# Hardware- and Software-Based Collective Communication on the Quadrics Network

Fabrizio Petrini*, Salvador Coll*[†], Eitan Frachtenberg* and Adolfy Hoisie*

*CCS-3 Modeling, Algorithms, and Informatics

Los Alamos National Laboratory

[†]Technical University of Valencia - SPAIN

scoll@lanl.gov

# Outline

- **Introduction**

- **Quadrics network design overview**

  - Hardware

  - Communication/programming libraries

- **Collective communication on the QsNET**

- **Barrier synchronization**

- **Broadcast**

- **Performance analysis**

  - Experimental framework

  - Results

- **Conclusions**

# Introduction

- The efficient implementation of collective communication is a challenging design effort

- Very important to guarantee scalability of barrier synchronization, broadcast, gather, scatter, reduce, etc.

- Essential to implement system primitives to enhance fault-tolerance.

- Software or hardware support for multicast communication can improve the performance and resource utilization of a parallel computer

  - Software multicast: based on unicast messages, simple to implement, no network topology constraint, slower
  - Hardware multicast: require dedicated hardware, network dependent, faster

# Introduction

- Some of the most powerful systems in the world use the Quadrics interconnection network and the collective communication services analyzed in this job:

    - The Terascale Computing System (TCS) at the Pittsburgh Supercomputing Center – the second most powerful computer in the world
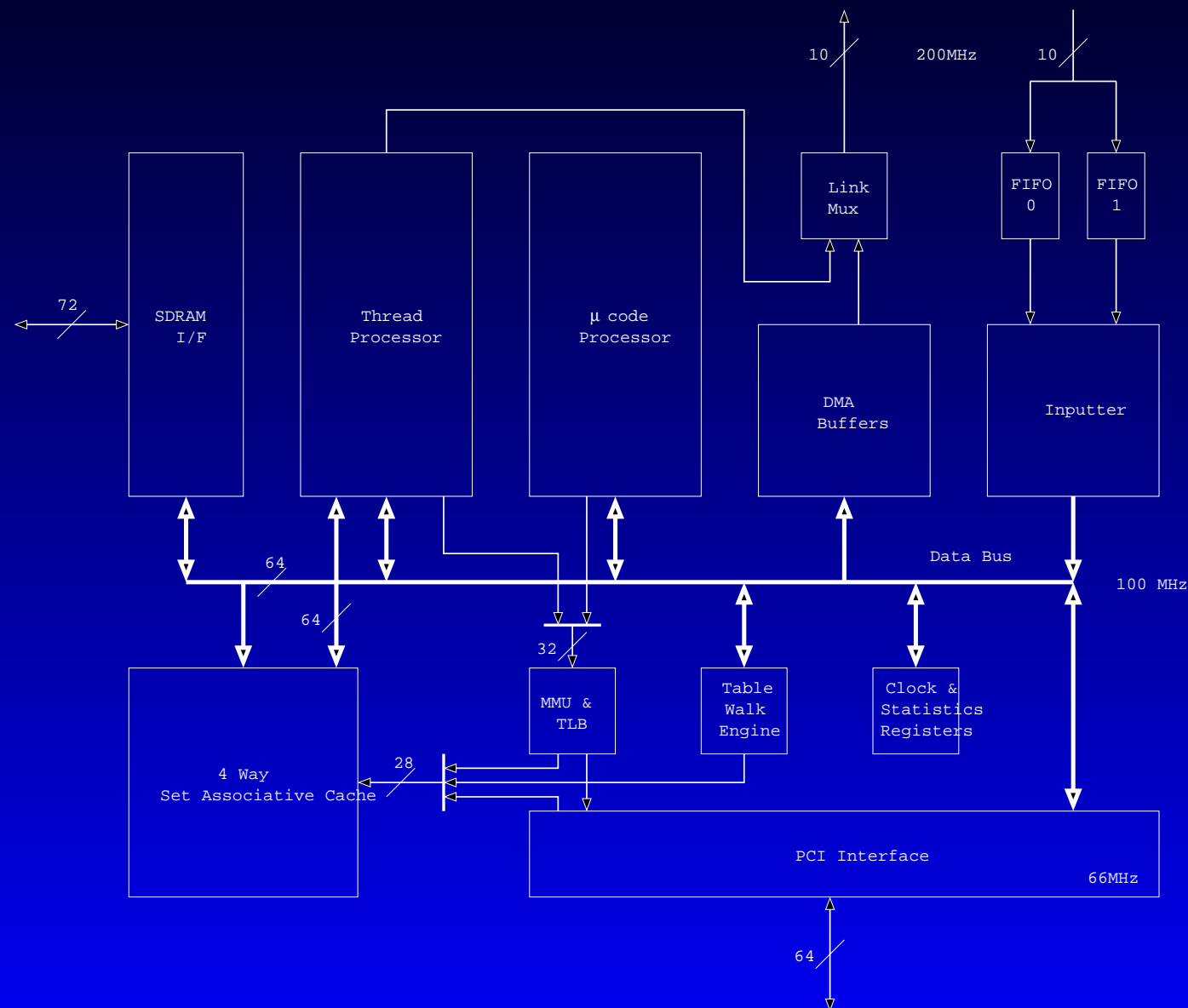
# Introduction

- Some of the most powerful systems in the world use the Quadrics interconnection network and the collective communication services analyzed in this job:

  - The Terascale Computing System (TCS) at the Pittsburgh Supercomputing Center – the second most powerful computer in the world



Barrier Test

# Introduction

- Some of the most powerful systems in the world use the Quadrics interconnection network and the collective communication services analyzed in this job:

    - The Terascale Computing System (TCS) at the Pittsburgh Supercomputing Center – the second most powerful computer in the world

    - ASCI Q machine, currently under development at Los Alamos National Laboratory (30 TeraOps, expected to be delivered by the end of 2002)

**Los Alamos**
NATIONAL LABORATORY

# Quadrics Network Design Overview

- QsNET provides an abstraction of distributed virtual shared memory

- Each process can map a portion of its address space into the global memory

- These address spaces constitutes the virtual shared memory

- This shared memory is fully integrated with the native operating system

- Based on two building blocks:

  - a network interface card called Elan
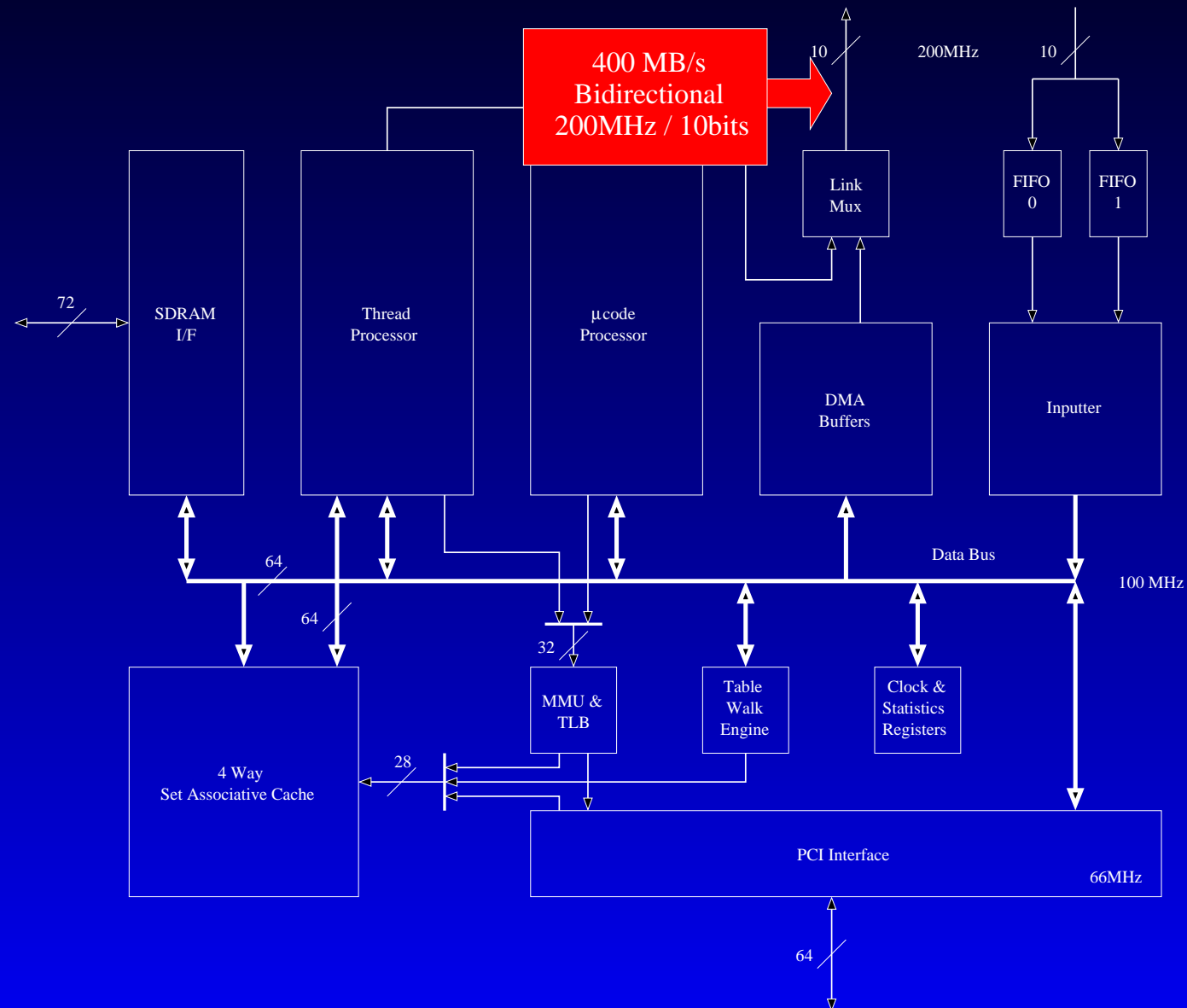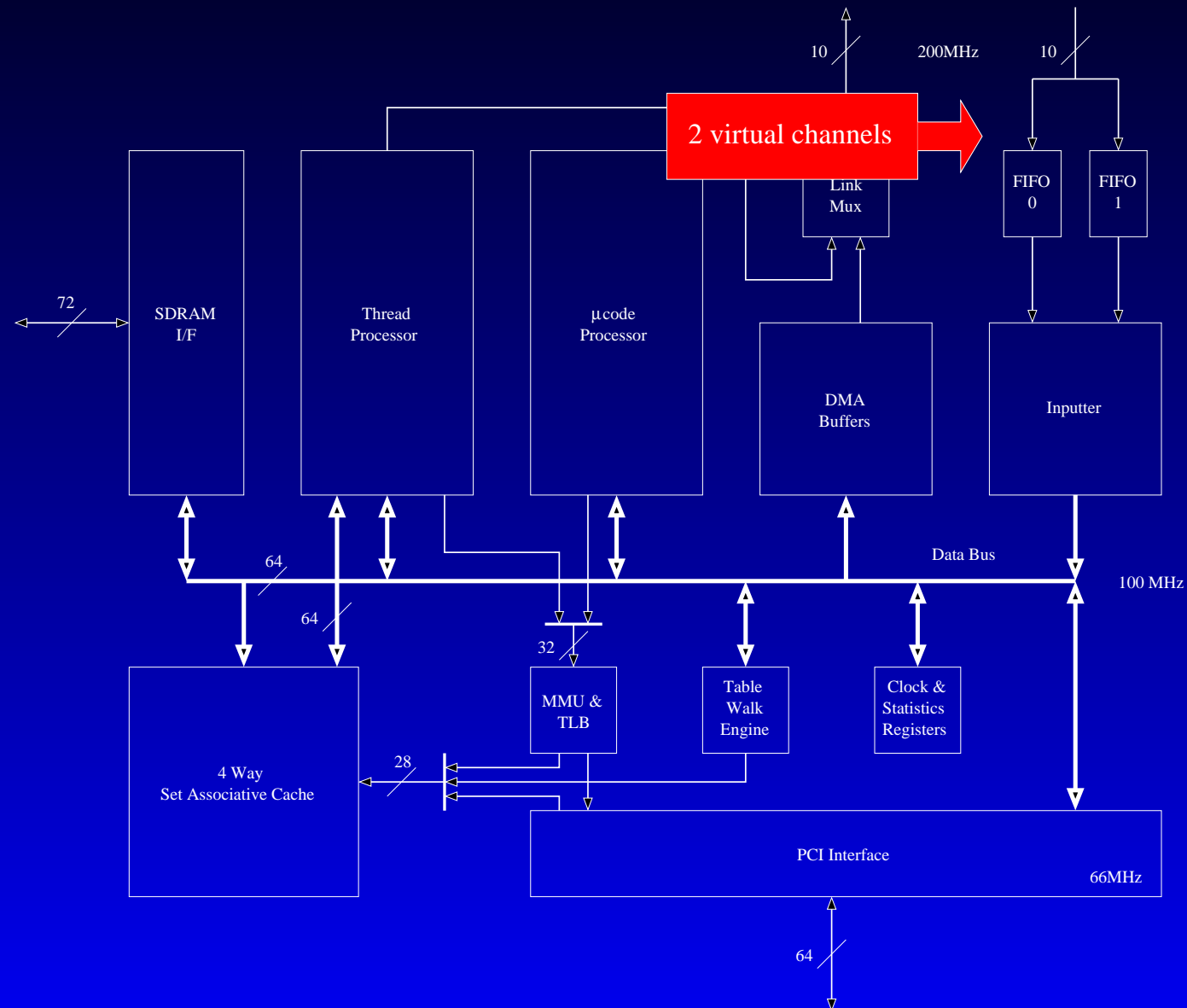
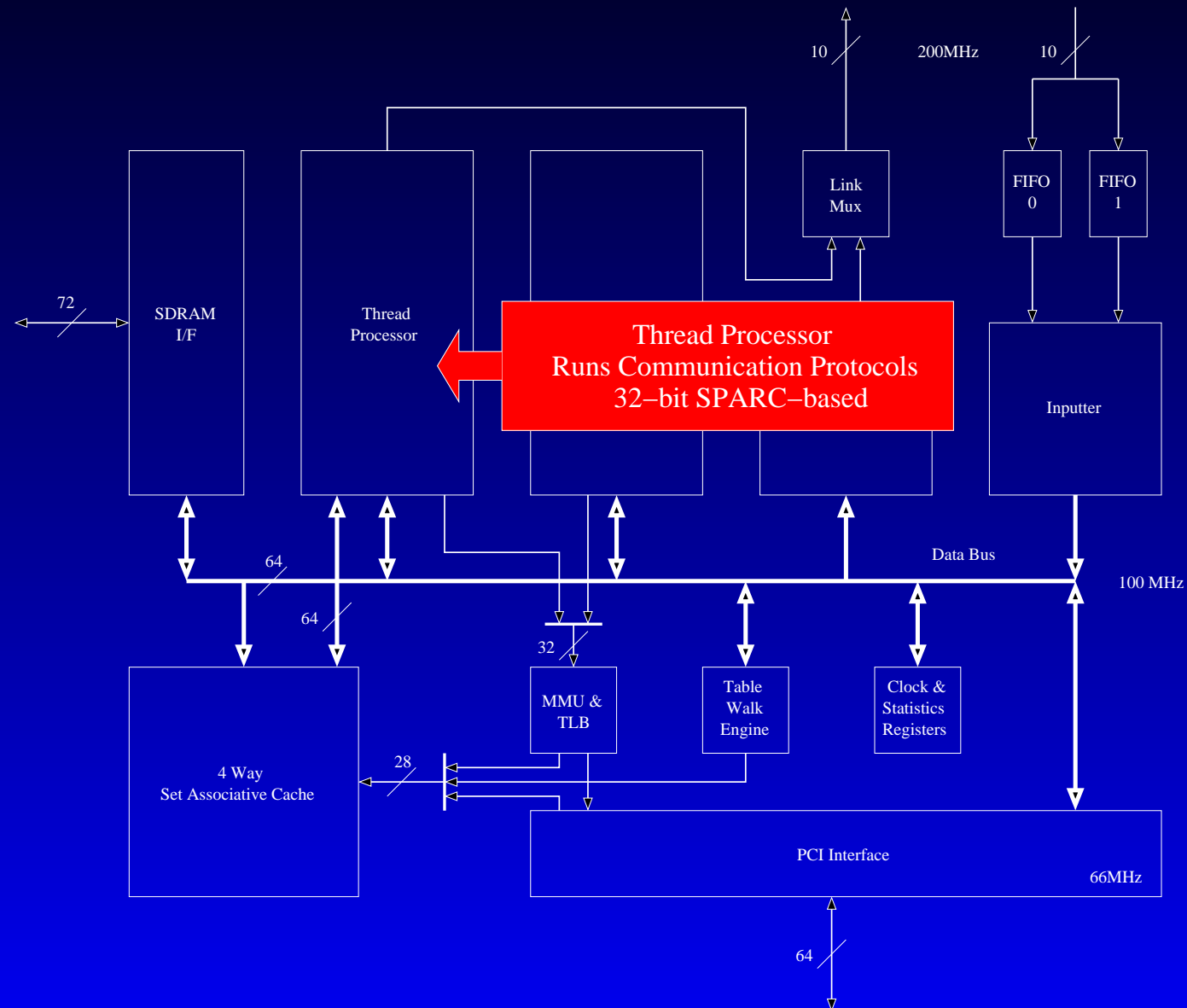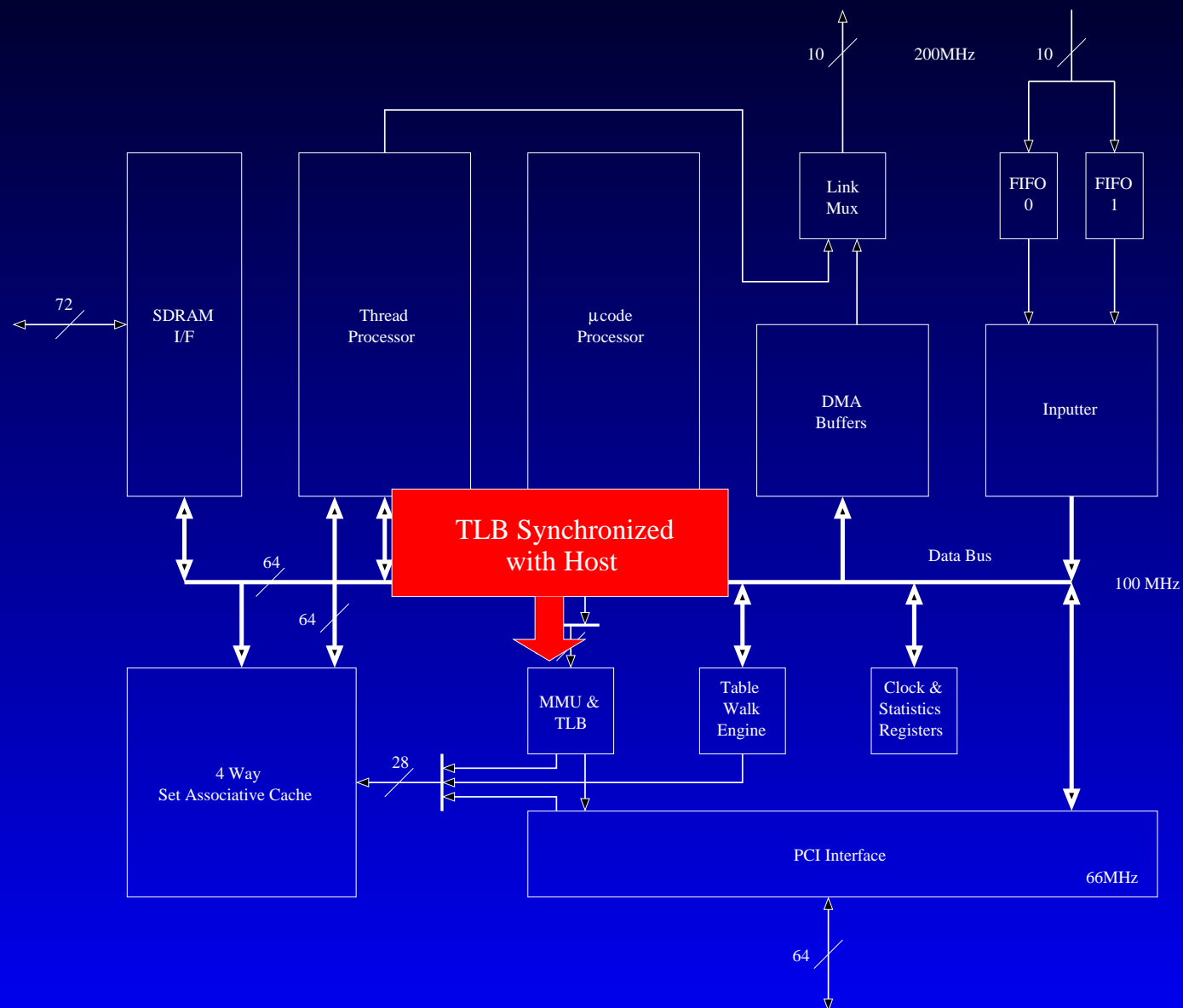  - a crossbar switch called Elite

Collectives

# Elan

# Elan

400 MB/s
Bidirectional
200MHz / 10bits

10    200MHz    10

Link
Mux

FIFO
0

FIFO
1

72

SDRAM
I/F

Thread
Processor

μ code
Processor

DMA
Buffers

Inputter

64    Data Bus

64    100 MHz

32

4 Way
Set Associative Cache

28

MMU &
TLB

Table
Walk
Engine

Clock &
Statistics
Registers

PCI Interface

66MHz

64

# Elan

# Elan

SDRAM
I/F

72

Thread
Processor

Thread Processor
Runs Communication Protocols
32–bit SPARC–based

Link
Mux

10          200MHz          10

FIFO
0

FIFO
1

Inputter

64

Data Bus

100 MHz

64

32

4 Way
Set Associative Cache

28

MMU &
TLB

Table
Walk
Engine

Clock &
Statistics
Registers

PCI Interface

66MHz

64

Los Alamos
NATIONAL LABORATORY

# Elan



10    200MHz    10

SDRAM I/F

Thread Processor

μcode Processor

Link Mux

FIFO 0

FIFO 1

72

DMA Buffers

Inputter

64

TLB Synchronized with Host

Data Bus

100 MHz

64

MMU & TLB

Table Walk Engine

Clock & Statistics Registers

4 Way Set Associative Cache

28

PCI Interface

66MHz

64

Los Alamos
NATIONAL LABORATORY

# Elan

# Elite

- 8 bidirectional links with 2 virtual channels in each direction

- An internal 16x8 full crossbar switch

- 400 MB/s on each link direction

- Packet error detection and recovery, with routing and data transactions CRC protected

