



# Cisco InfiniBand and Server Virtualization



**Bryan Bradley**

**Sr. Systems Engineer**

**Server Virtualization and InfiniBand**

# Agenda

- Introduction to InfiniBand Technology
- RDMA and InfiniBand Upper Layer Protocols
- High Performance Computing (HPC)
- Designing an InfiniBand Fabric
- Server Virtualization and VFrame
- Examples and Case Studies

# Server Networking Challenges

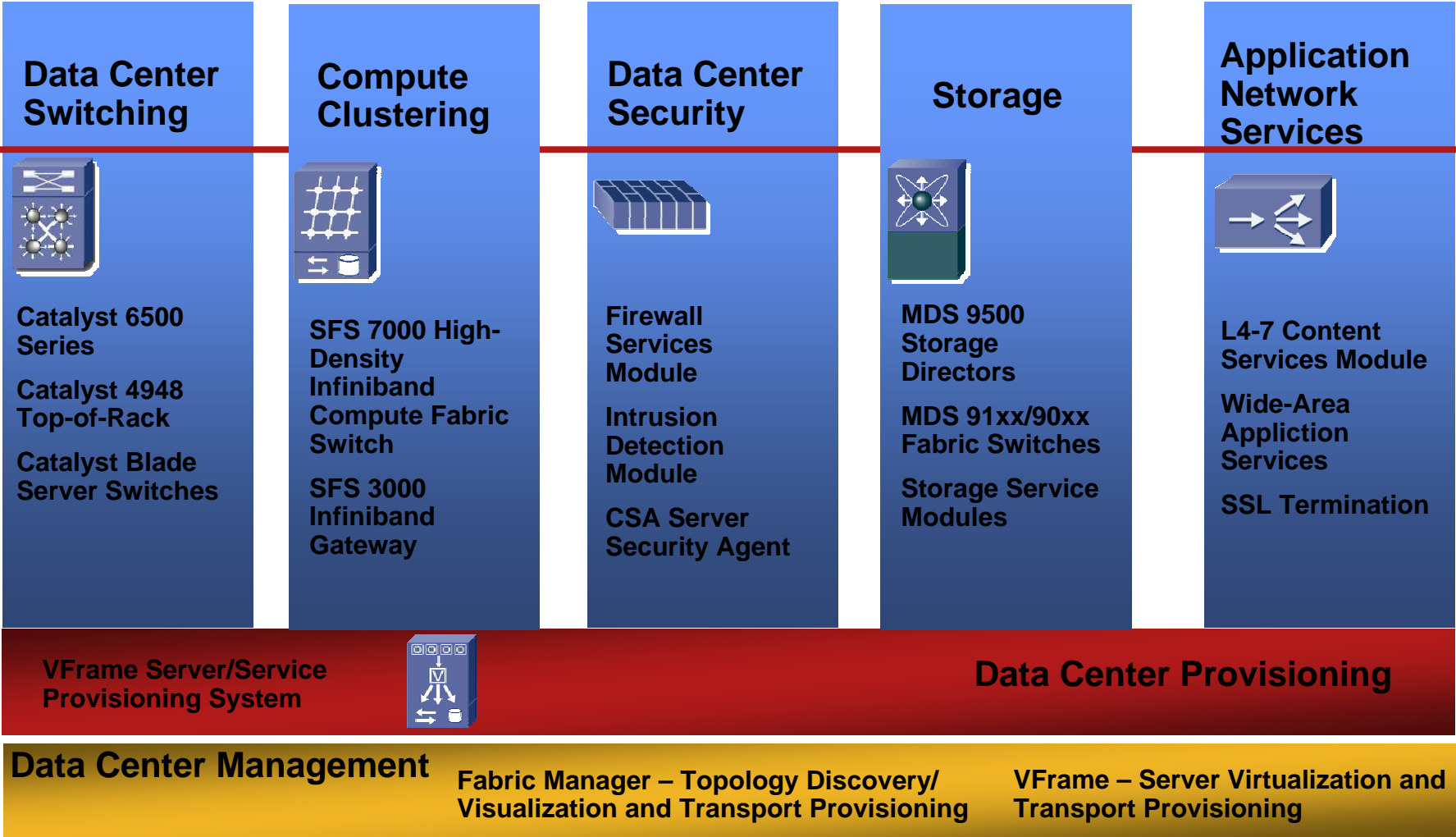
## Growing Bandwidth Demands

- **Highly distributed X86 architectures with increased I/O requirements**
- **Increased Inter-process communications between servers**
- **New X86 systems with channelized I/O**
- **Increased server utilization with server PCI-express / hypervisor and utility grid technologies**
- **Multi-CPU's and multi-core on exponential performance curves**

## Growing Network Complexities

- **Average server has 4-8 network interfaces including NICs, HBA's and HCA's**
- **Complex wiring with increased services requirements**
- **Strong division of labor between server, storage and networking teams**
- **Blurring of operational divisions within blade chassis**

# Cisco Data Center Products

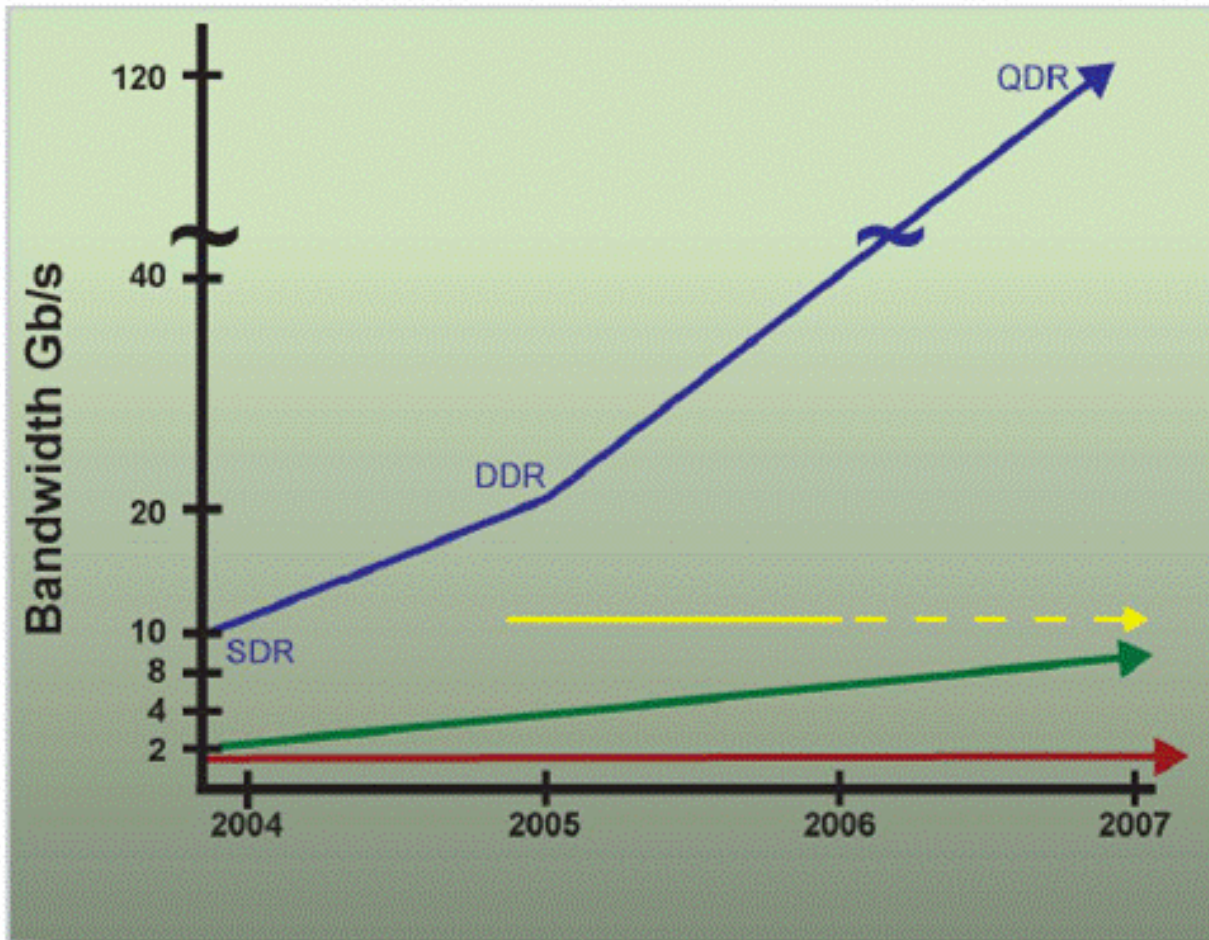


# InfiniBand Overview

- **Standards-based interconnect (since 2001)**
- **Channelized, connection-based interconnect optimized for high-performance computing**
- **Supports server and storage attachments**
- **Bandwidth capabilities (SDR/DDR)**
  - 1x—2.5/5 Gbps: 2/4 Gbps actual data rate (base rate for InfiniBand)**
  - 4x—10/20 Gbps: 8/16 Gbps actual data rate**
  - 12x—30/60 Gbps: 24/28 Gbps actual data rate**
- **Built-in RDMA as core capability for inter-CPU communication**



# InfiniBand Bandwidth



- InfiniBand
- 10GigE, 10G iSCSI, Proprietary
- Fibre Channel\*
- GigE, iSCSI









\*FCIA estimates 2007/8 for 8Gb/s FC

# Where to use Infiniband?

- Cost-Effective 10 Gbps to the Server/Desktop
  - 10 Gigabit Ethernet switches very costly
  - 10 Gbps Infiniband is established is cost effective
- Improve Time to Market (Low Latency/Application Acceleration)
  - Financial Markets
  - Manufacturing
  - Oil & Gas Industry
  - Life Sciences (Pharma/Biotech)
  - Higher Education
- I/O Consolidation
  - Single interface for network and storage
- Data Center Virtualization
  - VFrame 3.x (currently shipping)
  - VFrame Data Center 1.1 (available Fall 2006)

