
Partitioning on Origin 3000 and SNIA

Russ Anderson

RAS/Partitioning Project Tech Lead

rja@sgi.com

SUMMIT 2001

43rd CUG Conference on High Performance Computing & Visualization

Agenda

- Partitioning Overview
- Partitioning Successes
- Current Status
 - Origin 3000 (Mips/Irix)
 - SNIA (Intel/Linux)



SUMMIT 2001

43rd CUG Conference on High Performance Computing & Visualization

Agenda

- Partitioning Overview
 - What is partitioning?
 - Advantages of Partitioning
 - Partitioning Specifics



SUMMIT 2001

43rd CUG Conference on High Performance Computing & Visualization

Partitioning Overview

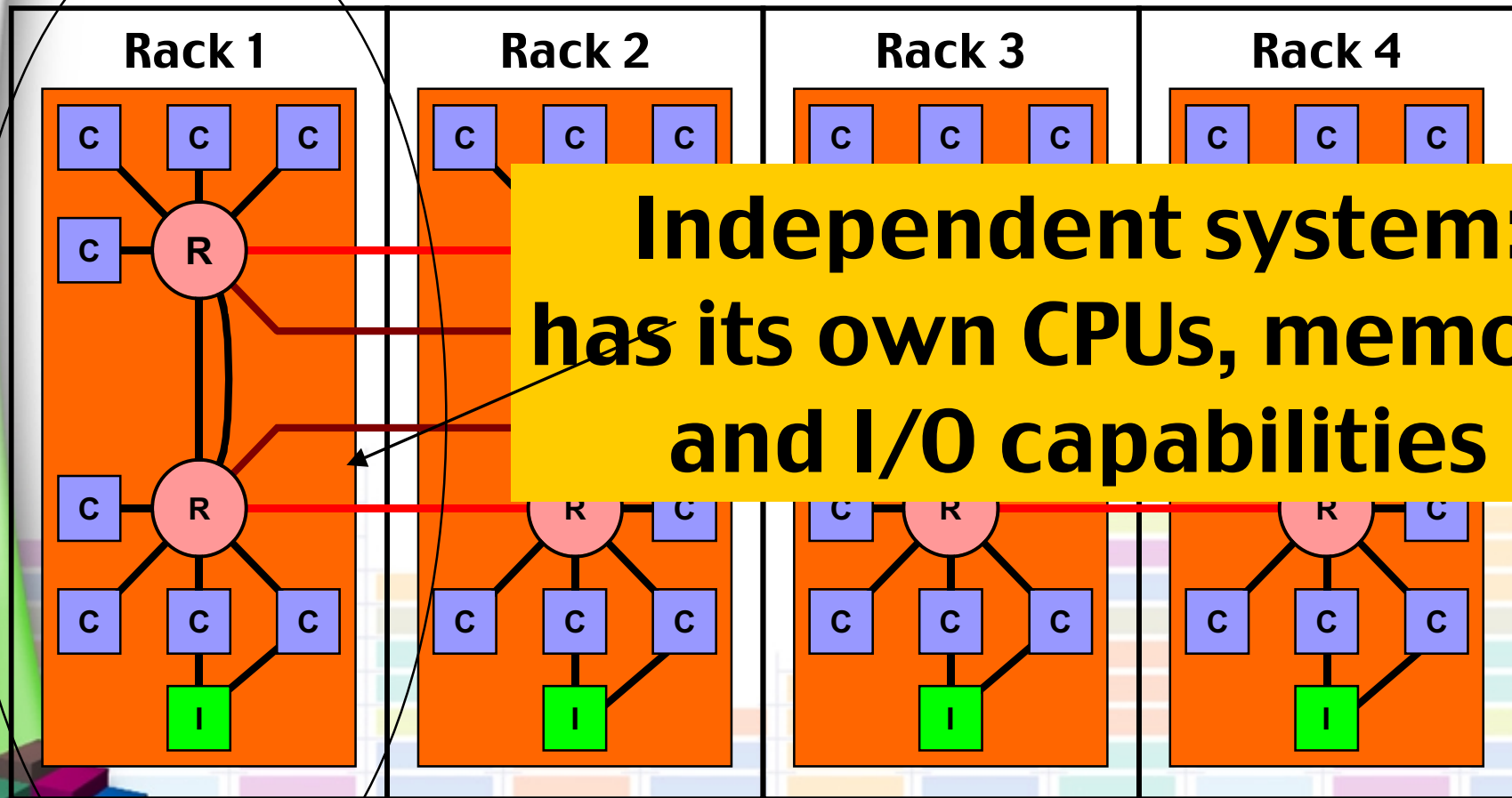


- What is Partitioning?
 - The ability to run a large system as multiple smaller independent systems
 - Use NUMALink as a fast interconnect between partitions (TCP/IP, MPI)



