



**thern Lights**  
computer.com



**≡ SCALI**

# Scalable Linux Systems

**Einar Rustad**

**Scali AS**

**einar@scali.com**

**<http://www.scali.com>**

**Junzo Tamada**

**Northern Lights Computers**

**Kåre Løchsen**

**Dolphin Interconnect Solutions**



---

## Scalable System Requirements

---

- **Balanced Hardware Resources**
  - High Processor Speed
  - Scalable Memory and Storage
  - High Bandwidth Interconnect
  - Low latency Communication
- **Efficient Middleware, Standard APIs**
- **Ease of Use**
- **Easy and Flexible System Administration**



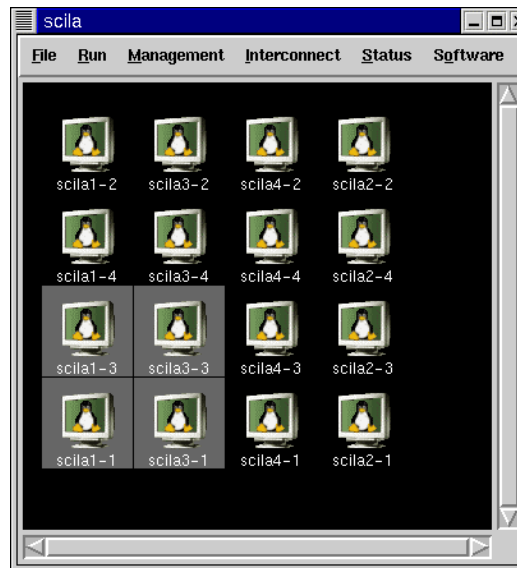
## The Value Chain

**Interconnect  
Hardware**



**WulfKit™**

**Interconnect Software  
(MPI, SCI-SAN)  
System Management-  
Tools**



**Complete-  
Integrated Systems  
Support**



NL Server 2000



## Dolphin Technology

- Based on SCI, the only bus extension technology
- Proven over years in applications
  - Clustering of SUN's high availability servers
  - Fujitsu Siemens large scale IO system
  - Data General AViiON ccNUMA servers
  - Mirage and Rafale flight computers
  - Scali and Northern Lights ISP, ASP and I
- Solutions for Servers and Embedded Computing
- 2µs ° Application Latency, 500 Mbytes/s Link Speed
- Dolphin sells Chips, Cards, Switches and licenses technology
- Dolphin makes the Scali and WulfKit card assembly and sells the WulfKit to system builders

**WulfKit™**





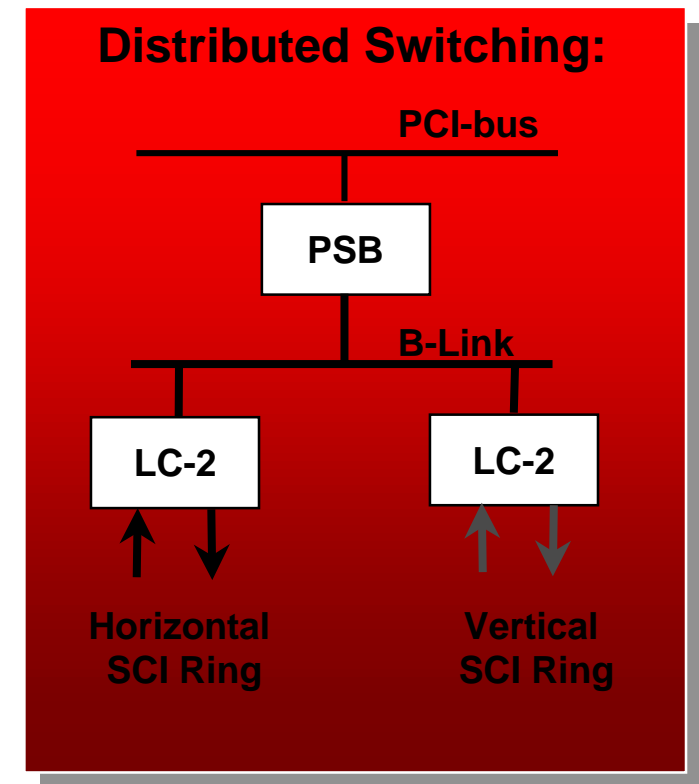
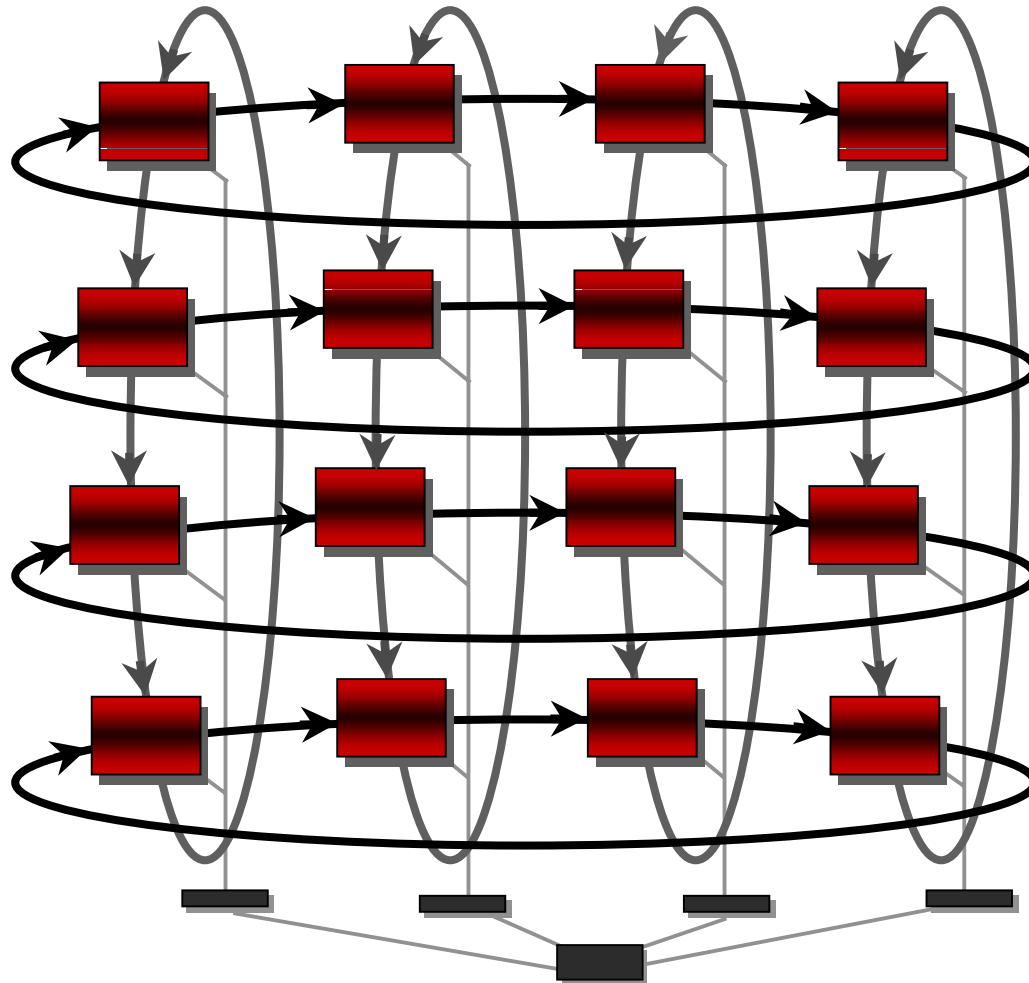
---

