# INFINIBAND ℠

TRADE ASSOCIATION

# Key Measures of InfiniBand™ Performance in the Data Center

Driving Metrics for
End User Benefits

# Benchmark Subgroup

## Benchmark Subgroup Charter

The InfiniBand™ Benchmarking Subgroup has been chartered by the InfiniBand Trade Association Marketing Work Group to propose and support InfiniBand Trade Association programs targeted at benchmarking and reporting InfiniBand Architecture related data. Subgroup deliverables will include: facilitation, collection, publication and promotion of InfiniBand Trade Association Member company InfiniBand benchmarks and related data.

## Member Companies

Agilent Technologies, Banderacom, Crossroads Systems, InfiniCon Systems, InfiniSwitch, Lane15 Software, Mellanox Technologies, Vieo

**INFINIBAND™**
TRADE ASSOCIATION

# Abstract

This presentation highlights the performance benefits of InfiniBand Architecture deployment in the data center. Key application performance drivers are introduced and InfiniBand Architecture's role in optimizing these drivers discussed. Actual InfiniBand technology performance measurements are presented by individual InfiniBand Trade Association member companies. The session concludes with a general Q&A with the presenters.

# Agenda

Introduction to InfiniBand Performance

Chris Petty (Banderacom)

Member Company Performance Presentations

Chris Petty (Banderacom)

Chris Eddington (Mellanox)

Brad Benton (Lane15)

Thomas Dippon (Agilent)

Lawrence Didsbury (Auspex)

General Q&A

INFINIBAND™
TRADE ASSOCIATION

# InfiniBand Performance Measurement

*Driving Metrics for End User Benefits*

# Agenda

## Application and system level performance issues

- Where IB will be deployed in the data center
- Performance attributes

## Performance impact of IB in the data center

- Theoretical impact
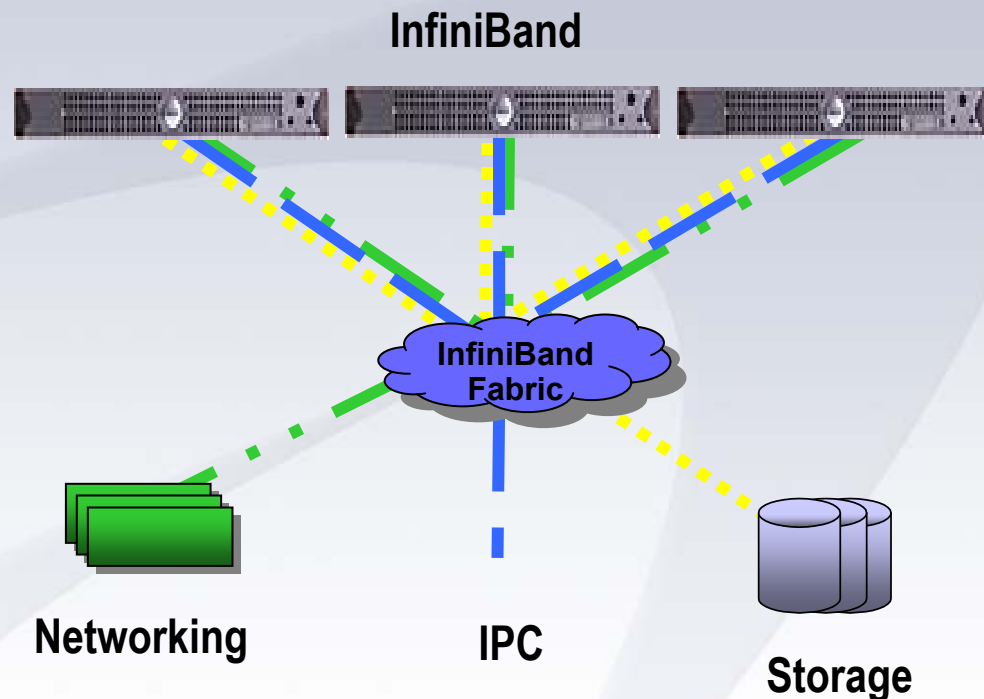- Performance impact of real applications

# Why a New I/O Interconnect

**Current Implementations: A Separate Fabric for *Each* Communication type -** *Costly, Harder to Manage/Administer & With Considerable Host Processor Overhead*

**InfiniBand offers *One Integrated* Fabric for all types of Communication - *Converges:* IPC, Storage I/O & Network I/O**

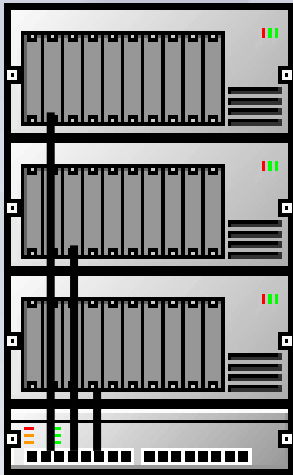*InfiniBand Silicon Solutions*

# Data Center Performance Impact

- All three types of I/O exist in a server and data center

- I/O fabric must provide performance for each type

- InfiniBand delivers on performance metrics for all I/O types
  - CPU utilization
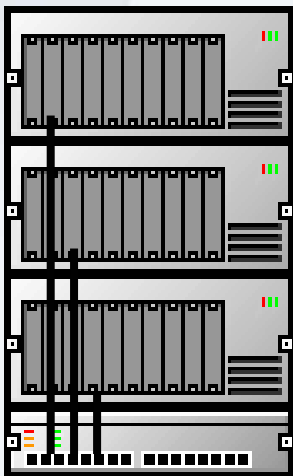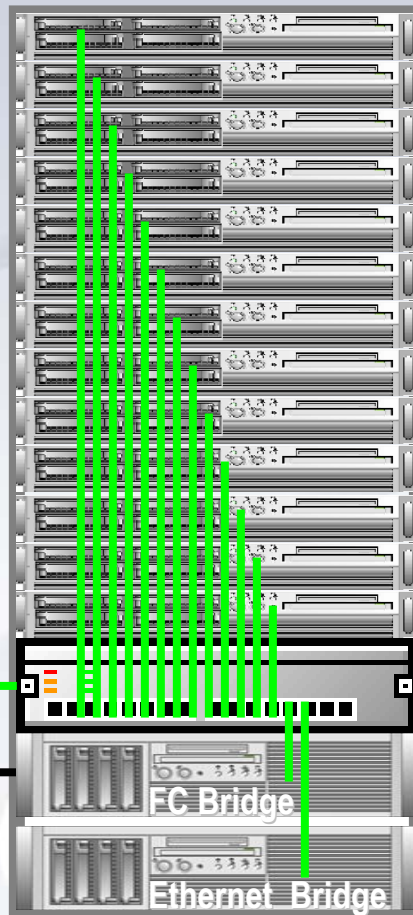  - Scalable throughput
  - Extremely low latency

**InfiniBand**

**InfiniBand Fabric**

**Networking**

**IPC**

**Storage**

*InfiniBand Silicon Solutions*

# InfiniBand Architecture in the Data Center

**Fibre Channel SAN**

**InfiniBand Storage**

**InfiniBand* Rack**

**Network Storage (NAS / iSCSI)**

**Fibre Channel SAN**

**LAN – MAN - WAN**

FC Bridge

Ethernet Bridge

*InfiniBand Silicon Solutions*

# With InfiniBand ...

- Data center scales for all I/O
  - Bandwidth ready for large installs
  - Multiple 10GigE and 10GigFC easily handled
- Applications receive more CPU MIPS
  - CPU usage returned to application from I/O
- High performance clustering
  - Low latency enables large cluster size
  - Service partitions divide I/O and IPC on the same fabric