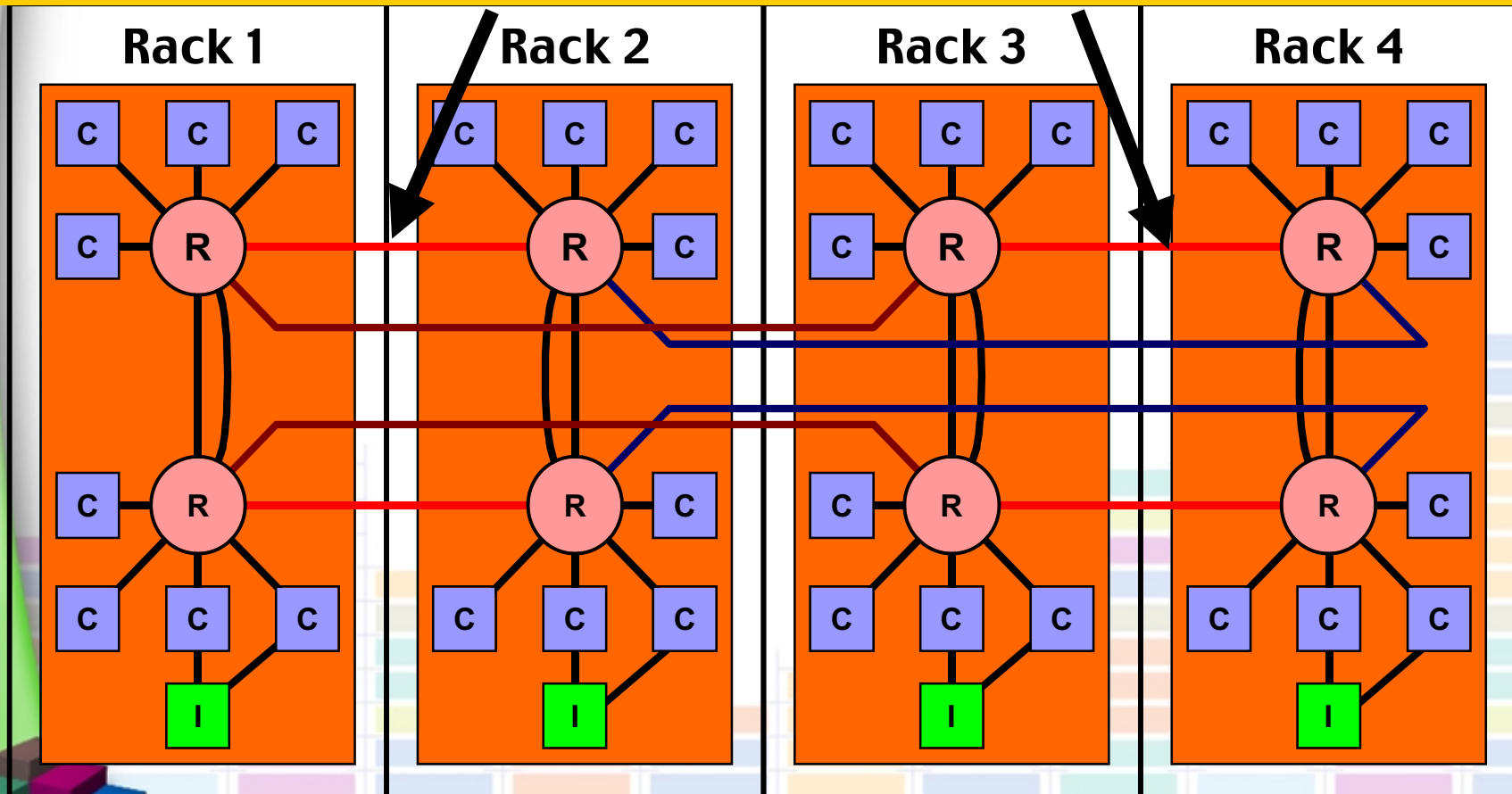
SUMMIT 2001

Partitioning Overview



Partitioning Overview

NUMAlink: Memory to memory communication.



SUMMIT 2001

43rd CUG Conference on High Performance Computing & Visualization

Partitioning Overview



- Advantages of Partitioning
 - Flexible configuration
 - Use NUMALink as a fast interconnect between partitions
 - Fault containment
 - Enhanced HW serviceability



SUMMIT 2001

43rd CUG Conference on High Performance Computing & Visualization

■ Advantages of Partitioning

■ Flexible configuration

Change configurations without re-cabling

Rolling OS Upgrades

Small kernel image & Large kernel image
in single system

Change partition size (Linux)

■ Advantages of Partitioning

- Use NUMALink as a fast interconnect between partitions

TCP/IP communication (NFS, CXFS)

