# Creating an agile infrastructure with Virtualized I/O

Richard Croucher

May  2009

*Smart Infrastructure Solutions*

London ⚭ New York ⚭ Singapore

*www.citihub.com*

- Majority of current servers virtualized have been low utilization windows servers

- To push further and make Virtualised the default we need to be able to economically accommodate more I/O intensive workloads



- Issues which virtualized I/O help to address:

  - Lack of intrinsic I/O capacity management within existing Virtualized OS environments

  - Plethora of I/O cards and cables which need to purchased, fitted and maintained

  - Multiplicity of different I/O configurations on physical servers which limit mobility options

- ## Block level solutions

  - PVSCSI and VM Direct Path I/O

  - iSCSI

  - FibreChannel over Ethernet (FCoE)

  - InfiniBand

There are also NAS solutions, ie. Logical level solutions which we won't go into in this session.  These are commonly used today but  do not provide the same level of throughput as block level solutions and make diskless support non-standard (Windows) or complex (UNIX)

**Acitihub**™

- VMware continue to add capabilities to their flagship ESX platform

- Paravirtualized SCSI

  - Separate virtualised adapter used by Guest for high performance block level access over existing physical I/O paths

  - Intended for Guests demanding high performance I/O – still dependent on underlying physical connectivity

  - Requires Guest OS support – currently limited to WinSvr2003/2008 and RH Ent5

- VMDirect path I/O for storage

  - Maps a physical HBA to a single guest

  - One for one, no sharing

  - Limited to specific physical adapters

- Enables block level access to storage across TCP/IP networks

- Initiators supported by most OS environments out of the box, and can be used for guests via IP. Supported by VMware ESX.

- Free to use on servers with existing Ethernet ports and uses existing Ethernet infrastructure. Fully routable protocol. Relies on TCP for data integrity and in-order delivery

- Complexity of setup, particularly setting up and maintenance of the CHAP authentication

- No intrinsic boot support in Windows. Boot disk support dependent on INT13 provided with HW iSCSI cards and not accessible to Guests. Guests can attach to storage via iSCSI after booting or boot from an iSCSI disk presented to it via it's host.

- High TCP/IP overhead. Can be handled using TOE cards but drives cost up. Bigger problem on storage arrays (targets) where it limits maximum I/O throughput compared with other protocols.

- Storage vendors typically only target SMB and avoid Enterprise due to the overhead issues on arrays and increased CPU load on servers

- A server running iSCSI target and with FibreChannel can act as a gateway, however adding in resiliency increases complexity significantly.

- Initiative led by Cisco and others under ANSI T11.3  (see www.fcoe.com) to deliver FibreChannel protocol directly over Ethernet.

- Creates no TCP load on server and storage arrays

- Encapsulates FC frame with Header and SCSI command/data inside Ethernet data field. Simplifies de-encapsulation onto physical FC networks

- Requires guaranteed, in-order delivery demanding priority based flow control be added to the Ethernet standard

- All Ethernet components in path need to be FCoE compliant, i.e. Switches need to support Converged Enhanced Ethernet (CEE) to prevent problems.

- Relevant  proposed standards are:

    - IEEE 802.1Qbb Priority Flow control

    - IEEE 802.1Qaz Enhance Transmission Selection

    - IEEE 802.1Qau Congestion Notification

- Issues:

    - Standards are still draft, expect to be ratified late 2009 or early 2010

    - Non-routable, server and storage need to be in same subnet, fabric extenders create a STP free L2 fabric
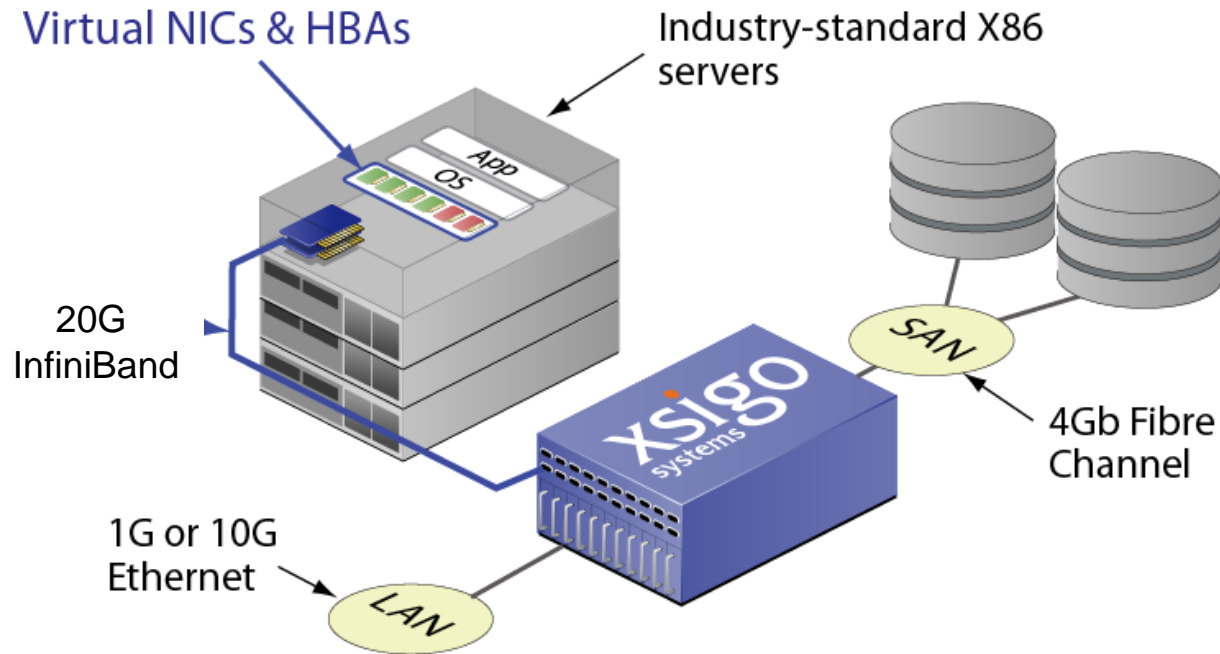
**acitihub**

## Converged Network Adapters

- Qlogic QL8100  2x 10Ge *note

- Emulex LP2100 – 2x 10Ge *note

- Brocade 1020 – 2x 10Ge

- ServerEngine – 2x 10Ge *note

## Converged Enhanced Ethernet Switches (pre standard products)

- Cisco Nexus 5010  with 20x 10Ge + 1 expansion slots

- Cisco Nexus 5020  with 40x 10Ge + 2 expansion slots

  - Expansion cards

    - 8x 4G FC, 6x 10Ge, 4x10Ge+ 4x 4G FC

- Brocade 8000 24x 10Ge+ 8x 8G FC

- Blade Network Technologies 24x 10Ge SFP+

*Note: iSCSI TOE included for backward compatibility

- Virtualized I/O over 20Gb/s InfiniBand

- Guests run vHBA and vNIC, with CIR/PIR based QoS, no need for InfiniBand drivers

- VMware ESX runs InfiniBand stack and allocates vHBA's/vNICs to guests

- I/O director provides physical connectivity into LAN/SAN environments

- 1 physical connector/cable (2 for resiliency) per physical server providing 40Gb/s

