# Chelsio Communications

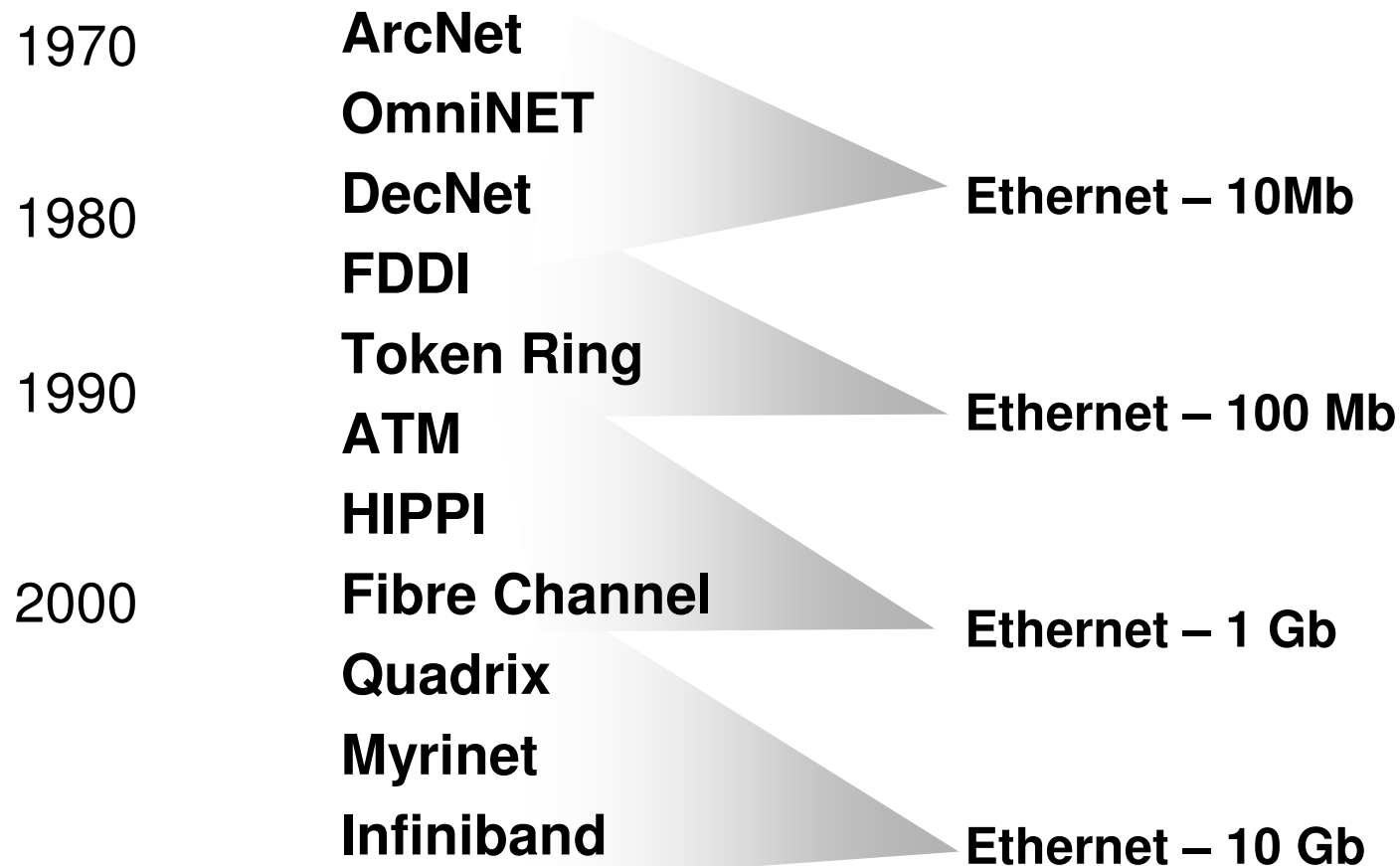## The State-of-the-Art of TOE Technology

*Michael Chen, PhD*

*PFLDnet2006 Presentation, Feb. 2006*

# Agenda

- **Technology Trend**

- **10GbE TOE Architecture**

- **TOE Support of LFN**

- **10GbE TOE Performance**
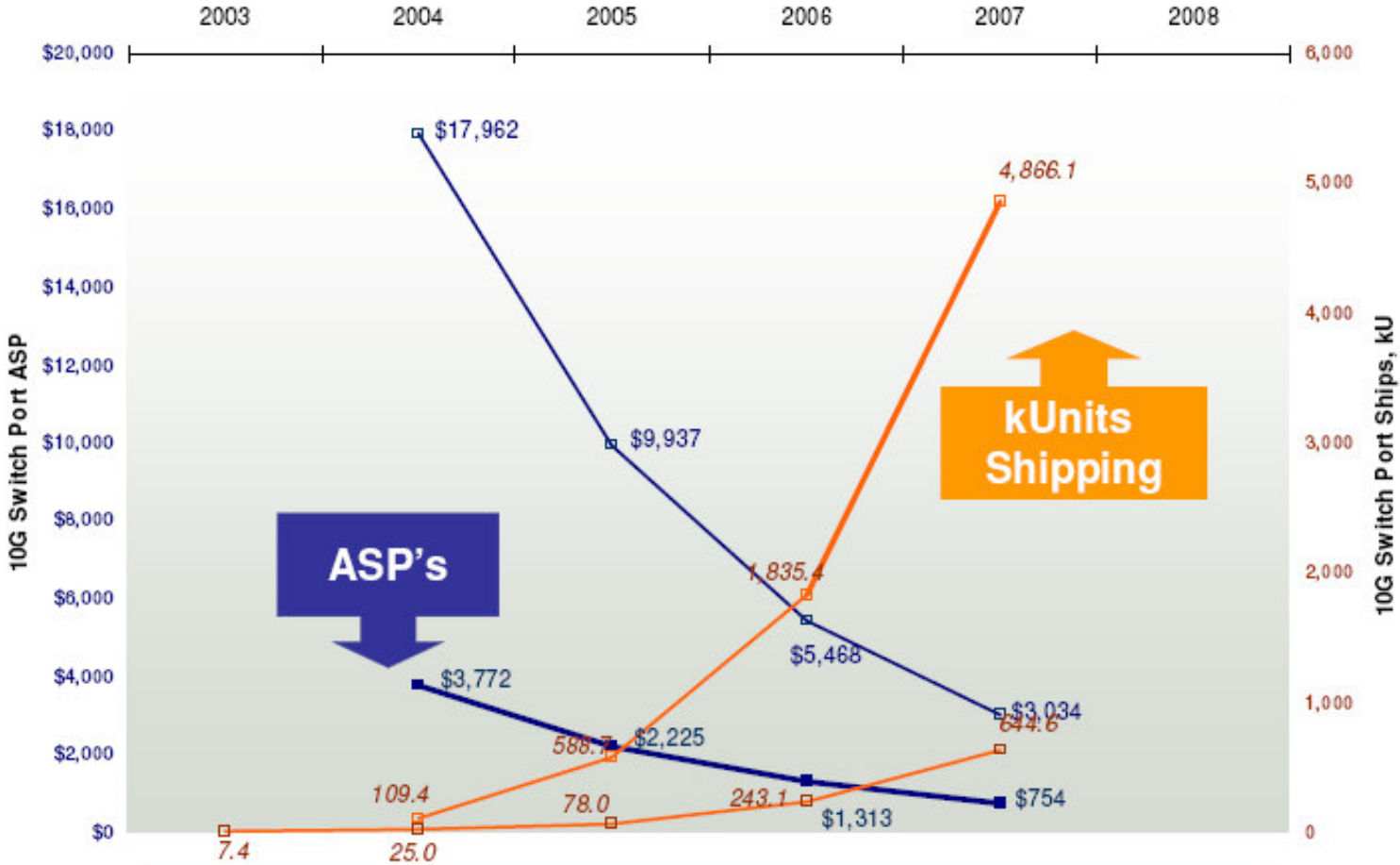
- **Network Convergence and ULP Acceleration**

# Ethernet's History of Absorbing Proprietary Networking Technologies

1970  ArcNet
      OmniNET
      DecNet                          Ethernet – 10Mb
1980
      FDDI
      Token Ring
1990                                  Ethernet – 100 Mb
      ATM
      HIPPI
      Fibre Channel
2000                                  Ethernet – 1 Gb
      Quadrix
      Myrinet
      Infiniband                      Ethernet – 10 Gb