# Define Drivers

- CPU Load
  - Time spent by application processors on I/O protocol
  - Current GigE NIC and FC HBA evaluation
- IOPS
  - Number of discrete I/O's the system can drive through the I/O device
  - Current FC HBA evaluation
- Bandwidth
  - The total throughput of data an application can drive through the I/O
  - Current GigE NIC and Switch and FC HCA and Switch evaluation
- Latency
  - The time taken from an input of a packet or I/O to the output of that packet or I/O
  - Current GigE Switch and FC Switch evaluation

*InfiniBand Silicon Solutions*

# Application Performance Drivers

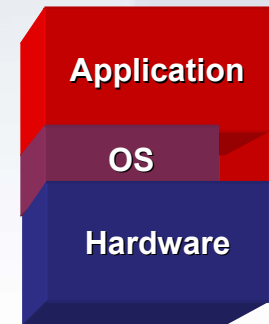|  | Storage | Networking | IPC |
|---|---|---|---|
| CPU Load<br>*processing metric* | Moderate | Important | Important |
| IOPS<br>*processing metric* | Important | Important | Moderate |
| Bandwidth<br>*data metric* | Important | Moderate | Low |
| Latency<br>*data metric* | Low | Low | Important |

# Impact of InfiniBand on Processing

- Full Offload of reliable communication
  - No SW required to ensure data delivery
  - Linear growth of CPU MIPS to I/O capability
- User mode and Virtual Addressing
  - No OS protection required to access fabric
  - Zero copy for data structures
- Multiple Channels and Priorities
  - Independent resources eliminate sharing
  - "Fire and Forget" for applications

**Traditional HW/SW Stack**

Application

OS Overhead

Hardware

**InfiniBand HW/SW Stack**

Application

OS

Hardware

**(Not to Scale)**
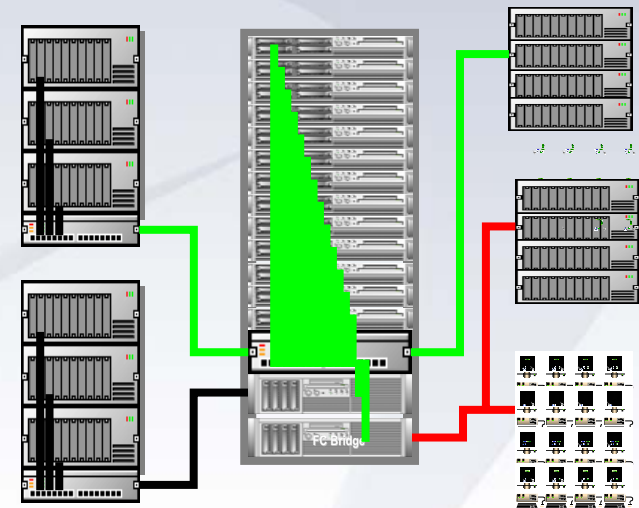
*InfiniBand Silicon Solutions*

# Impact of InfiniBand on Data Movement

- User Mode Access to Fabric
    - Applications move data without Kernel transition
    - Multiple Applications access in parallel
- Independent resources for Priority Levels
    - Different traffic streams have fabric paths
    - Cluster traffic not blocked by I/O traffic
- Low Hop Delay through Fabric Components
    - Stacked switches don't introduce large delay
    - Fabric level delays constant as fabric grows
- Scalable Link Speed
    - Addition of link width transparent to fabric
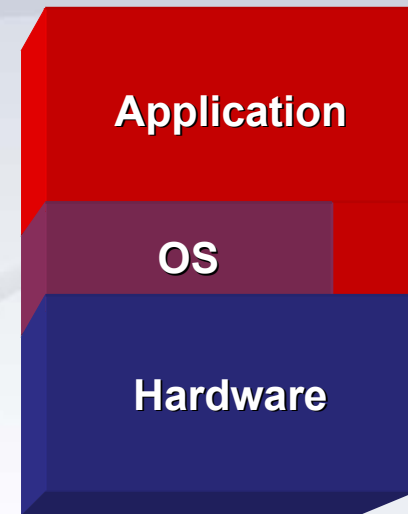    - Seamless addition of bandwidth anywhere

# Application Benefits of InfiniBand Performance

- Storage
  - Increased bandwidth yields faster access to data
  - Increased IOP allows more simultaneous accesses

- Networking
  - Lower processing overhead increases throughput to LAN clients
  - Increased IOP allows more efficient small packet web traffic

- Clustering
  - Lower processing overhead gives more CPU to the application
  - Lower latency allows clusters to scale efficiently

**Data Center**

*InfiniBand Silicon Solutions*
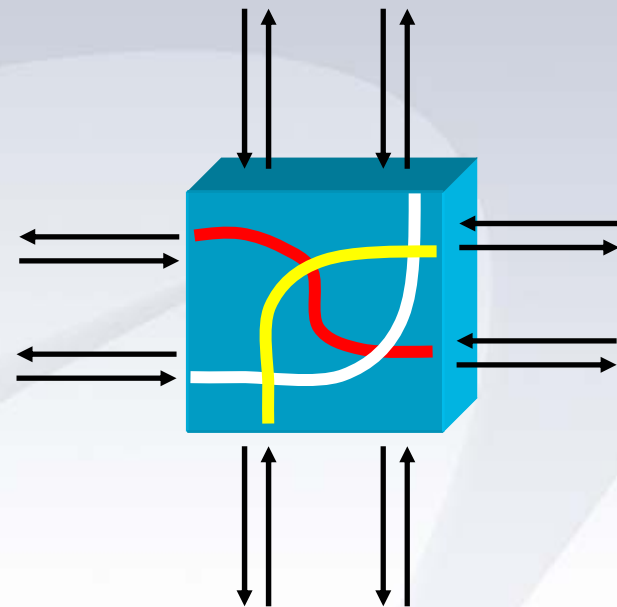
# Metrics for InfiniBand Performance

- CPU Utilization
  - Channel Services interface measurement
  - Price paid for accessing the fabric

- Messages Per Second
  - Channel Adapter Send/Receive work measurement
  - Transport message processing capability

**Application**

**OS**

**Hardware**

**InfiniBand Software**

# Metrics for InfiniBand Performance

- Bandwidth
  - Bus and wire level Gbit rate measurement
  - Efficiency of interfaces and internal structure
- Latency
  - Bus and wire level measurement in nano seconds
  - Speed of internal structures in "cut-through" mode



**8p4x InfiniBand Switch**

# Company Examples

*InfiniBand Silicon Solutions*

# **Banderacom**
## *InfiniBand Silicon Solutions*

InfiniBand Solutions Conference

April 2002

Chris Pettey, CTO

# **Metrics for InfiniBand**

| | Storage | Networking | IPC |
|---|---|---|---|
| **CPU Load** *processing metric* | Moderate | Important | Important |
| **IOPS** *processing metric* | Important | Important | Moderate |
| **Bandwidth** *data metric* | Important | Moderate | Low |
| **Latency** *data metric* | Low | Low | Important |