# Scalable Linux Systems Advantages

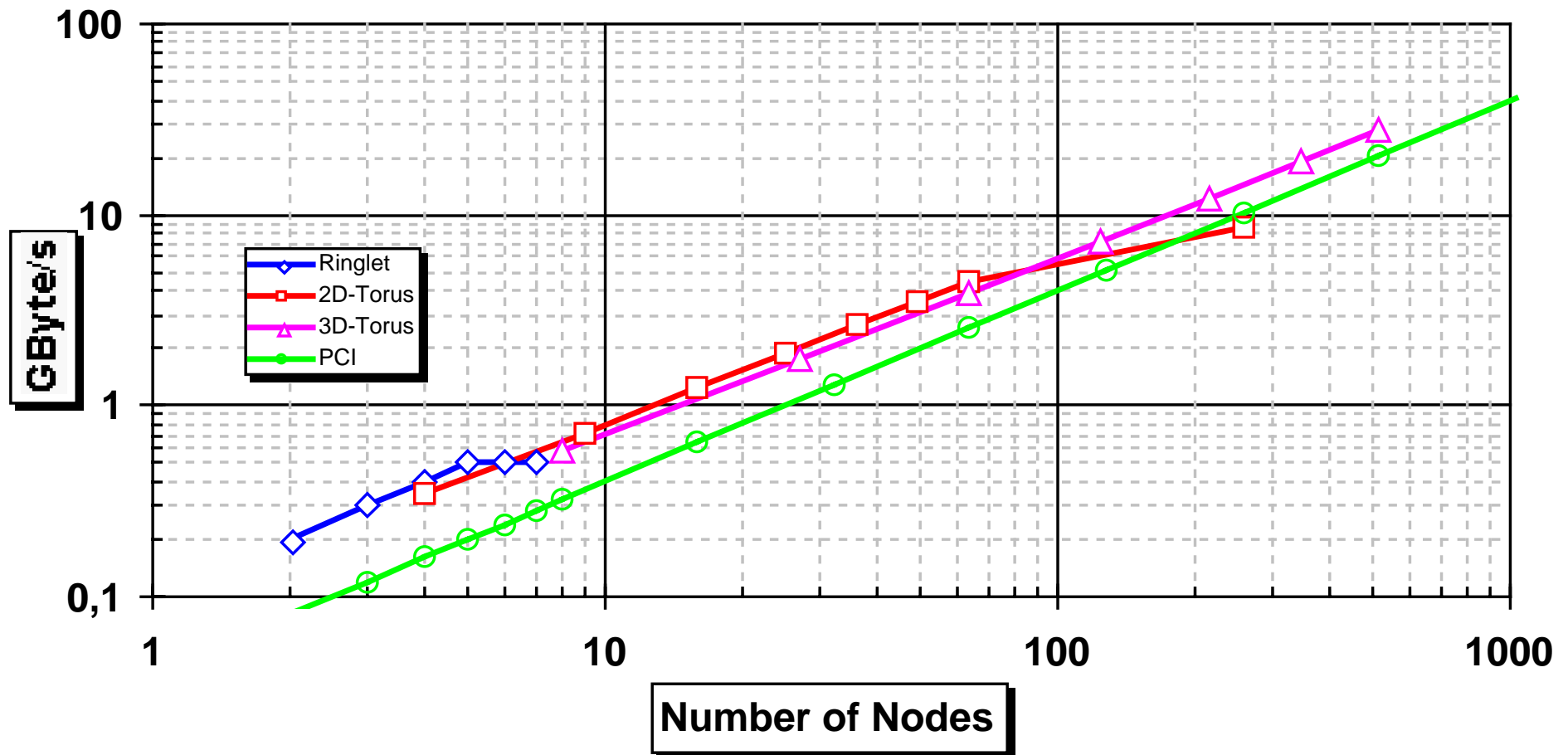
---

- **Industry Standard Programming Model - MPI**
  - Porting = Recompile
- **Lower Cost**
  - COTS based Hardware = lower system price
  - Lower Total Cost of Ownership
- **Better Performance**
  - Always "Latest & Greatest" Processors
  - Superior Standard Interconnect - SCI
- **Scalability**
  - Scalable to thousands of Processors
- **Redundancy**
- **Single System Image to users and administrator**
- **Choice of Front-End OS**
  - Linux
  - Solaris
  - Windows NT