MPI / XPMEM



SUMMIT 2001

Partitioning Overview

■ Advantages of Partitioning

■ Fault containment

Increase MTTI

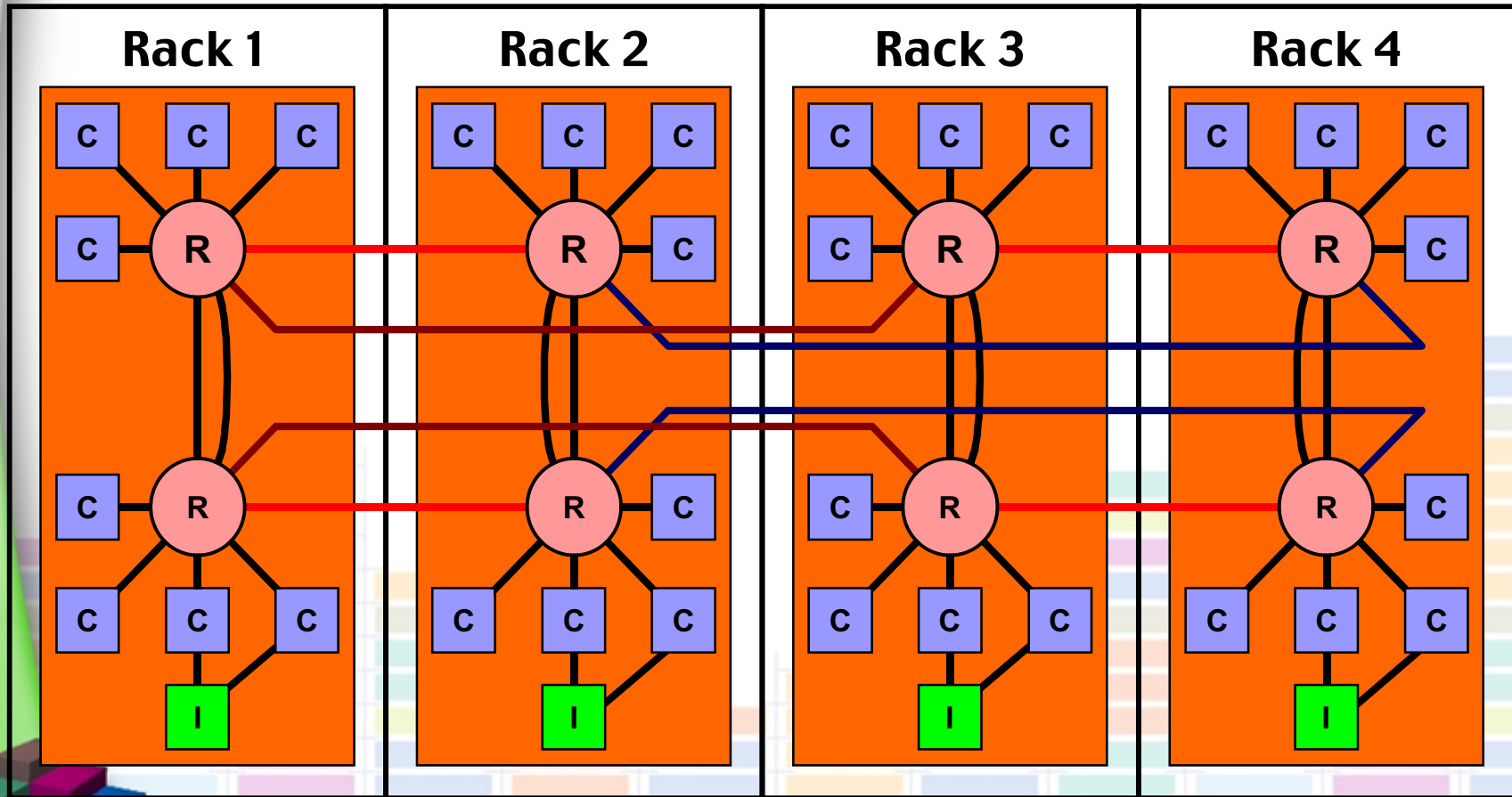
Maximize availability

■ Enhanced HW serviceability

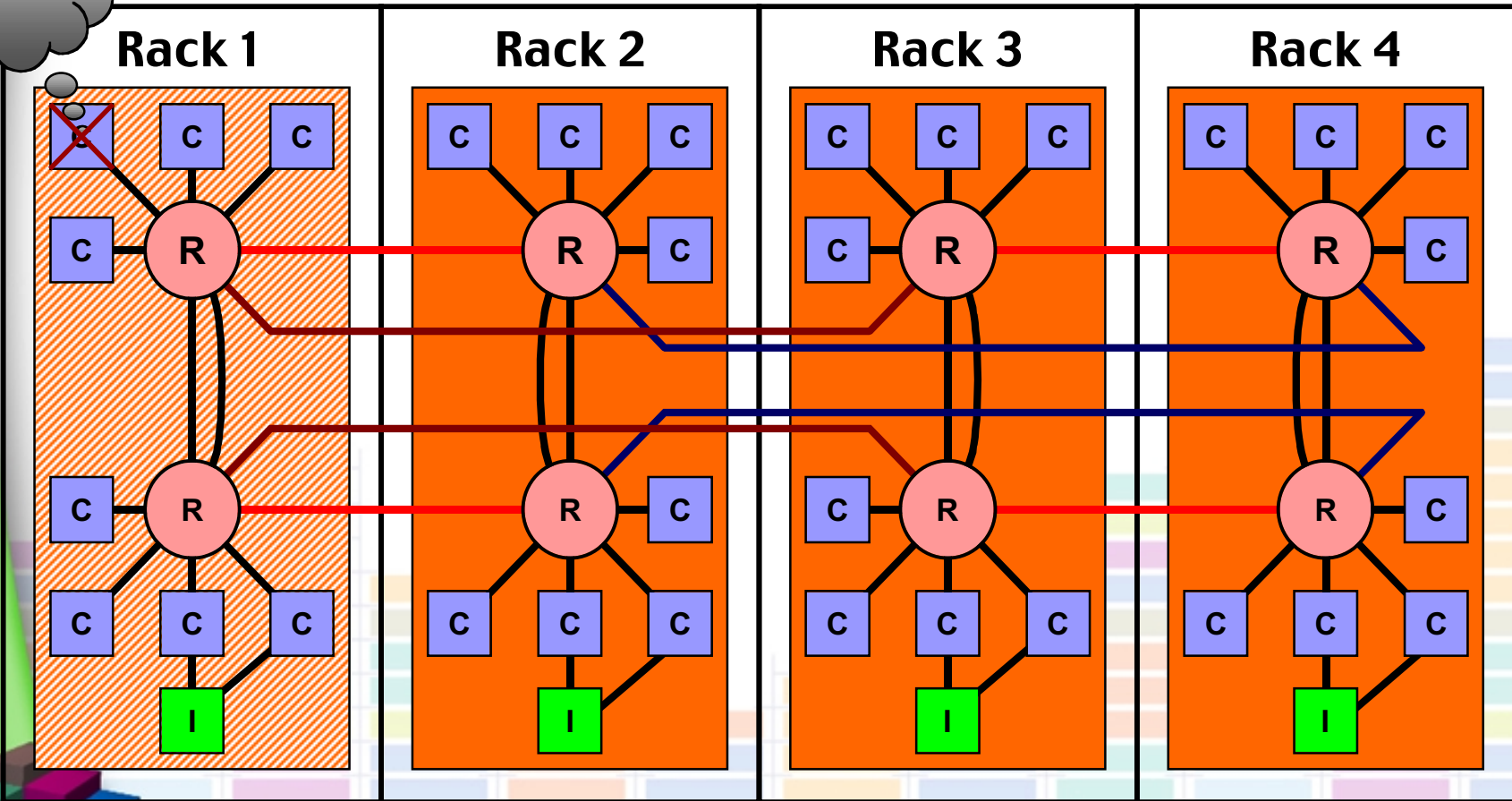


SUMMIT 2001