*InfiniBand Silicon Solutions*

# Example Requirements for IB Performance - IOPs

### IB-FC Bridge

- SRP protocol
  - $3^+$ IB IOPs per FC IOP
  - 2 SEND + $1^+$ RDMA
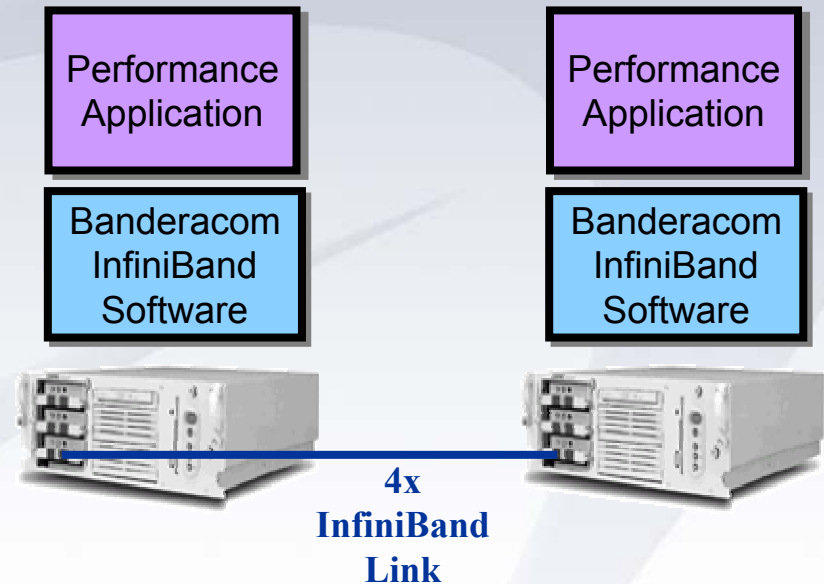- 50K to 100K FC IOPs
- 150K to 300K IB IOPs

### IB-Ethernet Bridge

- vNIC protocol
  - 1 IB IOPs per packet
  - Either SEND or RDMA
- 100K to 500K GigE packets per second
- 100K to 500K IB IOPs
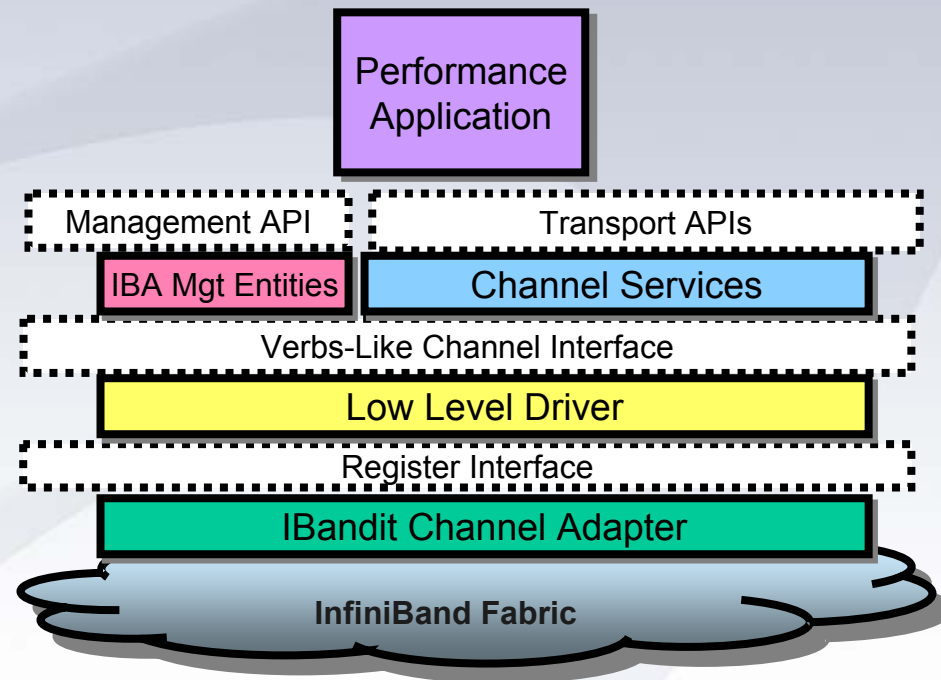
# Performance Application

### Hardware

- *Prototype PCI-X system at 133MHz*
- *Single Pentium III at 1 GHz*
- *IBandit-TCA operating in 4x mode and point-to-point interconnect*

Performance Application

Banderacom InfiniBand Software

Performance Application

Banderacom InfiniBand Software

**4x InfiniBand Link**

*InfiniBand Silicon Solutions*

# Performance Application

## Software

- *Performance application stressing at the channel services interface using RC*
- *Banderacom software running on Windows 2000*

| Performance Application |
|---|

| Management API | Transport APIs |
|---|---|
| IBA Mgt Entities | Channel Services |

Verbs-Like Channel Interface

Low Level Driver

Register Interface

IBandit Channel Adapter

**InfiniBand Fabric**

*InfiniBand Silicon Solutions*

# IBandit-TCA Performance

- IBandit-TCA saturates 133MHz PCI-X bus
  - > 6 Gbps bi-directional transfer rate
  - 70-80% PCI-X bus utilization using prototype PCI-X chipset
- IBandit-TCA is capable of driving >350 KIOPS
  - Currently limited by bus and CPU capabilities, IBandit-TCA can scale even higher
- IBandit-TCA performance scales linearly across QPs
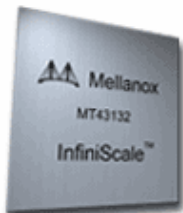
# INFINIBAND℠

## TRADE ASSOCIATION

# Mellanox Technologies, Inc
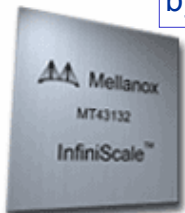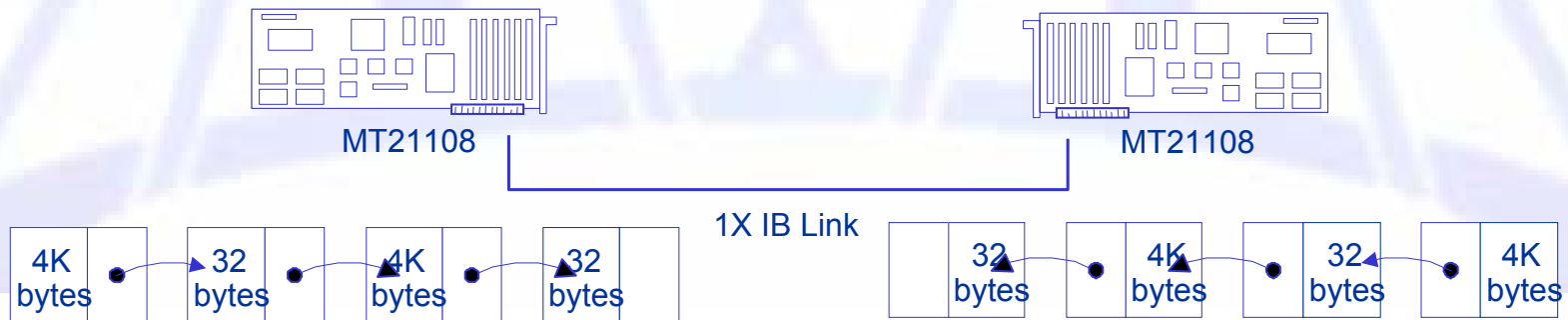## InfiniBand Solution Conference
## April, 2002