- Xsigo Management Server integrated into VMware Virtual Center

- Resold by Dell

**Hardware-based architecture**
- Fully non-blocking fabric
- 780 Gb/s aggregate bandwidth

**Custom silicon**
Line rate throughput

**24 Server ports**
- Expansion switch available for connection to hundreds of servers

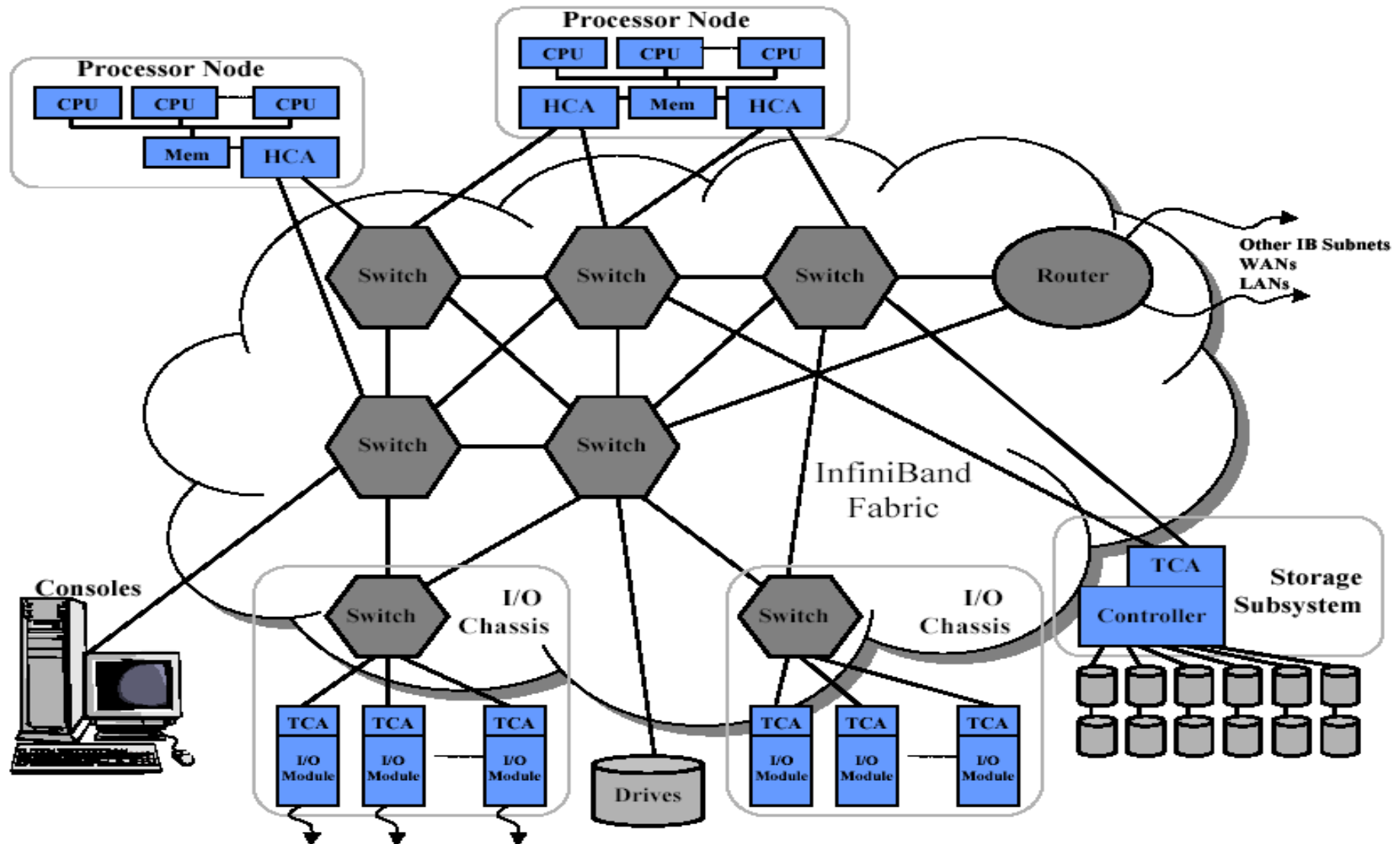**4U height**

**15 I/O module**

**4x 1Ge**

**10x 1Ge**

**1x 10Ge**

**2x 4G FC**

InfiniBand is a I/O protocol designed to provide high bandwidth, low-latency interconnect for clustering. It has been designed to offload CPU overhead by incorporating powerful RDMA (Remote Direct Memory Access) engines.

Mature standard with proven vendor interoperability

Historically deployed for HPC Grids but now entering Enterprise to support Virtualization rollouts

IB Switch functionality allows cut-through packets at wire speed with link and end-to-end data integrity