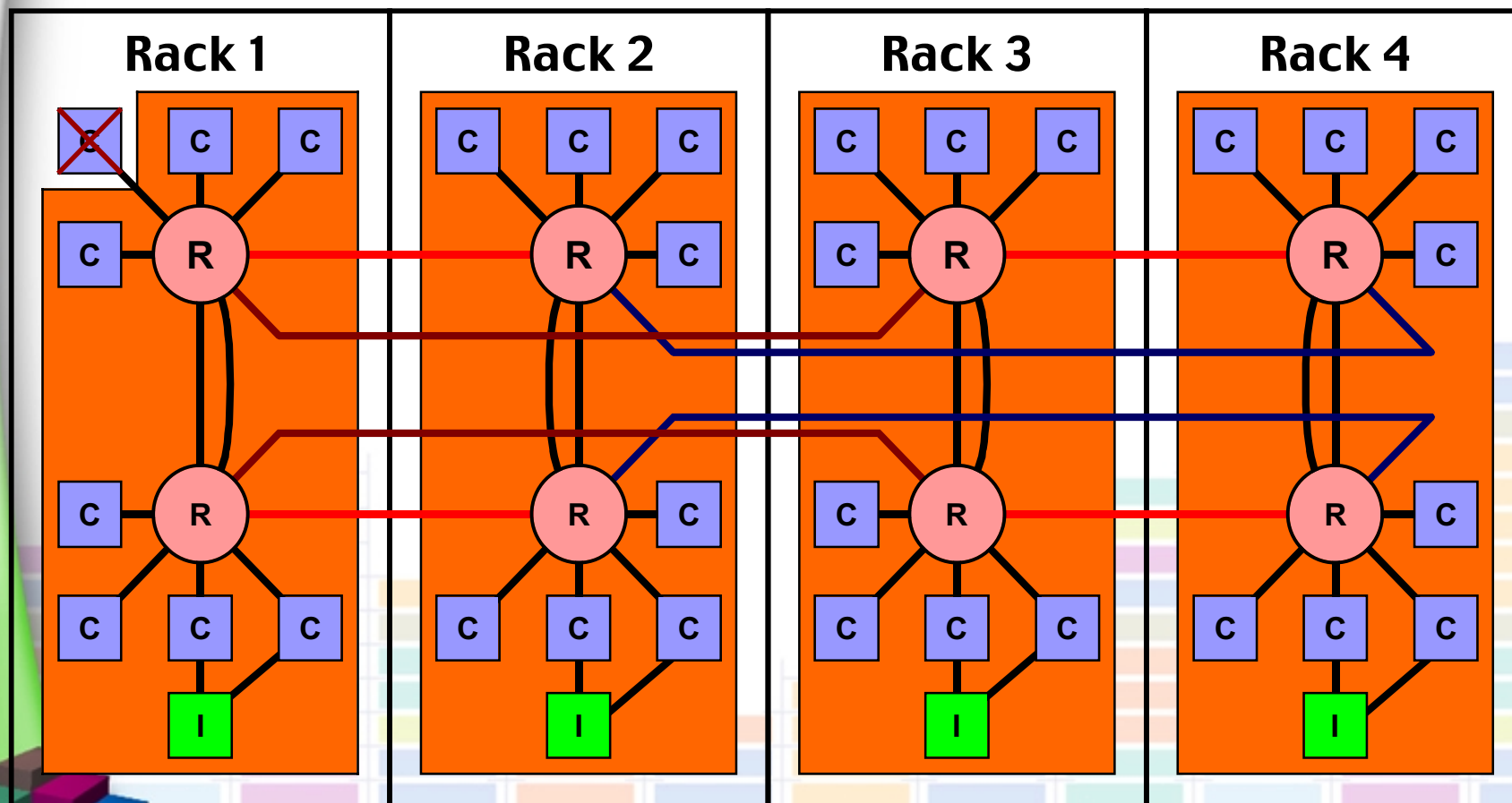
Partitioning Example



Partitioning Example

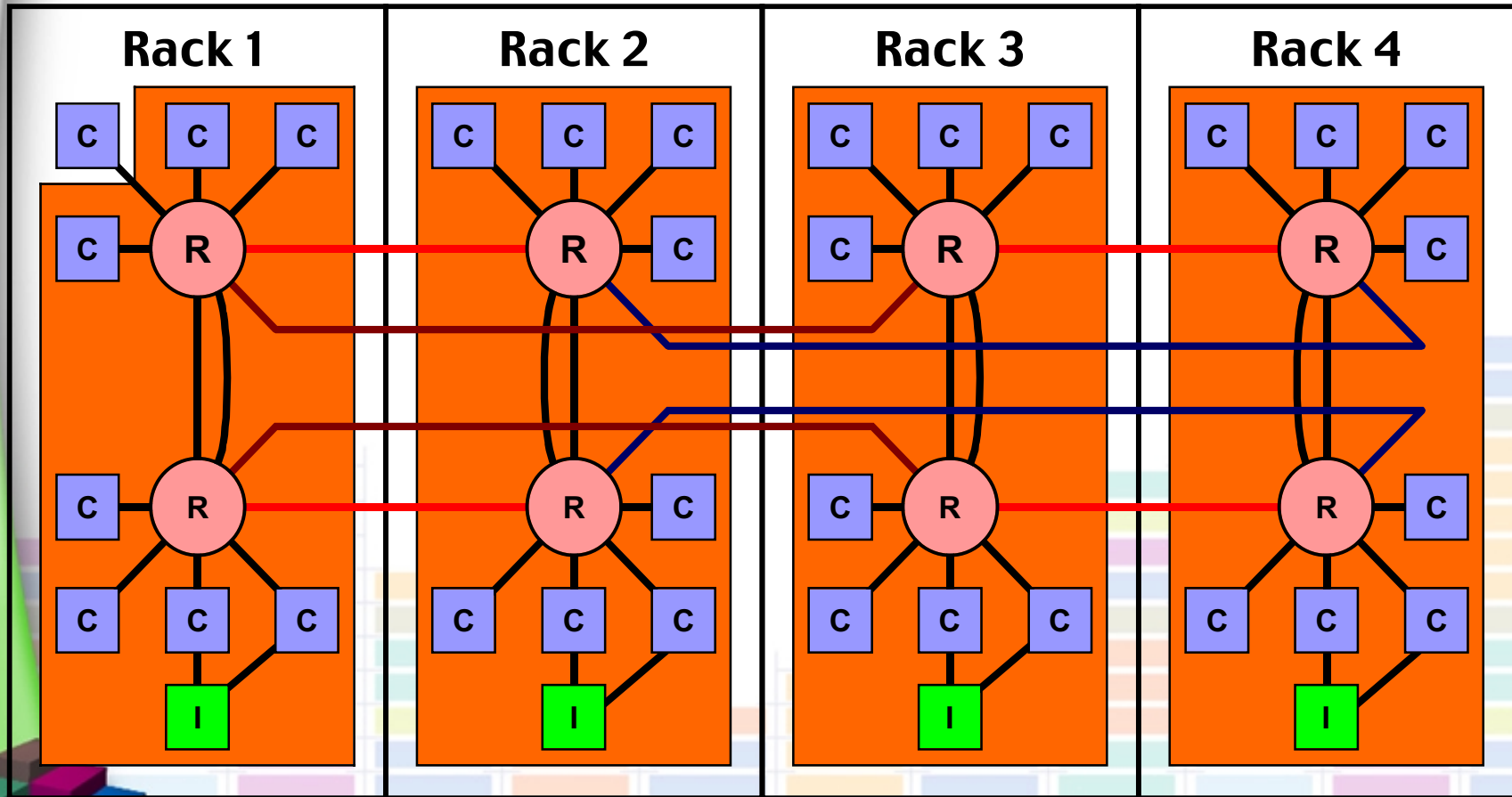


Partitioning Example



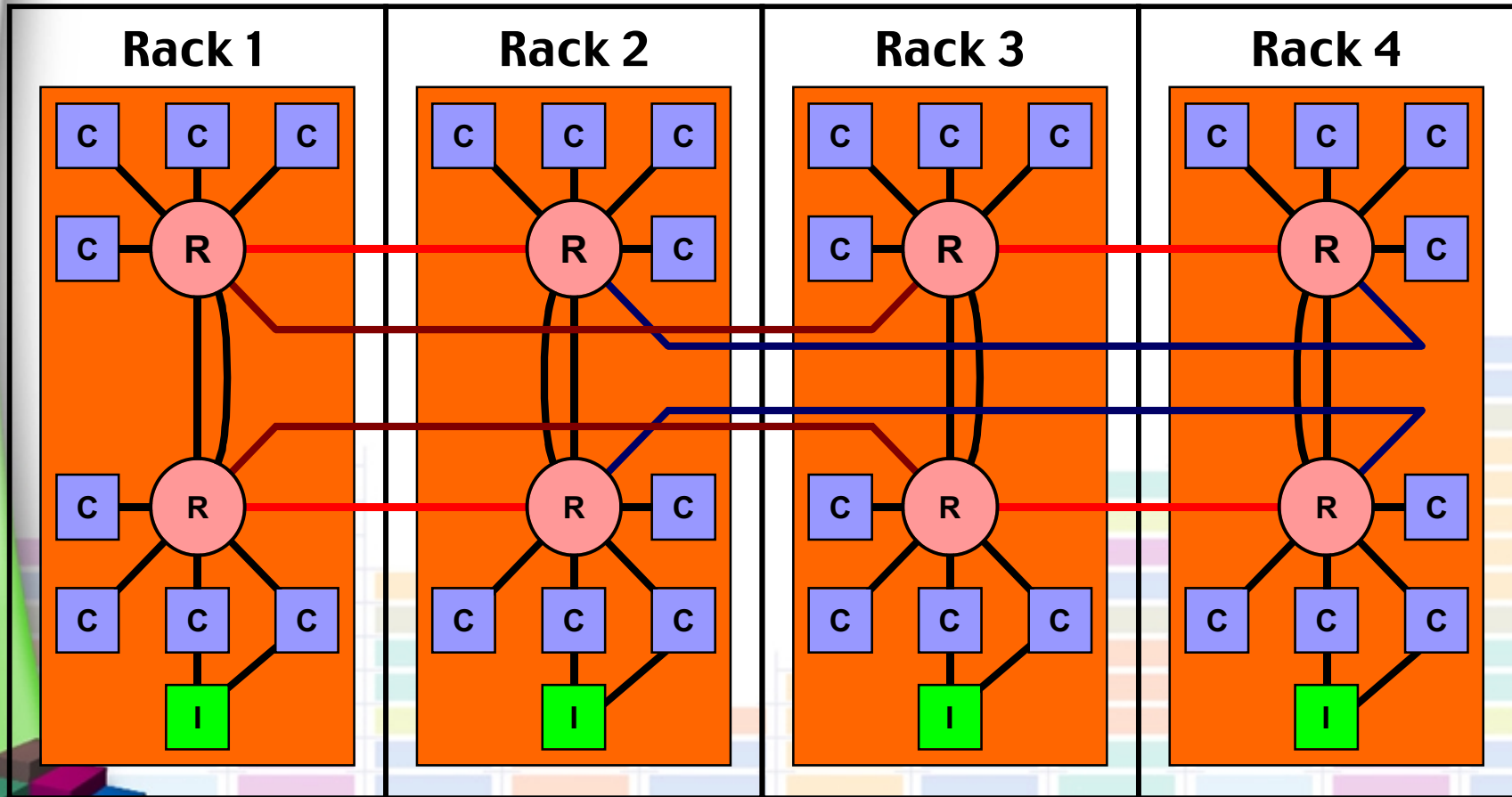
SUMMIT 2001

Partitioning Example



SUMMIT 2001

Partitioning Example



SUMMIT 2001

■ Partitioning Specifics

- Origin 3000 Architecture Overview
- Hardware support for partitioning
- Partitioning Configuration Rules
- Brick replacement



SUMMIT 2001

Origin 3000 Partitioning



■ NUMAflex technology

■ Bricks

C-Brick - Compute (MIPS)

I-Brick - Basic I/O slots/Root disk/Etc.