- 2 priority levels plus an aging mechanism

- Adaptive routing

- Hardware support for broadcast

# Network Topology: Quaternary Fat-Tree
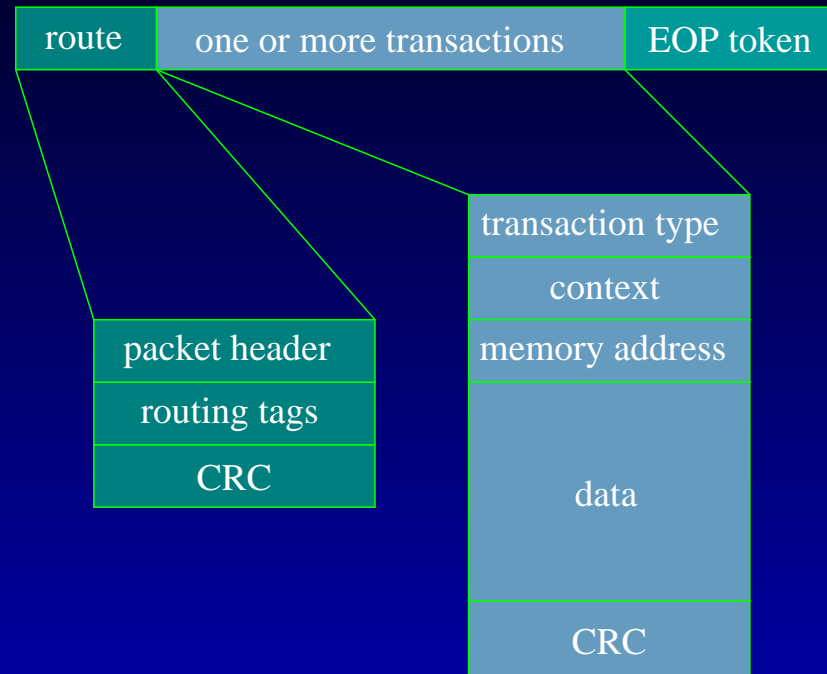
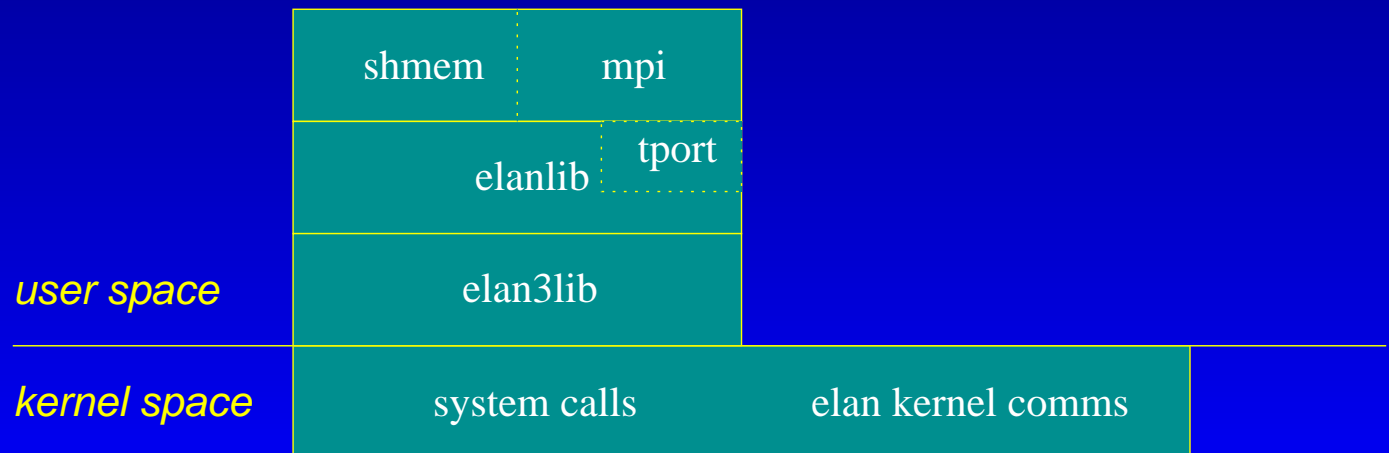# Network Topology: Quaternary Fat-Tree
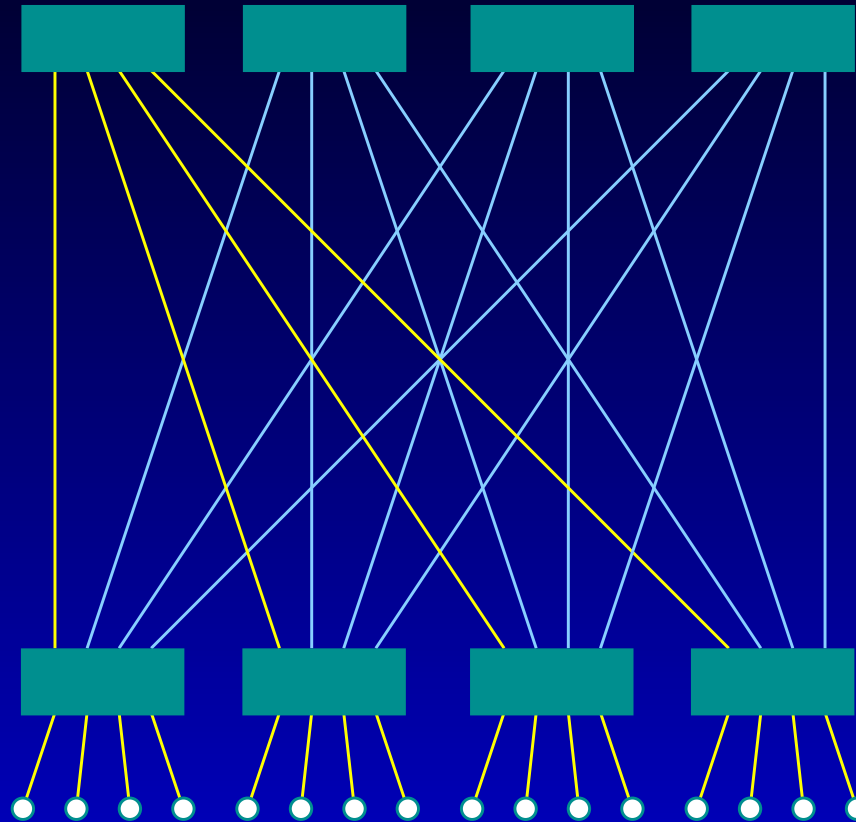
# Network Topology: Quaternary Fat-Tree

# Packet Format

| route | one or more transactions | EOP token |
|---|---|---|

| packet header |
|---|
| routing tags |
| CRC |

| transaction type |
|---|
| context |
| memory address |
| data |
| CRC |

- 320 bytes data payload (5 transactions with 64 bytes each)
- 74-80 bytes overhead

# Programming Libraries

- Elan3lib
  - event notification
  - memory mapping and allocation
  - remote DMA