Implicit trunking across multiple serial lanes provides protocol independent speed improvements
1x = 2.5 Gbps  10B/8B = 2.0 Gbps data
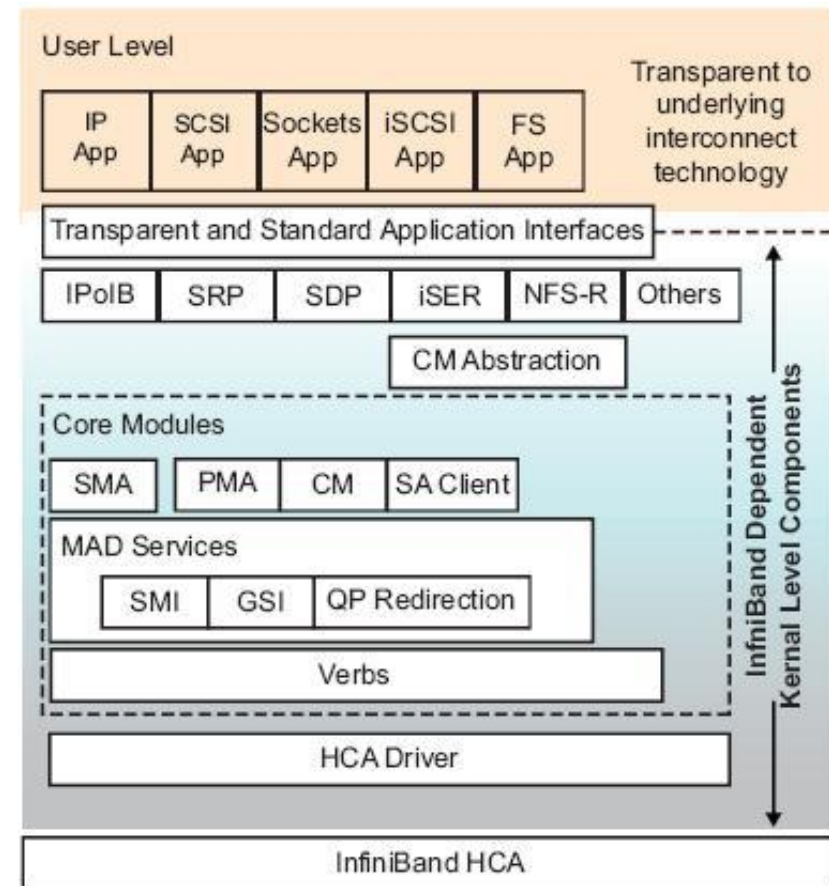4x = 10 Gbps  == 1G Byte/second data transfer
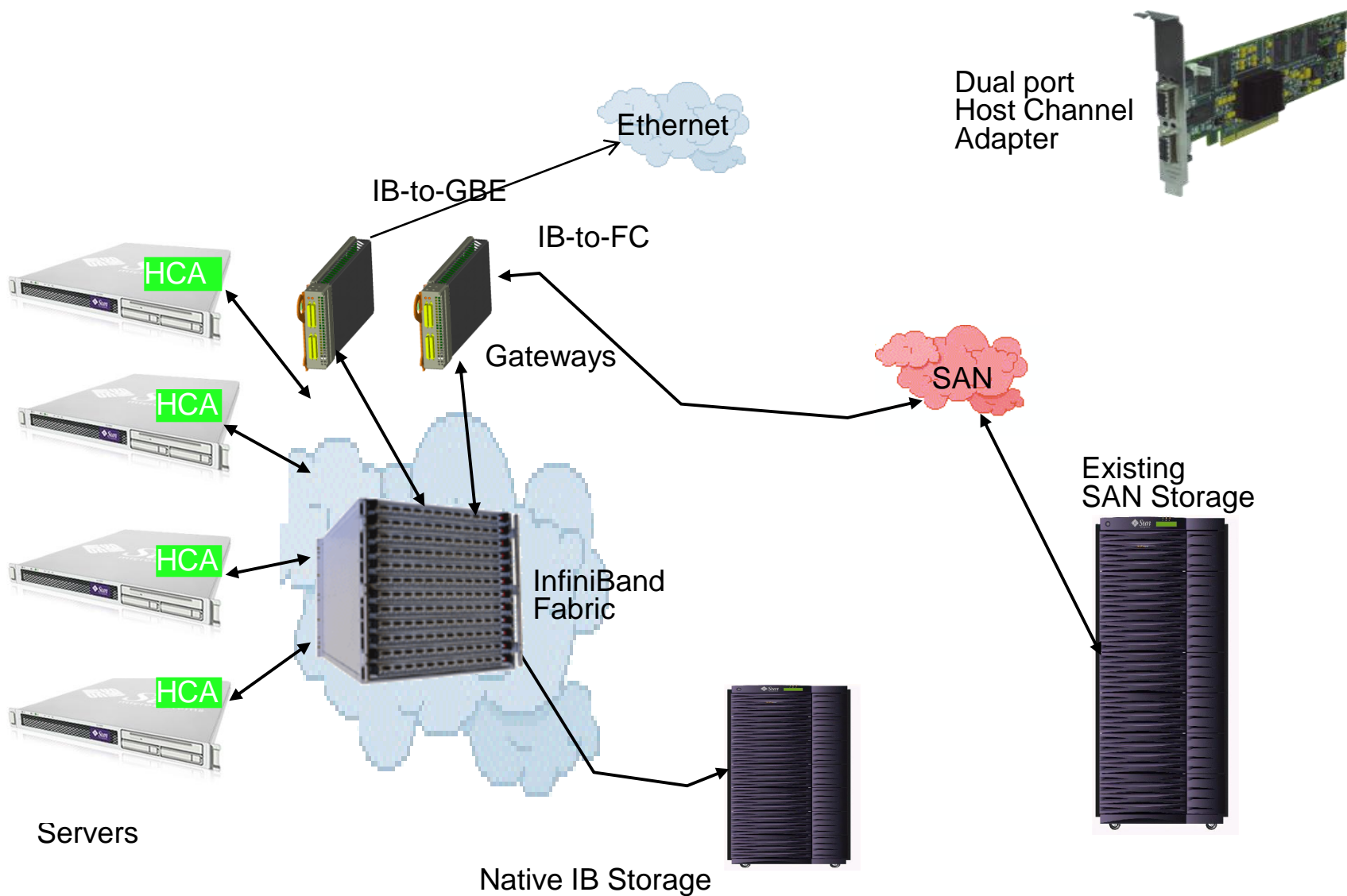12x = 30 Gbps == 3G Byte/second data transfer
Double Data Rate = 5 Gbps
4xDDR = 20 Gbps == 2G Byte/second data transfer
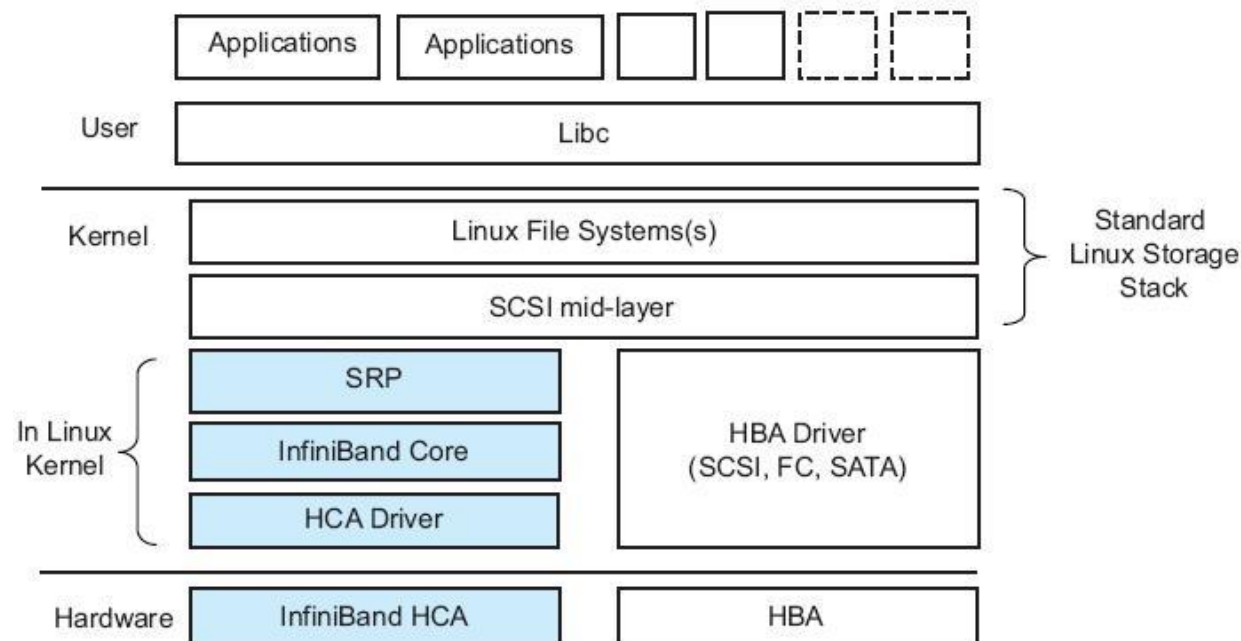Quad Data Rate = 10 Gbps
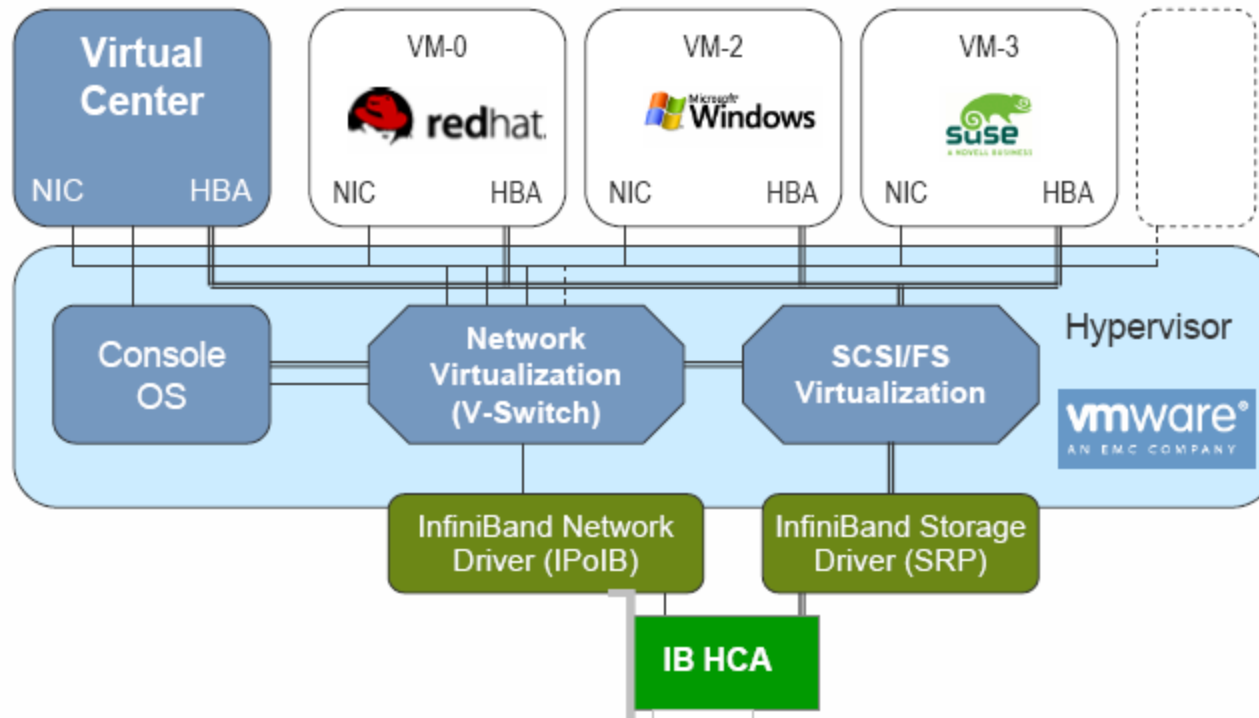4xQDR == 4G Byte/second transfer