# Torus Topology - Distributed Switching

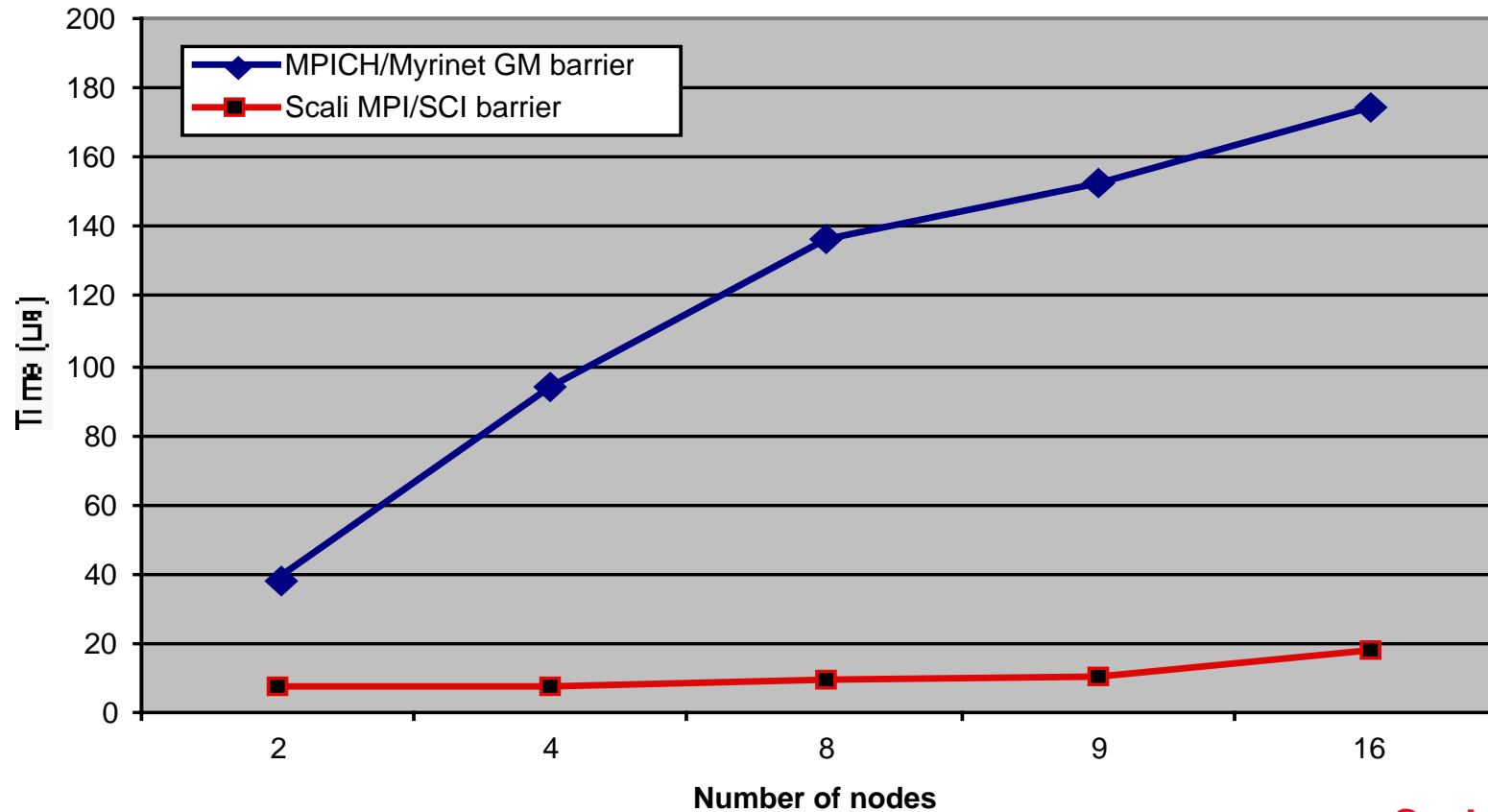


# Theoretical Scalability



# Versus