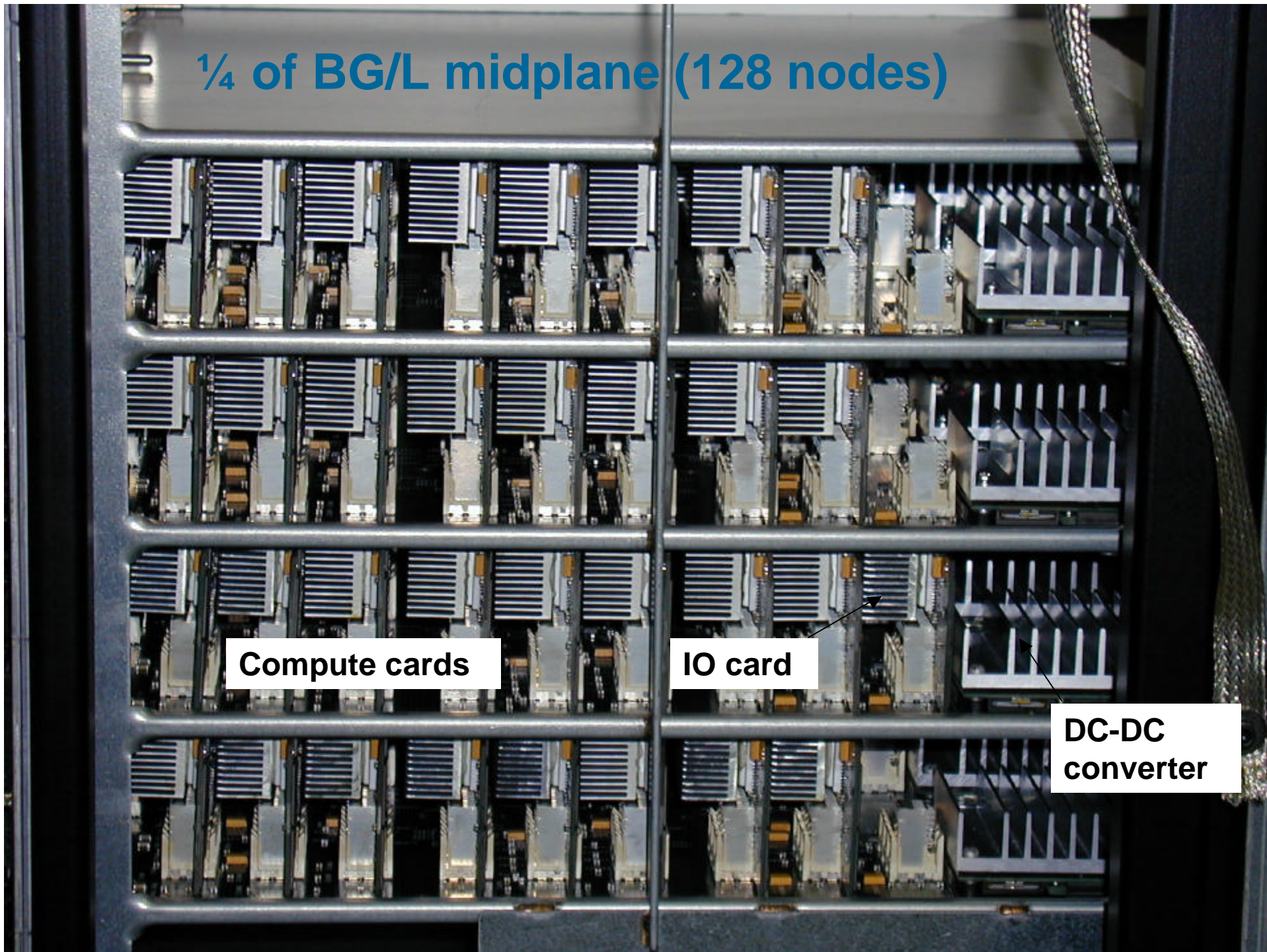32- way (4x4x2) node card

Gb Ethernet connectors