# 10G Ready for Prime Time

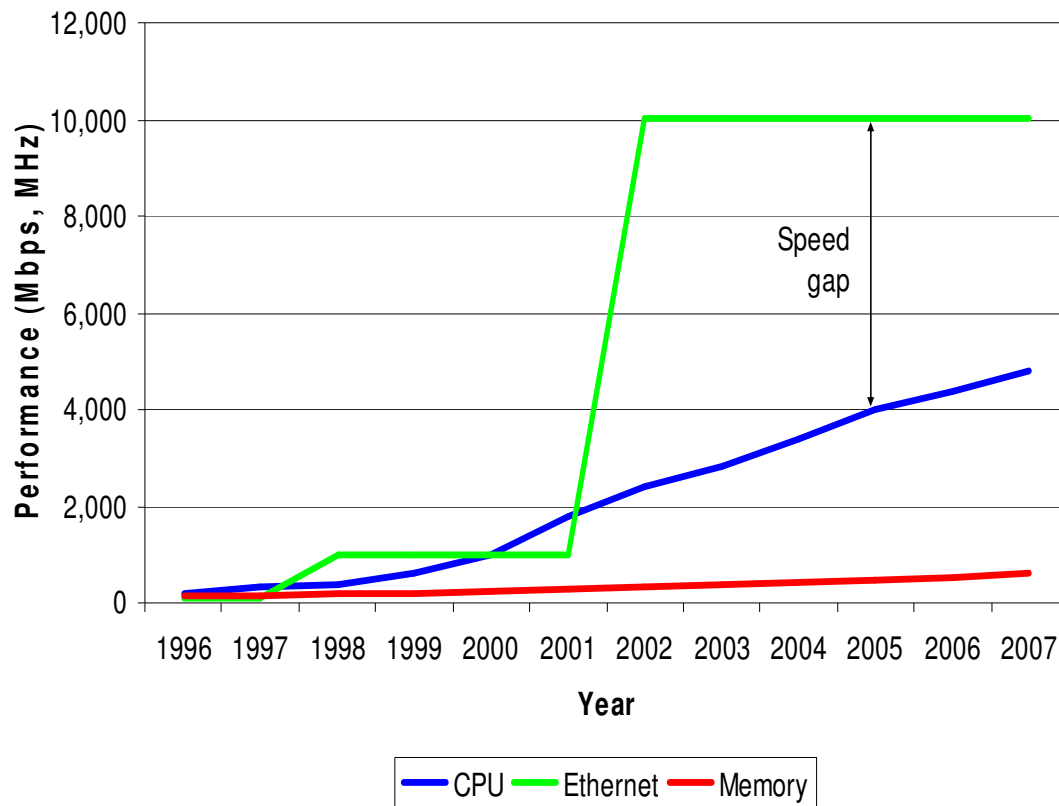| Criteria | Market Drivers / Enablers |
|---|---|
| **Units growth** | ▪ 3x volume growth in the 10GbE NIC market in 2006 (synergy report)<br>▪ 3x volume growth in 10GbE switch market in 2006 (Dell'Oro)<br>▪ iSCSI market growing at 50% per quarter<br>▪ 55% of top 500 HPC installations are Gb Ethernet |
| **Infrastructure** | ▪ High density 10GbE optical switches available now<br>▪ High density 10GbE CX4 switches available now<br>▪ Low latency switch chips available from at least 4 vendors<br>▪ UTP PHY chips available from at least 3 vendors |
| **Prices** | ▪ XFP over 12 months has dropped by more than 100%<br>▪ CX4 switch port at $700/port list price now<br>▪ CX4 adapter pricing dropping past the knee<br>▪ 10GbE HBA pricing has been halving every 12 months |
| **Standards** | ▪ 10G CX4 (copper media) introduced and shipping<br>▪ 10G-baseT silicon expected by year-end |

# 10GbE is Beating Forecast



**ASP's**

**kUnits Shipping**

10G Switch Port ASP: $20,000 — $18,000 — $16,000 — $14,000 — $12,000 — $10,000 — $8,000 — $6,000 — $4,000 — $2,000 — $0

10G Switch Port Ships, kU: 6,000 — 5,000 — 4,000 — 3,000 — 2,000 — 1,000 — 0

Years: 2003, 2004, 2005, 2006, 2007, 2008

ASP values: $17,962, $9,937, $5,468, $3,772, $2,225, $1,313, $3,034, $754

kUnits values: 7.4, 25.0, 78.0, 109.4, 243.1, 588.7, 644.6, 1,835.4, 4,866.1

**With CX4, HBA prices are at 2007 levels**

# The Speed Gap

**The Network/System Speed Gap**



- **Rule-of-thumb: 1GHz of CPU needed to process 1Gbps data rates**

- **At 10Gbps, today's highest performance CPUs lag by 2.5x**

- **In 2006, 10GbE full-duplex (20Gbps) will further widen the gap**

- **Memory speeds lag even further behind and will become main obstacle in the future**

- **The SOLUTION is Protocol Offload**

# Ethernet Popular HPC Deployments

| 2003 Top Supercomputer Cluster Interconnects | | | | | |
|---|---|---|---|---|---|
| | Clusters | % | Servers | % | Avg Server/Cluster |
| Ethernet | 88 | 55% | 15,112 | 53% | 172 |
| Myrinet | 57 | 35% | 8,890 | 31% | 156 |
| InfiniBand | 3 | 2% | 1,484 | 5% | 495 |
| Quadrics | 9 | 6% | 2,608 | 9% | 290 |
| SCI | 4 | 2% | 310 | 1% | 78 |
| Total | 161 | 100% | 28,404 | 100% | 176 |
| Source Top500.org Nov 2003 (161 new clustered systems were added to the list) | | | | | |

**Ethernet is the *dominant*
High Performance Cluster Interconnect
Today!**

# Chelsio Product Family

**N210: 10GbE** Server Adapters

**T210: 10GbE** Protocol Engines – Fiber & Copper
*Server Adapter + TCP + iSCSI + RDMA*

**T204: 4-port 1GbE** Protocol Engines

Protocol software and drivers

8

# 10GbE PHY Technologies

- **Fiber**
  - 10GBase-SR          85 m                    shipping
  - 10GBase-LR          10,000 m                shipping

- **Copper**
  - 10GBase-CX4         15 m                    shipping
  - 10GBase-T           55-100 m

- **Backplane**
  - 10GBase-KX4         0.5-1 m
  - 10GBase-KR          0.5-1 m

# Chelsio
# Communications

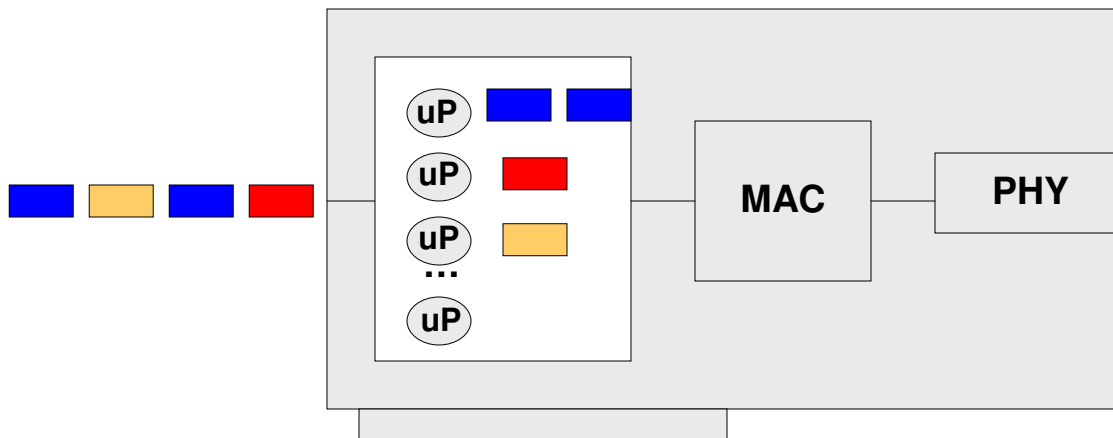## TOE Architecture

# Alternative 10G Solutions

## Basic 10G NIC



- Layer-2 protocols only
- No protocol offload intelligence

- Saturates host CPU
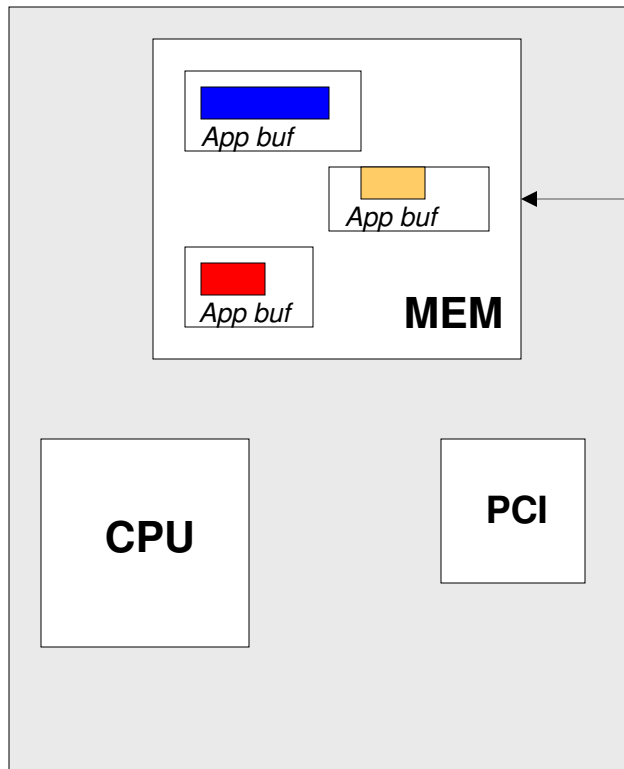- Inadequate for high-performance

## Multi-RISC based Architecture



- Offload engine consists of multiple RISC cores
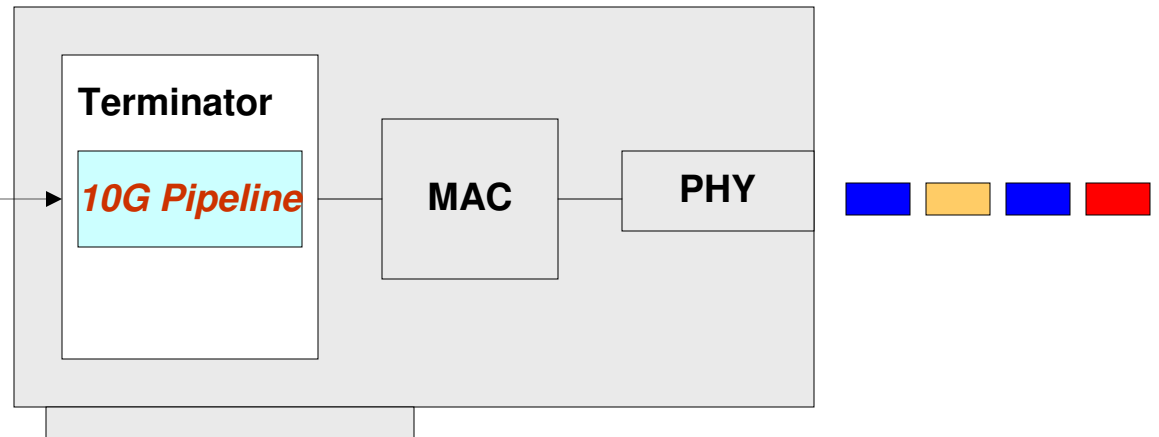- Each TCP connection's bandwidth limited by the core frequency