Myrinet (1)

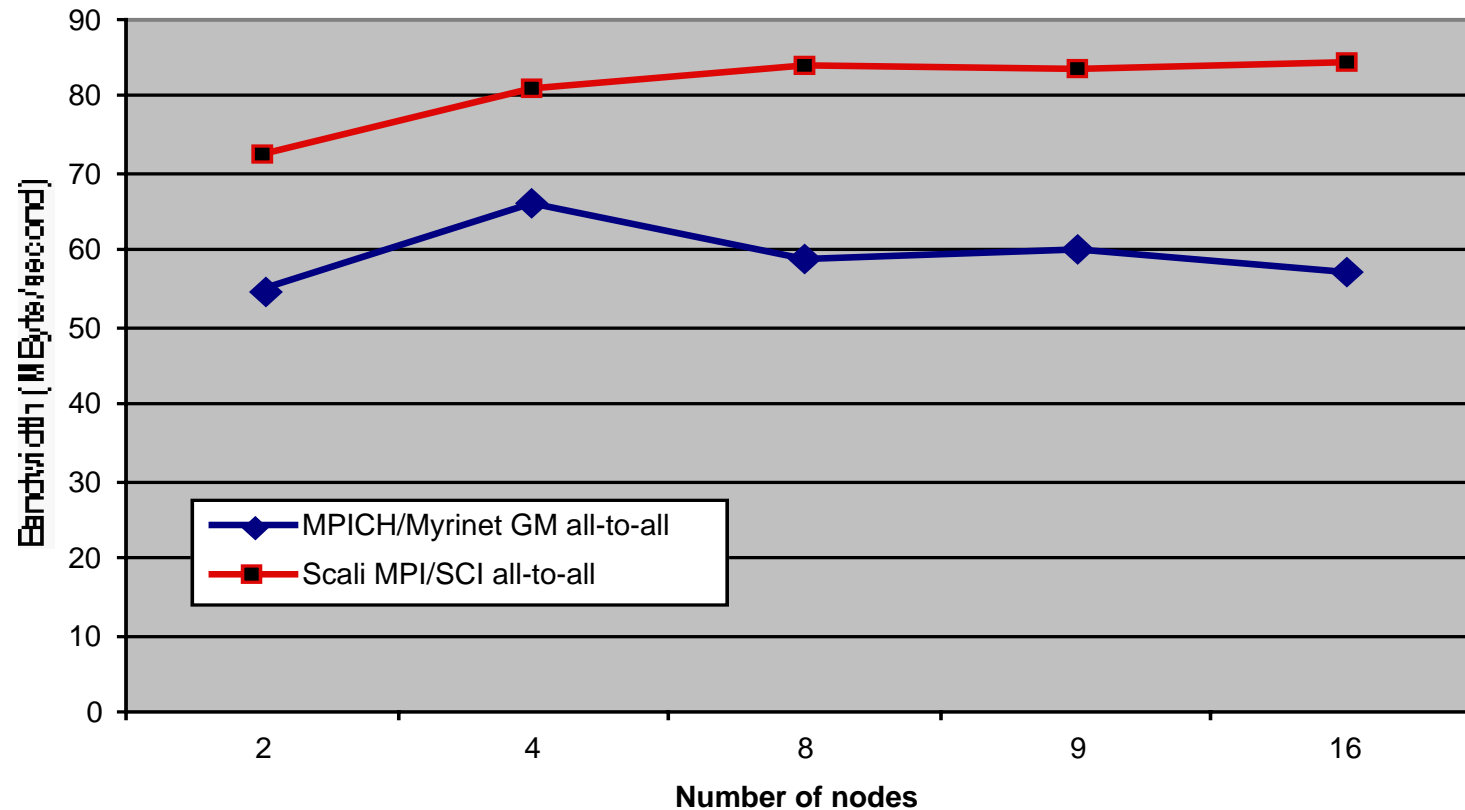
Barrier synchronization





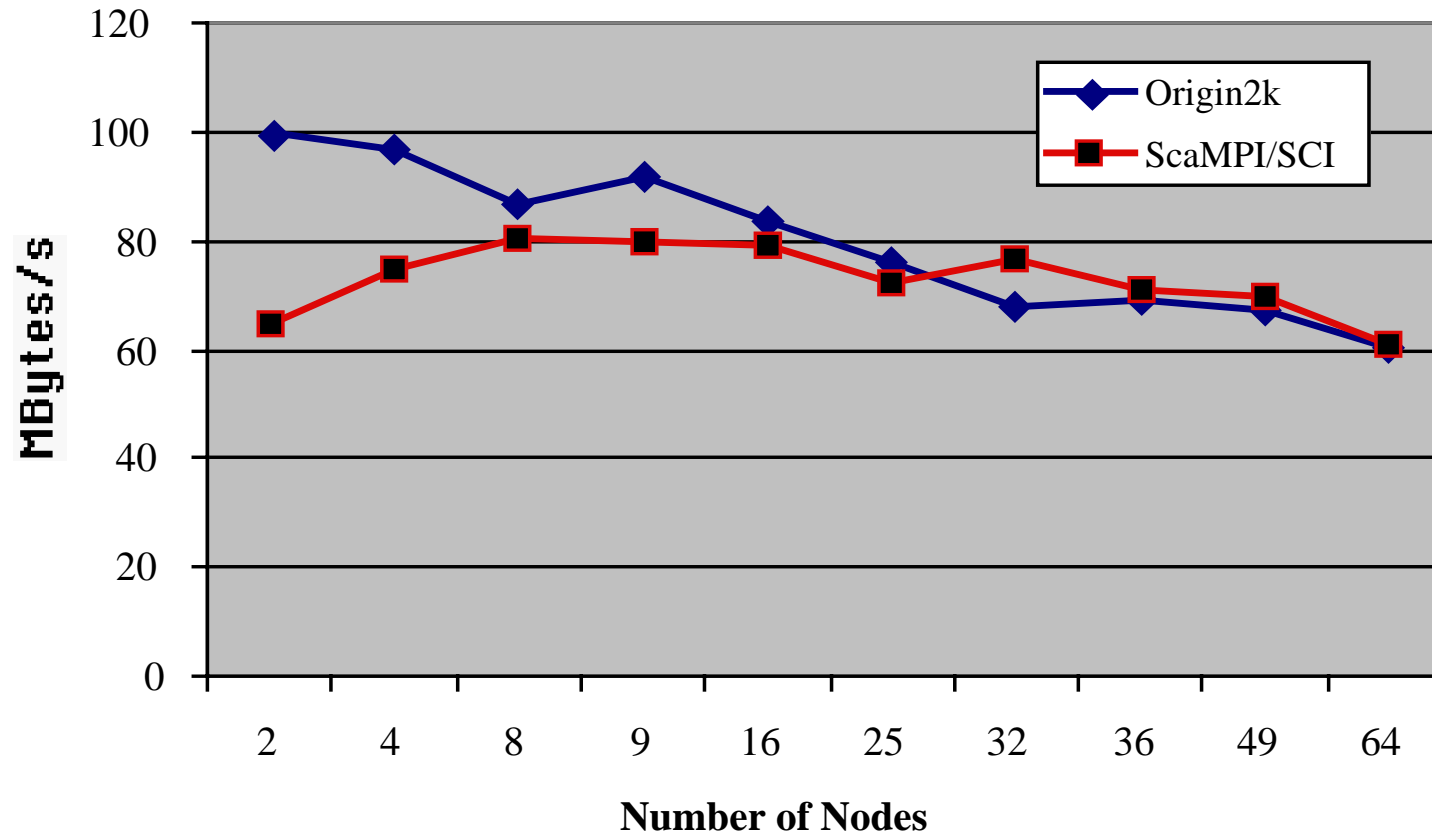