X-Brick - XTALK I/O Brick

P-Brick - PCI I/O Brick

R-Brick - Router Brick

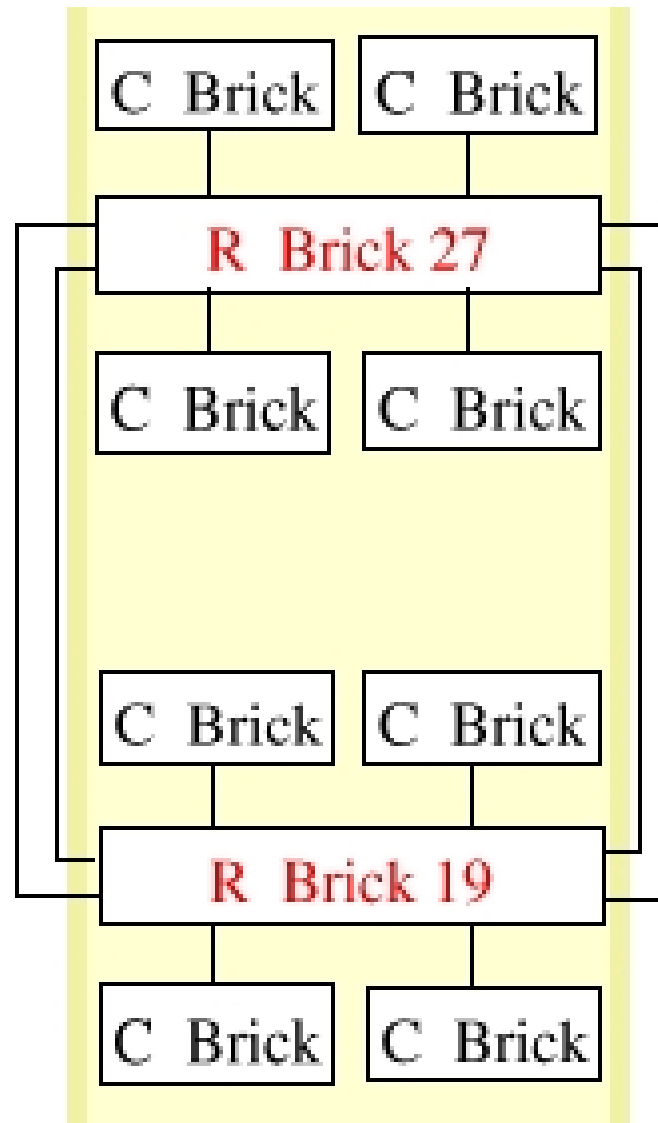
■ Building Block Approach



SUMMIT 2001

43rd CUG Conference on High Performance Computing & Visualization

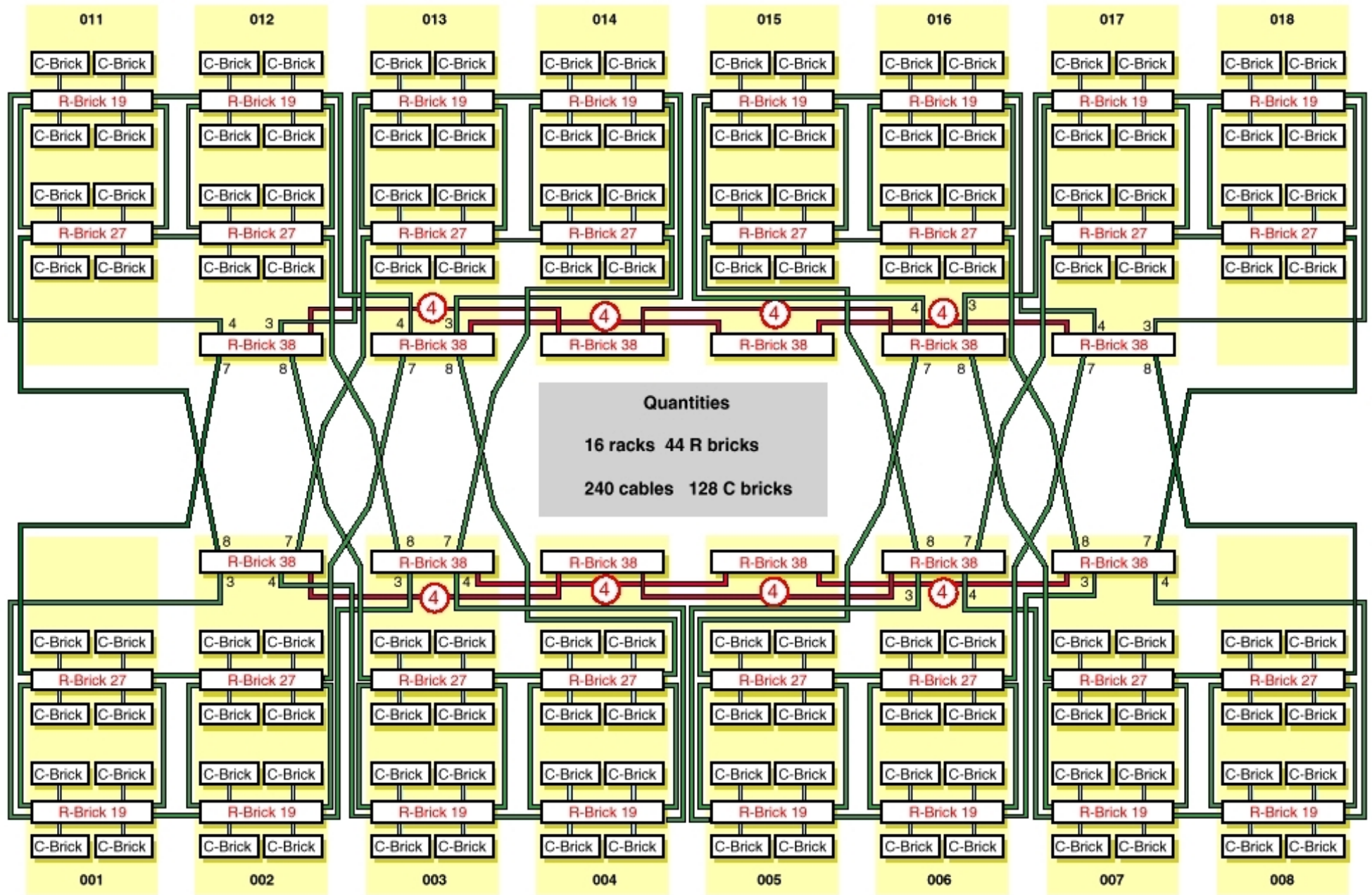
■ Rack



SUMMIT 2001

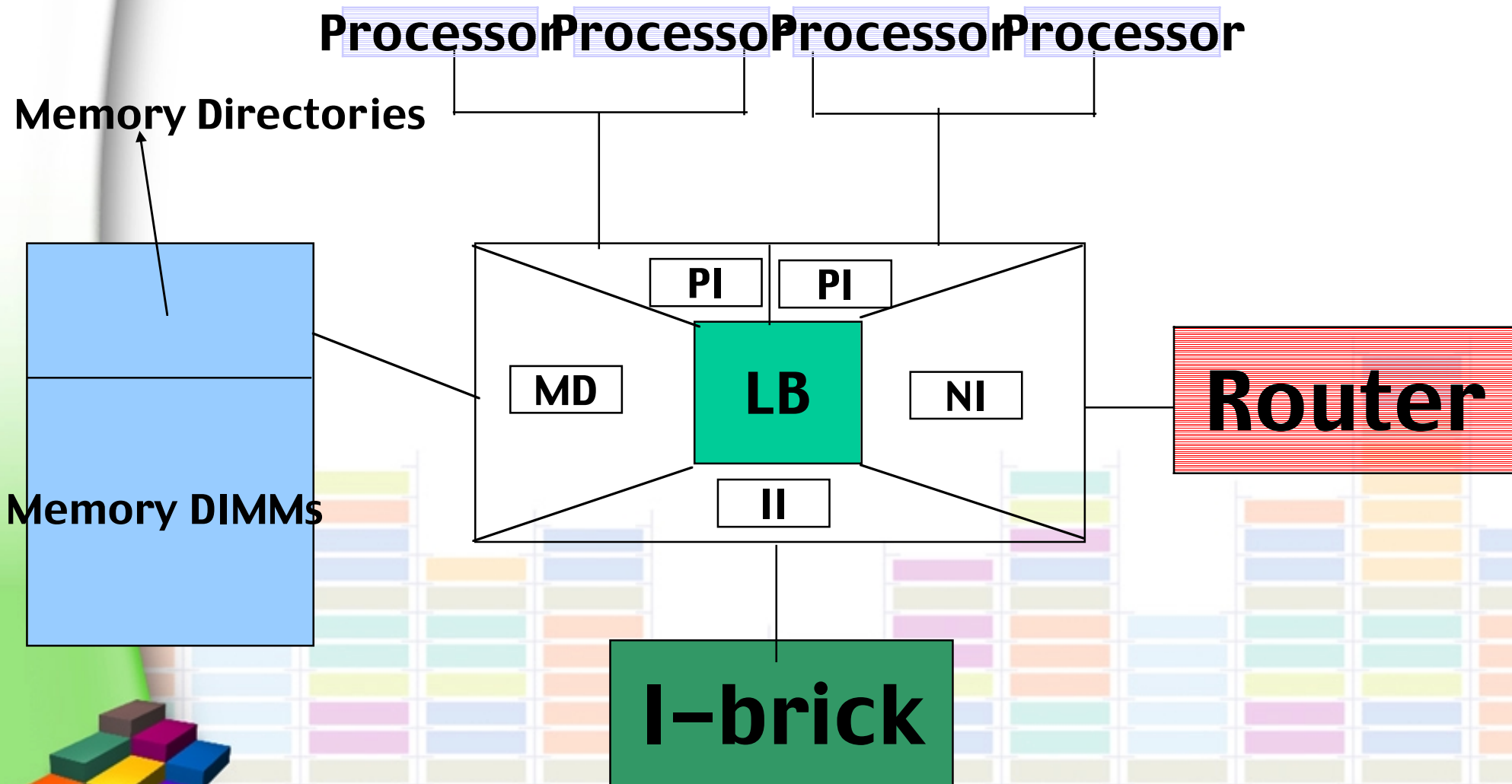
Origin 3000 512p

sgi



SUMMIT 2001

SN1 Architecture



Partitioning Hardware Support

sgi

- Partitioning Hardware Support
 - Memory Protection
 - Reset Fences
 - BTE and NUMALink