| User Level | | | | | Transparent to underlying interconnect technology |
|---|---|---|---|---|---|
| IP App | SCSI App | Sockets App | iSCSI App | FS App | |

Transparent and Standard Application Interfaces

| IPoIB | SRP | SDP | iSER | NFS-R | Others |
|---|---|---|---|---|---|

CM Abstraction

**Core Modules**

| SMA | PMA | CM | SA Client |
|---|---|---|---|

MAD Services

| SMI | GSI | QP Redirection |
|---|---|---|

Verbs

HCA Driver

InfiniBand HCA

*InfiniBand Dependent Kernal Level Components*

Dual port
Host Channel
Adapter

Ethernet

IB-to-GBE

IB-to-FC

HCA

HCA

HCA

HCA

Gateways

SAN

Existing
SAN Storage

InfiniBand
Fabric

Servers

Native IB Storage

# InfiniBand Ethernet   and TCP/IP Integration

- IP over IB

  – Included with all InfiniBand implementations

  – Runs full TCP/IP stack including TCP/UDP and multicast

  – Connects InfiniBand attached nodes or through gateway to Ethernet

  – Does not leverage RDMA or support raw Ethernet

- vNIC

  – Layer 2 interface support raw Ethernet

  – Connects via compatible gateway to Ethernet

  – Does not leverage RDMA

- Sockets Direct Protocol

  – Bypasses TCP stack/overhead

  – Supports existing TCP socket applications

  – Preload libraries avoid recompilation

  – Connects InfiniBand attached nodes

- SCSI RDMA Protocol  (SRP)

  – Defined by ANSI T.10

  – Bypasses TCP/IP stack

  – All InfiniBand switch vendors offer gateways to FibreChannel

  – Driver support for Linux, Solaris, Windows, VMware ESX

- iSCSI extension for RDMA (iSER)

  – Initially defined by Voltaire, now IETF draft

  – Now open sourced and accepted by OpenIB into Linux implementation

  – Bypasses TCP/IP stack

- vHBA

  – Non standardized, Gateway specific, available for Qlogic, Mellanox and Xsigo

**Acitihub**



- Provides block level support

- Interposes below existing SCSI driver, just appears as a block level device

- Uses RDMA, requiring less CPU than FibreChannel

- Supported by VMware ESX, Windows, Linux, Solaris

- Host can connect through SRP and present block level device to Guest

InfiniBand: Enhancing Virtualization ROI with new Data Center Efficiencies, Sujal Das, Mellanox