## Versus Myrinet (2)

All-to-all performance



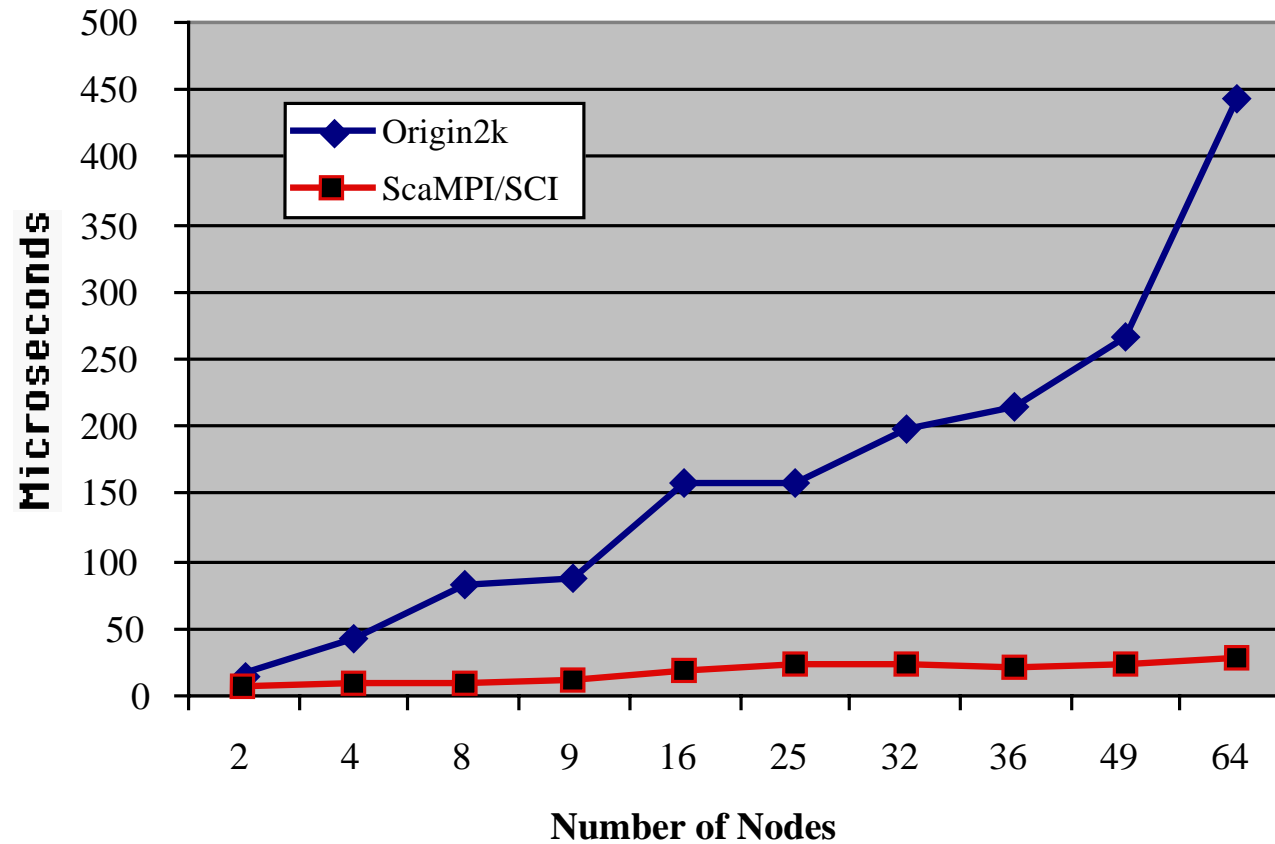
## Versus Origin 2000 (1)