- Complex internal software for control & management of multiple cores
- Inadequate single channel performance & scalability
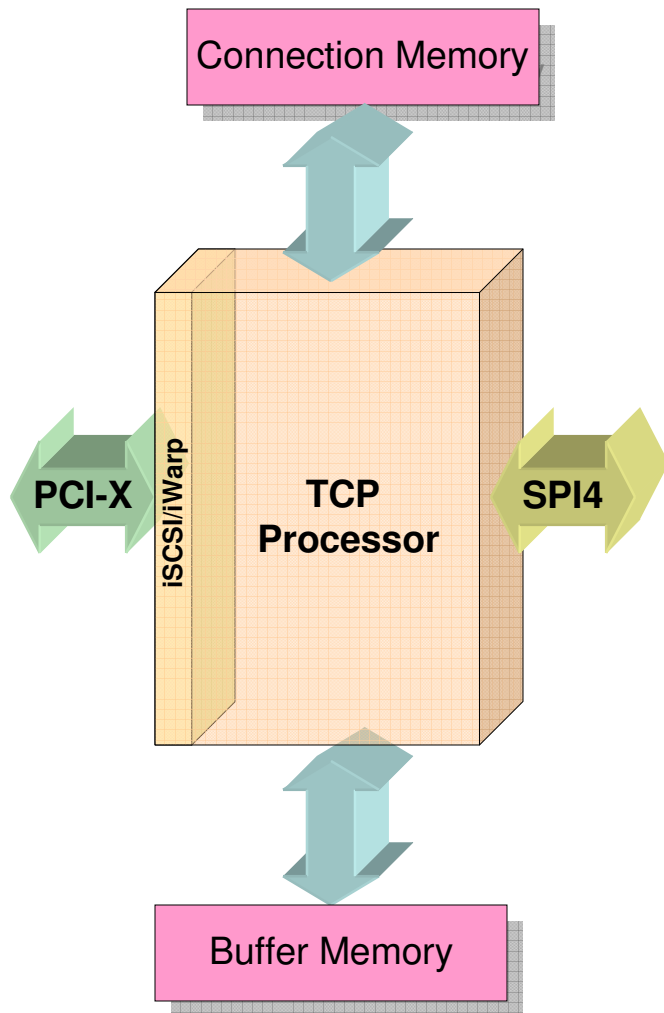
# Chelsio's Unique Architecture

**Host**

**Chelsio Terminator Architecture**



- Optimal partitioning of functions between hardware, firmware and software

- 10G VLIW processor delivers highest performance from 1 to 1000s of connections

- Pipelined architecture uses cut-through processing for low latency

- Direct data placement into application buffers eliminates copy overhead

# "Terminator" Processor ASIC

- Cut-through, wire-speed architecture
- Scalable from 10G to 1G line speeds
- TCP, iSCSI, RDMA, DDP acceleration
- 400+ configuration registers
- Programmable TCP rules per connection

**Connection Memory**

PCI-X | iSCSI/iWarp | **TCP Processor** | SPI4

**Buffer Memory**

| > L4 | iSCSI | | iWARP | | DDP | | *Accelerated* |
|---|---|---|---|---|---|---|---|

| L4 | UDP/TCP Checksum Offload | Full TCP Protocol Conformance | | | Per-Connection, per-Server Control | |
|---|---|---|---|---|---|---|
| | | Congestion Avoidance and Control | High Performance Extensions | | Increased Initial Window Size | |
| | | TCP Traffic Filtering | SYN Cookies | | Wildcard Sockets | |

*Offloaded*

| IP | IP Checksu\m Offload | Route Table Cache | QoS, DiffServ, DSCP | ARP Cache | Path MTU Discovery |
|---|---|---|---|---|---|
| MAC | Multiple MAC Address per Port | Jumbo Frames | VLAN Tagging *802.1q* | Prioritization *802.1p* | Ethernet-II and 802.3 Frames |
| Ether | 1 Gbps, 10 Gbps Ethernet | | | Copper / LR and SR Fiber | |

13

# 10GE TCP Processing

- **Current generation NPU not good match for 10GE TCP processing**

- **TCP at 10Gb characteristics and requirements**
  - Stateful protocol -> efficient RMW
  - Large number of connections -> scalable architecture
  - Jumbled byte stream -> intelligent memory system
  - < 10us latency requirement -> cut-through processing
  - Still an evolving protocol -> programmability

# 10GE VLIW TCP Processor Innovation

- **The TCP protocol is stateful**
  - At 10Gbps, 1500 byte packets are 1us apart
  - TCP state has poor cache locality
  - Wire speed needs to be attainable for 1 connection
  - Wire speed needs to be attainable for 1000s connections

  - Requires
    - An efficient pipelined pre-fetch of the TCP state
    - Single engine that can process 10Gbps traffic

- **The TCP protocol provides a FIFO stream abstraction to the end points**
  - Packets can partially overlap and/or arrive out of order
  - Latency requirement rules out store-and-forward

  - Requires
    - Specialized memory subsystem to unravel the packet jumble @10Gbps speeds

# 10GE VLIW TCP Processor Innovations

- **There are stringent end-to-end latency requirements (< 10us) in addition to high BW requirement (10Gbps)**
  - This requires cut-through processing
  - Cut-through processing refers to the packet arriving on one terminal, being processed and the being forwarded out the other terminal without ever being stored in off-chip memory

- **Measured end-to-end latency < 10us with L2 switch**
  - DMA engine used interrupts but could push number lower by using polling mode

# Optimized Architecture

| | 10G VLIW | Multi-RISC |
|---|---|---|
| **Scalability** | **Unlimited** | **Limited by # of CPUs** |
| **Firmware Complexity** | **Low** | **Typically 1+ year firmware debug** |
| **Cache Capacity** | **Unlimited** | **Limits maximum # of accelerated connections** |
| **Performance Profile** | **Linear, uniform bandwidth per connection** | **Falls off once IPC becomes significant** |
| **Roadmap** | **Low-risk upgrade path** | **Complex firmware more difficult to scale** |

# Chelsio Communications

## TOE Support of LFN

# Congestion Control in LFN

- **RFC 3649: Highspeed TCP**

- **In Congestion Avoidance,**
  - for each ACK, increase the window by
    - $w = w + a(w) / w$

    Note: in standard TCP
      a(w)=1                    // when w is maintained by #segs
      a(w) = mss * mss        // when w is maintained by bytes

  - For each congestion event, decrease the window by
    - $w = (1-b(w)) * w$,   // where $0 < b(w) <= 0.5$

# Congestion Control in LFN

- **Table-driven Implementation**

  - For each ACK received, using current *w* as index to lookup a table for $a(w)$

  - For each congestion event, using current *w* as index to lookup a table for $b(w)$

  - The lookup table is SW configurable which provides the max flexibility for various LFN environments.

## AB-Table

| W | A(w) | B(w) |
|---|---|---|
| < 38 | 1 | 0.5 |
| 38 (56k) | 1 | 0.5 |
| 118 (172k) | 2 | 0.44 |
| 221 (322k) | 4 | 0.41 |
| ... | ... | ... |
| 5610 (11M) | 21 | 0.24 |
| ... | ... | ... |
| 83000 (120M) | 70 | 0.09 |
| ... | ... | ... |

# Traffic Pacing and Shaping

- **Researches [Hiraki-SC04, etc] indicated the importance of pacing TCP streams across LFN to reduce the traffic burstness**