**Chris Eddington**
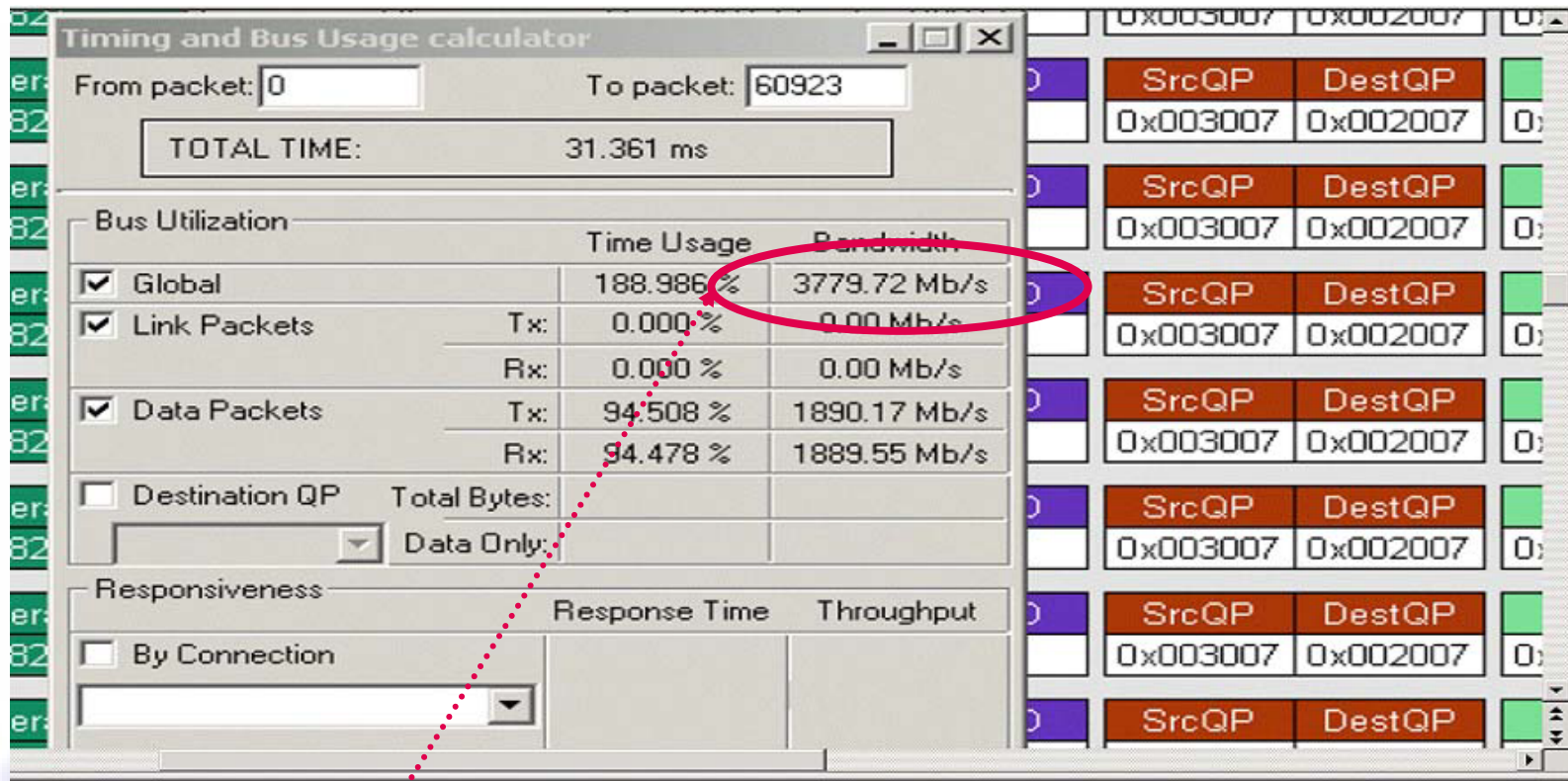
**Director of Technical Marketing**
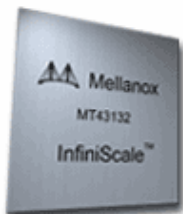
# Test Configuration

- InfiniBridge™ MT21108 in a storage application that pushes for the highest InfiniBand 1X bandwidth

- Throughput measurements from storage application under development

- Systems: 800 MHz Pentium III $^®$ CPU 64 bit 66 MHz PCI bus

- Even mix of large (4KB) and small (32B) packets

- Fully pipelined descriptor posting and completion

- Takes advantage of completion aggregation



MT21108                    MT21108

1X IB Link

| 4K bytes | 32 bytes | 4K bytes | 32 bytes |   | 32 bytes | 4K bytes | 32 bytes | 4K bytes |

# InfiniBridge™ Link Performance Update



- Result: ~3.8 Gb/s or 94% of maximum bandwidth achieved with only 7% processor utilization