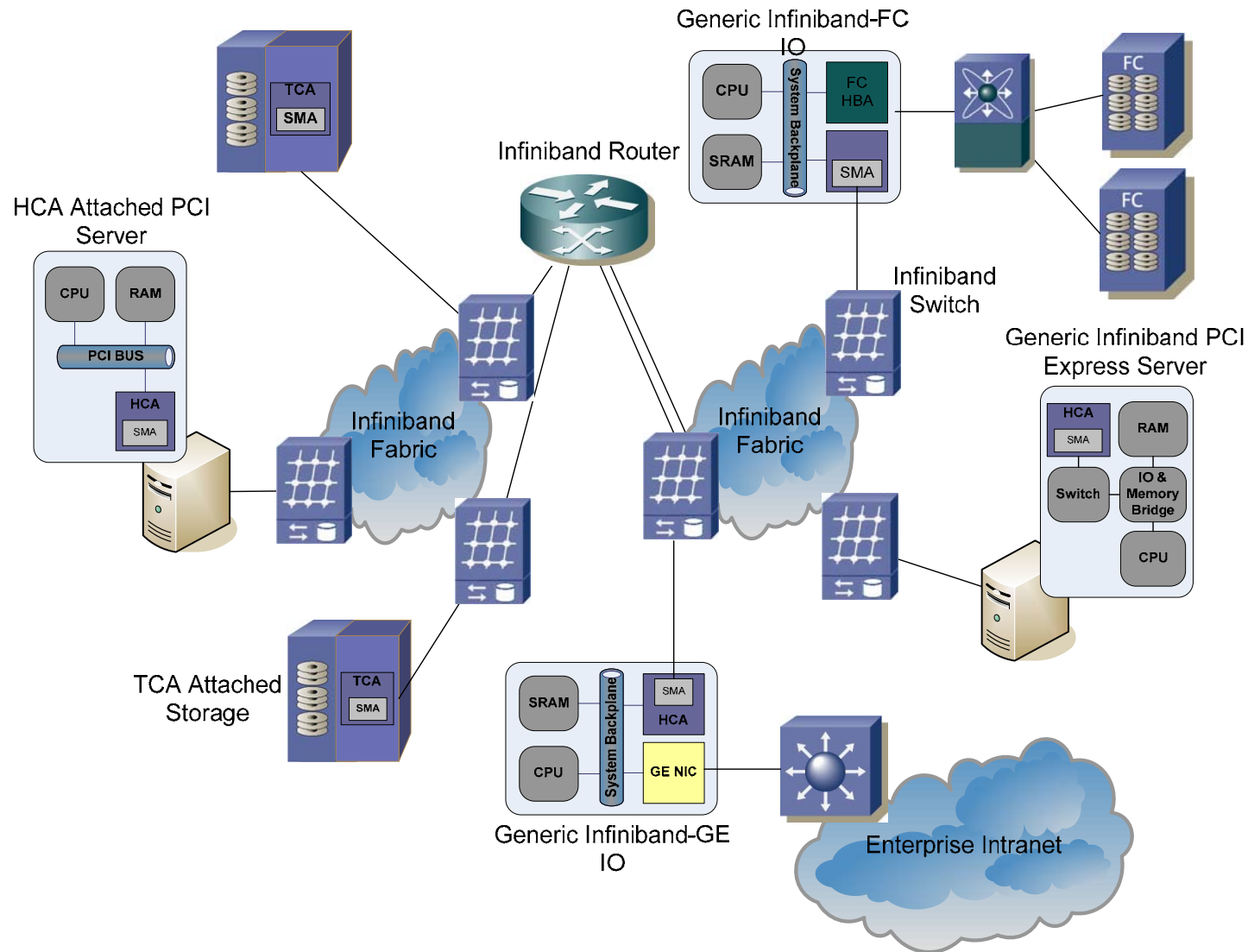
# Distributed Applications & Systems Networking

- **High-Performance Computing and distributed application processing**
  - Extract meaningful trends and business metrics
  - Reduces business risk
  - Reduce R&D costs
- **Low-latency data distribution between application systems**
  - Process & Systems Integration
  - Autonomic trading
- **Distributed Databases & High-volume Storage**
- **Distributed Computing requires high-performance Networking to reduce inter-process latency**

High Performance Computing	
	Computational Fluid Dynamics
	Finite Element Analysis
	Computational chemistry
	Monte Carlo simulations
Distributed Systems	
	Market Data Distribution
	Application Integration
Distributed Database & Application Tier	
	



# Architecture of an InfiniBand Solution



# Sources of overhead in Datacenter Servers

Sources of Overhead in Server Networking	CPU Overhead
Transport Processing	40%
Intermediate Buffer Copying	20%
Application Context Switches	40%

## Solutions for Overhead in Server Networking

### Transport Offload Engine (TOE)

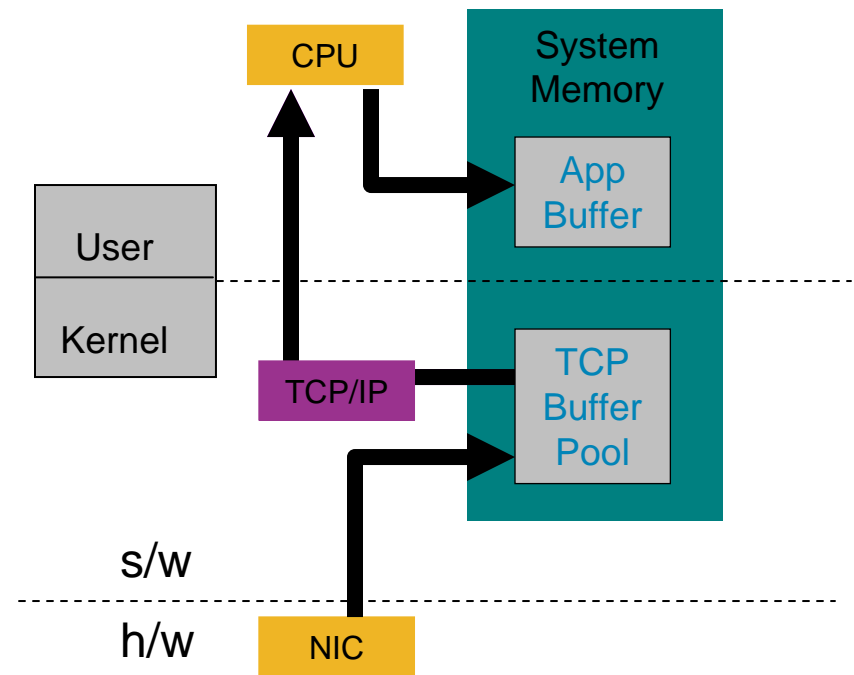
- Moves Transport processor cycles to the NIC
- Moves TCP/IP protocol stack buffer copies from system memory to the NIC memory

### RDMA

- Eliminates intermediate and application buffer copies (memory bandwidth consumption)

### Kernel Bypass – direct user-level access to hardware

- Dramatically reduces application context switches



# RDMA: Remote Direct Memory Access

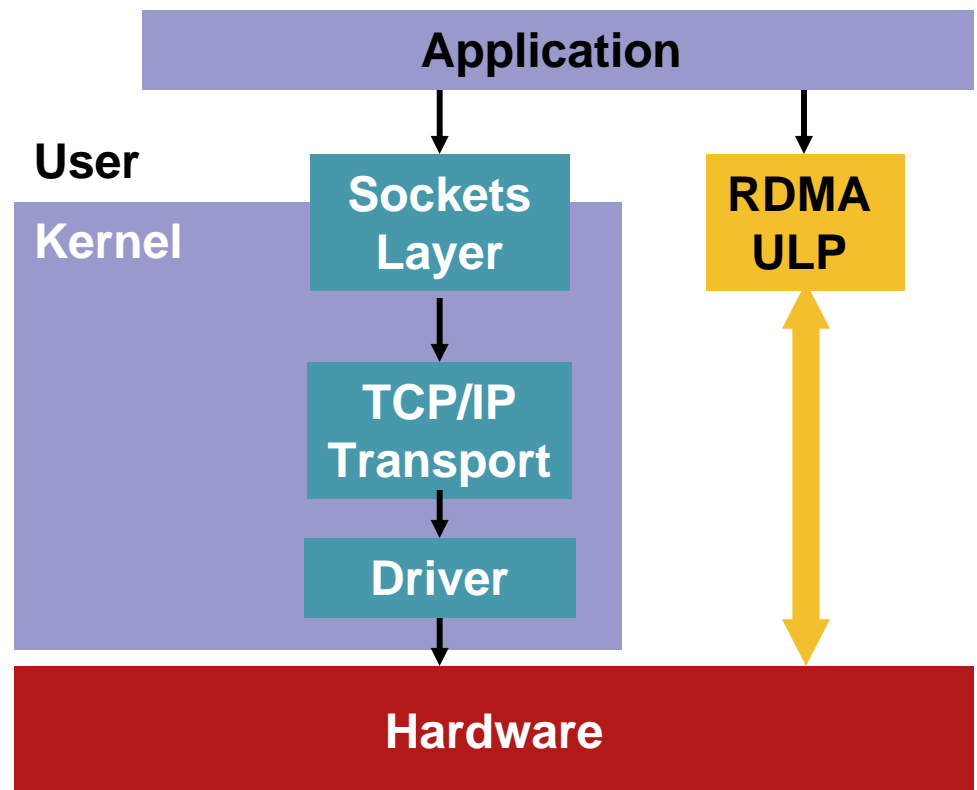
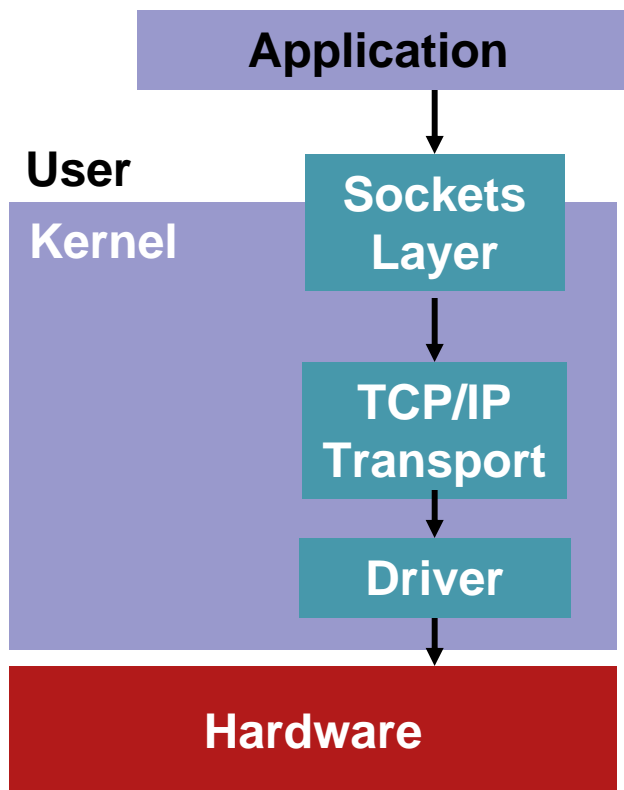
- RDMA enables data to be moved between application (user) memory space without CPU intervention
- RDMA is transport agnostic
  - InfiniBand has native support for RDMA
  - Ethernet RDMA and TOE NICs can support RDMA
- RDMA can significantly increase application and transport performance
  - Kernel bypass
  - Zero copy data transfer
  - No CPU intervention