- **SW traffic pacing/shaping at 10Gb rate is CPU intensive**

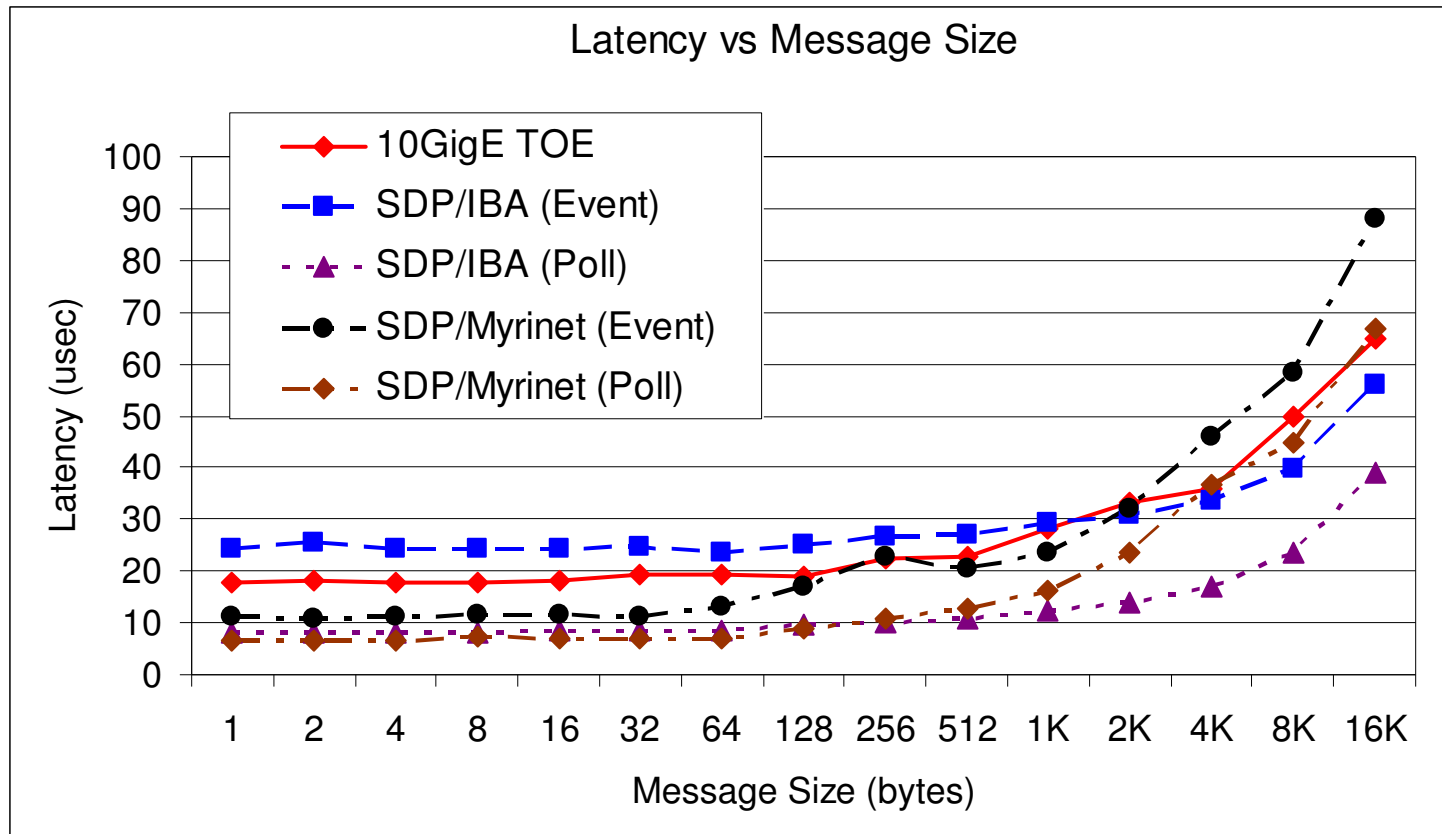- **TOE enables HW traffic pacing at TCP level**

# Chelsio Communications

## TOE Performance

# OSU/LANL HPC Benchmarks 10GbE TOE vs IB & Myrinet
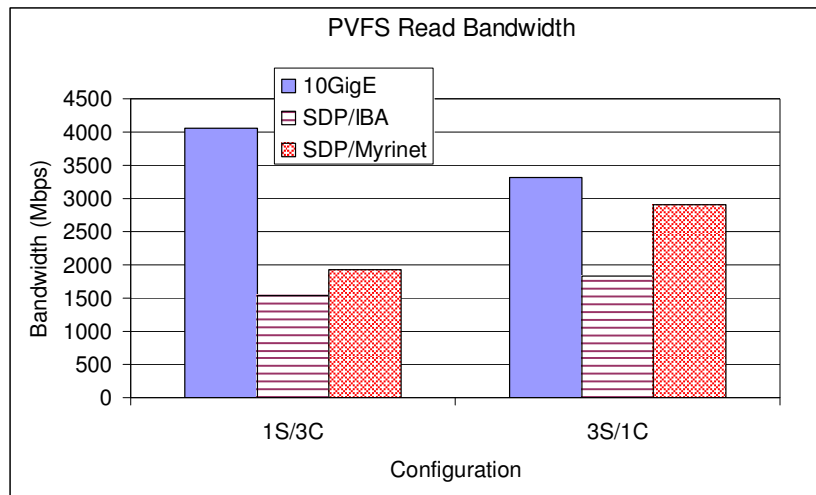


**Bandwidth vs Message Size**

**Source:** *Head-to-TOE Evaluation of High-Performance Sockets over Protocol Offload Engines* by DK Panda et al
**Test configuration:** 4-node cluster connected through 10GbE switch running single connection
**System configuration:** Dual 32-bit Intel Xeon 3.0GHz processors running Red Hat 9.0 Linux kernel 2.4.25smp
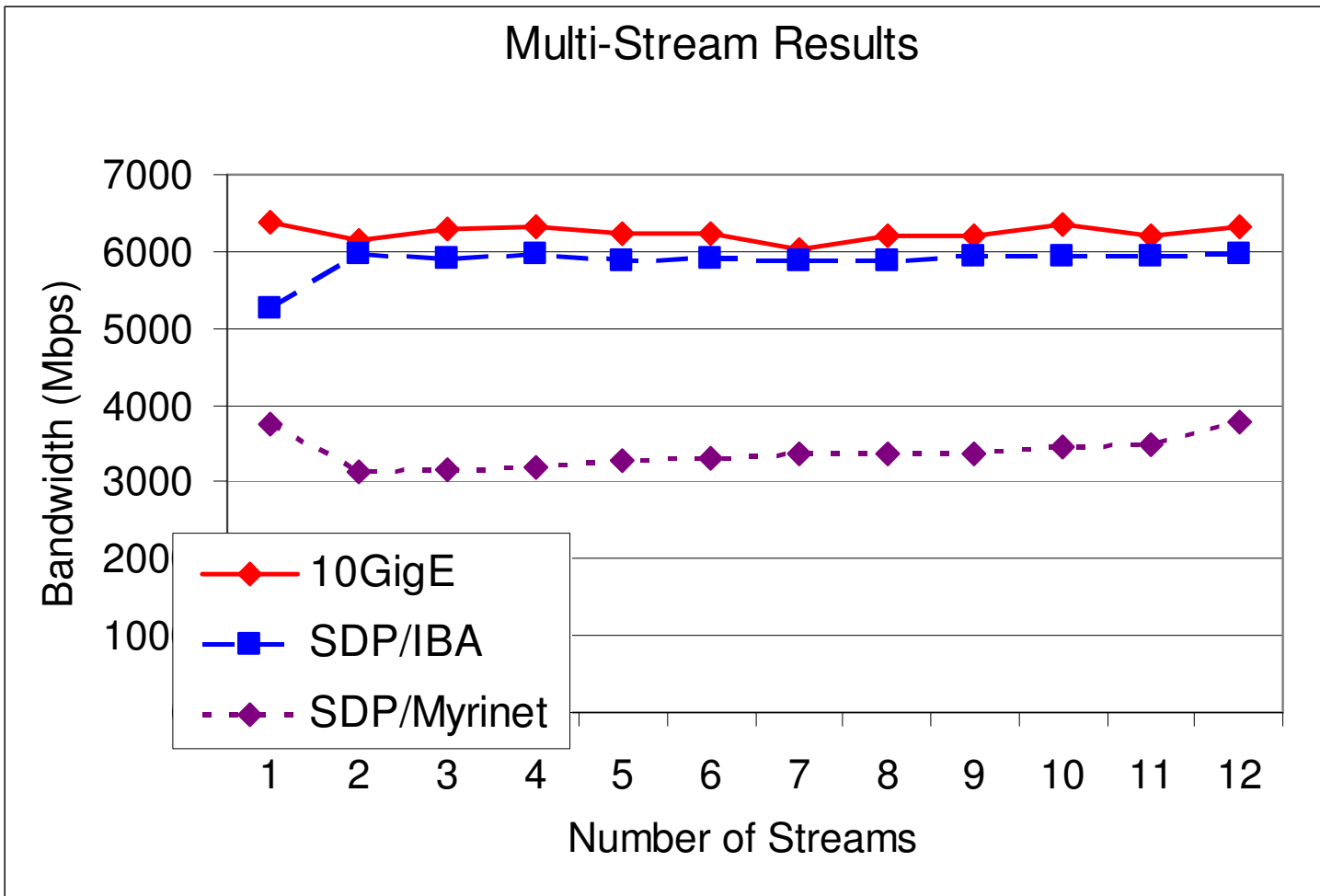
# OSU/LANL Benchmarks
# 10GbE TOE vs IB & Myrinet



Latency vs Message Size

**Source**: *Head-to-TOE Evaluation of High-Performance Sockets over Protocol Offload Engines* by DK Panda et al
**Test configuration**: 4-node cluster connected through 10GbE switch running single connection
**System configuration**: Dual 32-bit Intel Xeon 3.0GHz processors running Red Hat 9.0 Linux kernel 2.4.25smp