**Ocitihub**



InfiniBand: Enhancing Virtualization ROI with new Data Center Efficiencies, Sujal Das, Mellanox

**Ocitihub**



128KB Read benchmarks from four VMs

## Same as four dedicated 4Gb/s FC HBAs

**VMWORLD** 2006

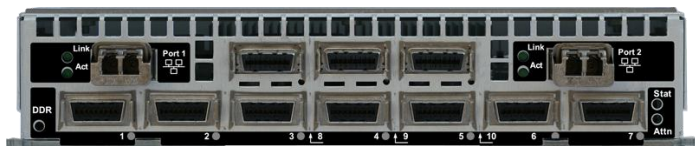InfiniBand: Enhancing Virtualization ROI with new Data Center Efficiencies, Sujal Das, Mellanox

**Acitihub™**

## Fibre Channel Leaf Module

## Ethernet Leaf Module

Fibre Channel VIC   (vHBA)
- 10-port  20G InfiniBand
- 8-port 1/2/4 Gb FC

10 Gigabit Ethernet VIC (vNIC)
- 10 port 20G InfiniBand
- 2-port 10 GbE

- Modules fit inside Qlogics director class InfiniBand switches

- Directors for 4,8,12, 24 modules

- Also 12 port 20G InfiniBand modules

**QLOGIC®**

- Latest range is 36 – 864x 40G InfiniBand switches

- 40G gateway planned for 2010

**citihub**

BridgeX BX4000
4x 40G InfiniBand
  12x 10Ge (vNIC and IPoIB)
  Or
  16x 8G FC (vHBA)

**Mellanox**®
**TECHNOLOGIES**

ConnectX
2x 40G IB or 2x 10Ge

MTS3600
32x 40G InfiniBand

10Ge ports FCoE (draft)  compatible
Full driver stack for VMware, Linux and Windows
Boot over InfiniBand for diskless host servers

**SR4G**



FibreChannel gateway
4x 4G FC + 2x 20G InfiniBand
1U server profile
iSER + iSCSI support

**sRB-20210G**



Ethernet gateway
2x 10Ge + 22x 20G InfiniBand
Module option for Director Class switches
Directors for 4, 6, 12 modules
Also
24 port 20G InfiniBand modules
32 port 40G standalone switch