through tailstock
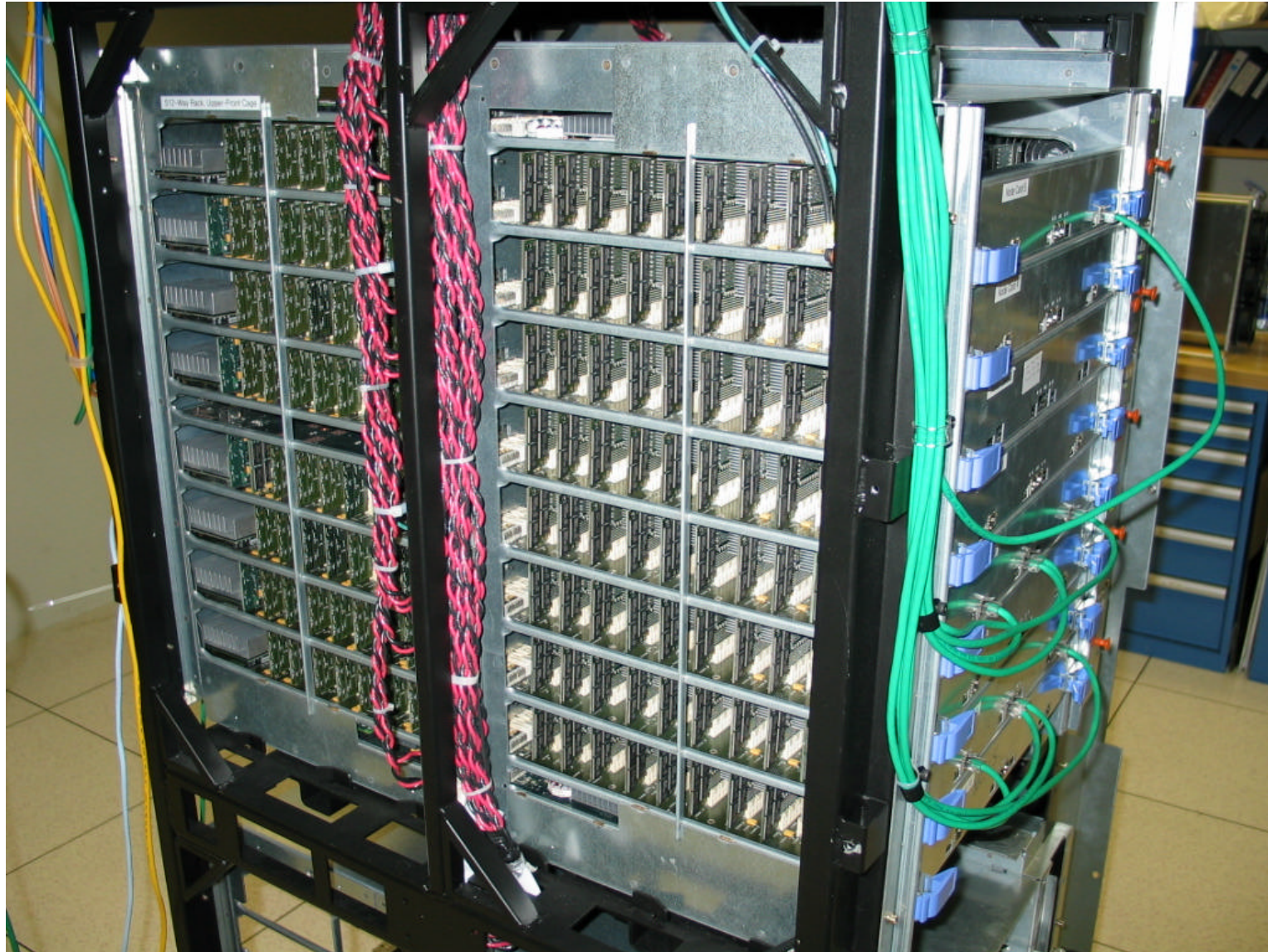
Latching and retention