- Parallel Virtual File System (PVFS) – concurrent read/write performance

**PVFS Read Bandwidth**

**PVFS Write Bandwidth**



**Source:** *Head-to-TOE Evaluation of High-Performance Sockets over Protocol Offload Engines* by DK Panda et al
**Test configuration:** 4-node cluster connected through 10GbE switch running single connection
**System configuration:** Dual 32-bit Intel Xeon 3.0GHz processors running Red Hat 9.0 Linux kernel 2.4.25smp
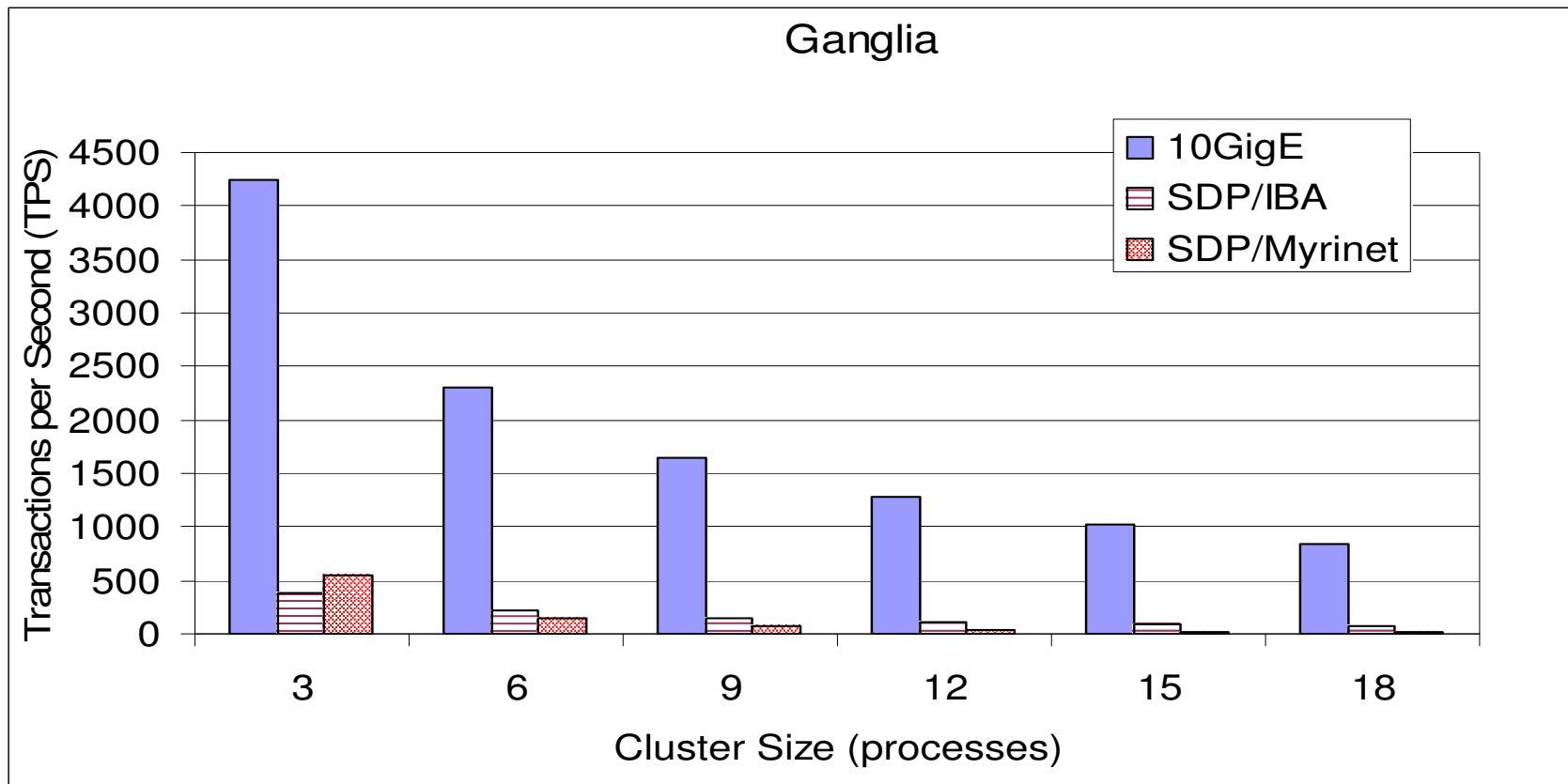
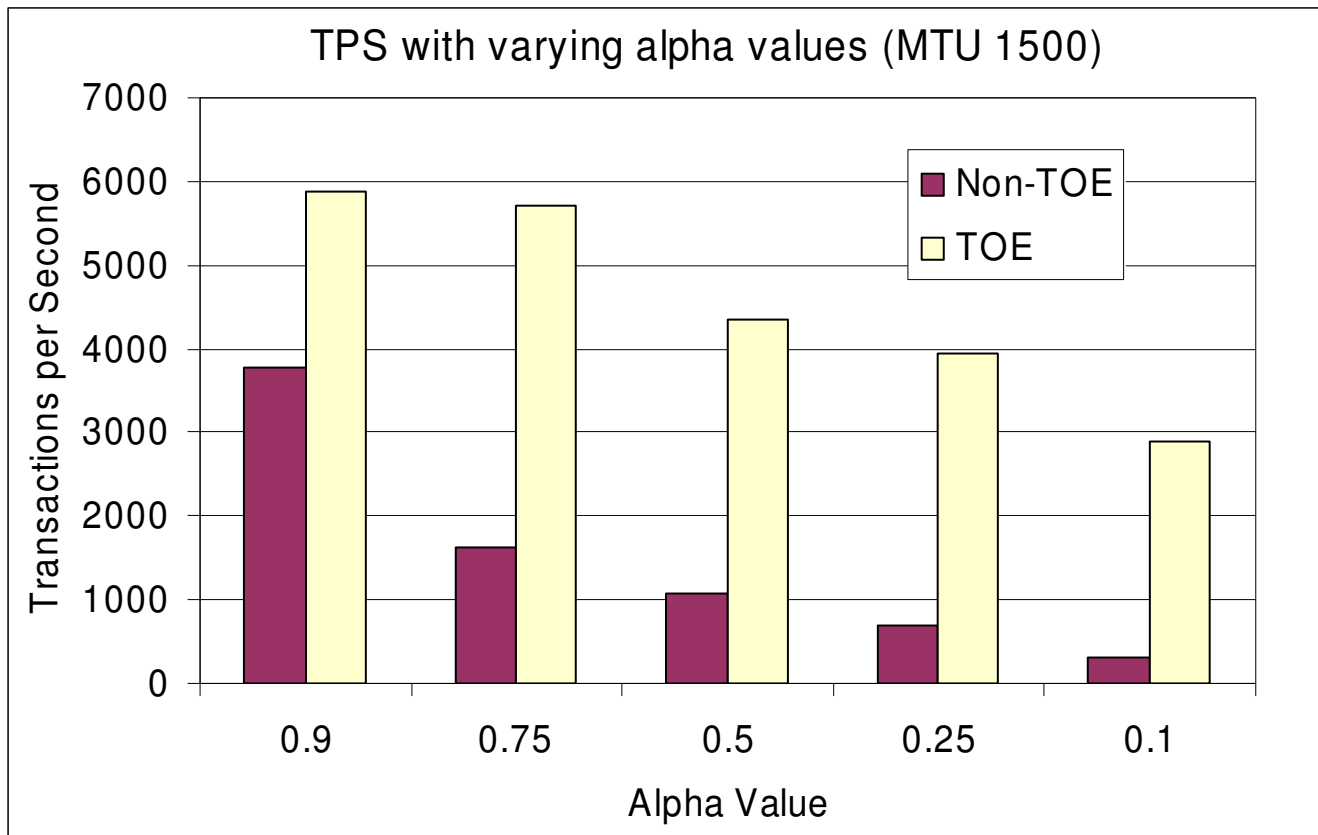# OSU/LANL Benchmarks 10GbE TOE vs IB & Myrinet



Source: *Head-to-TOE Evaluation of High-Performance Sockets over Protocol Offload Engines* by DK Panda et al
Test configuration: 4-node cluster connected through 10GbE switch running single connection
System configuration: Dual 32-bit Intel Xeon 3.0GHz processors running Red Hat 9.0 Linux kernel 2.4.25smp

- ## Ganglia Monitoring

# OSU/LANL Benchmarks 10GbE TOE vs 10GbE NIC

- Apache Web Server



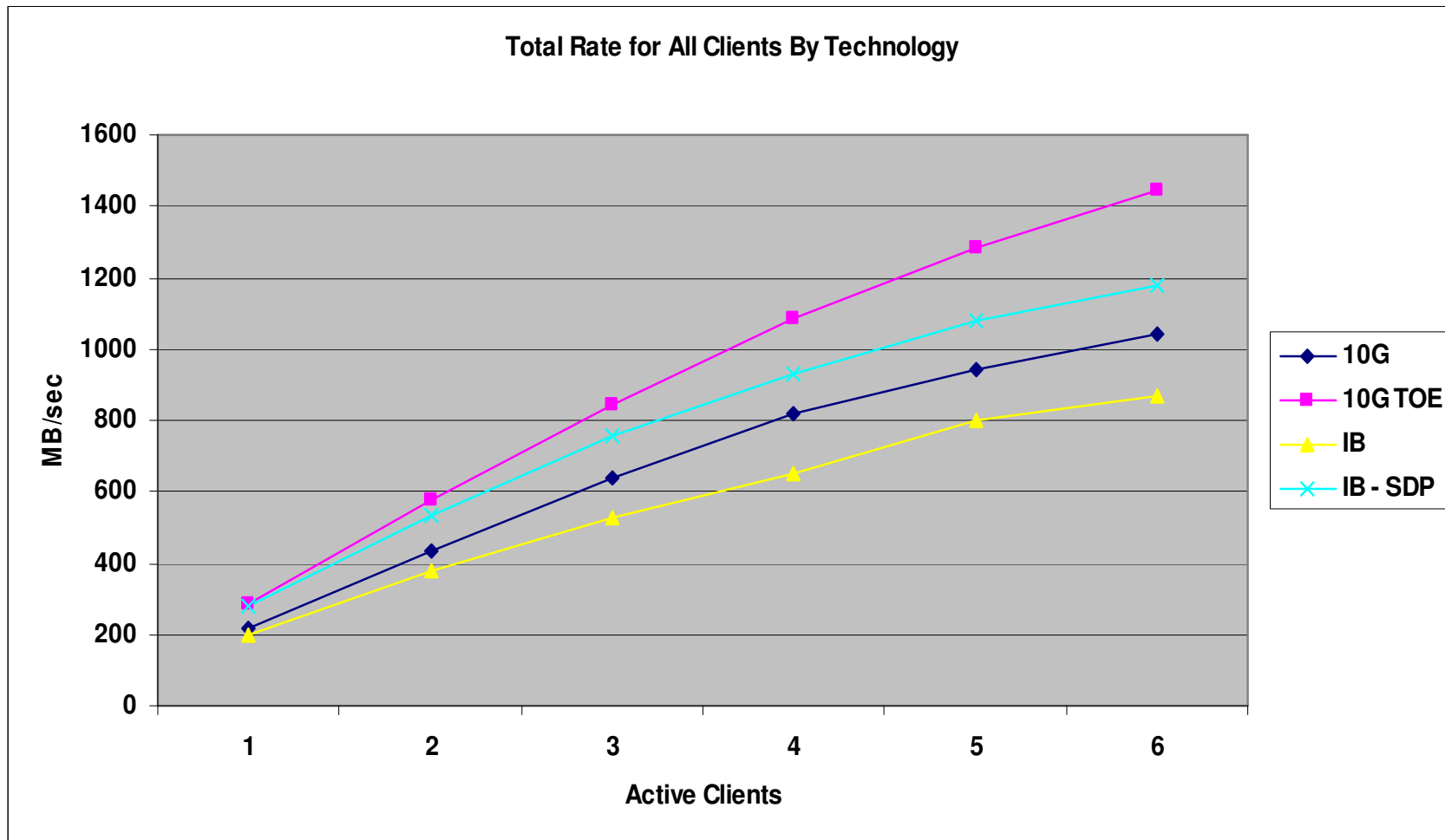TPS with varying alpha values (MTU 1500)