All-to-all Bandwidth per Node



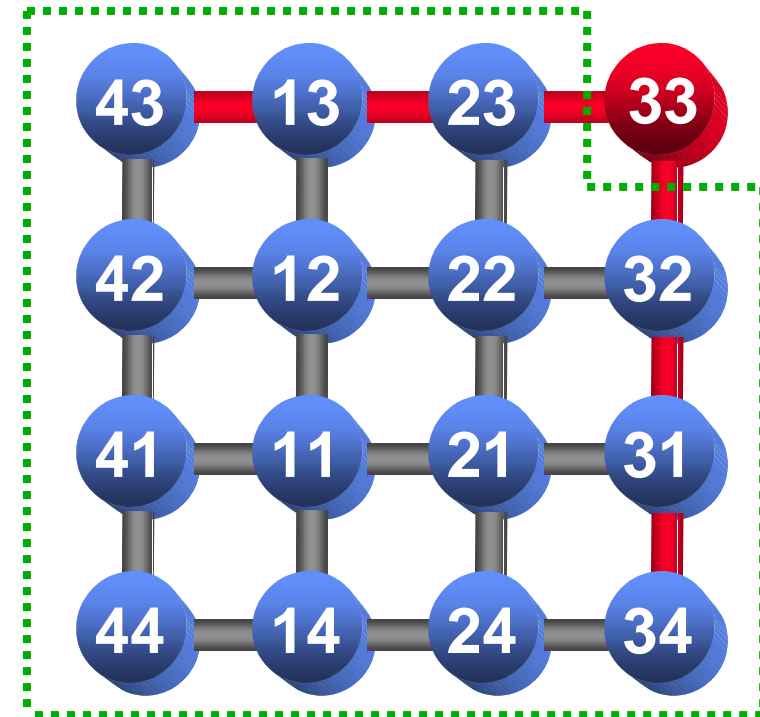
## Versus Origin 2000 (2)

### Barrier Synchronization

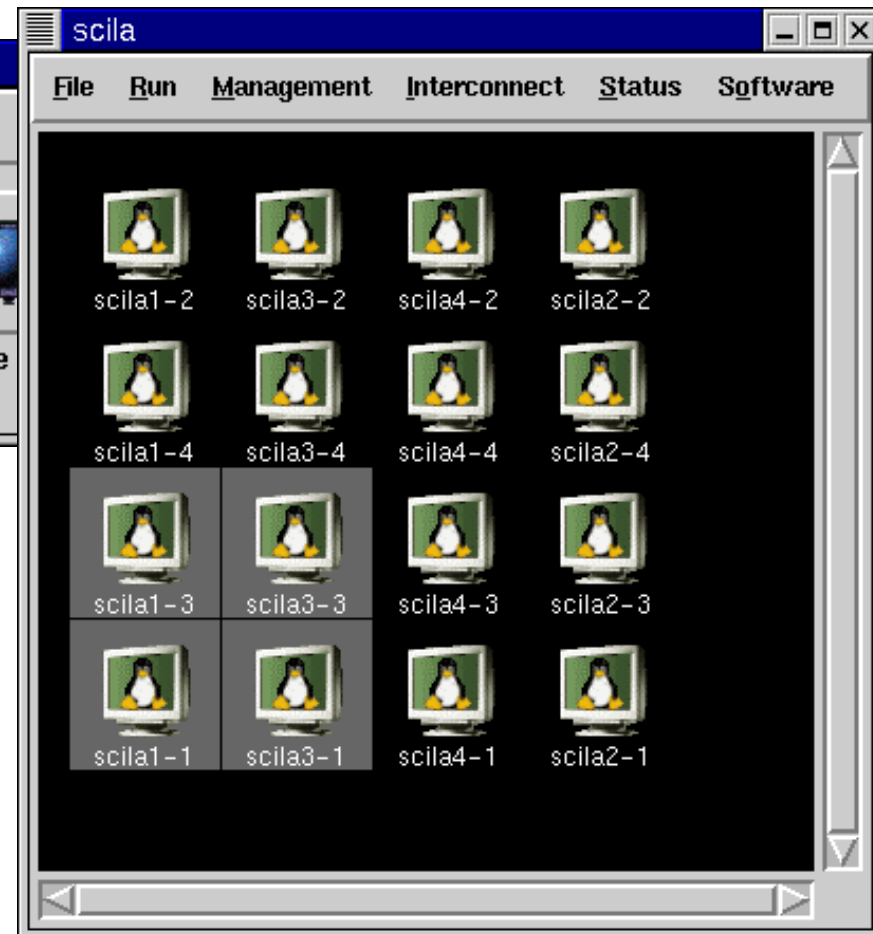
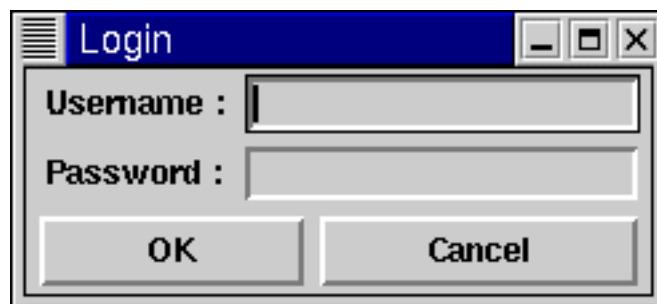
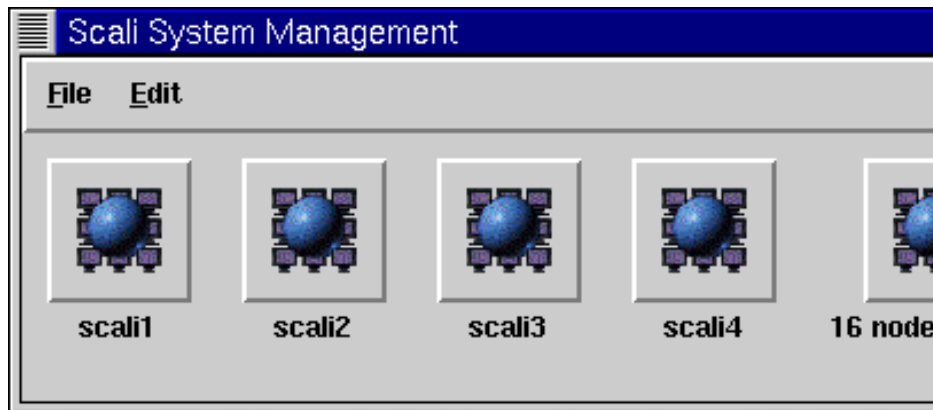