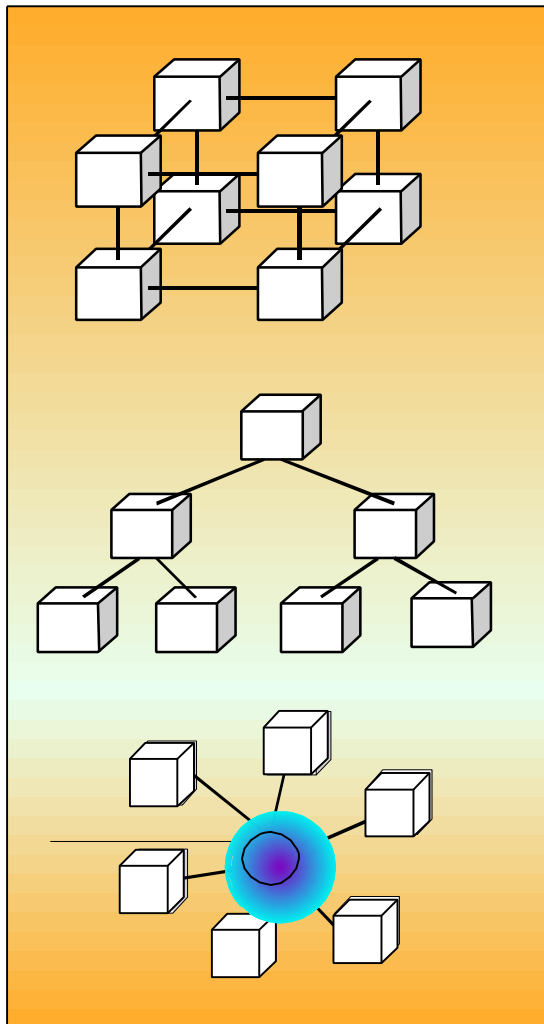¼ of BG/L midplane (128 nodes)

Compute cards

IO card

DC-DC converter

# 512 Way BG/L Prototype

# BlueGene/L Interconnection Networks

**3 Dimensional Torus**

- Interconnects all compute nodes (65,536)
- Virtual cut-through hardware routing
- 1.4Gb/s on all 12 node links (2.1 GB/s per node)
- Communications backbone for computations
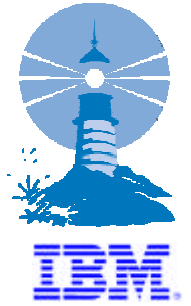- 0.7/1.4 Tb/s bisection bandwidth, 67TB/s total bandwidth

**Global Tree**

- One-to-all broadcast functionality
- Reduction operations functionality
- 2.8 Gb/s of bandwidth per link
- Latency of tree traversal 2.5 µs
- ~23TB/s total binary tree bandwidth (64k machine)
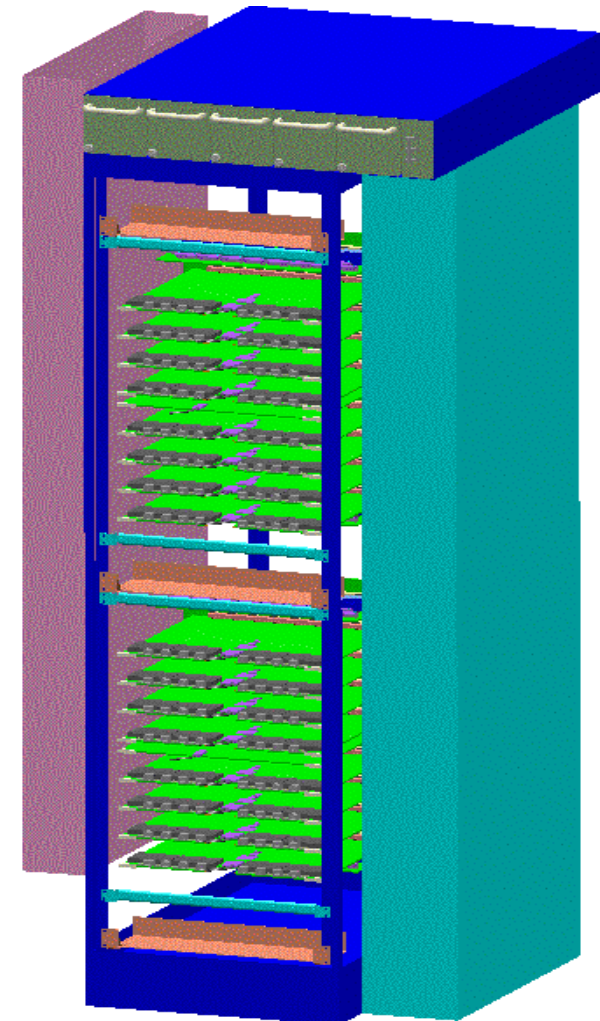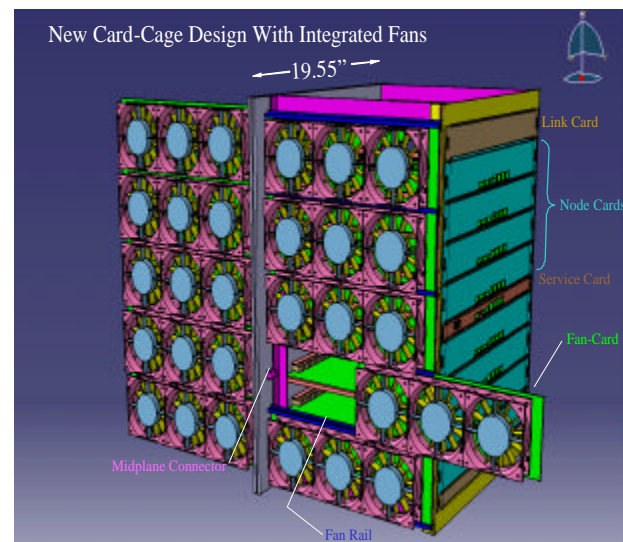- Interconnects all compute and I/O nodes (1024)