**Source:** *Performance Characterization of a 10-Gigabit Ethernet TOE* **by Wu Feng et al**
**Test configuration:** **4-node cluster connected through 10GbE switch running single connection**
**System configuration:** **Quad AMD Opteron 2.0GHz processors running Suse Linux with 2.6.6 stock kernel**
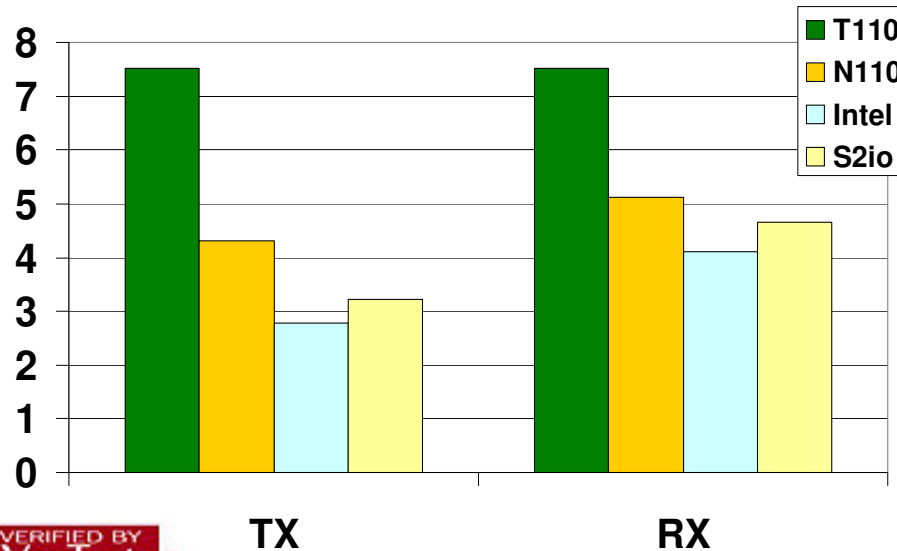
# Sandia Benchmarks
# 10GbE TOE vs IB & 10GbE NIC



**Total Rate for All Clients By Technology**

Legend:
- ◆ 10G
- ■ 10G TOE
- ▲ IB
- ✕ IB - SDP

Y-axis: MB/sec (0, 200, 400, 600, 800, 1000, 1200, 1400, 1600)
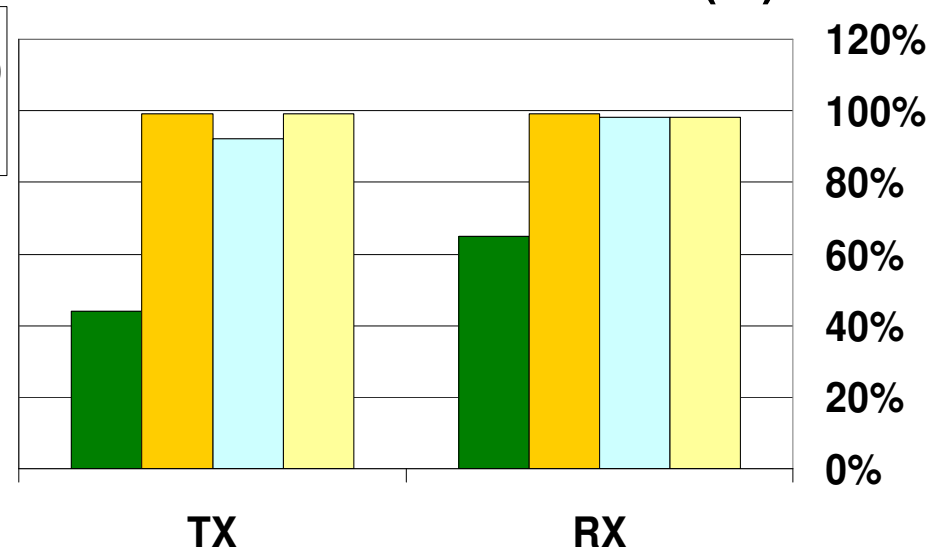
X-axis: Active Clients (1, 2, 3, 4, 5, 6)

**Source**: *Infniband and 10-Gigabit Ethernet for I/O in Cluster Computing* by Helen Chen et al
**Test configuration**: 8-node cluster connected through 10GbE switch running IOzone
**System configuration**: Dual AMD Opteron 2.2GHz processors running Linux kernel 2.4.25smp

# Chelsio Competitive Advantage

**Chelsio** Communications — *Accelerate*

## Higher Throughput (Gb/s)

## Lower CPU Utilization (%)

Legend:
- T110
- N110
- Intel
- S2io



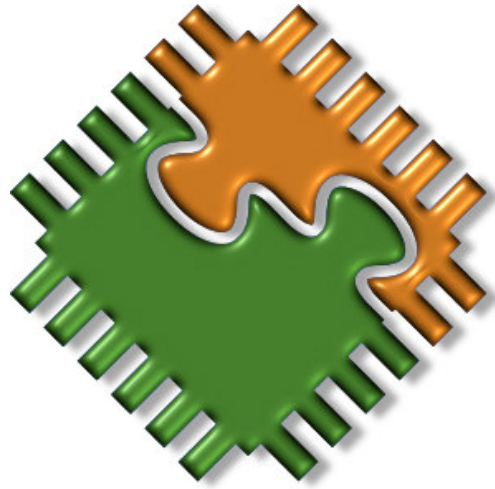VERIFIED BY VeriTest. Testing by the World's Leading Independent Lab

**Source:** Independently verified by VeriTest, Inc.
**Test tool:** netperf
**Test configuration:** 2 systems connected through 10GbE switch running single TCP channel with *1500-byte Ethernet frames*
**System configuration:** AMD Opteron 248 2.2GHz uniprocessor running Linux kernel 2.6.6
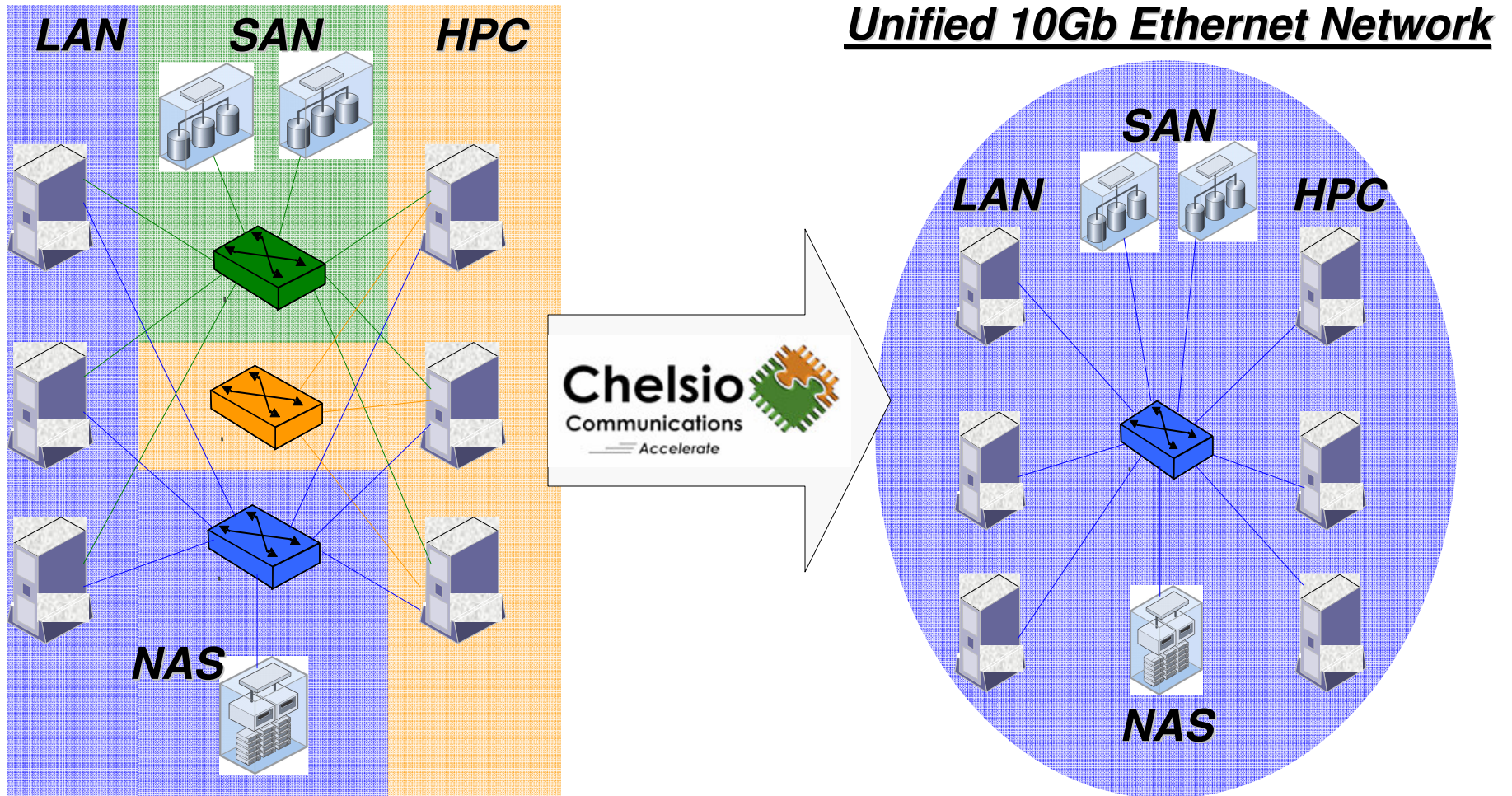
- T110 achieves 2x network throughput vs basic 10GbE NICs
- T110 utilizes only ½ x CPU resources vs basic 10GbE NICs
- <u>RESULT:</u> T110 delivers <u>4x performance efficiency</u> vs NICs