- Elanlib and Tports
  - collective communication
  - tagged message passing

- MPI, shmem

User  Applications

| | |
|---|---|
| shmem | mpi |
| elanlib | tport |
| elan3lib | |

*user space*

---

*kernel space*

| system calls | elan kernel comms |

# Collective communication on the QsNET



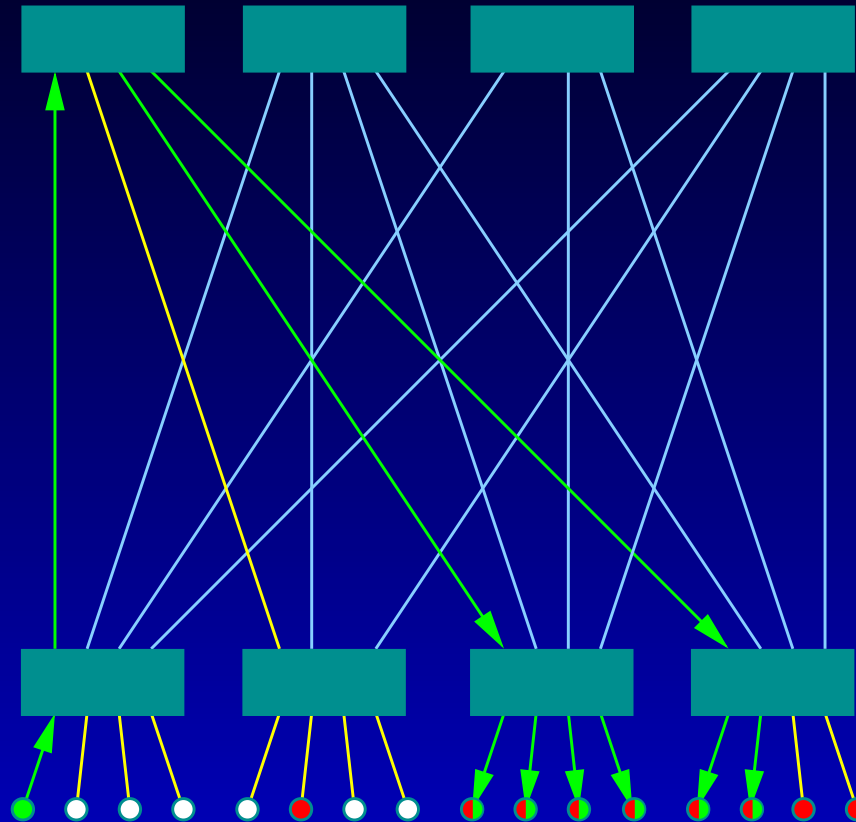Broadcast tree for a 16-node network

# Collective communication on the QsNET
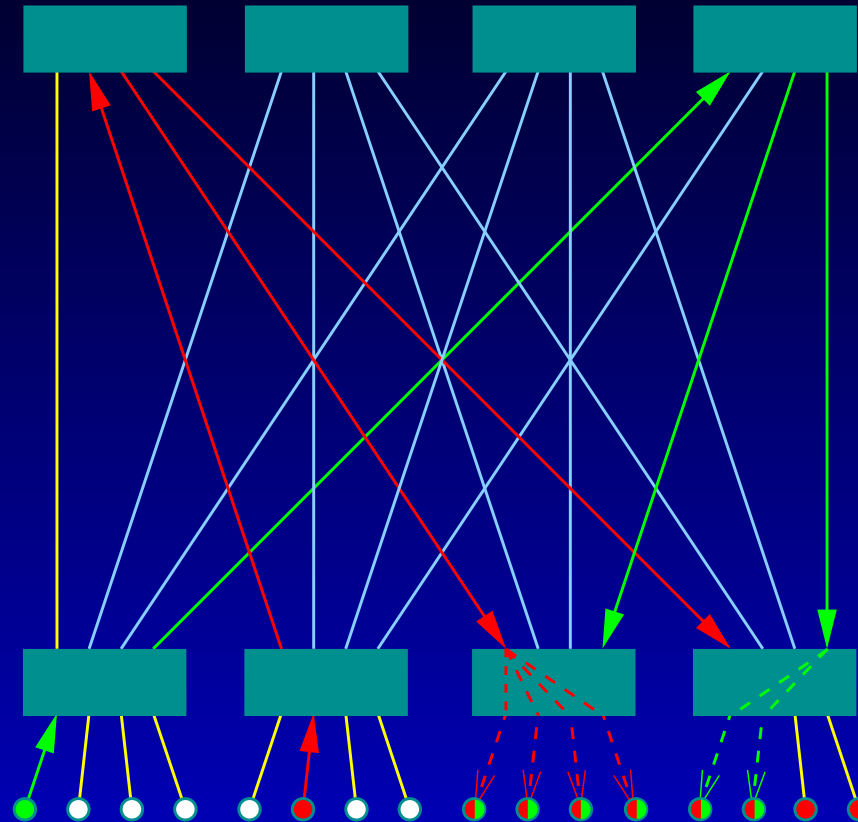
# Collective communication on the QsNET



Serialization through the root switch to avoid deadlocks

# Collective communication on the QsNET

# Collective communication on the QsNET
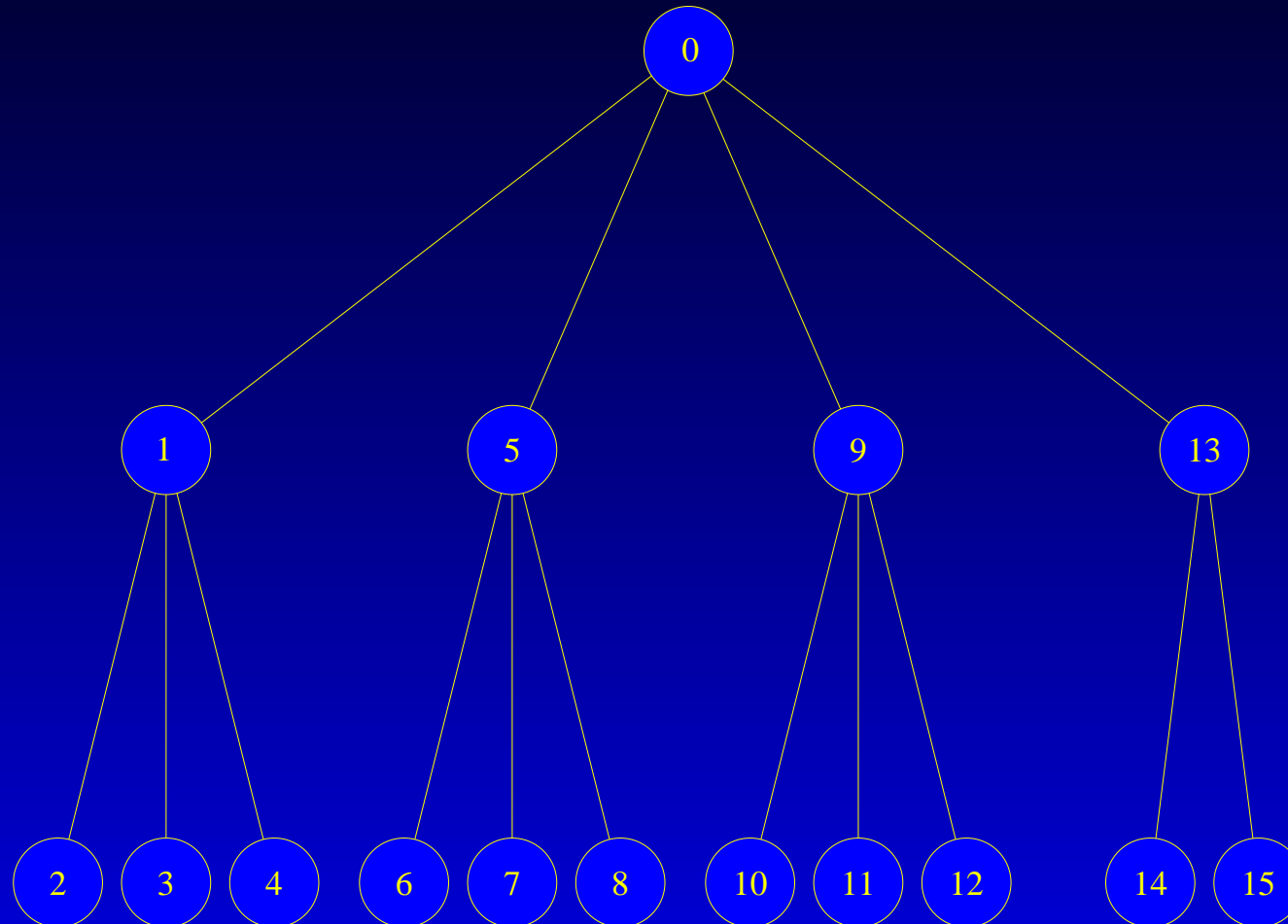
Deadlocked situation

# Barrier Synchronization

QsNET implements two synchronization primitives:

- Software-based: it uses a balanced tree and point-to-point messages

    - `elan_gsync()`

- Hardware-based: it uses the hardware multicast support

    - `elan_hgsync()`: busy-wait
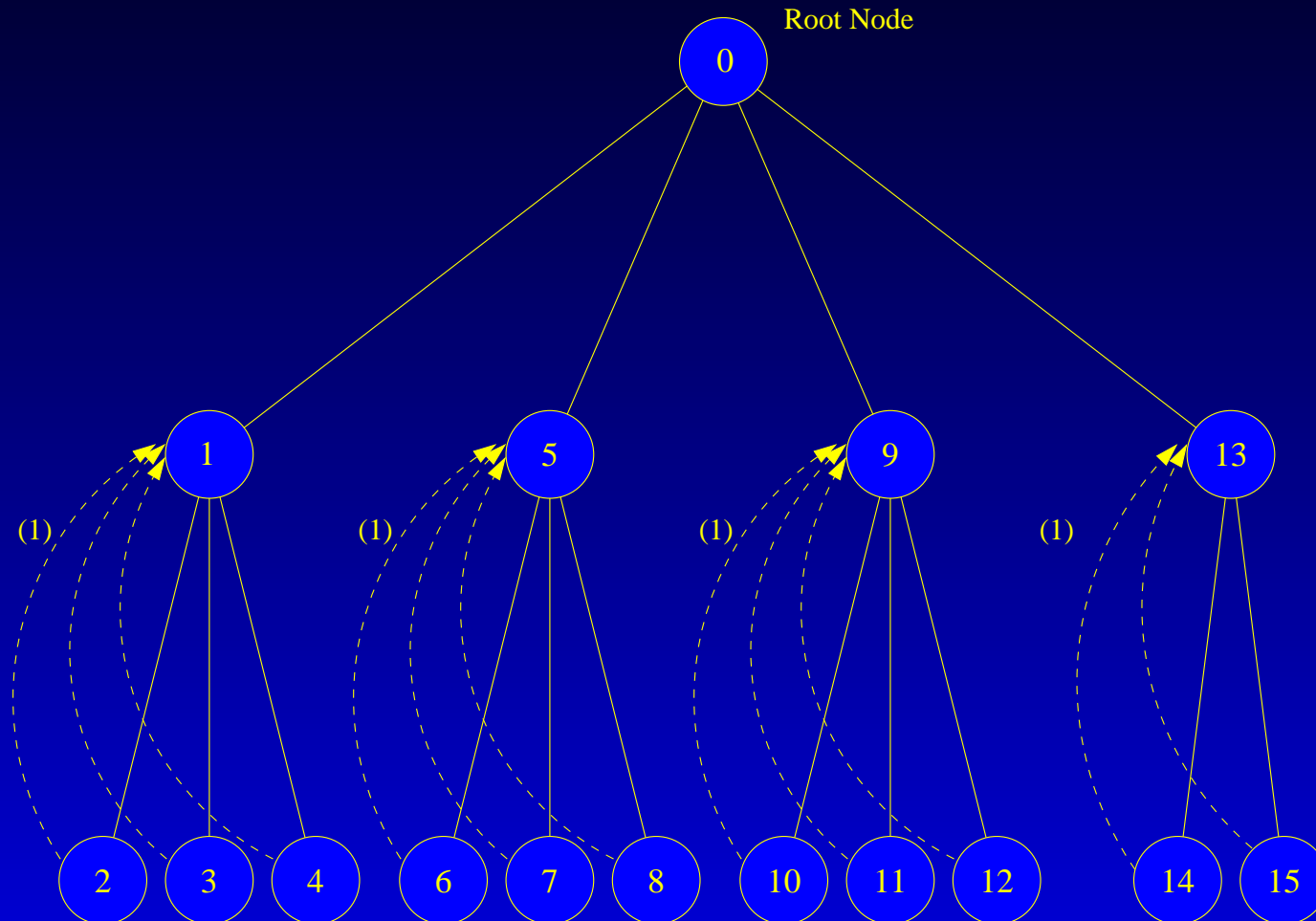
    - `elan_hgsyncevent()`: event-based

# Software-Based Barrier

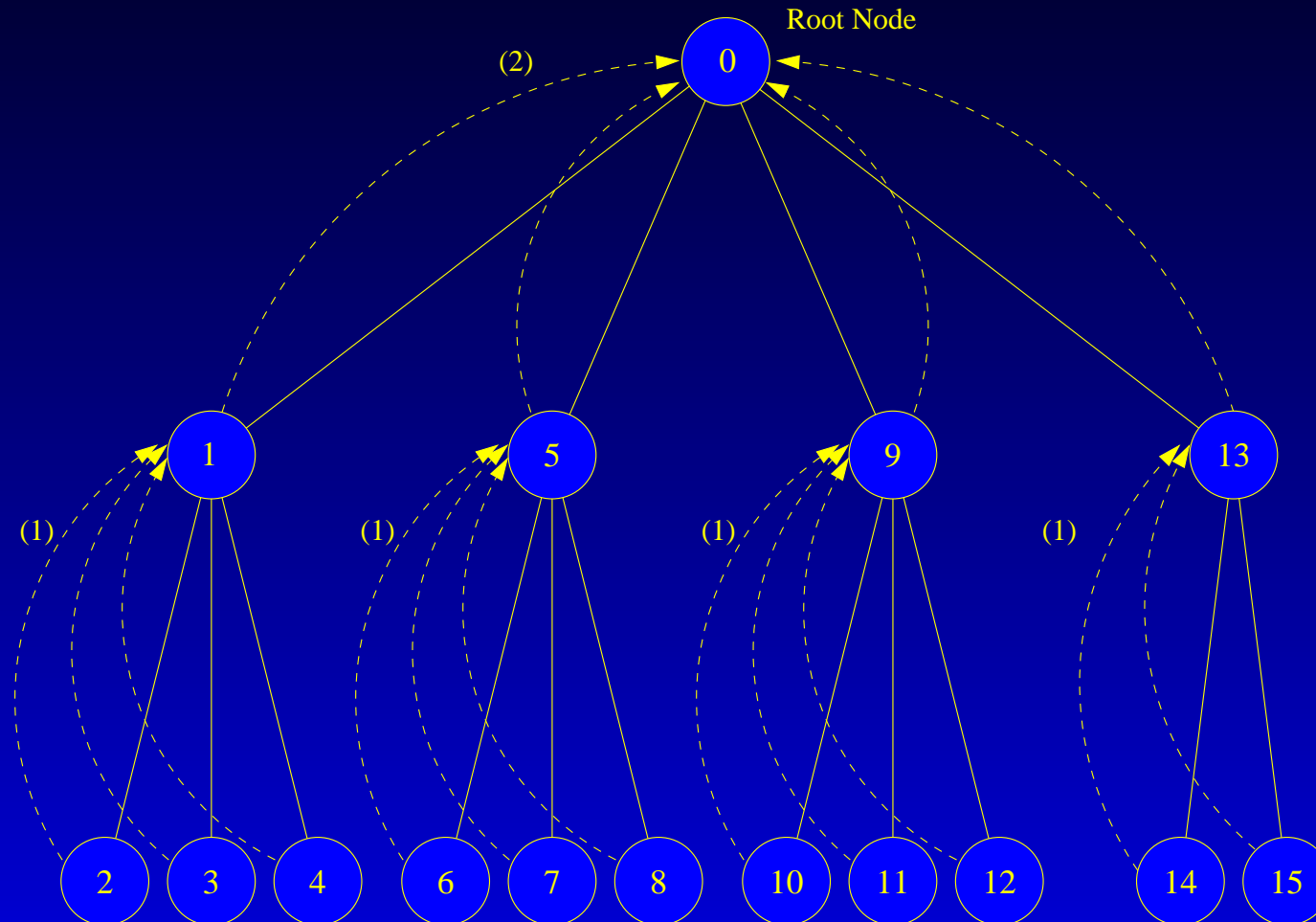Each process waits for 'ready' signals from its children