**Ethernet**

- Incorporated into every node ASIC
- Active in the I/O nodes (1:64)
- All external comm. (file I/O, control, user interaction, etc.)

# Rack, without cables

- 1024 compute nodes
  - ƒ 1024 compute proc
  - ƒ 1024 communication proc
  - ƒ 256 GB DRAM
  - ƒ 2.8TF peak

- 16 I/O nodes
  - ƒ 8 GB DRAM
  - ƒ 16 Gb Ethernet channels

- ~20KW, air cooled
  - redundant power
  - redundant fans

- ~36inch wide,
  ~36 inch deep,
  ~80 inch high (40-42U)



New Card-Cage Design With Integrated Fans
19.55"
Link Card
Node Cards
Service Card
Fan-Card
Midplane Connector
Fan Rail

10/30/2003

# System Organization (conceptual)

Top View of system

Host System:
- Diagnostics, Booting, Archive
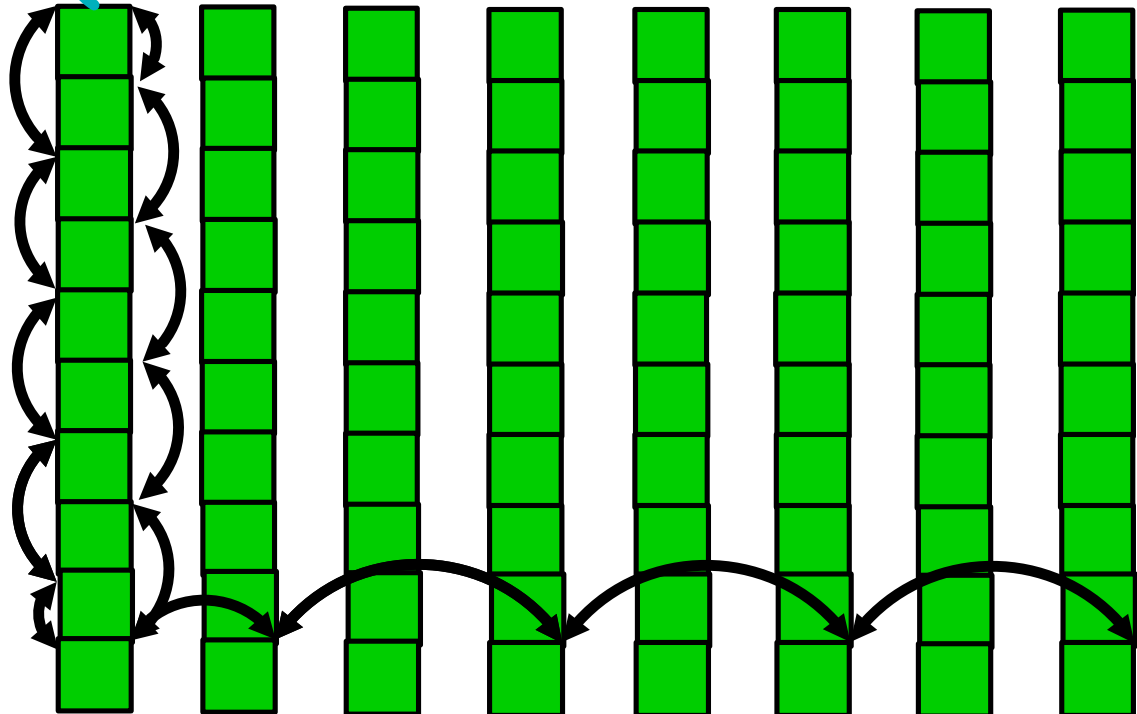- Application dependent requirements