# Remote Direct Memory Access (RDMA) and Kernel Bypass

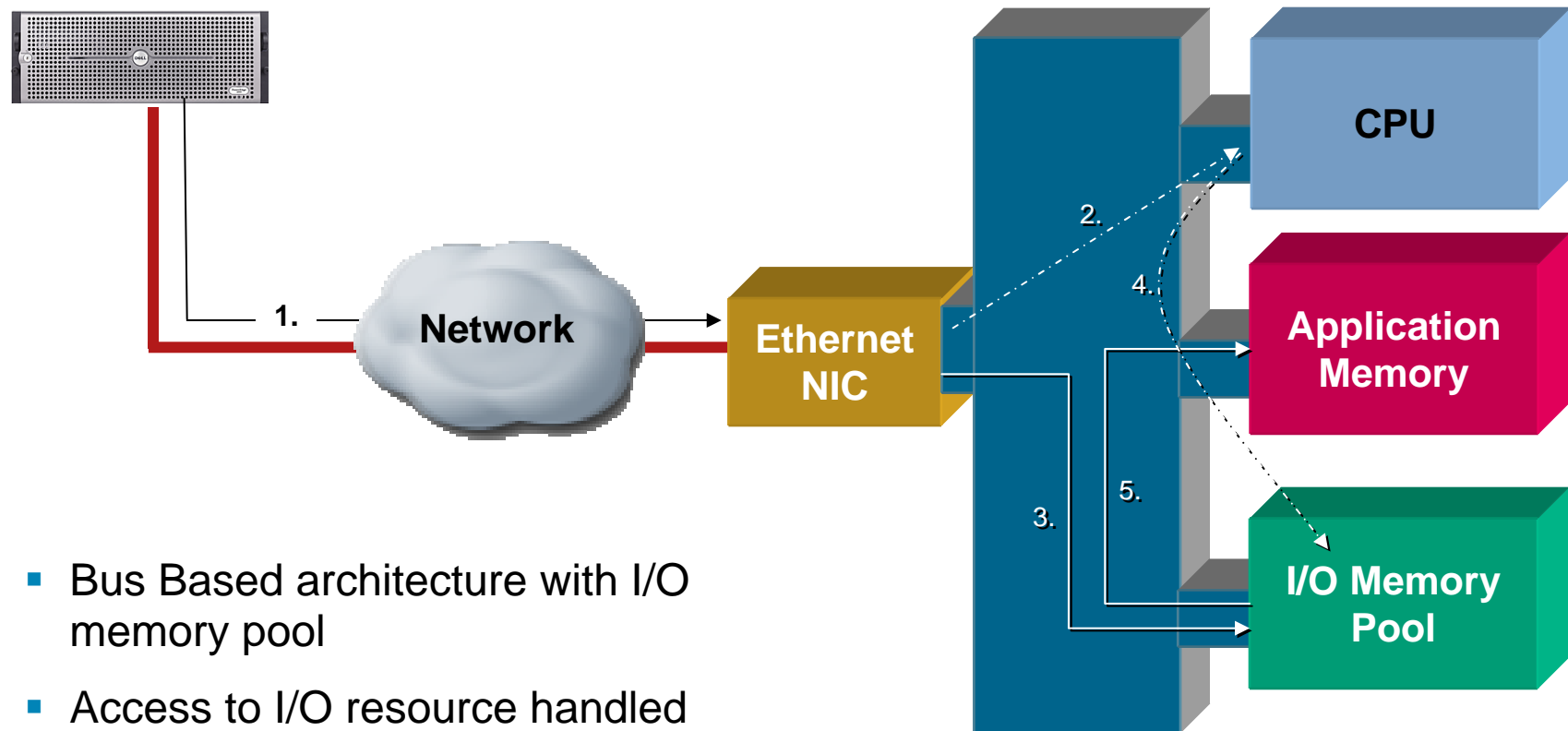
**Traditional Model**



**Kernel Bypass Model**

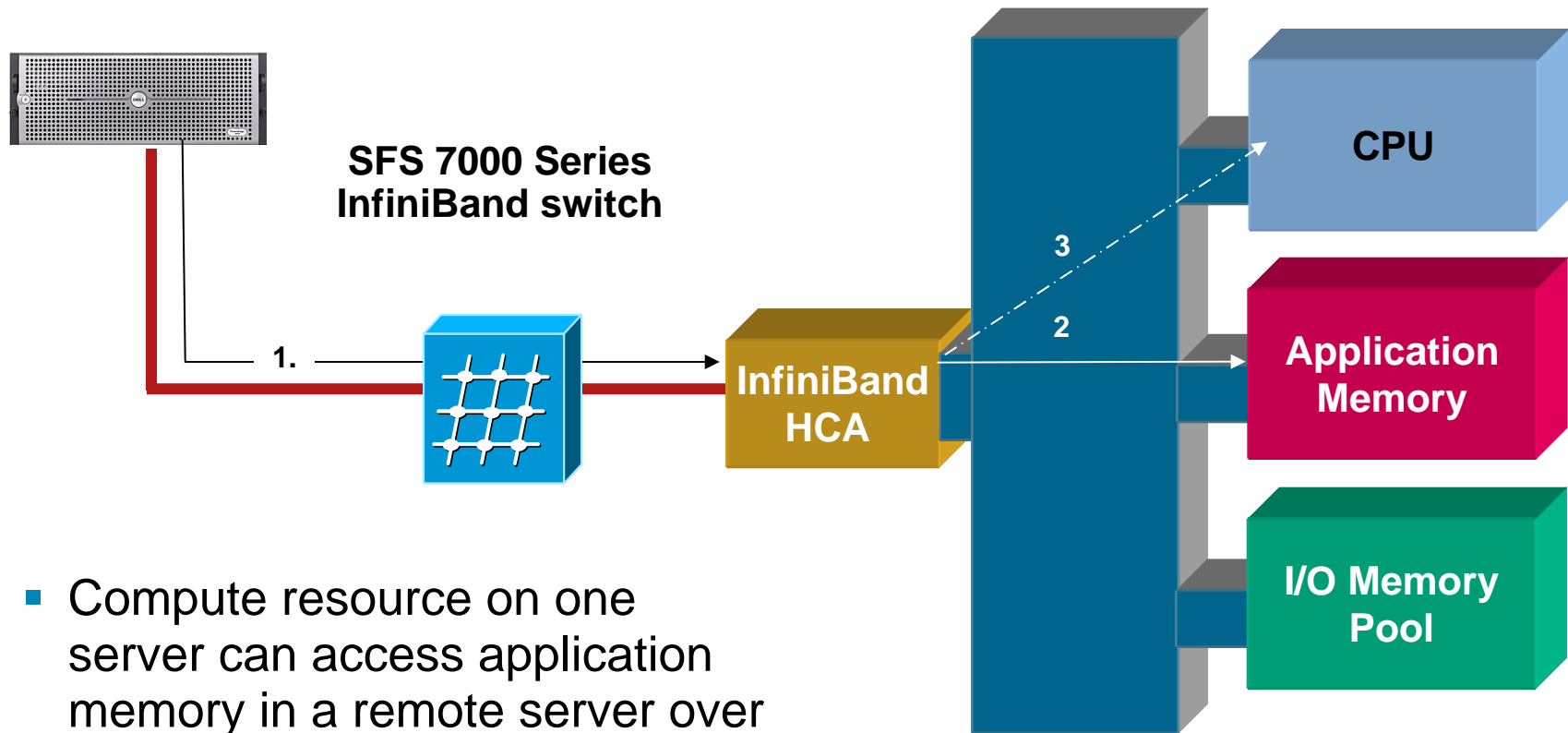


# Traditional Server I/O architecture



- Bus Based architecture with I/O memory pool
- Access to I/O resource handled by BIOS
- A data packet is typically copied across the bus three times  
CPU Interrupts, Bus bandwidth constrained, Memory bus constrained

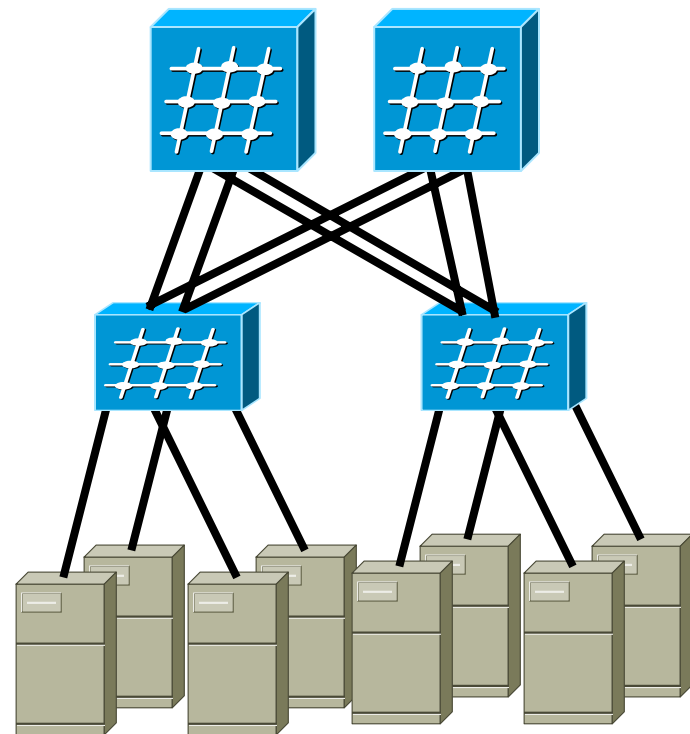
# RDMA I/O Architecture



- Compute resource on one server can access application memory in a remote server over InfiniBand connection

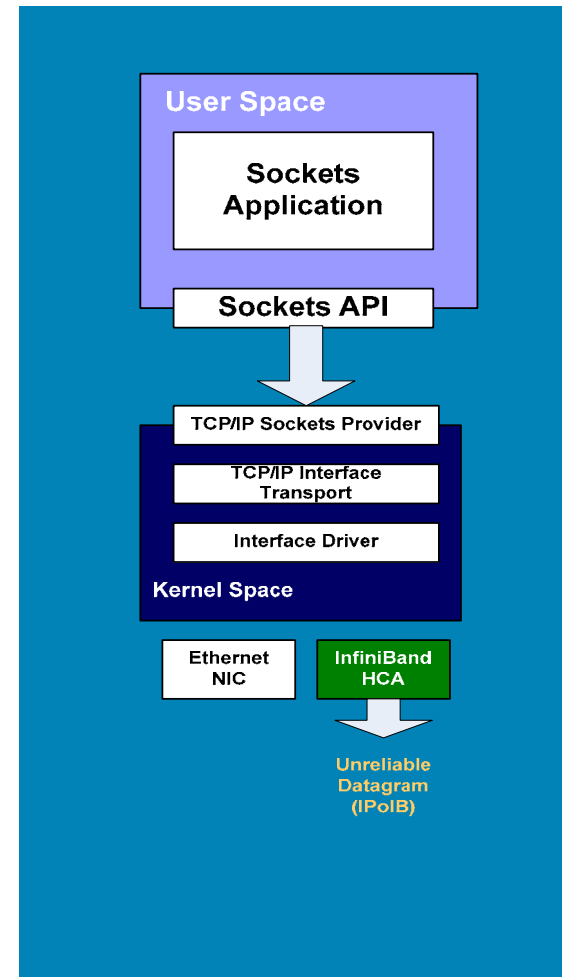
# InfiniBand Subnet Manager

- InfiniBand Fabric is called an InfiniBand Subnet
  - All devices under the control of a single Master Subnet Manager (SM)
  - May have multiple slaves with replicated SM database state
- At system startup, all devices register with the SM
  - Central Routing function
  - Shortest Path First Routing
  - Equal Paths Load-balanced with static round robin distribution
  - Connection endpoint lookup
- Built within the SFS switches and also offered externally