# Software-Based Barrier

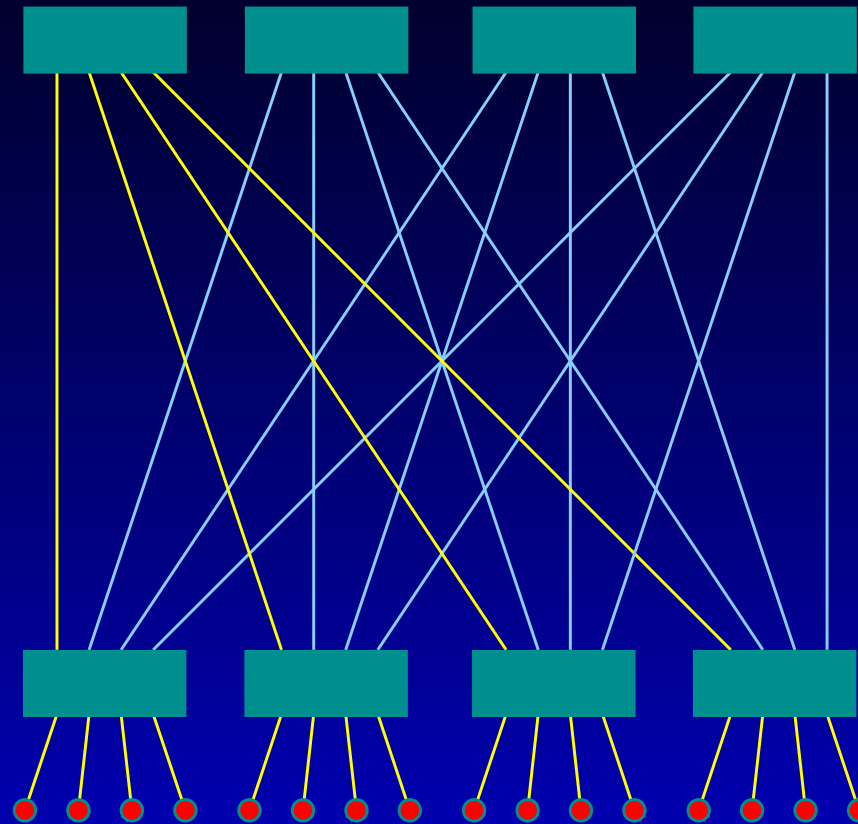Each process waits for 'ready' signals from its children (1) ...

# Software-Based Barrier

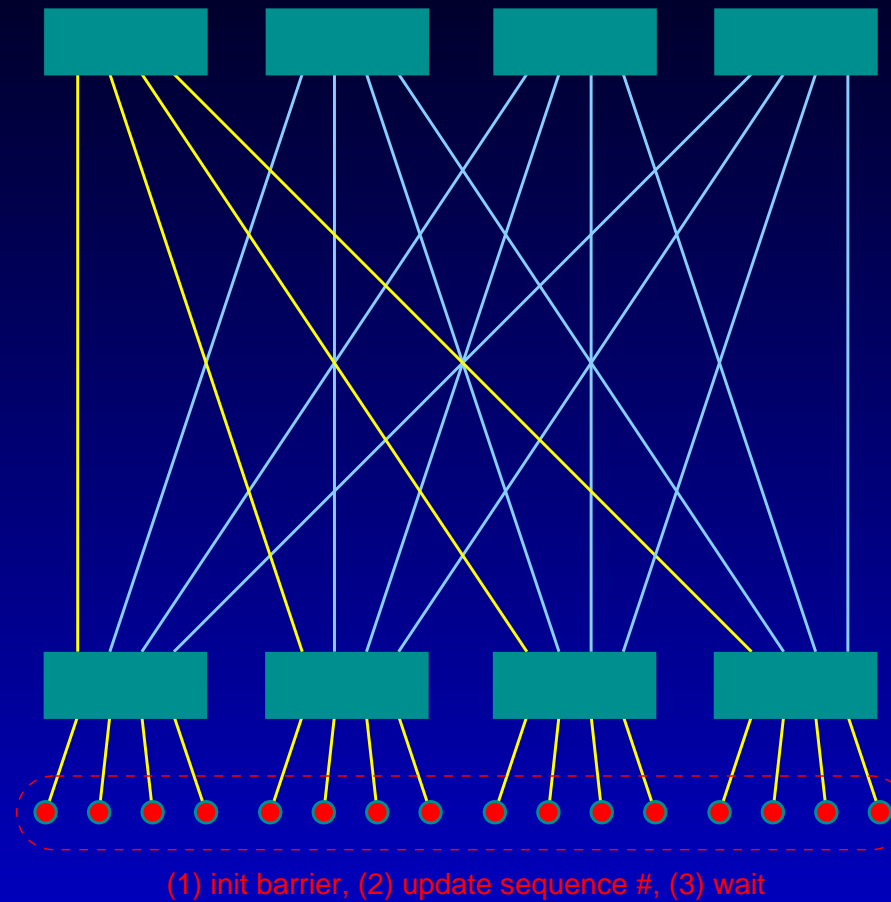... and sends its own signal up to the parent process (2)
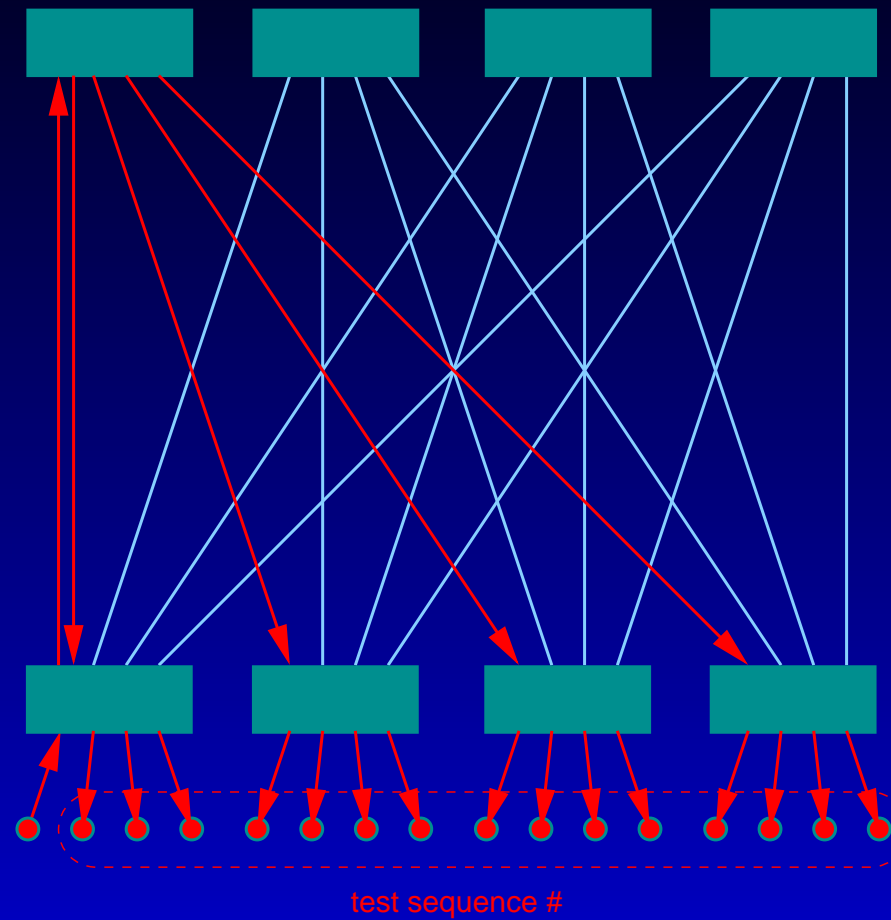
# Hardware-Based Barrier



Example for 16 nodes
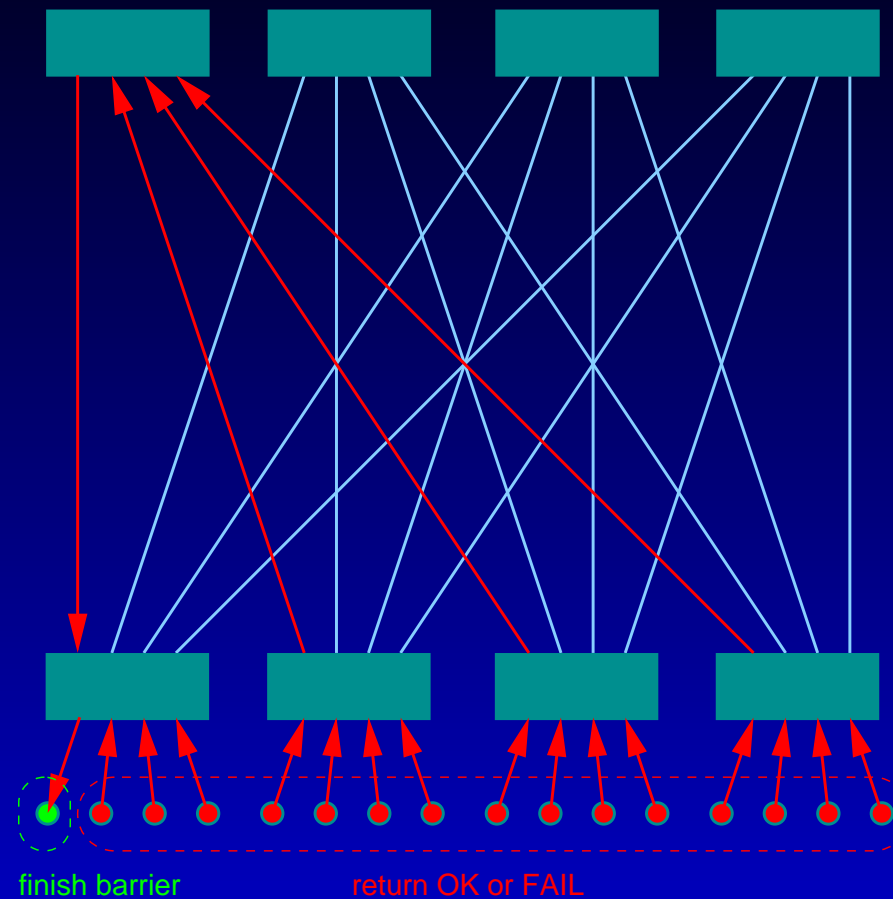
# Hardware-Based Barrier