Host Computer

File Server Array
- ~ 500 RAID PC servers
- Gb Ethernet and/or Infiniband
- Application dependent requirements
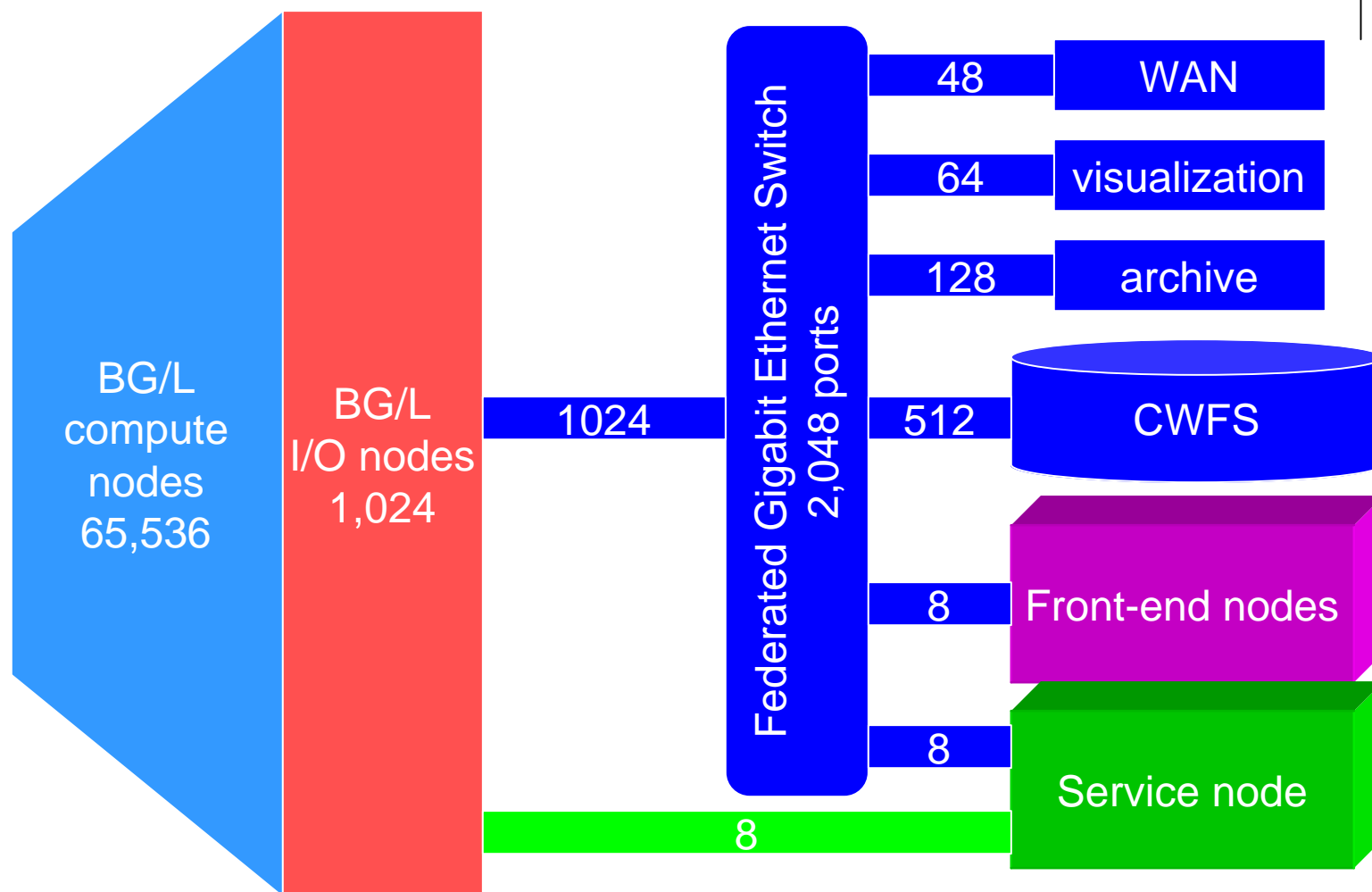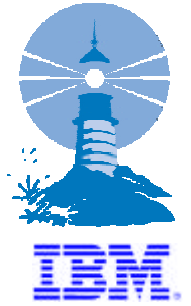
BlueLight Processing Nodes
- 81920 Nodes
  Two major partitions
  - 65536 Nodes Production Platform (256 TFlops Peak)
  - 16384 Nodes partitioned into code development Platforms