# IP over InfiniBand (IPoIB)

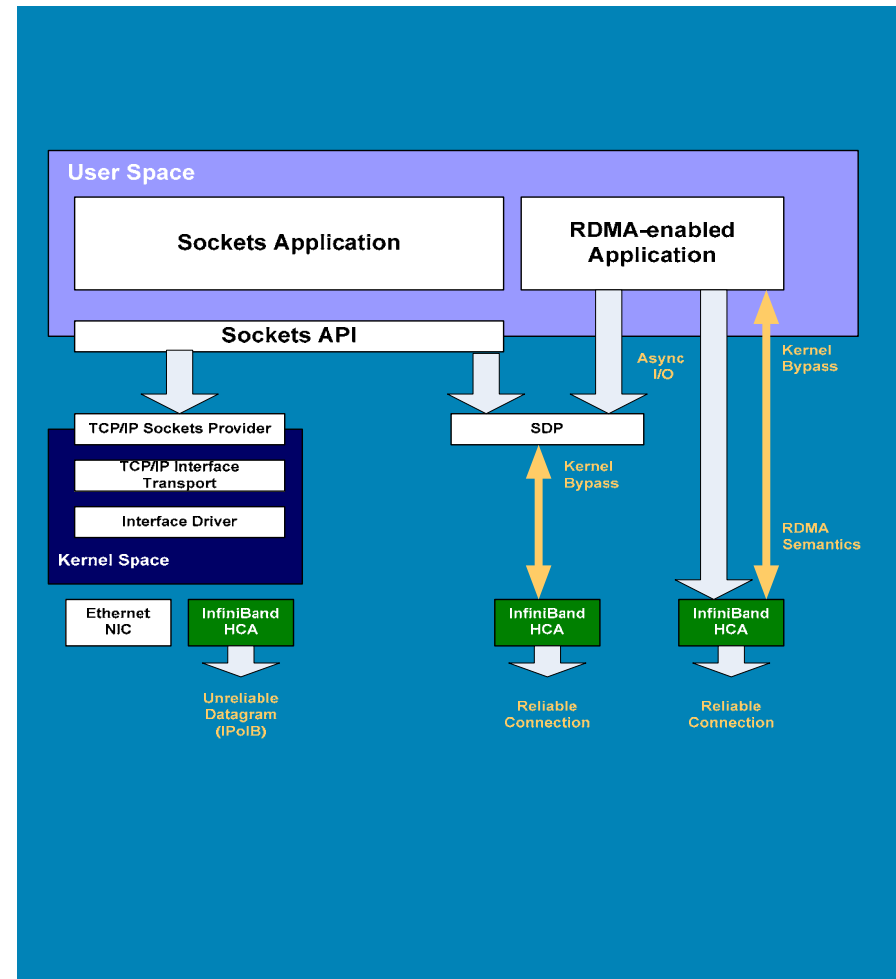
- IETF Standard (RFC 4931 / 4932)
- IP over InfiniBand provides the TCP / UDP socket interface for InfiniBand
  - Uses IB as a transport for IP
  - Support for IP Multicast over IB
  - No RDMA available in IPoIB
  - IPoIB is also used for address resolution for other protocols such as RDS, SDP, iSER
- Highest level of application compatibility no application change necessary
- Supported under Linux, Solaris, AIX, HP-UX, Windows
- SFS 3000 Ethernet Gateway may be used to bridge IPoIB traffic from InfiniBand to Ethernet





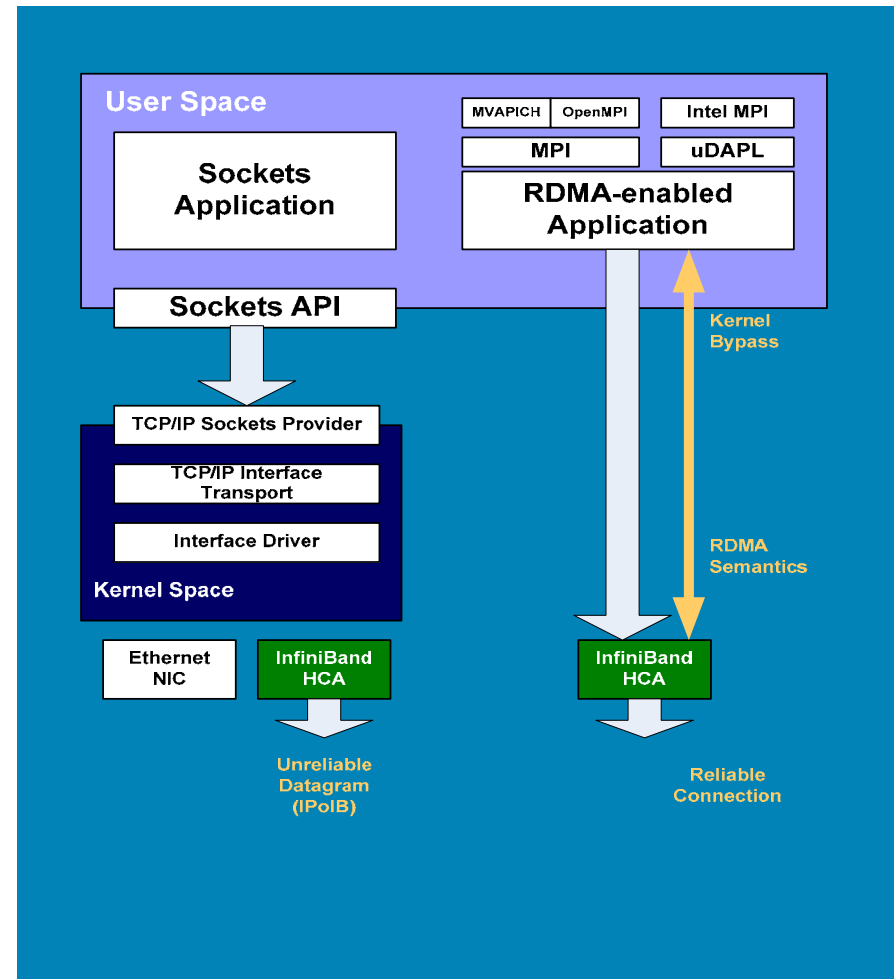
# Sockets Direct Protocol (SDP)

- Provides a compatible sockets interface that takes advantage of RDMA features
- Runs socket based TCP traffic with TCP and copy offload
- Application doesn't require any changes but host-side configuration is required
  - Configurable by process name, port numbers, destination address, etc.
- Zero copy SDP is possible but needs application / OS kernel support



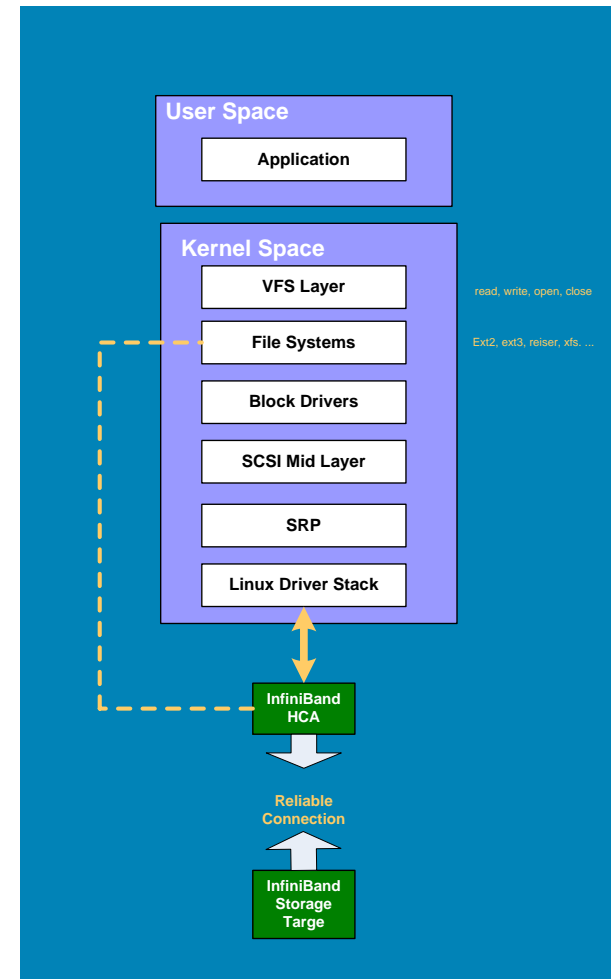
# Message Passing Interface (MPI)

- MPI is *message-passing* middleware
  - An extended message-passing model
  - Not a language or compiler specification
  - Not a specific implementation or product
  - Not a driver
- For parallel applications running on multiprocessor computers and clusters
- Supports heterogeneous compute environments
  - CPU, Memory and Interconnect agnostic
- Feature Rich protocol library: >300 functions
- Used extensively in HPC clusters
- MPI is not a JMS based message passing API



# SCSI RDMA Protocol (SRP)

- Intended to run SCSI protocol to run over InfiniBand for SAN usage
  - T10 specification, similar to FCP
  - Transactions use RDMA for data movement from target to initiator
- Host drivers tie into standard SCSI/Disk interfaces in kernel/OS
- Linux, Windows, Solaris implementations
- SFS 3000 FC Gateway is also a SRP Target
- Native SRP storage targets available today
- Not IB specific (no iWARP implementation yet)