# INFINIBAND℠

TRADE ASSOCIATION

# LANE 15 SOFTWARE

*The Leader in InfiniBand™ Management Software*

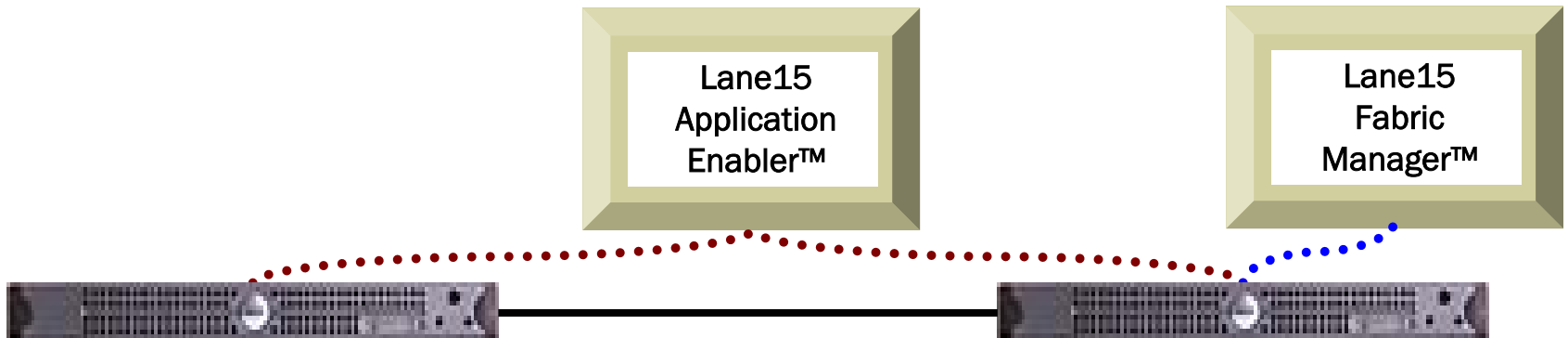*IBSC Performance Presentation*

*Brad Benton*

*Technical Relationship Manager*
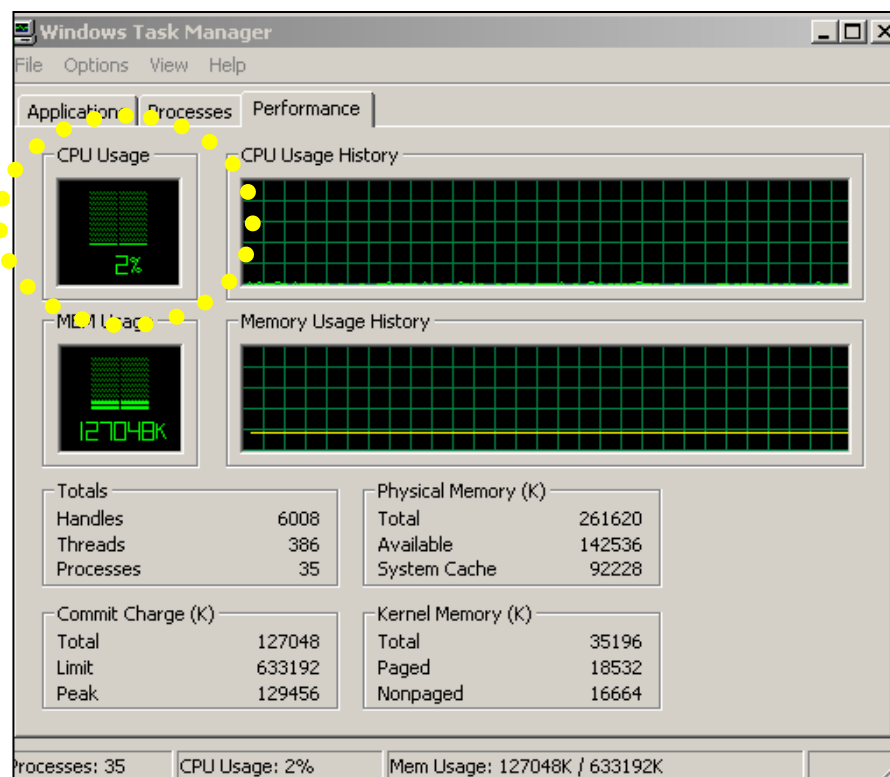
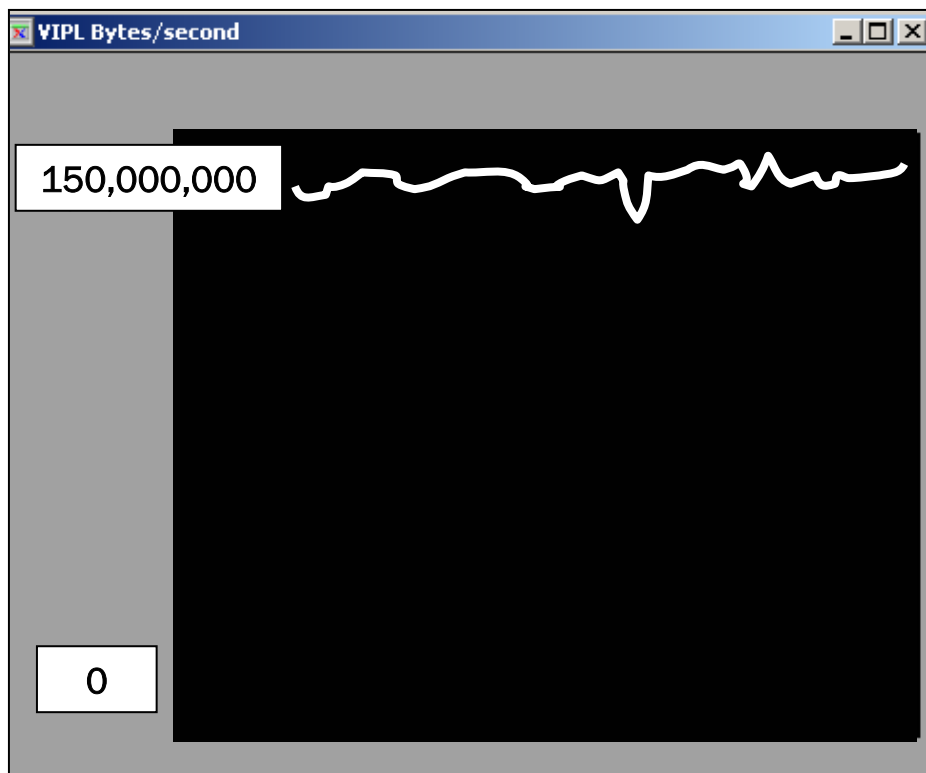# Lane15 Application Performance

## Demonstration

- VIPL-based data pump
- High-speed data transfers
- Low protocol overhead
- OS bypass technology
- Yield high data rates
- Minimal CPU overhead

**150,000,000 bytes per second with < 2% CPU Utilization!**
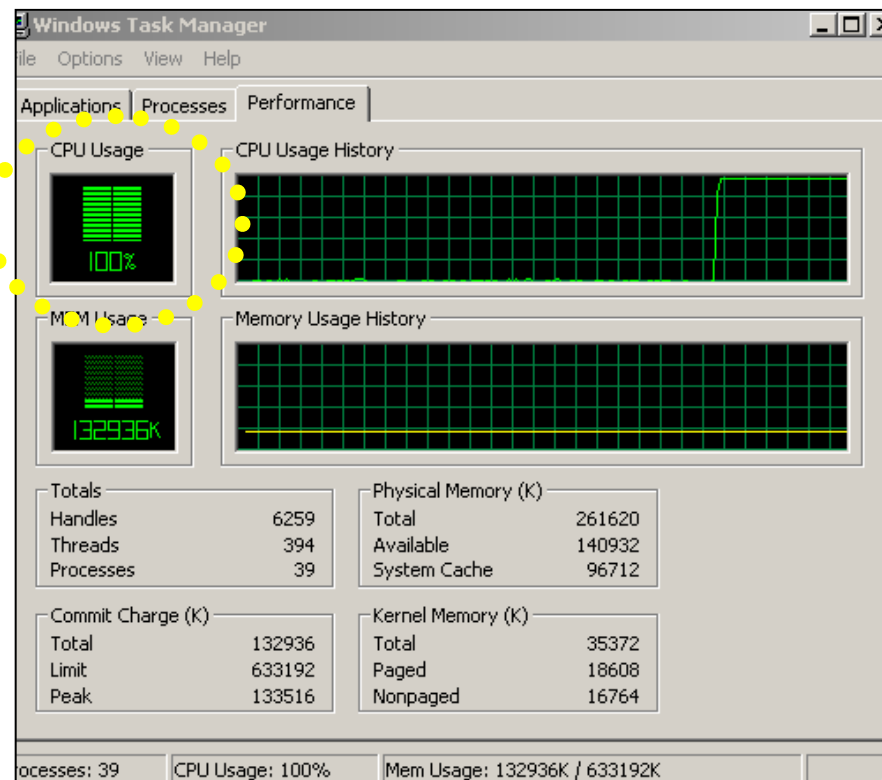
Lane15
Application
Enabler™