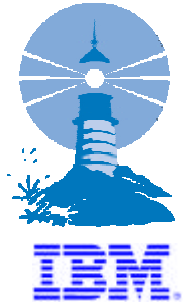50 Feet

# Complete BlueGene/L System at LLNL



BG/L compute nodes 65,536

BG/L I/O nodes 1,024

Federated Gigabit Ethernet Switch 2,048 ports

1024

48 — WAN

64 — visualization

128 — archive

512 — CWFS

8 — Front-end nodes

8 — Service node

8

# Software Design Summary

- **Familiar software development environment and programming models**
  - Fortran, C, C++ with MPI
    - Full language support
    - Automatic SIMD FPU exploitation
  - Linux development environment
    - User interacts with system through FE nodes running Linux – compilation, job submission, debugging
    - Compute Node Kernel provides look and feel of a Linux environment – POSIX system calls (with restrictions)
  - Tools – support for debuggers, hardware performance monitors, trace based visualization

- Scalability to *O(100,000)* processors

# Summary of performance results

- DGEMM:
  - 92.3% of dual core peak on 1 node
  - Observed performance at 500 MHz: 3.7 GFlops
  - Projected performance at 700 MHz: 5.2 GFlops (tested in lab up to 650 MHz)
- LINPACK:
  - 77% of peak on 1 node
  - 69% of peak on 512 nodes (1413 GFlops at 500 MHz)
- sPPM, UMT2000:
  - Single processor performance roughly on par with POWER3 at 375 MHz
  - Tested on up to 128 nodes (also NAS Parallel Benchmarks)
- FFT:
  - Up to 508 MFlops on single processor at 444 MHz (TU Vienna)
  - Pseudo-ops performance (5N log N) @ 700 MHz of 1300 Mflops (65% of peak)
- STREAM – impressive results even at 444 MHz:
  - Tuned:    Copy: 2.4 GB/s, Scale: 2.1 GB/s, Add: 1.8 GB/s, Triad: 1.9 GB/s
  - Standard: Copy: 1.2 GB/s, Scale: 1.1 GB/s, Add: 1.2 GB/s, Triad: 1.2 GB/s
  - At 700 MHz: Would beat STREAM numbers for most high end microprocessors
- MPI:
  - Latency –  < 4000 cycles (5.5 $\mu$s at 700 MHz)
  - Bandwidth – full link bandwidth demonstrated on up to 6 links

# Applications

- BG/L is a general purpose technical supercomputer
- N-body simulation
  - ƒ molecular dynamics (classical and quantum)
  - ƒ plasma physics
  - ƒ stellar dynamics for star clusters, galaxies
- Complex multiphysics code
  - ƒ Computational Fluid Dynamics (weather, climate, sPPM...)
  - ƒ Accretion
  - ƒ Raleigh-Jeans instability
  - ƒ planetary formation and evolution
  - ƒ radiative transport
  - ƒ Magnetohydrodynamics
- Modeling thermonuclear events in/on astrophysical objects
  - ƒ neutron stars
  - ƒ white dwarfs
  - ƒ supernovae
- Radiotelescope
- FFT

# Summary

- Embedded technology promises to be an efficient path toward building massively parallel computers optimized at the system level.

- Cost/performance is ~20x better than standard methods to get to TFlops.

- Low Power is critical to achieving a dense, simple, inexpensive packaging solution.

- Blue Gene/L will have a scientific reach far beyond existing limits for a large class of important scientific problems.

- Blue Gene/L will give insight into possible future product directions.

- Blue Gene/L hardware will be quite flexible. A mature, sophisticated software environment needs to be developed to really determine the reach (both scientific and commercial) of this architecture.

10/30/2003