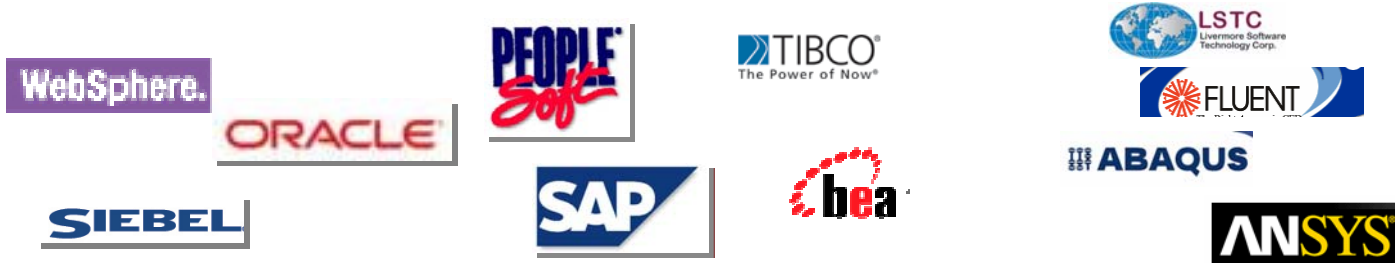


# InfiniBand Protocol Summary

Protocol / Application	Summary	Application Example
IPoIB (IP over InfiniBand)	Allows TCP/IP applications to run over the InfiniBand transport. Provides server to server and in-band management traffic from mgt station to switch and HCAs.	Standard IP-based applications. When used in conjunction with Ethernet Gateway, allows connectivity between IB network and LAN.
uDAPL (Direct Access Programming Library)	Allows application to take maximum advantage of RDMA benefits through flexible programming API. Requires custom development.	Used for IPC communication between cluster nodes for Oracle RAC.
SDP (Sockets Direct Protocol)	Adds RDMA benefits transparently to sockets-based applications. Can configure for all sockets applications or on a per port or application basis.	Communication between database nodes and application nodes, as well as between database instances.
SRP (SCSI RDMA Protocol)	Allows InfiniBand-attached servers to utilize block storage devices.	When used in conjunction with the Fibre Channel gateway, allows connectivity between IB network and SAN.
MPI (Message Passing Interface)	Low latency protocol used widely in HPC environments.	HPC applications.








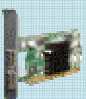
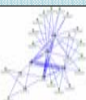
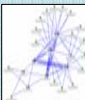
# InfiniBand Performance

## Measured Results

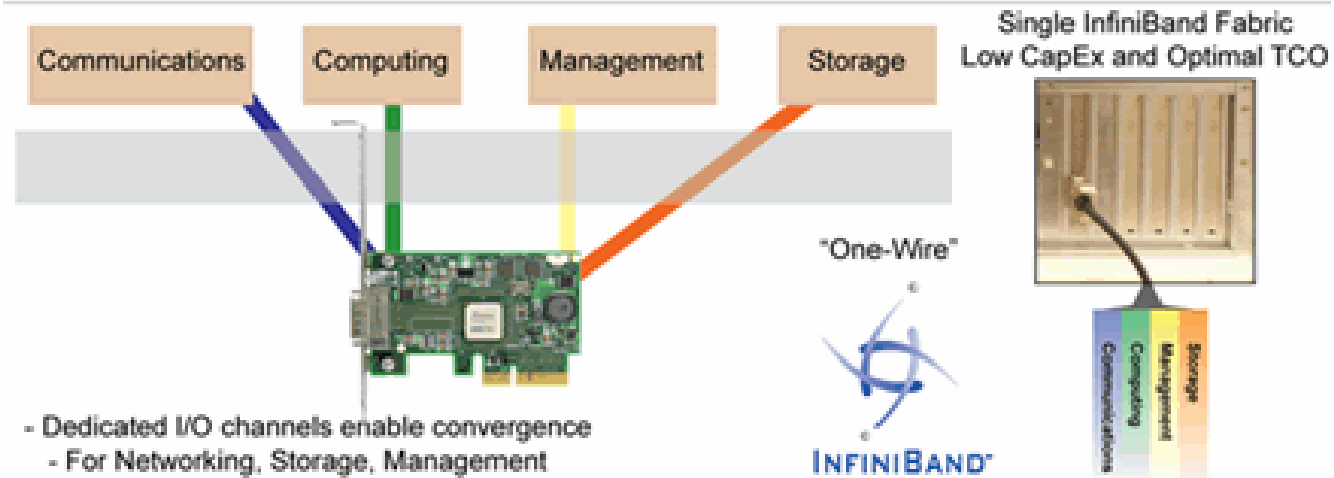
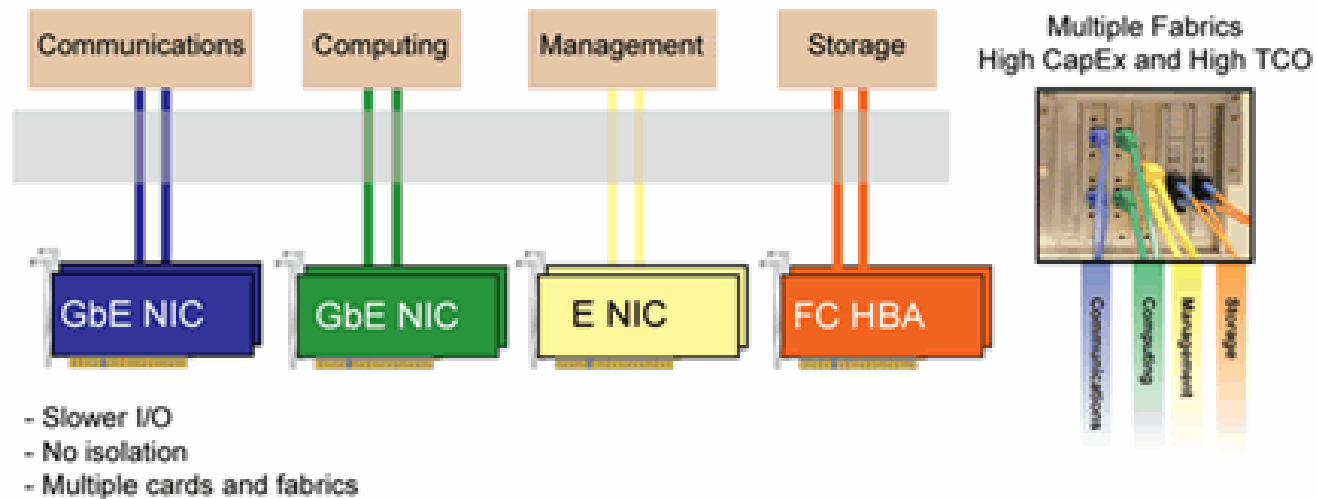


	Sockets API						MPI	
	TCP				SDP		MVAPICH	
	IP		IPoIB					
	Gigabit Ethernet	10 GE	SDR IB	DDR IB	SDR IB	DDR IB	SDR IB	DDR IB
<b>Latency (us)</b>	45.68	25.8	20.3	14.79	10	8.8	3.64	3.17
<b>Bandwidth MB/s</b>	118	1214	560	584	896	1033	960	1350
<b>CPU</b>	9%	25%	23%	26%	27%	28%	25%	25%

# The Cisco SFS Product Line: InfiniBand for High Performance Computing

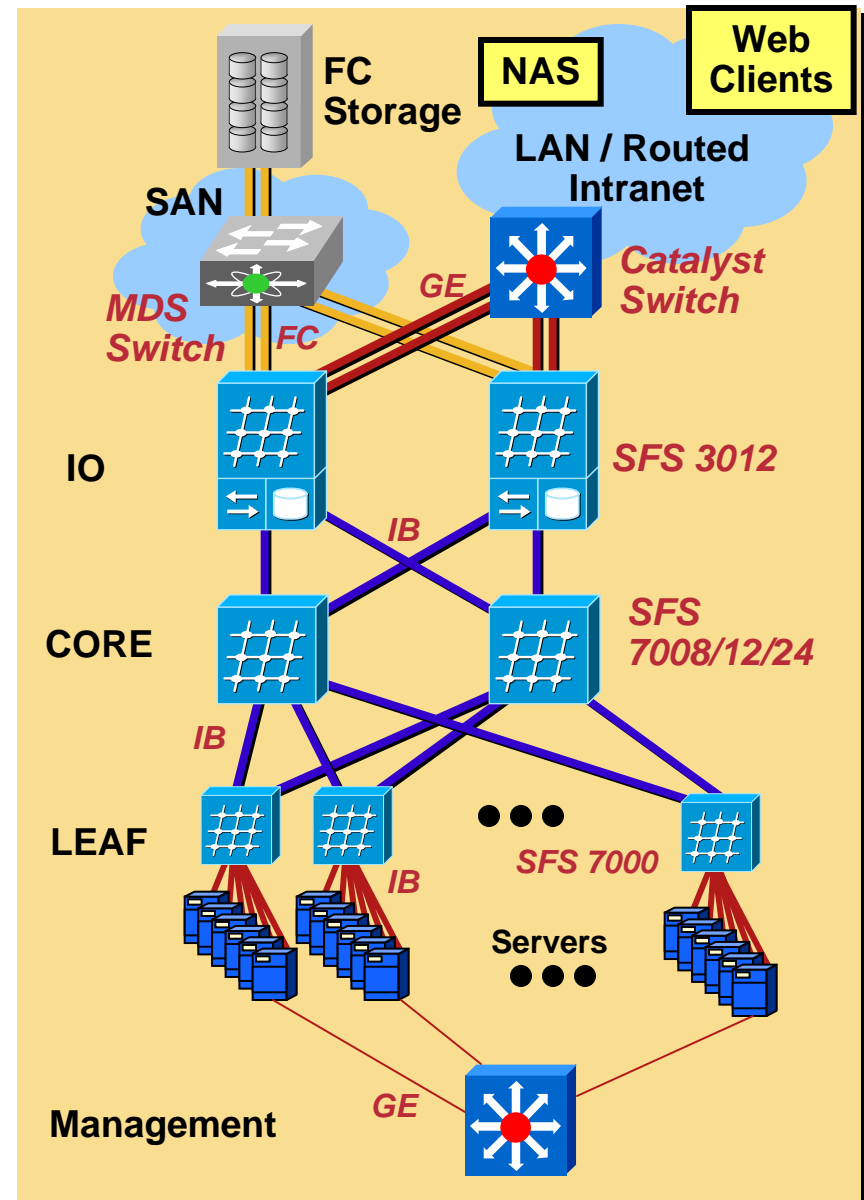
Server Fabric Switch	InfiniBand	<b>SFS 7000P</b>  (24) 4X IB	<b>SFS 7008P</b>  (96) 4X IB	<b>SFS 7012P</b>  (144) 4X IB	<b>SFS 7024P</b>  (288) 4X IB
	Multifabric	<b>SFS 3001</b>  (12) 4X IB + 1 GW	<b>SFS 3012</b>  (24) 4X IB + 12 GWs	 (2) 2-Gbps Fibre Channel Gateway (6) Gigabit Ethernet Gateway	
Blade Server	<b>IBM BladeCenter H</b> <ul style="list-style-type: none"> <li>• HCA (1) 4XIB PCI-Express</li> <li>• Embedded switch (14) 4X IB (internal) + (2) 4X IB and (2) 12X IB (external)</li> </ul>		<b>Dell PowerEdge 1855</b> <ul style="list-style-type: none"> <li>• HCA (2) 4X IB PCI-ex</li> <li>• Passthru Module (10) 4X IB</li> </ul>		
HCA	 <ul style="list-style-type: none"> <li>• (2) 4XIB PCI-X (Tall and Short Bracket)</li> <li>• (2) 4XIB PCI-ex (Tall and Short Bracket)</li> <li>• (1) 4XIB PCI-E (0 Mem, Tall and Short)</li> <li>• (2) 4XIB PCI-E (0 Mem, Tall and Short)</li> </ul>	<ul style="list-style-type: none"> <li>• Remote Boot</li> <li>• Linux Host Driver</li> <li>• Windows Host Driver</li> </ul>			
Subnet Mgmt	 <ul style="list-style-type: none"> <li>• High Performance Subnet Manager Software</li> </ul>	 <ul style="list-style-type: none"> <li>• Embedded Subnet Manager</li> <li>• SFS Element Manager / CiscoWorks LMS</li> </ul>			
Wire	<ul style="list-style-type: none"> <li>• 24 – 28 AWG Standard IB CX4</li> <li>• 30 AWG IB CX4 - SuperFlex</li> </ul>		<ul style="list-style-type: none"> <li>• 1, 3, 5, 7, 10, and 15 meters</li> <li>• 1, 3, 5 meters</li> </ul>		