Lane15
Fabric
Manager™

# High data rate with low CPU overhead
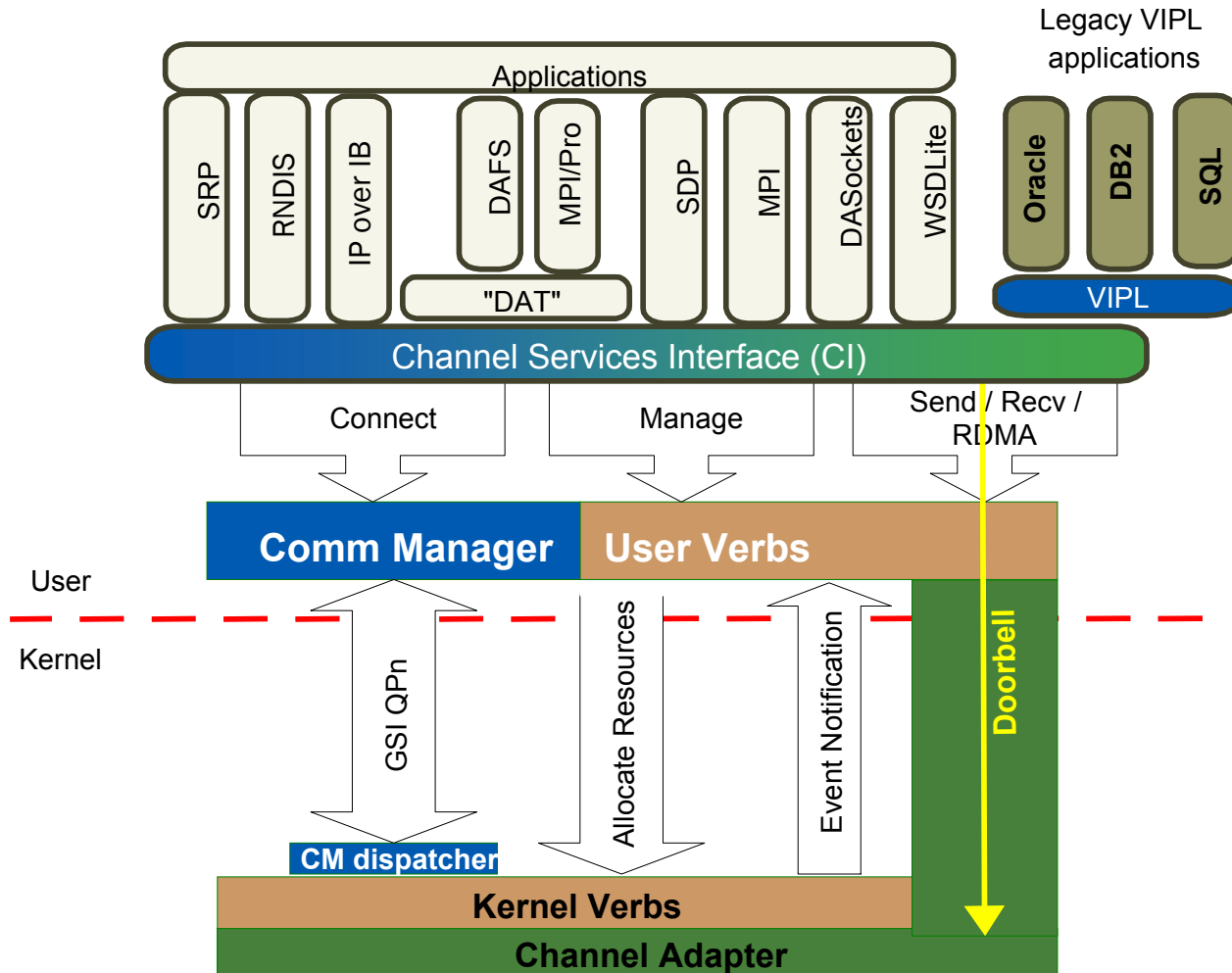


**CPU Utilization < 2%**

# Data rate with high application load



## Data rates are the same!

# High Performance Architecture

# Summary

- InfiniBand™ Performance is <u>REAL</u>

- VI based applications run <u>unmodified</u> on InfiniBand Fabrics today

- Native InfiniBand Applications are coming!

> InfiniBand Architecture gives you your CPU back!

# INFINIBAND℠

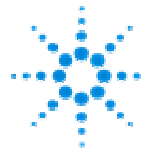## TRADE ASSOCIATION

# Contents

- **Overview**

- **Definition of Terms**

- **Test method and setup**

- **Measurement results**

- **Summary**

IB switch performance
14 March, 2002

Roland Scherzinger /
Thomas Dippon

**Agilent Technologies**

Page 3

# Overview

- **The performance parameters of switches are key to overall performance of fabric topologies such as InfiniBand™ Architecture**

- **Due to the number of input factors and the complexity of switch cores it is hard to predict actual performance under real-life load conditions**

- **Therefore measurement of switch performance parameters is critical for manufacturers to verify simulated results under real-life and worst-case load conditions**

- **End-users can benefit from measurements to compare various switch implementations and verify claims of manufacturers**

- **The goals of this presentation are**

  - **To show a method for measuring InfiniBand switch performance**

  - **To show how well InfiniBand switches work - even in this early stage of the technology**

  - **NOT a competitive analysis**

IB switch performance
14 March, 2002

Roland Scherzinger /
Thomas Dippon
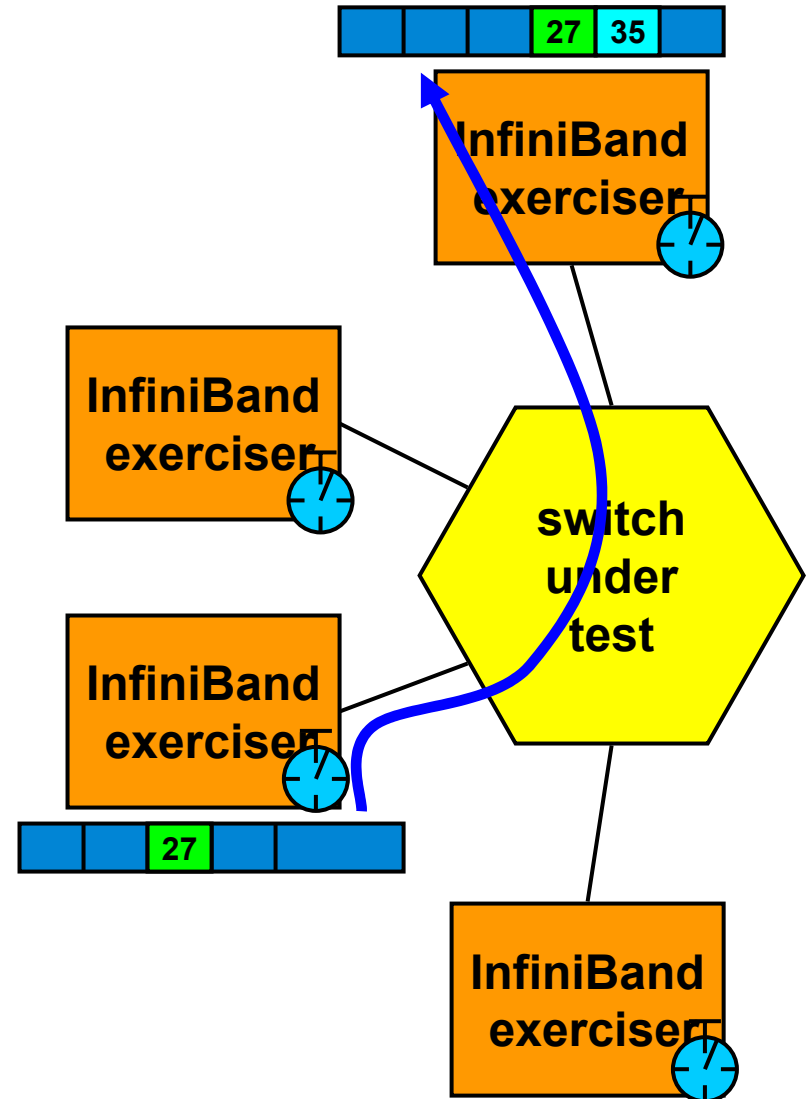
**Agilent Technologies**