SUMMIT 2001

43rd CUG Conference on High Performance Computing & Visualization

Partitioning Hardware Support

sgi

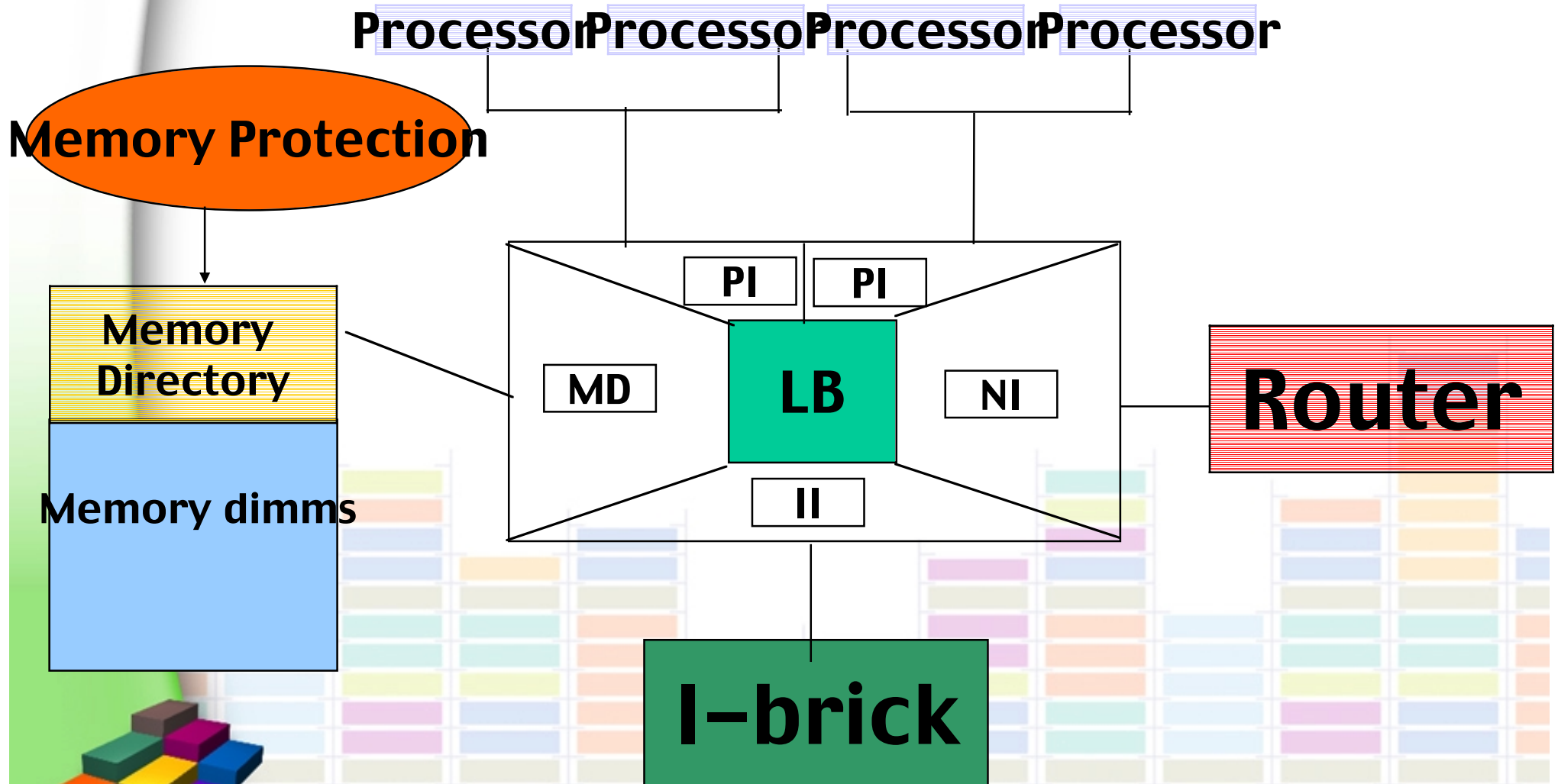
- Memory Protection
 - Built into Hub chip hardware
 - Protect a partition from unexpected writes from another partition (Fault Containment)
 - Protection can be modified to allow access to specific memory pages