# Virtual I/O Basics



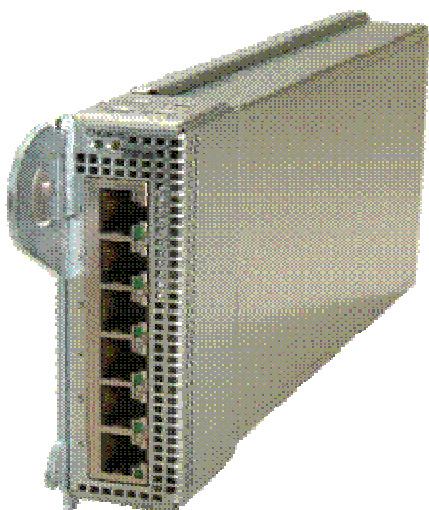
# InfiniBand – Accelerating Applications in the Data Center

- InfiniBand is a standards based technology
  1. Performance Improvement – 4 axis
    - Low Latency from RDMA
    - Improved Throughput
    - Message Rate
    - Lower CPU utilization for network transmission
  2. Multi-Fabric IO and Flexibility (CapEx & OpEx reduction benefit)
    - Inter-processor communication, Ethernet and Fibre-Channel traffic can all go over 1-wire
  3. Low cost 10G solution
    - Copper based solution; << \$1000 - \$1500 per port including switch port, HCA and cable
- Broad spectrum of support
  - Rack and Blade Servers (PCI-X, PCIe)
  - Major processors (Intel, AMD, IBM-PPC, Sun-Sparc)
  - Operating Systems (Linux, Windows, Solaris, HP-UX)
  - Storage & SAN (EMC, HDS, IBM, Cisco, Brocade, Multipathing support)





# SFS Ethernet Gateway Technology



- **3000 Series InfiniBand to Ethernet Gateway**
- **6 ports, 11.5M pps, 12Gbps Line Rate**
- **Ensures seamless integration with IP-based applications.**
- **IP bridge device**
- **Bridge group bridges one VLAN to one IB partition**
- **Ethernet bridge port can be tagged or untagged**
- **Ethernet bridge port can aggregate up to 6 ports**

# InfiniBand-to-Ethernet Gateway Features

- **IP-Only protocols**
- **802.1Q VLAN trunk support**
- **Link aggregation**
- **IPv4 multicast support**
- **Loop protection**
- **Ethernet jumbo frames up to 9k**
- **IP fragmentation**
- **High availability**

# SFS Fibre Channel Gateway



- **3000 Series InfiniBand to Fibre Channel Gateway,**
- **Two 2 Gbps Fibre Channel ports**
- **800 MBps throughput**
- **Supports**
  - SRP to FCP translation**
  - Dynamic load balancing and failover**
  - Load redistribution**
  - Global and individual ITL policies**
  - Topology Transparency**
    - Transparent support for zoning and LUN-based access controls**
  - Additional ITL security filters**

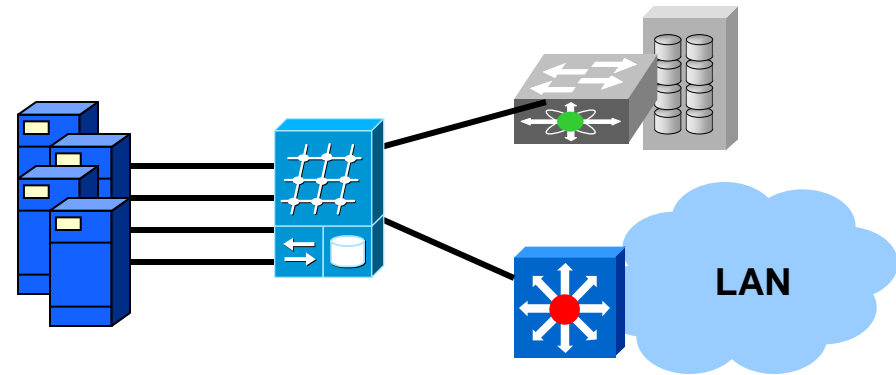
# InfiniBand-to-Fibre Channel Gateway

- **Ensures seamless integration with important SAN tools.**
  - Fabric-based Zoning**
  - LUN-based access controls**
  - Storage and host-based HA and load balancing tools**
- **Creates SAN network addresses on InfiniBand.**
  - SAN Management Tools must “see” each node.**
  - Creates “talk-through” mode with virtual WWNNs per server.**
- **Enables SAN Interoperability with InfiniBand.**
  - Appears as public AL-Port.**
  - Proven interoperability with Cisco, Brocade, McData, Qlogic, EMC, IBM, Hitachi, and more.**

# Physical vs. Logical View

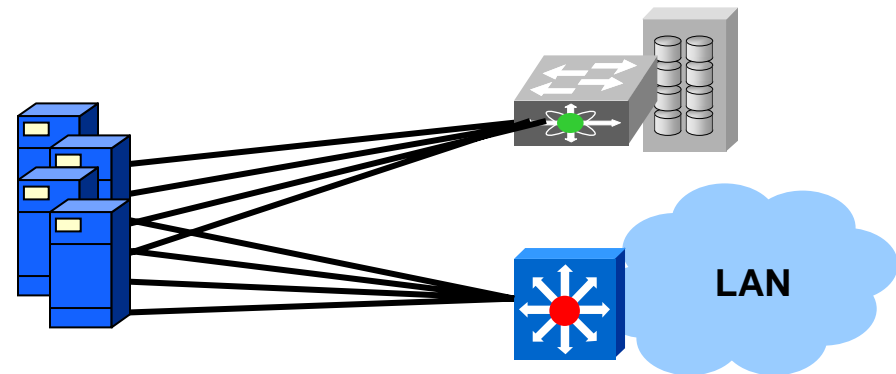
## Physical View

- Servers connected via IB
- SAN attached via public AL
- Ethernet attached via Gig Etherchannel



## Logical View

- Hosts present WWNN on SAN
- Hosts present IP address on VLAN



# InfiniBand Cabling Options

## Superflex and Optical Interconnect Options available



- **CX4 – Copper SDR – up to 15 meters, DDR – up to 9 meters**
- **Fiber Optics up to 200 meters – cross data center with pluggable optics modules on each end of fiber**
- **InfiniBand WAN capabilities have also been demonstrated over Optical Dense Wavelength Division Multiplexing (DWDM)**

# Cisco IOS Integration

## Consistent Configuration & Management

- Common CLI across all products
  - Command Syntax, scripting, etc.
- Consistent Security model
  - TACACS and RADIUS for Centralized Authentication
  - SSH/SSL/SNMPv3 for full management security
  - Multiple authorization levels
- File & Image Management
  - System image and configuration file libraries.
- Consistent Management Notification
  - Full SNMP v1/v2/v3 Support across all fabrics
  - Streaming Syslog: Integrates with Syslog Analyzer
  - Cisco Discovery Protocol (CDP)



# CiscoWorks LMS Support for Infiniband

- Single network management application for Ethernet and InfiniBand networks
- Resource Manager Essentials
  - Centralized Device, Software and Configuration Inventory Manager
- Dynamic Fault Manager
  - Diagnostic Tools and Syslog Analyze with centralized reporting
  - Device level fault analysis for network fabric, including high availability monitoring, pager/email/trap notification
- **Benefit: Eliminates administrative and usage barriers; identify and fix problems -> increased performance**

The screenshot displays the CiscoWorks LMS interface. At the top, there's a 'Cisco Systems Archive Mgmt' window showing configuration files. Below it is a 'Hardware Report' for a device at 10.64.158.162. The main window is 'Resource Manager Essentials', showing a 'Syslog Analyzer Standard Report' with a table of log entries. Below that is a 'Detailed Device Report' for the same device, which includes several tables: 'System Information', 'Chassis Information', and 'Interfaces'.