Page 4

# Definition of Terms

- **Stream**

    - **A continuous sequence of packets with certain characteristics (load, range of packet sizes, VL, SL, etc.) transmitted from one port of a switch to another**

- **Throughput**

    - **The number of bytes per second that are actually transferred through each port of a switch depending on various operating conditions (load on each port, packet sizes, size of forwarding tables, etc.).**

- **Latency**

    - **The difference between the time when the first byte of a given packet enters a switch and the time when the first byte of the same packet leaves the switch**
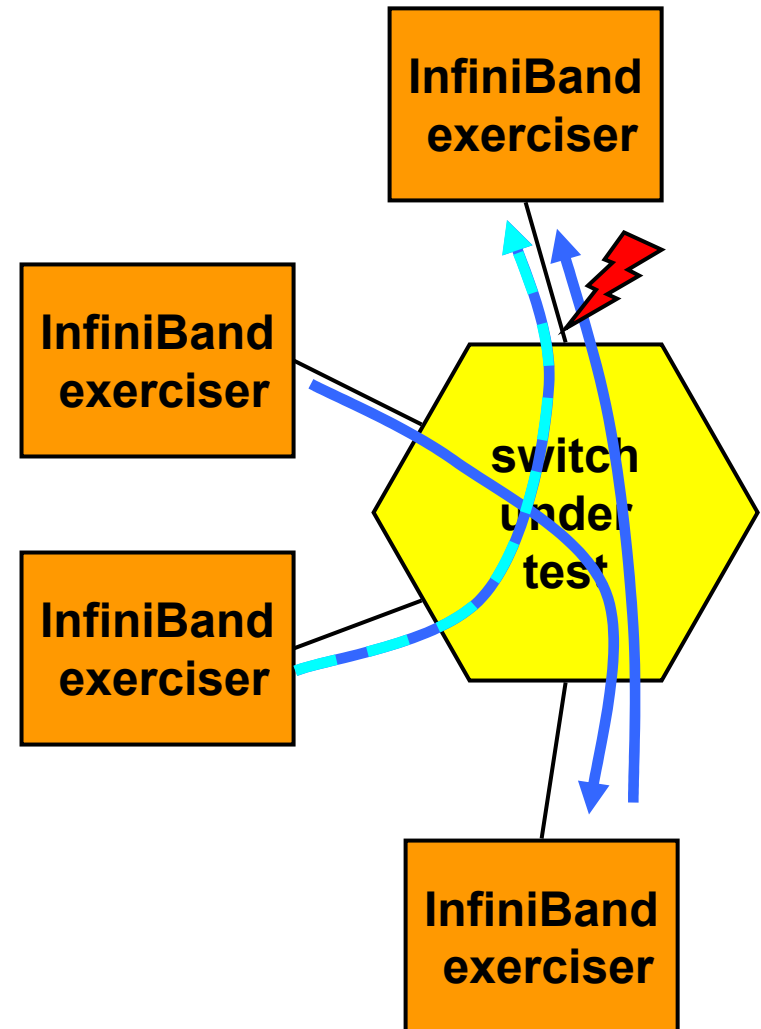
IB switch performance
14 March, 2002

Roland Scherzinger /
Thomas Dippon

Agilent Technologies

Page 5

# Test Method (1)

- **InfiniBand(TM) exercisers are connected to the ports of a switch**

- **All exercisers have precisely synchronized timestamp clocks (+/- few nanoseconds)**
- **Exercisers generate "instrumented packets" by inserting a timestamp in the payload of each packet just before transmission**
- **Upon reception, the difference between the receivers clock and the timestamp value in the packet indicates the latency**
- **In addition there are counters keep track of the number of transmitted/received bytes/packets**

| | | | 27 | 35 | |
|---|---|---|---|---|---|

**InfiniBand exerciser**

**InfiniBand exerciser**

**switch under test**

**InfiniBand exerciser**

| | | 27 | | |
|---|---|---|---|---|

**InfiniBand exerciser**

IB switch performance
14 March, 2002

Roland Scherzinger /
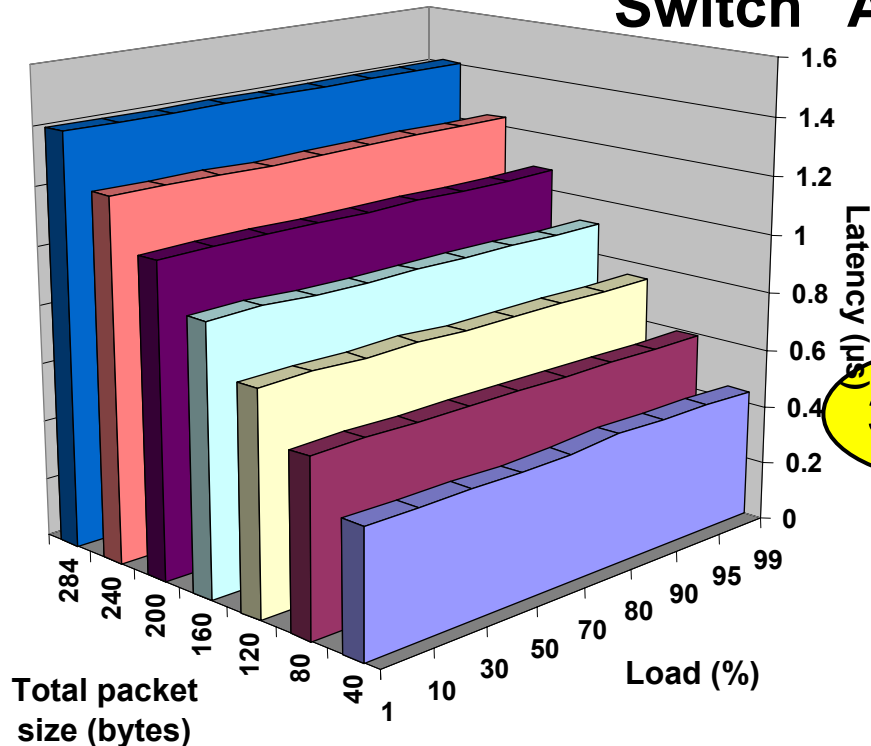Thomas Dippon

Agilent Technologies

Page 6

# Test Method (2)

- **Performance parameters are measured with one or more simultaneous streams**

    - **Single stream**
    - **With other background traffic**
    - **With traffic congestion (two or more streams destined for the same port)**