SUMMIT 2001

43rd CUG Conference on High Performance Computing & Visualization



Origin 3000 Architecture Memory Protection

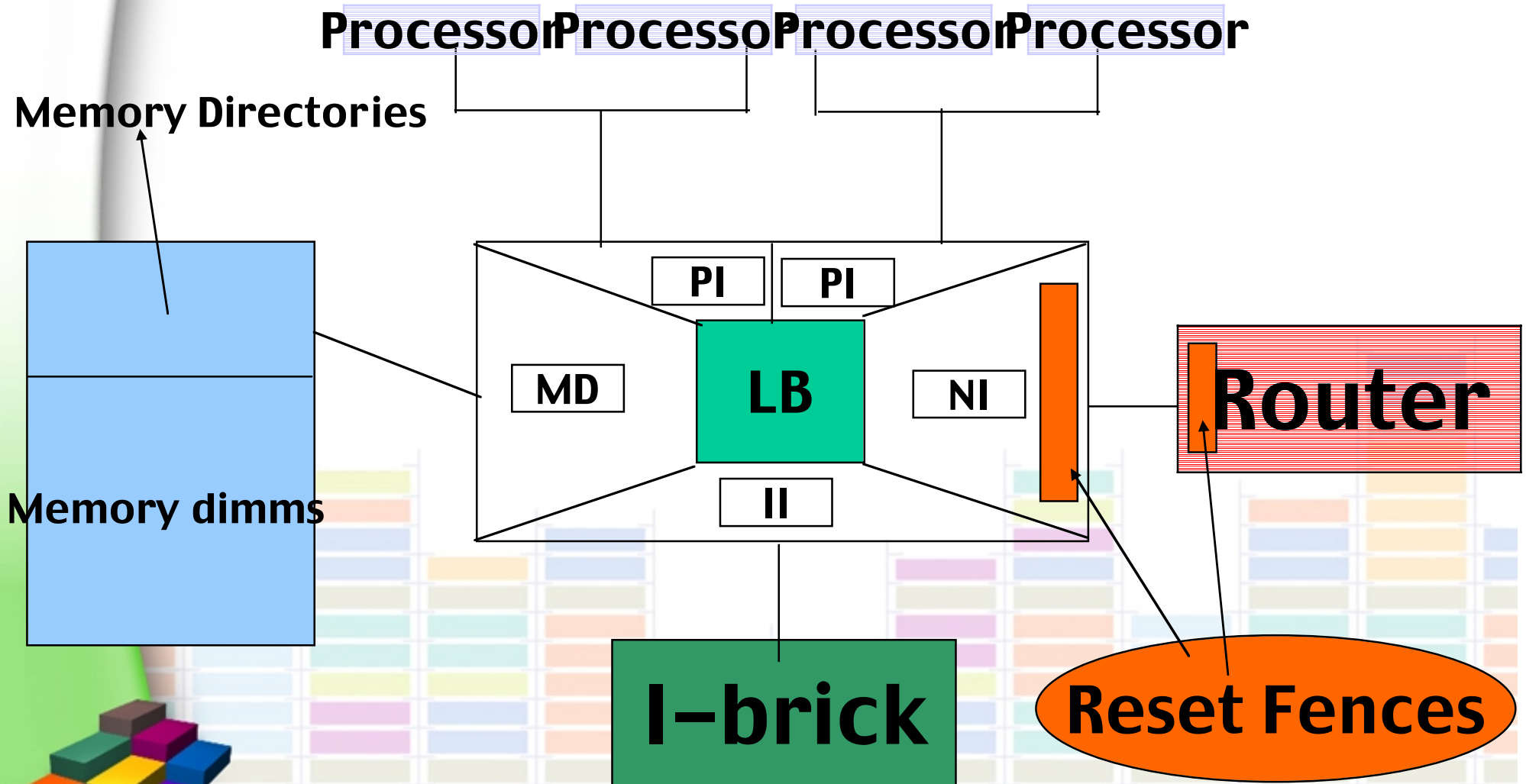


Partitioning Hardware Support

- Reset Fences
 - Built into Hub and Router chip hardware
 - Protect a partition from hardware resets in other partitions
 - Used to support concurrent brick replacement



Origin 3000 Architecture Reset Fences



SUMMIT 2001

Partitioning Hardware Support

sgi

- Block Transfer Engine
 - Built into Hub chip hardware
 - Transfers data between partitions (without changing memory protection)
 - Allows processors to do other work

SUMMIT 2001

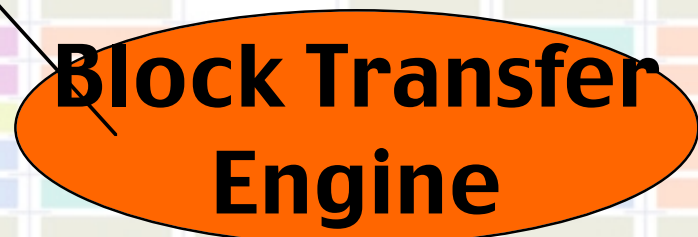
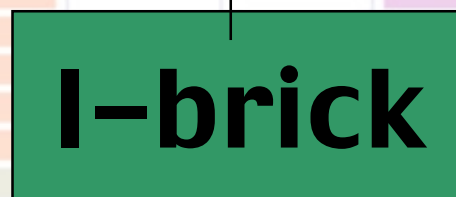
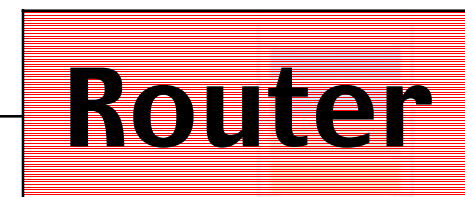
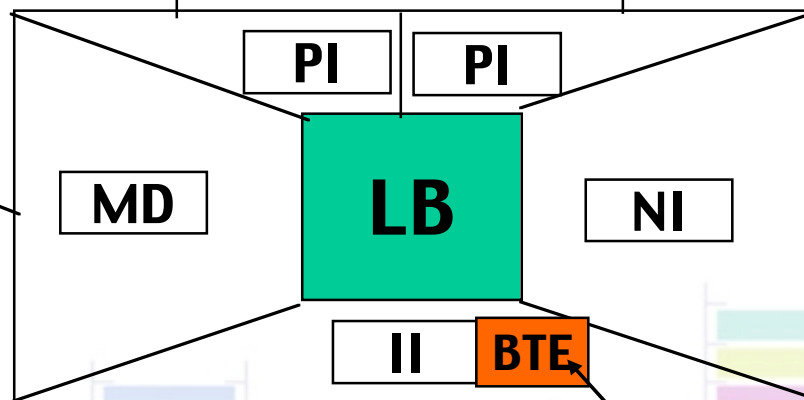
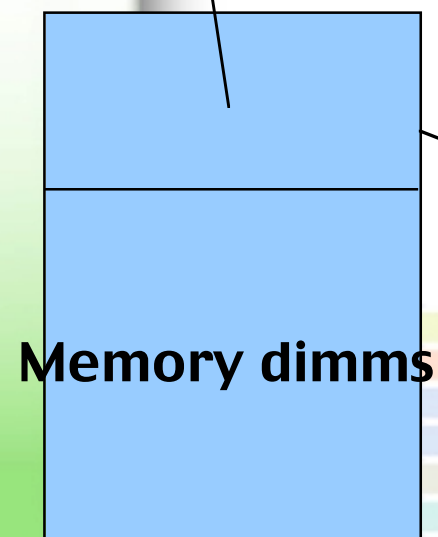
43rd CUG Conference on High Performance Computing & Visualization

Origin 3000 Architecture Block Transfer Engine

sgi

Processor Processor Processor Processor

Memory Directories



SUMMIT 2001

■ Current Partitioning Rules

- Partitions Run like stand alone systems
- Partition on rack & half-rack boundaries
- I/O Bricks belong to the partition of the attached C-Brick



SUMMIT 2001

- C-Brick Replacement
 - Power down independently
 - Diagnostic Testing
 - Re-integrate with partition Reboot



Agenda

- Partitioning Overview
- **Partitioning Successes**
- Current Status
 - Origin 3000 (Mips/Irix)
 - SNIA (Intel/Linux)



SUMMIT 2001

43rd CUG Conference on High Performance Computing & Visualization

Origin 3000 Partitioning Successes



- Partitioned system running at SARA
 - 128p partitioned into 2x64p
 - 512p partitioned 2x256p and 4x128p
 - 192p (6 rack) system partitioned (4x32p, 1x64p) in production
- Production + development partitioned system
- Partitioned systems with GFX & GSN



SUMMIT 2001

43rd CUG Conference on High Performance Computing & Visualization

Origin 3000 Current Status



- Irix 6.5.9 - Hard wall partitioning supported
- Irix 6.5.10 - Cross Partition communication (TCP/IP)
- Irix 6.5.11-12 - More supported configurations
- Irix 6.5.13 -
 - XPMEM (support for MPI)
 - BTE performance enhancements
 - Improved system administration



SUMMIT 2001

43rd CUG Conference on High Performance Computing & Visualization

SNIA Current Status

- Itanium/Linux SNIA development system running two partitions with TCP/IP over NUMALink
- SNIA low-level functionality working
 - Reset fences
 - Memory protection



SUMMIT 2001

43rd CUG Conference on High Performance Computing & Visualization

Origin 3000 Partitioning

sgi

Conclusion



SUMMIT 2001
43rd CUG Conference on High Performance Computing & Visualization