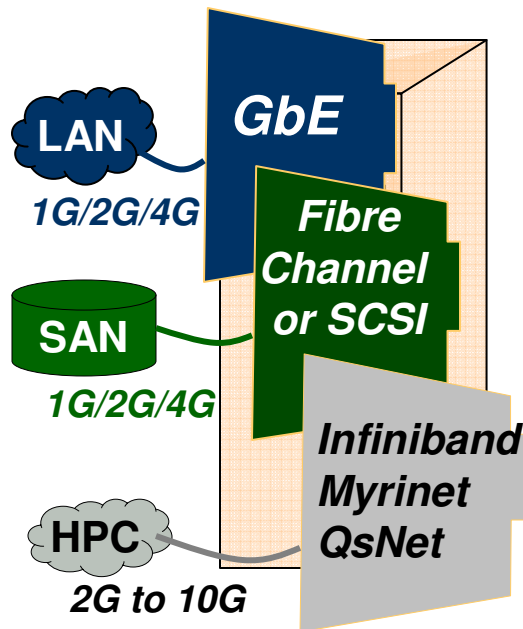# Chelsio Communications

Network Convergence & ULP Accelerations

# Network Fabric Convergence



LAN SAN HPC

Unified 10Gb Ethernet Network

SAN

LAN HPC

NAS

NAS

**Simplified network architecture – reduced operating costs**

# Simplified Server Architecture

**Chelsio Communications** — *Accelerate*

## Current Implementations

**GbE**

LAN
1G/2G/4G

**Fibre Channel or SCSI**

SAN
1G/2G/4G

**Infiniband Myrinet QsNet**

HPC
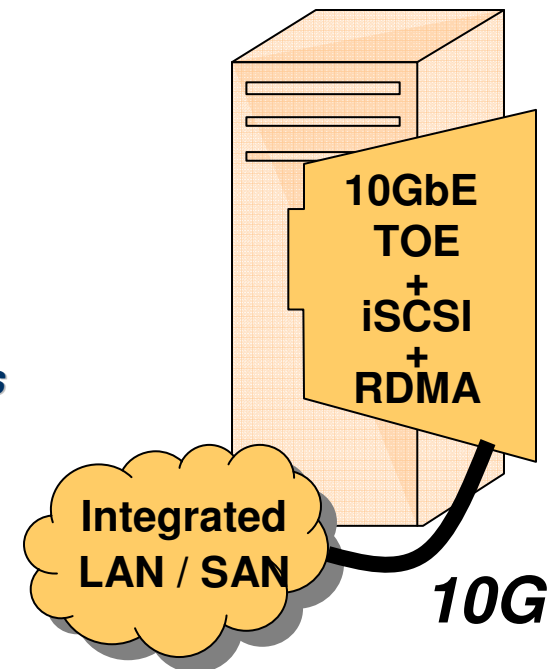2G to 10G

## Convergence Benefits

### Lower Total Cost of Ownership
- Improves CPU efficiency
- Minimizes software licenses
- Simplifies data center wiring
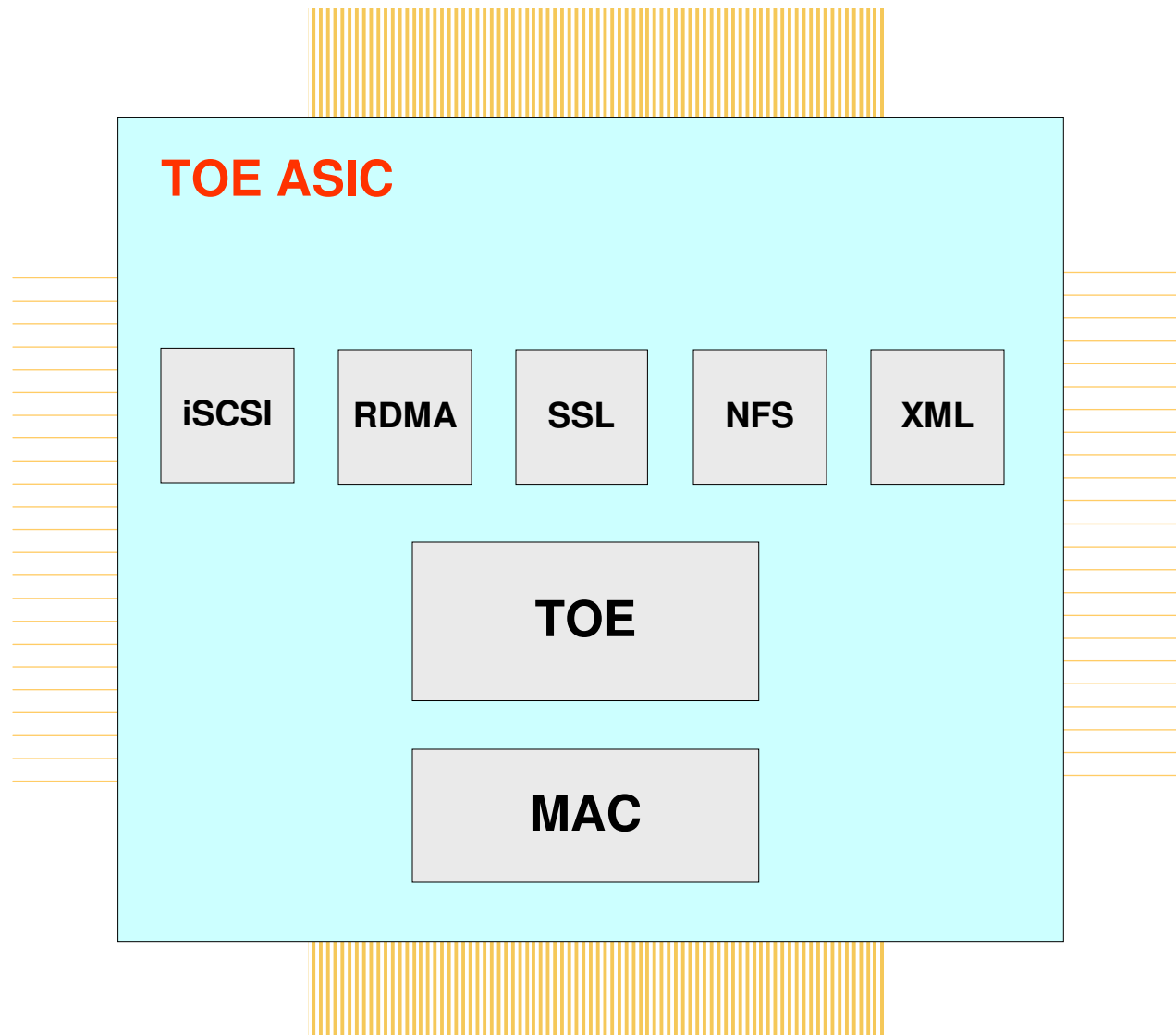- Leverages staffing skills & tools

### Higher Performance & New Apps
- Improves cluster performance
- Lowers application latency
- Faster backup and recovery
- Enables storage applications