(1) init barrier, (2) update sequence #, (3) wait

Init barrier

# Hardware-Based Barrier



test sequence #

Multicast transaction

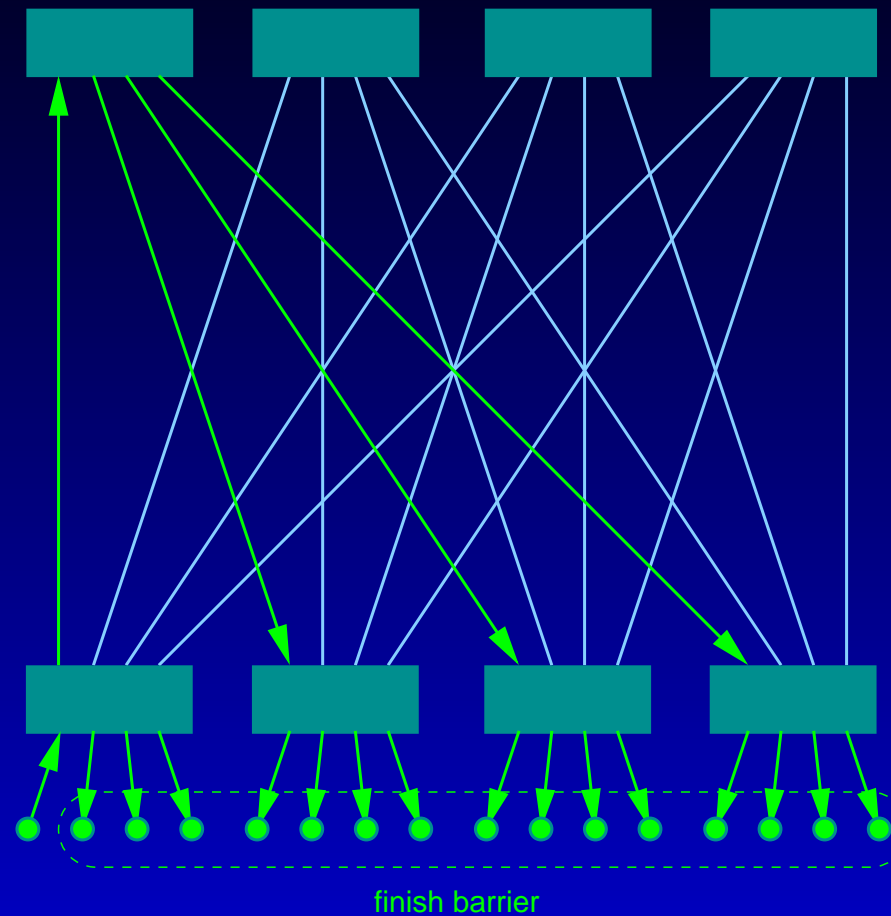# Hardware-Based Barrier



finish barrier          return OK or FAIL

Acknowledgment

# Hardware-Based Barrier



finish barrier

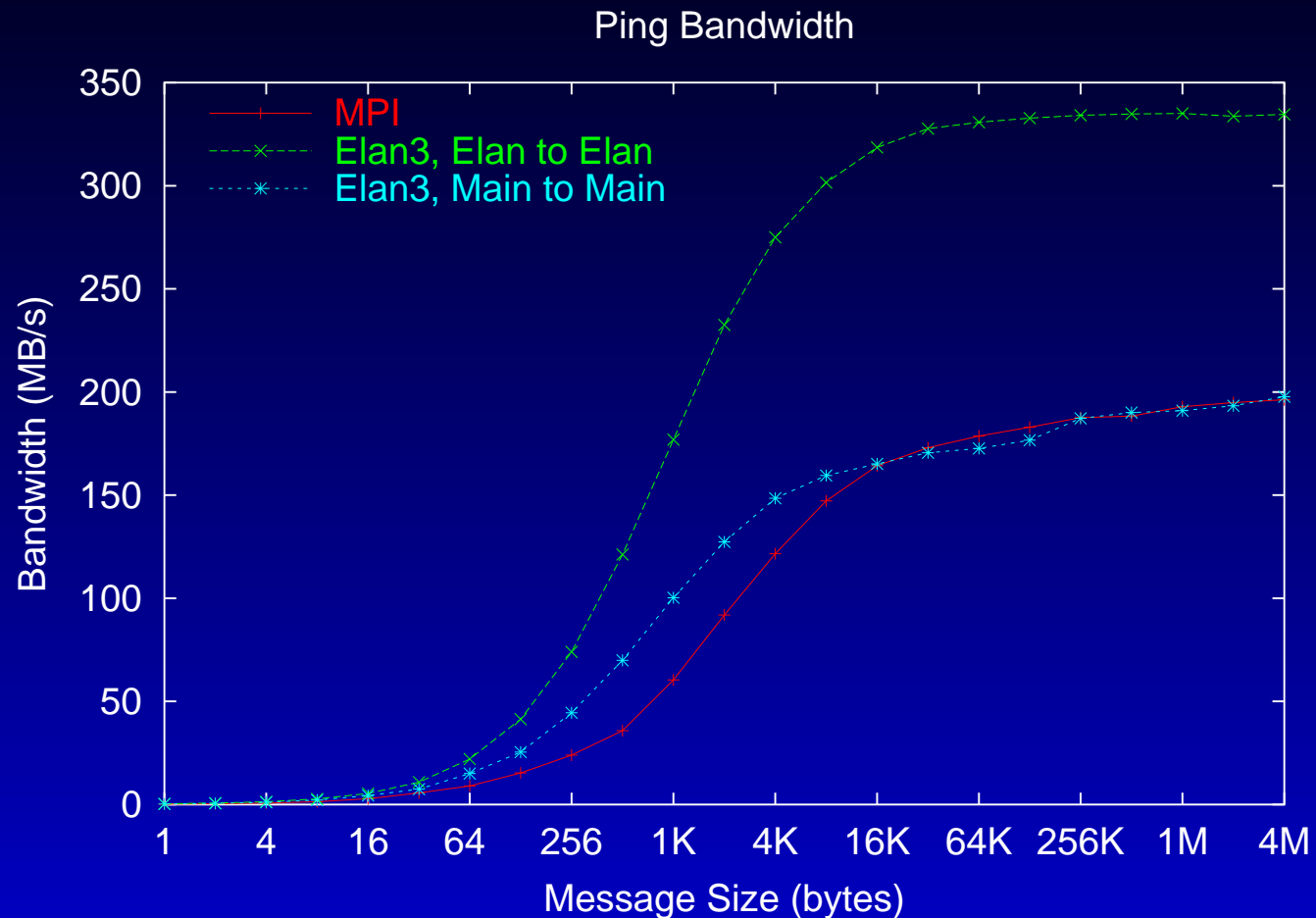Final 'EOP' (End-Of-Packet) token

# Broadcast

QsNET implements two broadcast primitives:

- Software-based: it uses a balanced tree and point-to-point messages

    - `elan_bcast()`

- Hardware-based: it uses the hardware multicast support

    - `elan_hbcast()`

- Both implementations perform an initial barrier to guarantee resources allocation
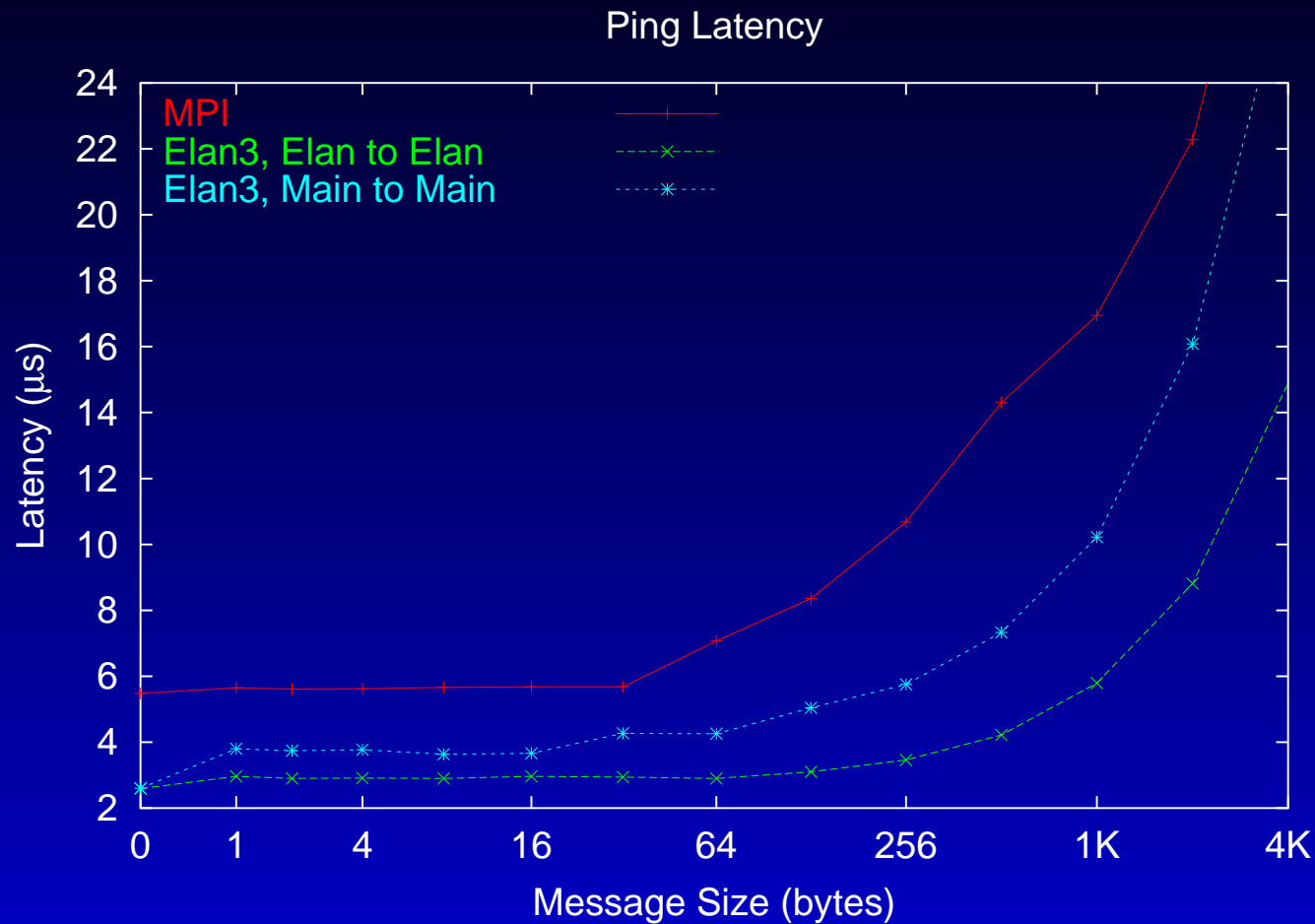
# Performance Analysis

- The experimental results are obtained on a 64-node cluster of Compaq AlphaServer ES40s running Tru64 Unix.

- Each Alpahserver is attached to a quaternary fat-tree of dimension three through a 64 bit, 33 MHz PCI bus using the Elan3 card.

- In order to expose the real network performance, we place the communication buffers in Elan memory.

- We present:

  - unidirectional ping results, as a reference, and
  - barrier and broadcast results, analyzing the effect of additional background traffic

Los Alamos
NATIONAL LABORATORY
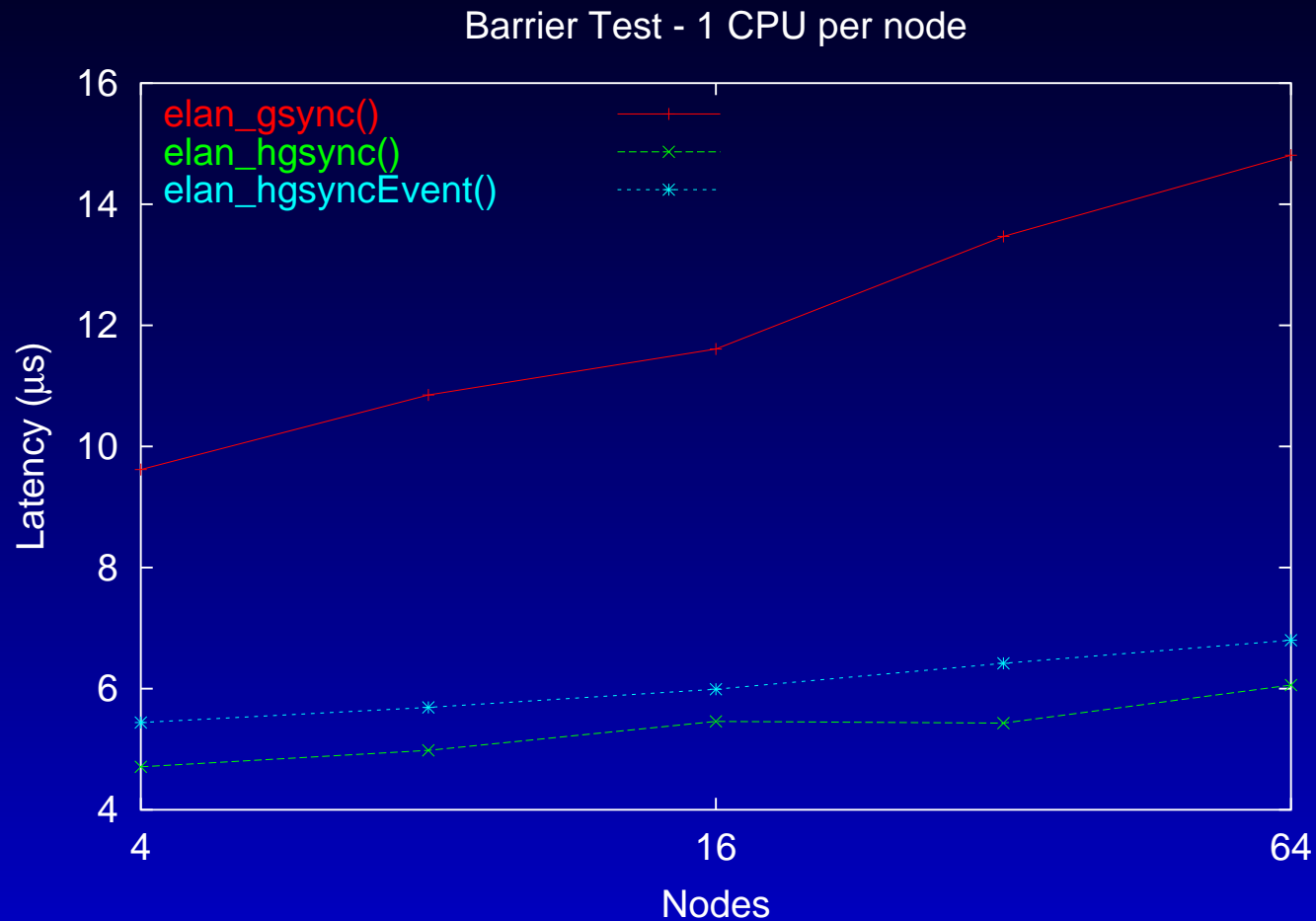
# Unidirectional Ping