Update Time	Device Name	Domain Name	User defined Serial No.	System Name	Description	Location	Contact
13 Aug 2001 03:00:35 GMT	10.64.158.162		69029236	CM5500	Cisco Systems WS-C5500 Cisco Catalyst Operating System Software Version 4.5(10) Copyright (c) 1995-2000 by Cisco Systems	chennai	rmuthua

Available Slots	Slot Capacity	Model	Power Supply1	Power Supply2	Backplane Type	Chassis Type	Mgmt Type	Network Mask	Broadcast Address
10	13	WS-C5500	none	wsc5508	giga3E	wsc5500	snmpV2cV1	255.255.255.240	10.64.158.175

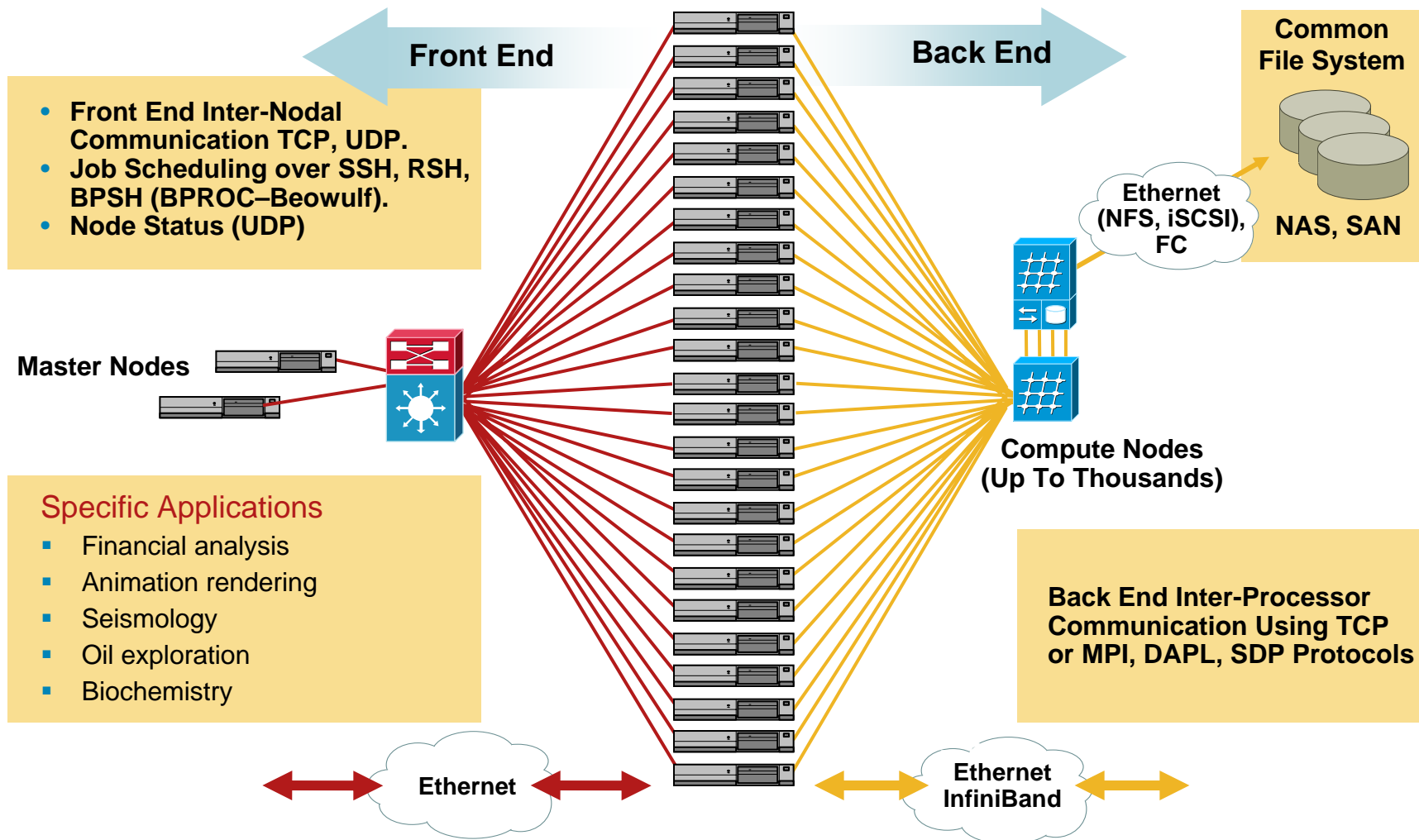
  

Slot No	Description	Type	Speed (bps)	Physical Address	Network Address	Status
N/A	sc0	ethernetCsmacd	10000000	00:90:5f:8b:2d:1f	10.64.158.162	up
N/A	sl0	slip	9600	00:00:00:00:00:00		down
N/A	VLAN 11	prop/Virtual	0	00:90:5f:8b:2a:0a		up
N/A	VLAN 12	prop/Virtual	0	00:90:5f:8b:2a:0b		up



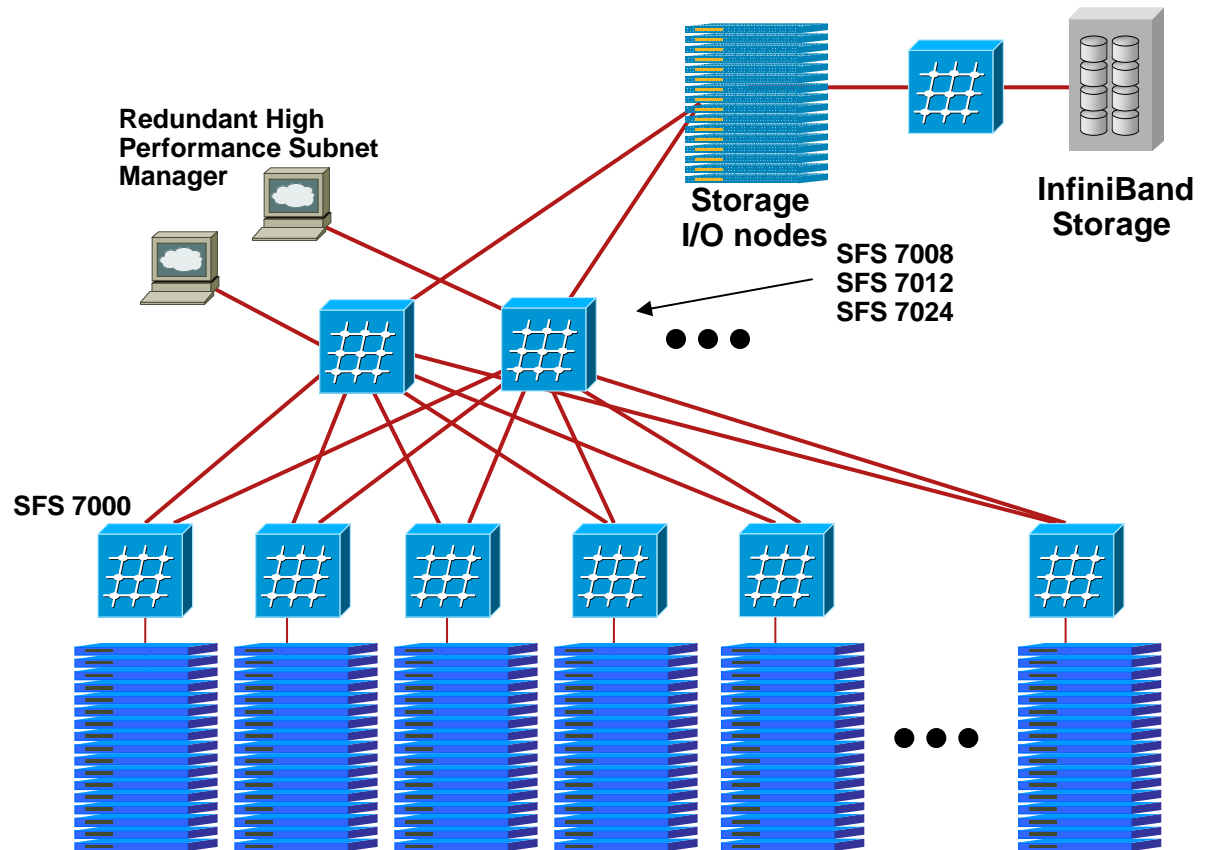
# Defining Clustered Servers

## High Performance Computing Clusters

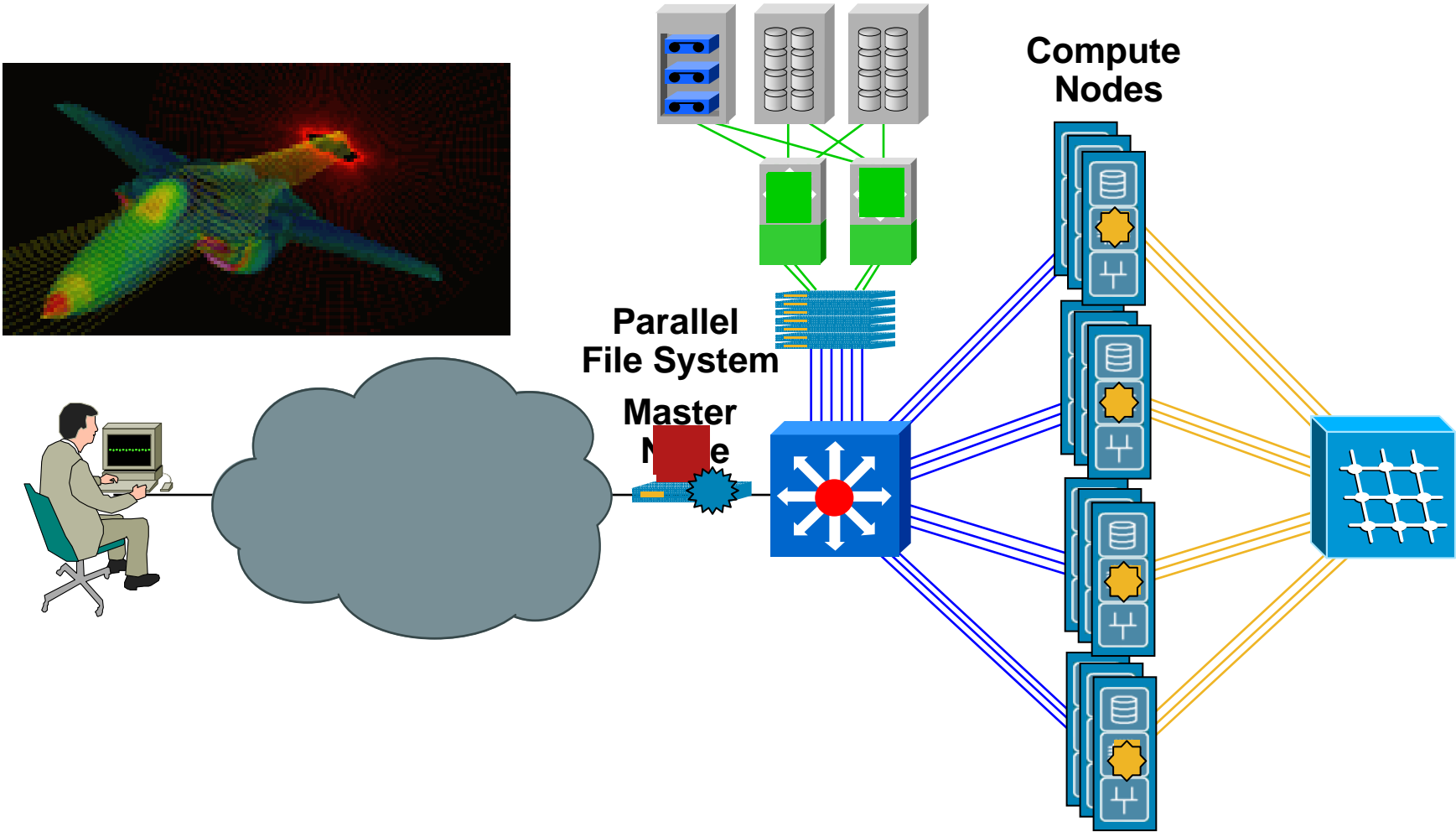


# HPC InfiniBand Networking

- Optimized for interprocessor communication for ultra-low latency applications
- Support SDR and DDR switching (10/20 Gbps)
- Scalable, Manageable, Modular design
- Integrated Subnet Manager in for “Plug-and-Play” operation
- High Performance Subnet Manager for the largest clusters
- InfiniBand 4X PCIe & PCI-X HCAs
- IB-attached storage for lower storage overhead