## Chelsio's 10GbE Solution

**10GbE TOE + iSCSI + RDMA**

**Integrated LAN / SAN**

*10G*

## Improved performance – reduced operating costs

# TOE Enables ULP Acceleration

**TOE ASIC**

| iSCSI | RDMA | SSL | NFS | XML |

**TOE**

**MAC**

# Software Architecture

# Chelsio Communications

## iSCSI over TOE

# 10GbE Best Suited for Storage

**Chelsio**
Communications
*Accelerate*

SQLServer
Windows Server

iSCSI Software
Initiator

SQL Server

*10 GbE*

1/10 GbE
Switch

Chelsio T210
iSCSI HBA

iSCSI Software
Initiator

MS Exchange
Windows Server

*GbE*

iSCSI Storage Array

iSCSI Software
Initiator

Other Apps Server
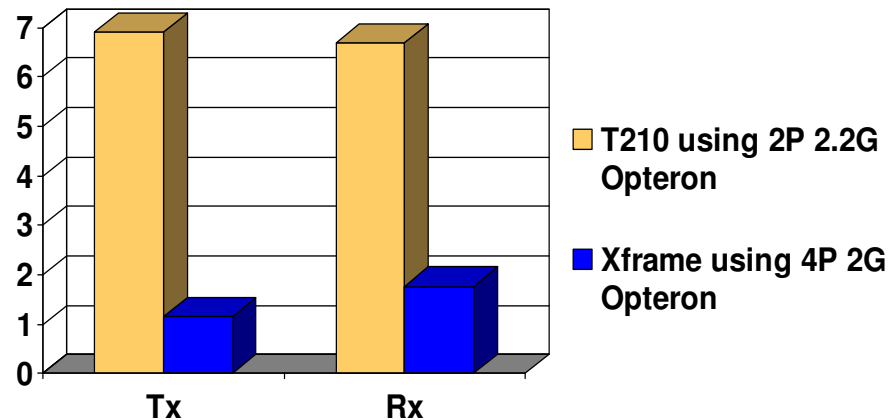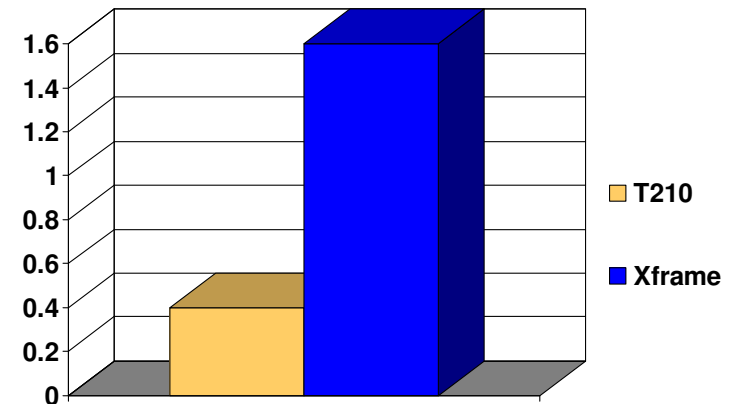
- Free client scalability
  - Free software initiators
  - Free GbE ports with ALL servers
  - HBAs not required for GbE initiators
  - GbE speed adequate for servers
- Similar Target costs to IB, FC
- No change to existing Apps
- Little change to infrastructure

38

# Chelsio iSCSI Performance vs basic 10GbE NICs

## Throughput
Gbps



- T210 using 2P 2.2G Opteron
- Xframe using 4P 2G Opteron

## Average CPU Utilization
x2.2 Ghz Opteron



- T210
- Xframe

**Sources:** T110 iSCSI Performance Analysis by Veritest and Xframe iSCSI Performance Analysis WP Published by Neterion
**Note:** Charts show performance at 4KB I/O size; iSCSI applications are transactional in nature using 2-4KB I/O sizes

- **T210 achieves 4x iSCSI network throughput vs basic 10GbE NICs**
- **T210 utilizes only ¼ x CPU resources vs basic 10GbE NICs**
- **RESULT:  T210 delivers 16x iSCSI performance efficiency vs NICs**

# T210 iSCSI Target Performance

| | Throughput | Avg. CPU |
|---|---|---|
| Read | 828MB (TOE) | 35% |
| Write | 857MB (TOE) | 46% |

| | IOPS | Avg. CPU |
|---|---|---|
| Read | 544k (TOE) | 88% |
| Write | 539k (TOE) | 99% |

Target Configuration:
- CPU: 2 x 2.2GHz Opteron
- SW:  Linux 2.4.25 and Chelsio Reference iSCSI stack
- IOmeter benchmark
- 28 GbE Microsoft Initiator to one 10GbE Target

# Chelsio Communications

**RDMA over TOE**

# The Benefit of RDMA
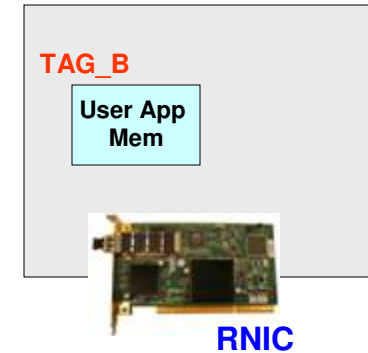
- **User space I/O**
- **OS bypass**
- **Direct Data Placement (DDP) and zero-copy**
- **Very low latency**
- **Very low CPU utilization**

# RDMA Operations

**Machine A**

TAG_A

User App Mem

RNIC

**Machine B**

TAG_B

User App Mem

RNIC

SEND ( "move 2MB from A to B, here is A's mem tag")

RDMA-READ ( "from TAG_A, off=0, to TAG_B, offset=0, len=64k")

**RNIC operation only. Host not gets Involved.**

RDMA-READ-RESP ( "here is the 64k data to TAG_B, offset=0")

. . .

SEND ( "done with the move 2MB from A to B")

# RDMA Protocol Stack

Chelsio Communications
*Accelerate*

| Oracle Parallel DB | NFS Over RDMA | iSER | ULP |

**RNIC HW**

| Layer | |
|---|---|
| **64k** | RDMAP |
| **64k** | DDP |
| **1.5k** | MPA |
| **1.5k** | TCP/IP |
| **1.5k** | Ethernet |

- RDMA Ops:
  RDMA Read, RDMA Read Response
  RDMA Write, Send

- ULP Message segmentation and reassembly
- Out-of-order placement
- In-order delivery

- Framing and CRC
- FPDU aligned with pkt, multiple FPDU in one pkt.
- Marker handling (Start from ISS, every 512B)
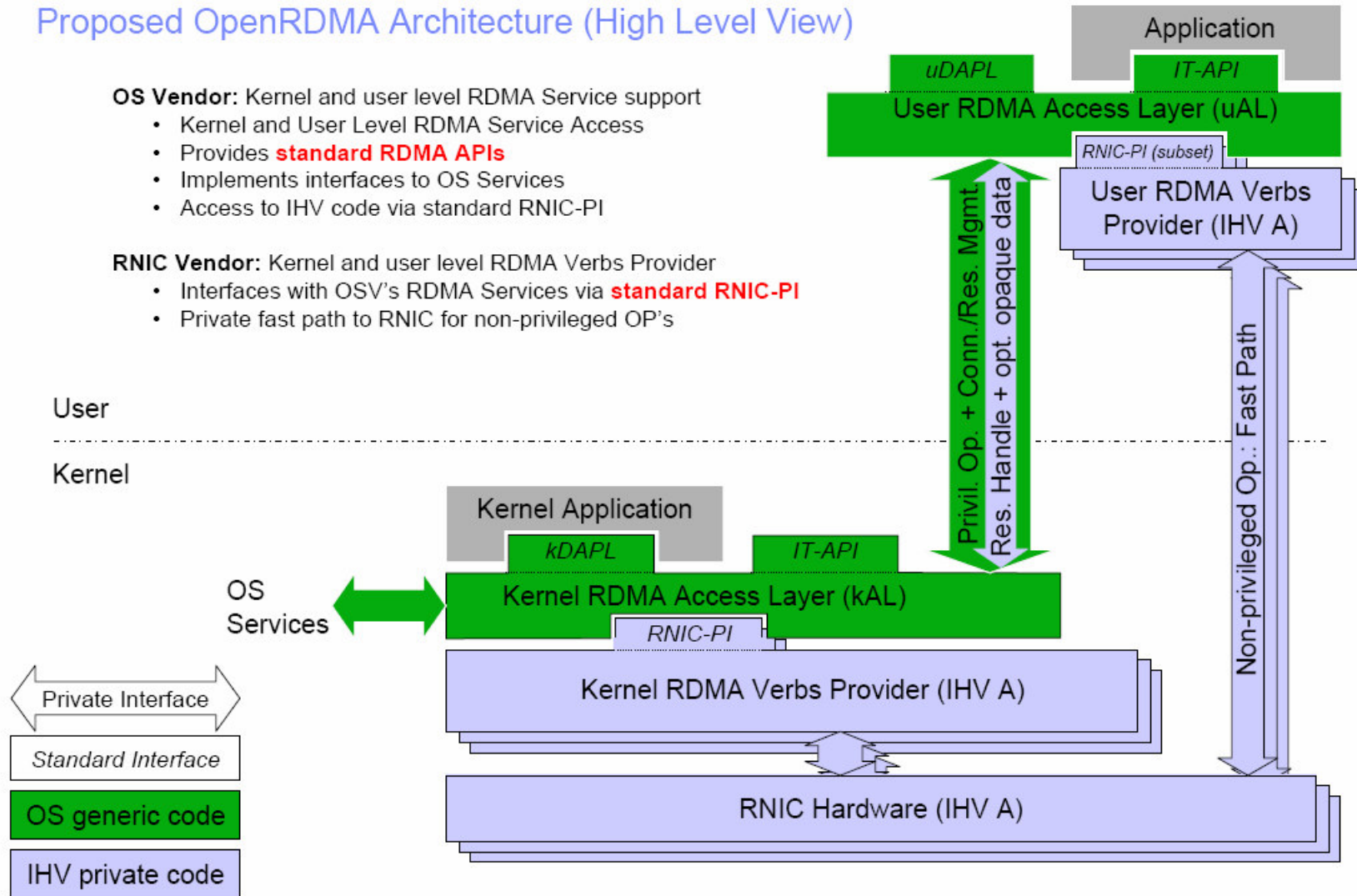
# RDMA Software Architecture



Proposed OpenRDMA Architecture (High Level View)

**OS Vendor:** Kernel and user level RDMA Service support
- Kernel and User Level RDMA Service Access
- Provides **standard RDMA APIs**
- Implements interfaces to OS Services
- Access to IHV code via standard RNIC-PI

**RNIC Vendor:** Kernel and user level RDMA Verbs Provider
- Interfaces with OSV's RDMA Services via **standard RNIC-PI**
- Private fast path to RNIC for non-privileged OP's

Application

uDAPL | IT-API
User RDMA Access Layer (uAL)
RNIC-PI (subset)
User RDMA Verbs Provider (IHV A)

Privil. Op. + Conn./Res. Mgmt. Res. Handle + opt. opaque data

Non-privileged Op.: Fast Path

User

Kernel

Kernel Application
kDAPL | IT-API
OS Services
Kernel RDMA Access Layer (kAL)
RNIC-PI
Kernel RDMA Verbs Provider (IHV A)

RNIC Hardware (IHV A)

Private Interface
Standard Interface
OS generic code
IHV private code

# Thank You!