    - **Etc…**
- **Stream parameters can be varied**

    - **Packet size (or range of packet sizes)**

    - **Inter-packet delay (used bandwidth)**

    - **Burst of packets or continuous stream**

IB switch performance
14 March, 2002

Roland Scherzinger /
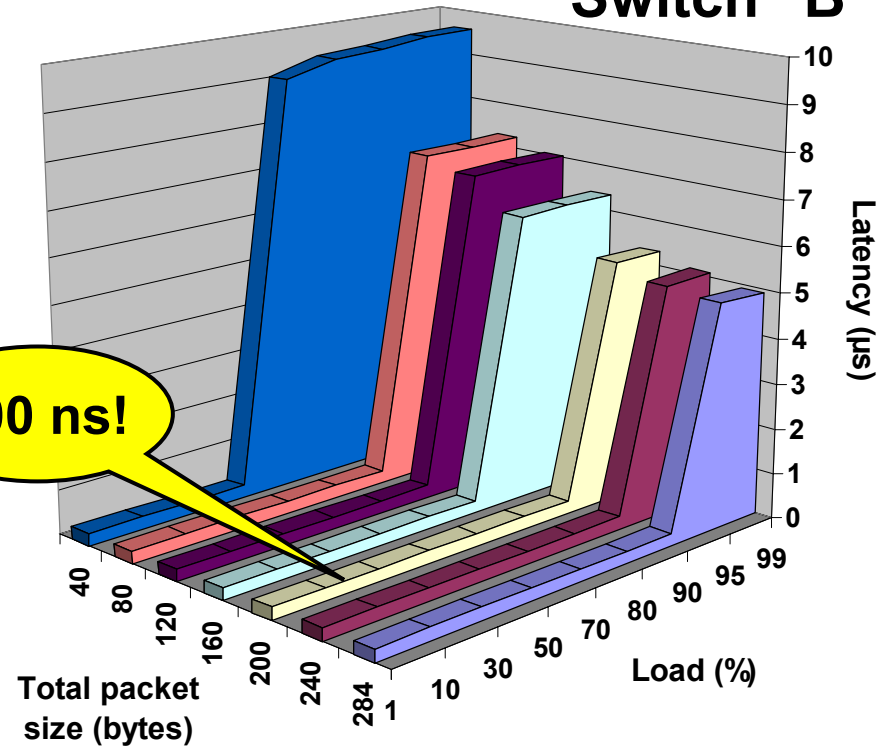Thomas Dippon

Agilent Technologies

Page 7

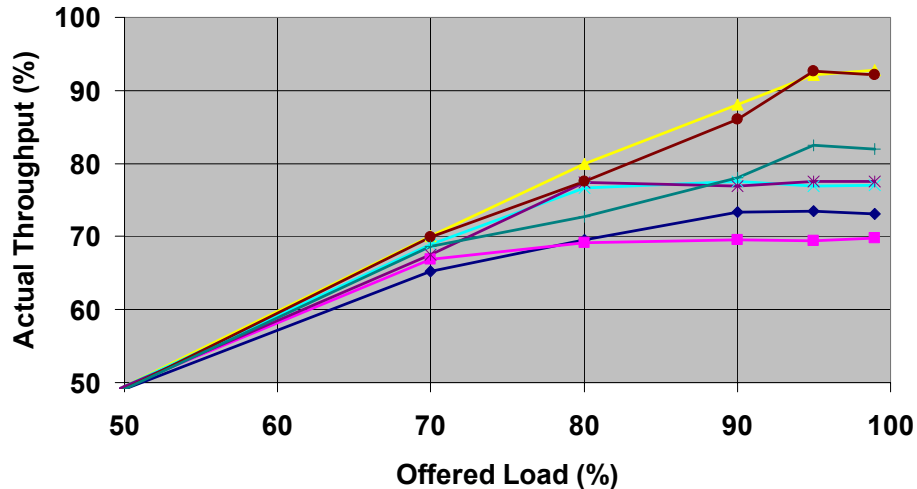# Results of Latency Measurements

## Switch "A"

## Switch "B"



- **Latency increases linearly with packet size: store & forward**
  - **~1μs per 256 byte packet size**
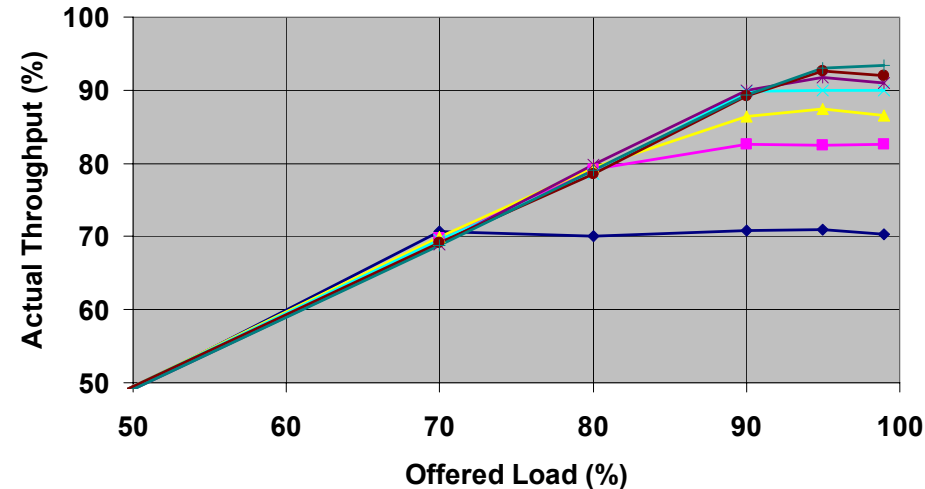- **No increased latency even at extreme load**

- **Latency in wide range independent of packet size: cut through switching**
- **Latency step function at a certain load level indicates throughput limitation on the output**

# Results of Throughput Measurements

## Switch "A"



## Switch "B"



- **Early testing shows excellent results up to ~70…80% load**

- **Throughput limitations are typically seen with small packets**

IB switch performance
14 March, 2002

Roland Scherzinger /
Thomas Dippon

Agilent Technologies

Page 9

# Summary

- **Even 1st generation InfiniBand switches show very impressive performance numbers**

    - **300 ns cut-through latency!**

- **It is important to understand the performance characteristics and limitations of switches**

    - **Even early switches perform excellent up to ~70…80% load**

- **Many other scenarios (multiple streams, multiple partitions, P_Key checking, SL to VL mapping, weighted VL priorities etc.) have to be taken into account to get a realistic picture**

- **InfiniBand 4x silicon and switches are now available - stay tuned for great performance numbers**

IB switch performance
14 March, 2002

Roland Scherzinger /
Thomas Dippon

Agilent Technologies

Page 10

How Data Gets Delivered

# Auspex InfiniBand™ Architecture Project

**Using an InfiniBand Architecture point solution to deliver even more enterprise data storage performance!**

# InfiniBand™ Architecture Meets NAS

- Improved NFS mounts and CIFS authentication by moving to an InfiniBand interconnect (higher number of mounts more rapidly)
  - Zero load on the CPU doing IO transfers

- Increase the scalability of the overall multi-node system immediately
  - from 3 to 4 I/O Nodes
  - eventually to N+1 nodes with an InfiniBand switch

- Increased interconnect bandwidth with InfiniBand Architecture
  - high-speed, point-to-point, bi-directional communication

# InfiniBand™ Architecture Meets Enterprise Storage

- Removal of existing software delays on both the network and host processors

- Significant load reduction on both network and host processors (processor load observed using PCI bus analyzer with various throughput quantities)

- Virtually unlimited I/O sizes on each DMA element

- NSc3000 provides NAS file services for SANs, extending the Infiniband benefit to the SAN data delivery

- Future interconnect to SAN? Within SAN? Other point solutions?

# Some Specific Numbers

*Before*
*InfiniBand™ Architecture*

*After*
*InfiniBand™ Architecture*

- 60 MB/sec PCI bus throughput causes 89% CPU utilization

- 60 MB/sec PCI bus throughput causes 12% CPU utilization

- 90 MB/sec PCI bus throughput causes 100% CPU utilization

- 90 MB/sec PCI bus throughput causes 18% CPU utilization