## Fault Tolerance

- **Automatic Rerouting**
- **Scali advanced routing algorithm:**
  - From the Turn Model family of routing algorithms
- **All nodes but the failing ones can be utilised as one big partition**

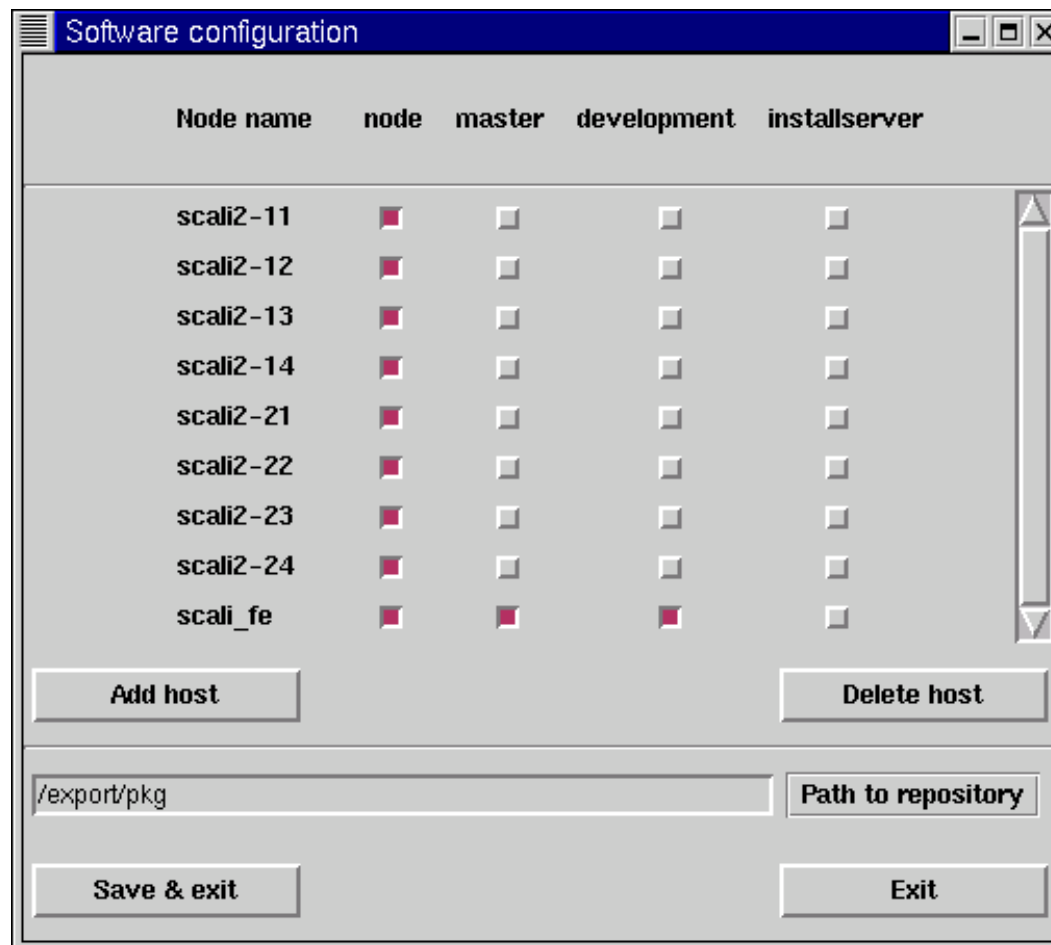


# The Scali Universe

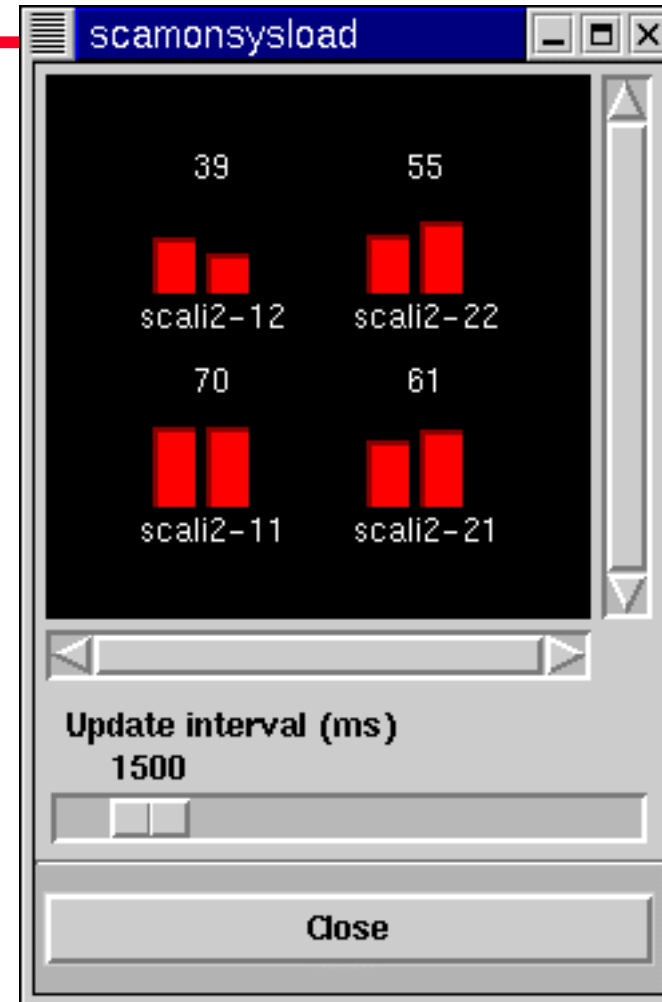
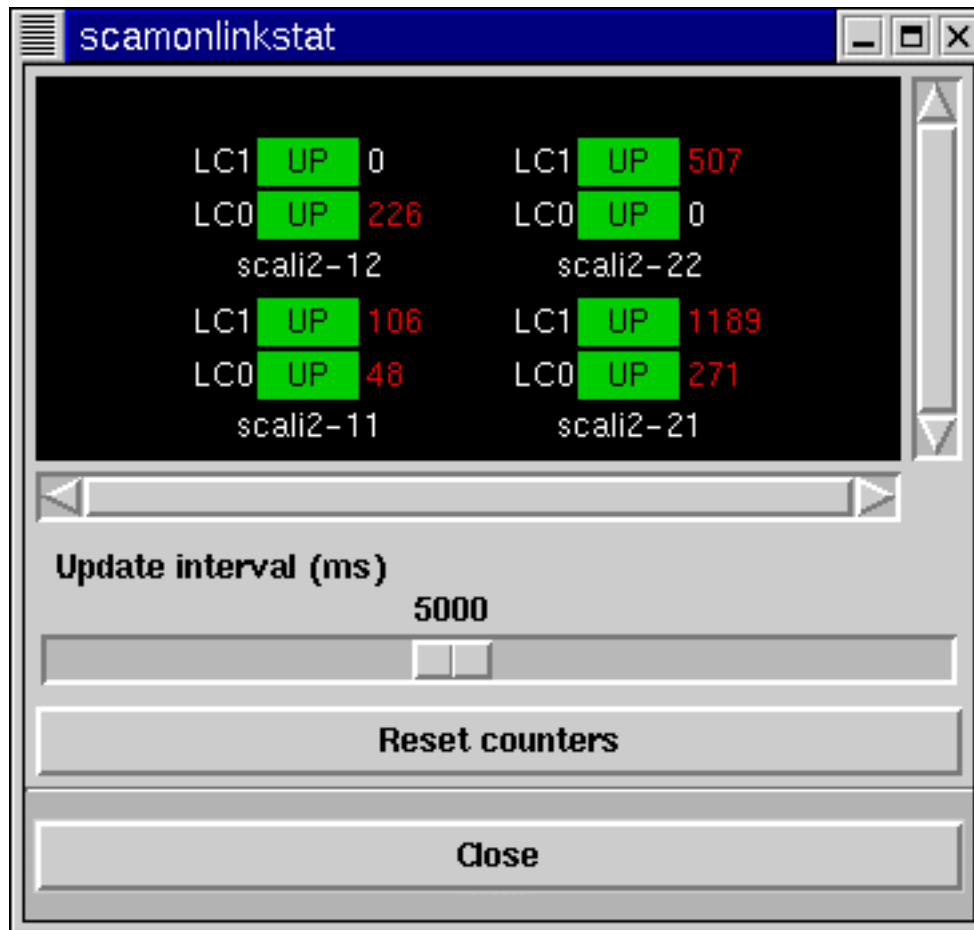




# Software Configuration Management



# System Monitoring



# Platform Attraction