Ping Bandwidth

- Peak data bandwidth (Elan to Elan) of 335 MB/s $\simeq$ 396 MB/s (99% of nominal bandwidth)
- Main to main asymptotic bandwidth of 200 MB/s

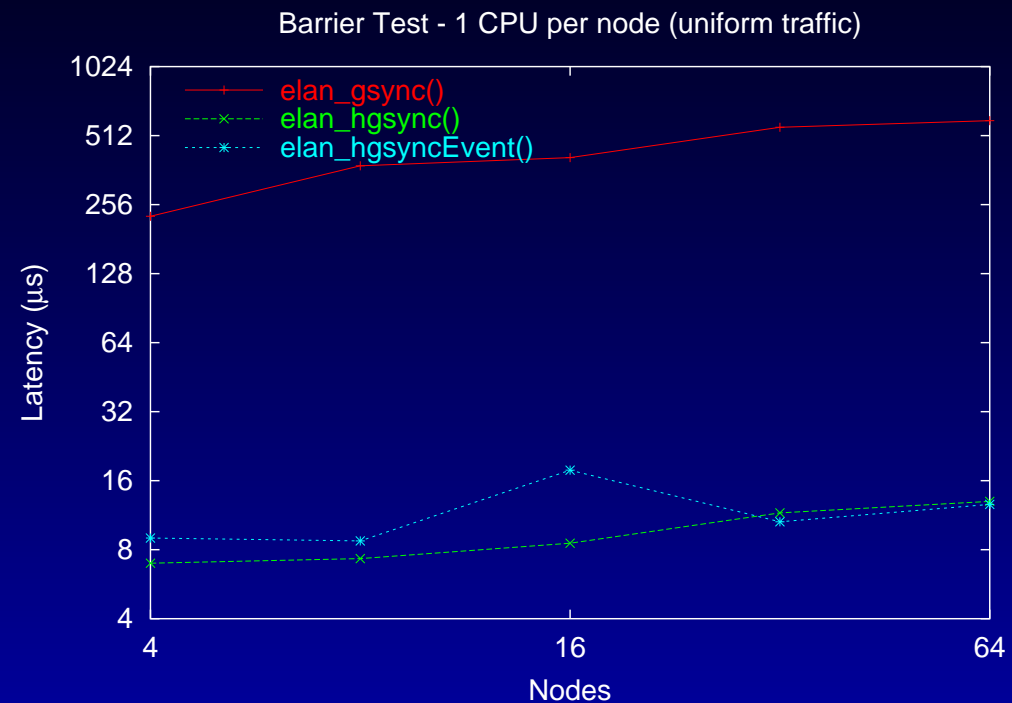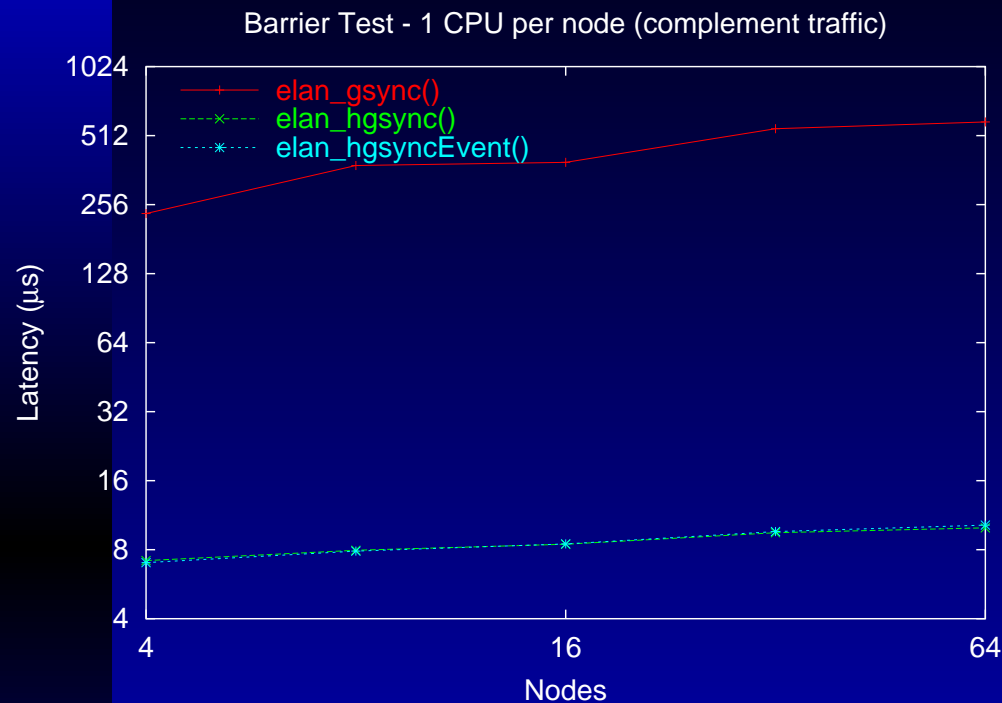# Unidirectional Ping



Ping Latency

- Latency of 2.4 $\mu$s up to 64-byte messages (Elan to Elan memory)
- Higher MPI latency due to message tag matching

# Barrier Synchronization



Barrier Test - 1 CPU per node

- **Good hardware barrier scalability**

# Barrier Synchronization with Background Traffic



Barrier Test - 1 CPU per node (complement traffic)

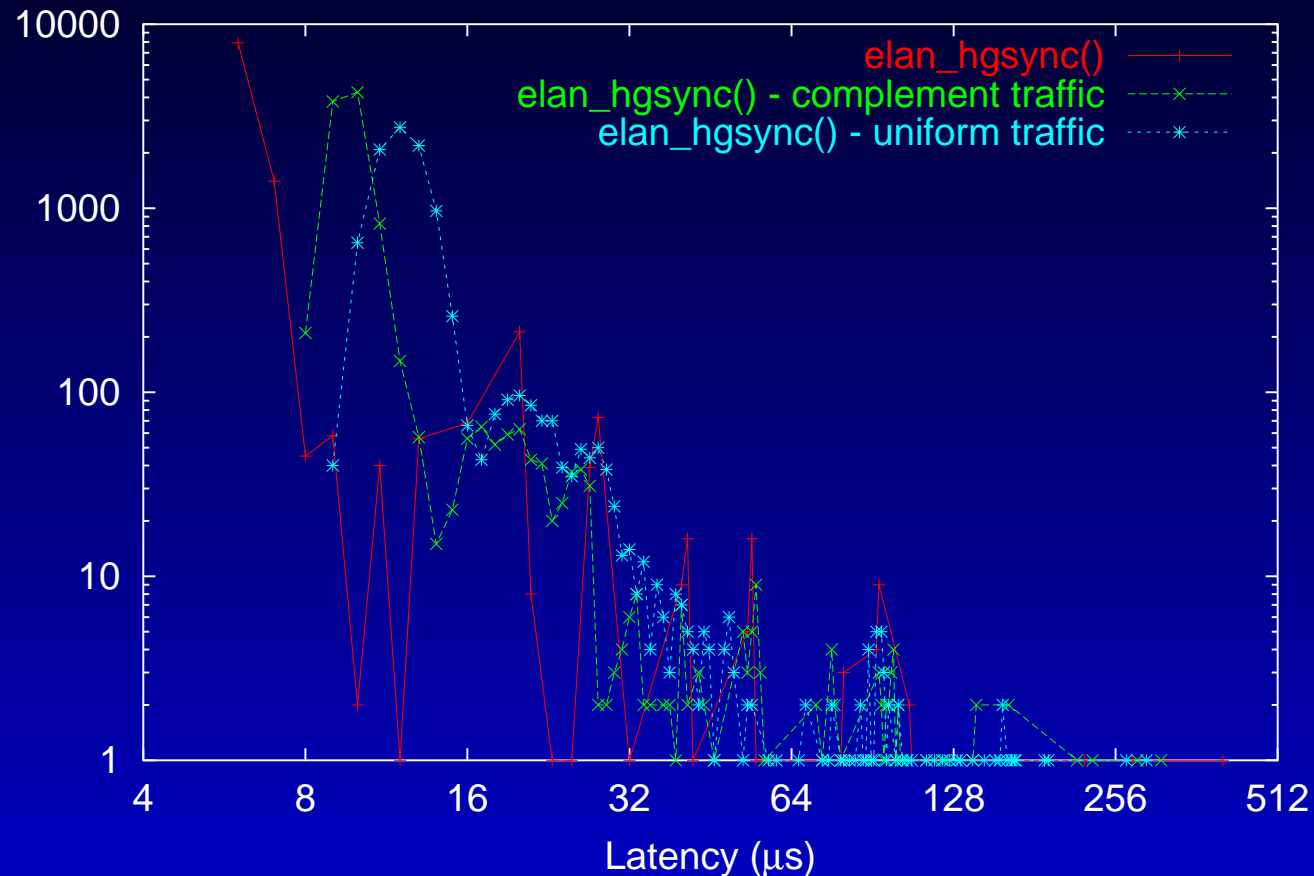Barrier Test - 1 CPU per node (uniform traffic)

- Software barrier significantly affected (the slowdown is 40 in the worst case)
- Little impact on the hardware barriers, whose average latency is only doubled

# Hardware Barrier with Background Traffic

Barrier Test - 64 nodes, 1 CPU per node (latency distribution)



- 94% of the operations take less than $9\mu$s with no bakground traffic
- 93% of the tests take less than $20\mu$s with uniform traffic