# Voltaire GridVision Virtualization Management



Automatically discover physical objects

**Switches, routers, storage, physical servers, …**
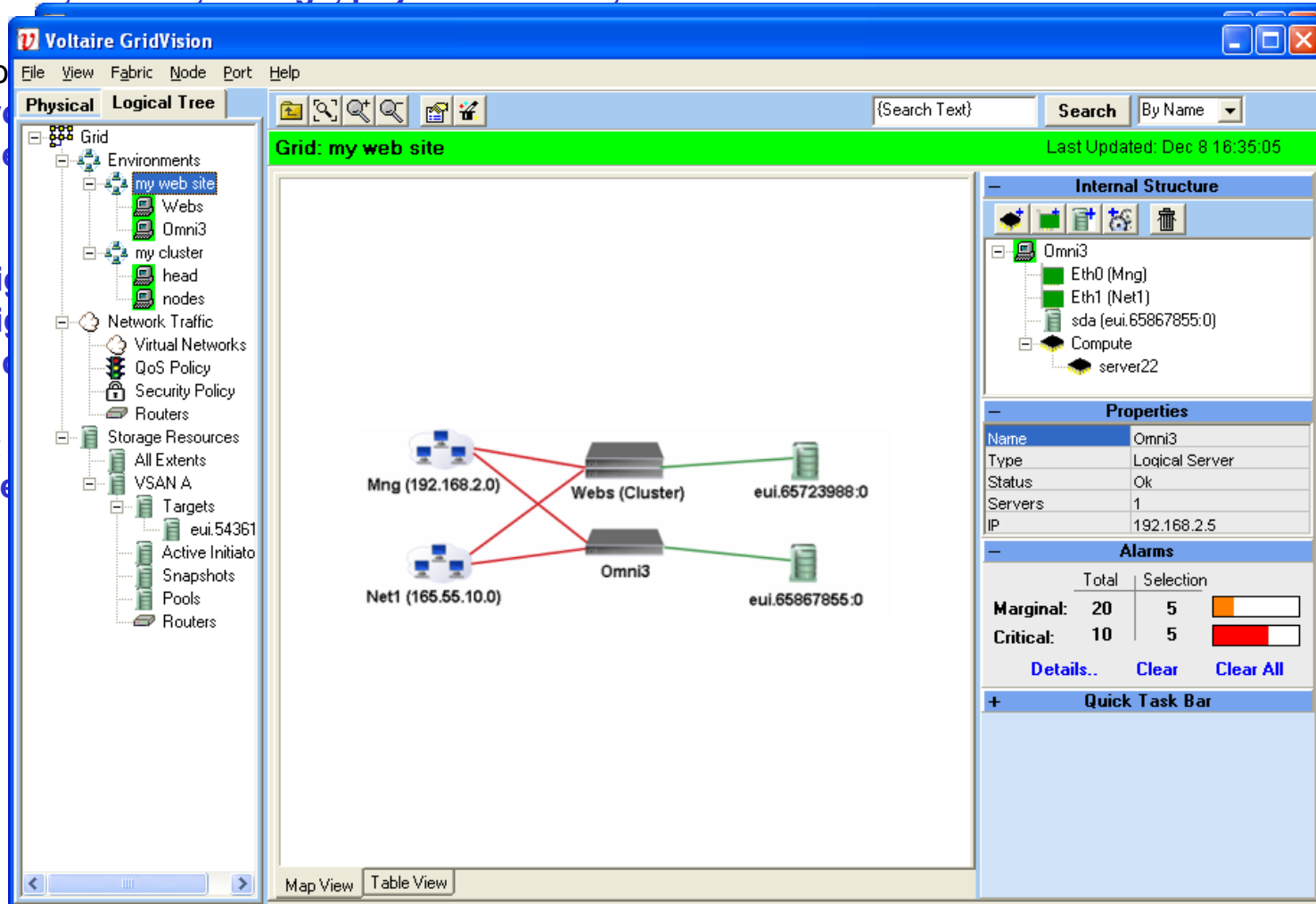
Define lo…

**Envir…**

**Define…**

Bind the …

**Config…**

**Config…**

**Provi…**

Abstract …

**Single…**

(c) – copyright Citihub, 2009

- Can use low profile servers since only one dual port PCIe card provides up to (80Gb/s)

- InfiniBand drivers to Manager only

- Guests use vNICs, vHBA's

- Connect multiple racks into a single virtual I/O cluster

  - Leverages high bandwidth

  - Shares cost of gateways across multiple servers particularly for resilient configurations

| | Traditional FibreChannel (assume 4 guests) |
|---|---|
| 2x server port (>=20G) | $ 5,200 |
| 2x switch ports for server | $ 4,000 |
| 4G B/W slice of FC gateway | $ - |
| FC switch ports (per guest) | $ - |
| 4G B/W slice of Ethernet gateway | $ - |
| Total (Physical Server) | $ 9,200 |
| Cost per Guest | $ 2,300 |

Excludes cost of server and storage

Your prices may vary

**A citihub**

| | iSCSI | |
|---|---|---|
| | **FC Storage** | **iSCSI storage** |
| 2x server port (>=20G) | $ 1,150 | $ 1,150 |
| 2x switch ports for server | $ 850 | $ 850 |
| 4G B/W slice of FC gateway | $ 1,000 | $ - |
| FC switch ports | $ 1,000 | $ - |
| 4G B/W slice of Ethernet gateway | | |
| Total (Physical Server) | $ 4,000 | $ 2,000 |
| Cost per Guest | $ 1,000 | $ 500 |

Excludes cost of server and storage
note increased storage cost of iSCSI due to TCP overhead

Your prices may vary

# Conclusions

**citihub**™

- I/O virtualization simplifies virtual server deployment and enables a wider range of applications to be supported

- iSCSI is low end solution due to  server CPU and storage array TCP/IP overheads

- FCoE is the strategic solution but:

  – ALL products are still pre-standard – demand free upgrade from suppliers

  – Recommend separate  storage Ethernet infrastructure

  – Still  missing routing and long distance solutions

  – Real savings only after storage migrated from FibreChannel to direct Ethernet

- InfiniBand is available today, well supported, lower cost, higher bandwidth server attach than both FibreChannel or 10Ge

  – Tactical solution until Ethernet standards settle.  May need to wait for next generation 40/100G Ethernet

  – Routing and long distance solutions available today for InfiniBand

  – Can be deployed without adding InfiniBand drivers to Guests

**Richard.Croucher@Citihub.com**

Independent advice for Data Centre Infrastructure
London – New York - Singapore

(c) – copyright Citihub, 2009