# Software Barrier with Background Traffic
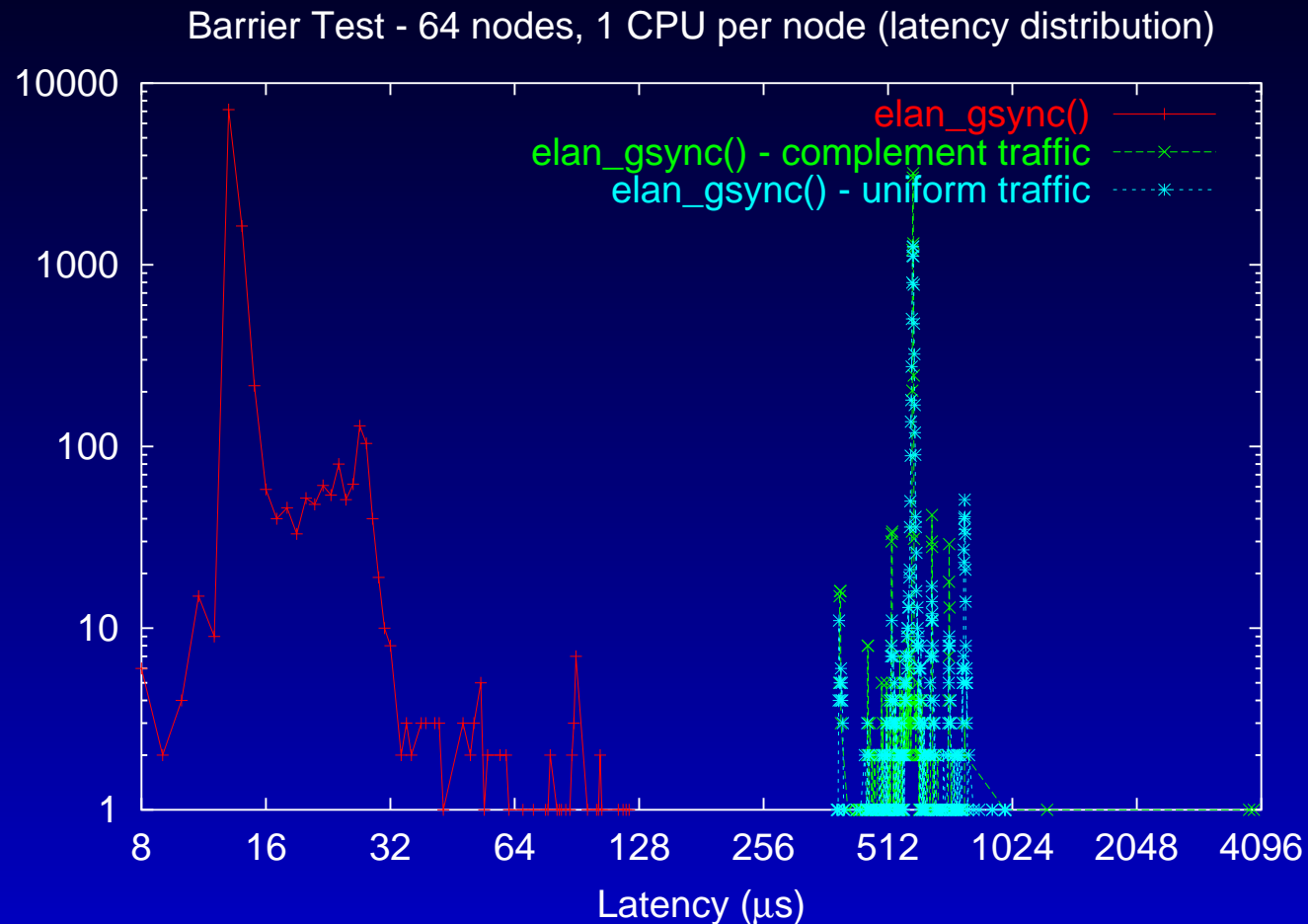


Barrier Test - 64 nodes, 1 CPU per node (latency distribution)

- 99% of the barriers take less than $30\mu$s with no bakground traffic

- 93% of the synchronizations complete with less than $605\mu$s with uniform traffic

# Broadcast Bandwidth

**Broadcast Test - 64 Nodes, 1 CPU per node**



- Asymptotic bandwidth of 288MB/s when using Elan memory for both implementations

# Broadcast Latency

Broadcast Test - 64 Nodes, 1 CPU per node



- Hardware latency with Elan buffers below 13$\mu$s for messages up to 256 bytes

- Software latencies are 3.5$\mu$s higher than hardware latencies

# Broadcast Scalability



Broadcast Test - 1 CPU per node (256k bytes)

Broadcast Test - 1 CPU per node (256k bytes)

- No significant effect when using buffers in main memory

- With buffers in Elan memory performance depends on the number of switch layers traversed

# Broadcast with Background Traffic



Broadcast Test - 64 Nodes, 1 CPU per node (complement traffic)

Legend:
- elan_bcast() - global - main
- elan_bcast() - global - elan
- elan_hbcast() - global - main
- elan_hbcast() - global - elan

Left plot: Bandwidth (MB/s) vs Message Size (bytes)
Right plot: Latency (μs) vs Message Size (bytes)

Los Alamos
NATIONAL LABORATORY

# Broadcast with Background Traffic

Broadcast Test - 64 Nodes, 1 CPU per node (uniform traffic)

Broadcast Test - 64 Nodes, 1 CPU per node (uniform traffic)



Legend (both plots):
- elan_bcast() - global - main
- elan_bcast() - global - elan
- elan_hbcast() - global - main
- elan_hbcast() - global - elan

Left plot: Bandwidth (MB/s) vs Message Size (bytes)

Right plot: Latency (µs) vs Message Size (bytes)

- Latency differences between hw and sw implementations increase

- Better performance with buffers in main memory (due to the background traffic application)

# Broadcast with Background Traffic



Broadcast Test - 1 CPU per node (256k bytes - uniform traffic)

- elan_bcast() - global - main
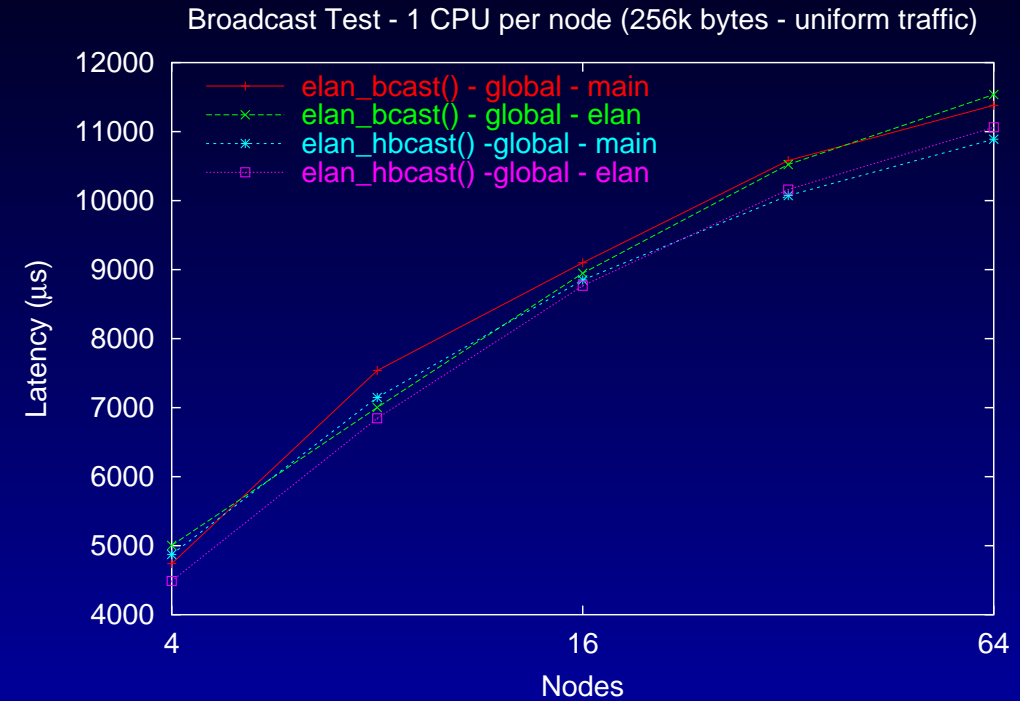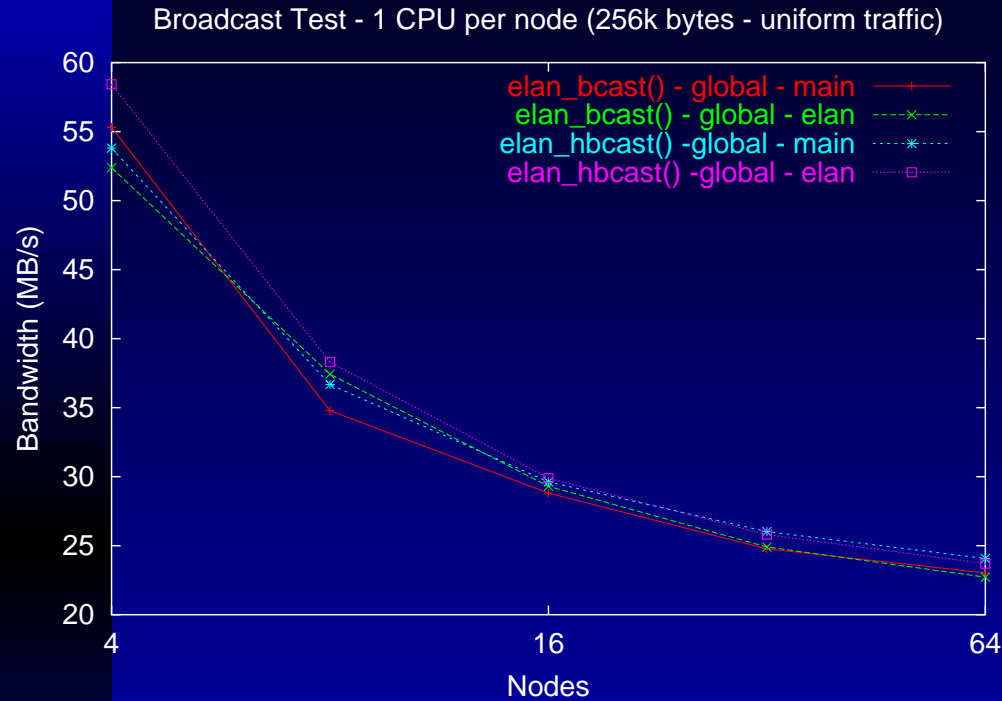- elan_bcast() - global - elan
- elan_hbcast() -global - main
- elan_hbcast() -global - elan

Bandwidth (MB/s) vs Nodes

Broadcast Test - 1 CPU per node (256k bytes - uniform traffic)

- elan_bcast() - global - main
- elan_bcast() - global - elan
- elan_hbcast() -global - main
- elan_hbcast() -global - elan

Latency (µs) vs Nodes

- Significant performance degradation for all the alternatives

# Conclusions

- Hardware-based synchronization takes as little as $6\mu$s on a 64-node Alphaserver cluster, with very good scalability.

# Conclusions

- Hardware-based synchronization takes as little as $6\mu$s on a 64-node Alphaserver cluster, with very good scalability.

- Good latency and scalability are achieved with the software-based synchronization too, which takes about $15\mu$s.

# Conclusions

- Hardware-based synchronization takes as little as $6\mu$s on a 64-node Alphaserver cluster, with very good scalability.

- Good latency and scalability are achieved with the software-based synchronization too, which takes about $15\mu$s.

- The hardware barrier is almost insensitive to background traffic, with 93% of the synchronizations completed in less than $20\mu$s.

# Conclusions

- Hardware-based synchronization takes as little as $6\mu$s on a 64-node Alphaserver cluster, with very good scalability.

- Good latency and scalability are achieved with the software-based synchronization too, which takes about $15\mu$s.

- The hardware barrier is almost insensitive to background traffic, with 93% of the synchronizations completed in less than $20\mu$s.

- With the broadcast, both implementations can deliver a sustained bandwidth of 288 MB/s Elan memory to Elan memory and 200 MB/s main memory to main memory.

Los Alamos
NATIONAL LABORATORY