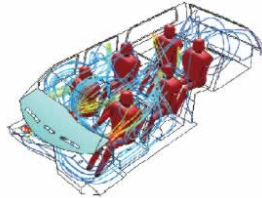
# How the HPC Network is Used



# HPC - Application Areas

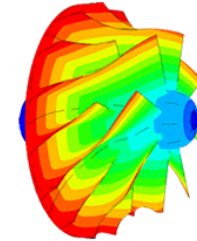
## Computational Fluid Dynamics

- Fluent
- Star-CD
- Exa PowerFlow
- Vectis



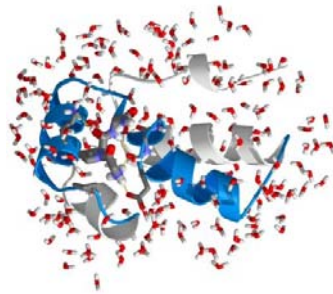
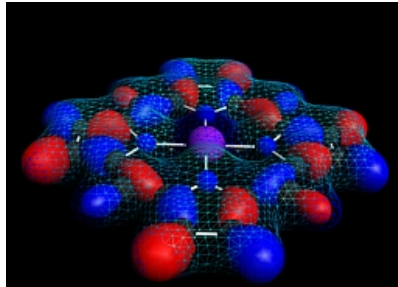
## Finite Element Analysis

- Abaqus
- Ansys
- LS-Dyna (3-D FEA)
- RADIOSS
- NASTRAN
- PAM-Crash



## Computational Chemistry Material Sciences:

- Amber
- Accelrys
- ADF
- GAMESS
- CHARMm
- Namd
- NWChem
- GROMACS



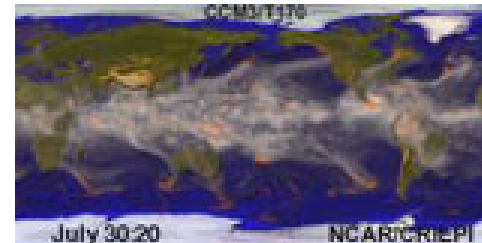
## Bioinformatics

- BLAST
- FASTA
- ClustalW
- EMBOSS



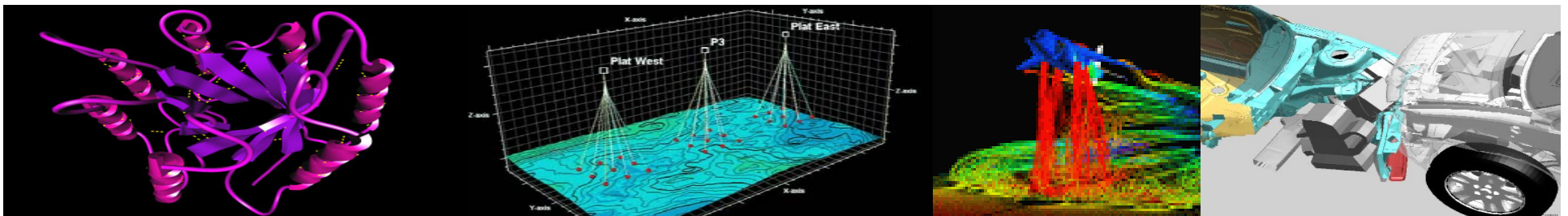
## Climate & Weather Simulation

- MM5
- WRF
- CCM3



# Application Classes and Characteristics

- **Tightly coupled applications**
  - Frequent IPC exchange
  - Latency sensitive, bursty traffic profiles
- **Loosely coupled**
  - Includes massively parallel, embarrassingly parallel or nearly-embarrassingly parallel
  - Little or no IPC traffic
  - Typically latency insensitive, Bandwidth may be a consideration for data set download
- **Parametric Execution Applications:**
  - no IPC traffic
  - latency insensitive, Bandwidth may be a consideration for data set download
  - Parametric Execution ~60% of HPC clusters
- **Application characteristics drive Network Technology & Design**



# HPC Inter-Process Communications Networking

- HPC cluster performance & efficiency is driven by IPC network characteristics
  - More time spent *communicating* is less time spent *processing*
- InfiniBand and Gigabit/10Gigabit Ethernet good IPC network technologies
  - Gigabit Ethernet ideal for smaller clusters and for loosely coupled applications
  - InfiniBand ideal for larger clusters and tightly coupled applications
- Understanding Application requirements and target CPU efficiencies are critical in technology selection for HPC

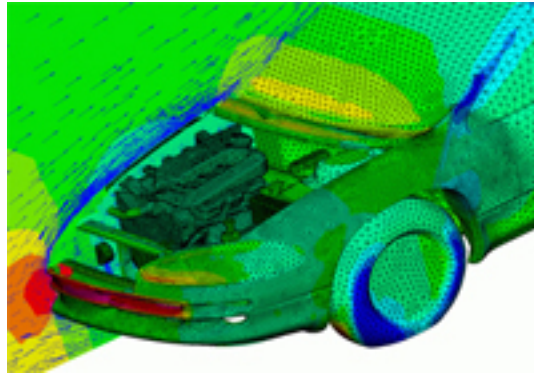


# HPC for Financial Sector

- Performance for Market Data Distribution
  - Low Latency, Throughput, Message Rate, CPU offload
  - Ecosystem: TIBCO, 29West, Wombat
- Traditional Applications
  - Monte Carlo simulations, Oracle RAC
  - Ecosystem: Platform, DataSynapse



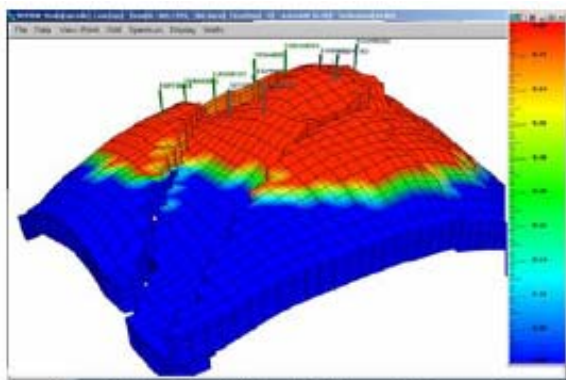
# HPC in Manufacturing



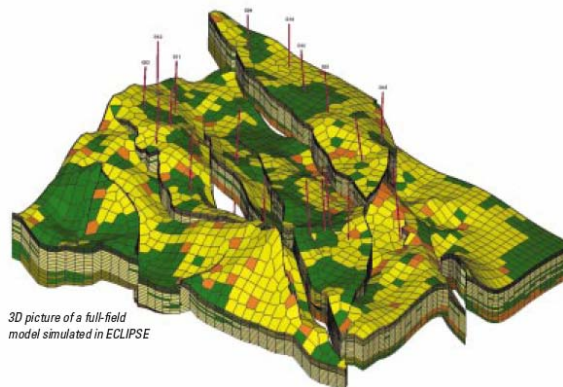
- Large clusters have been deployed with InfiniBand and Ethernet at major aerospace and manufacturing plants
- Reduce design cycles for automotive, aerospace, propulsion, mechanical for more rapid time to market
- Reduce time and cost of research and development



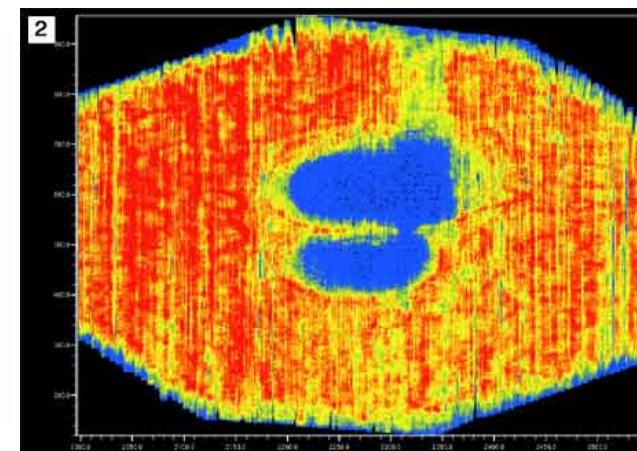
# HPC in Energy Sector



Landmark Graphics VIP: 3D View of Reservoir Simulation – Coarse Models for Risk Assessment



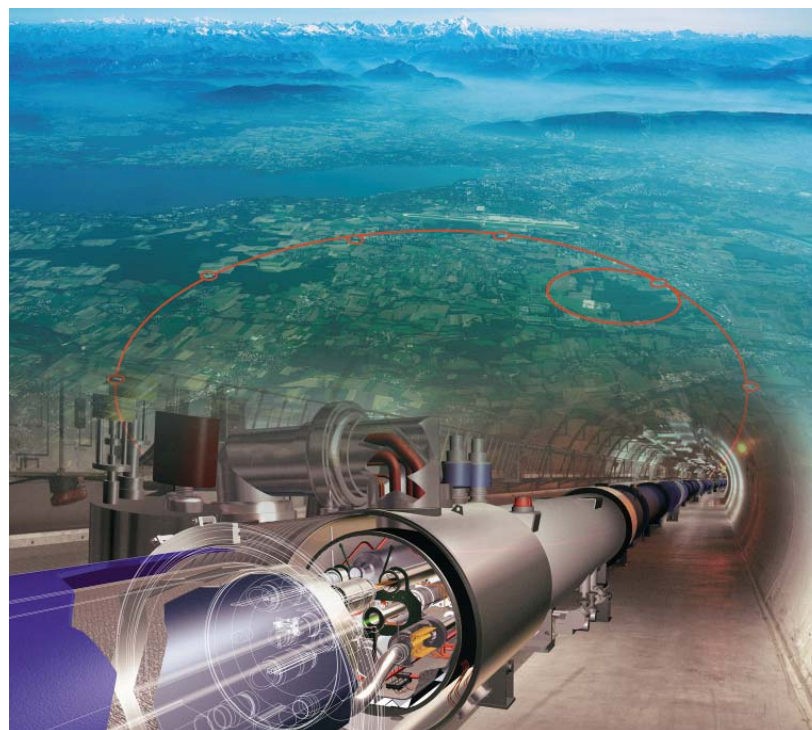
3D picture of a full-field model simulated in ECLIPSE



- Seismic Processing for Oil Exploration is a key driver for HPC
- Rapid processing of seismic data lead to more efficient drilling and “time-to-oil”
- Visualization of the reservoir for processing and analysis

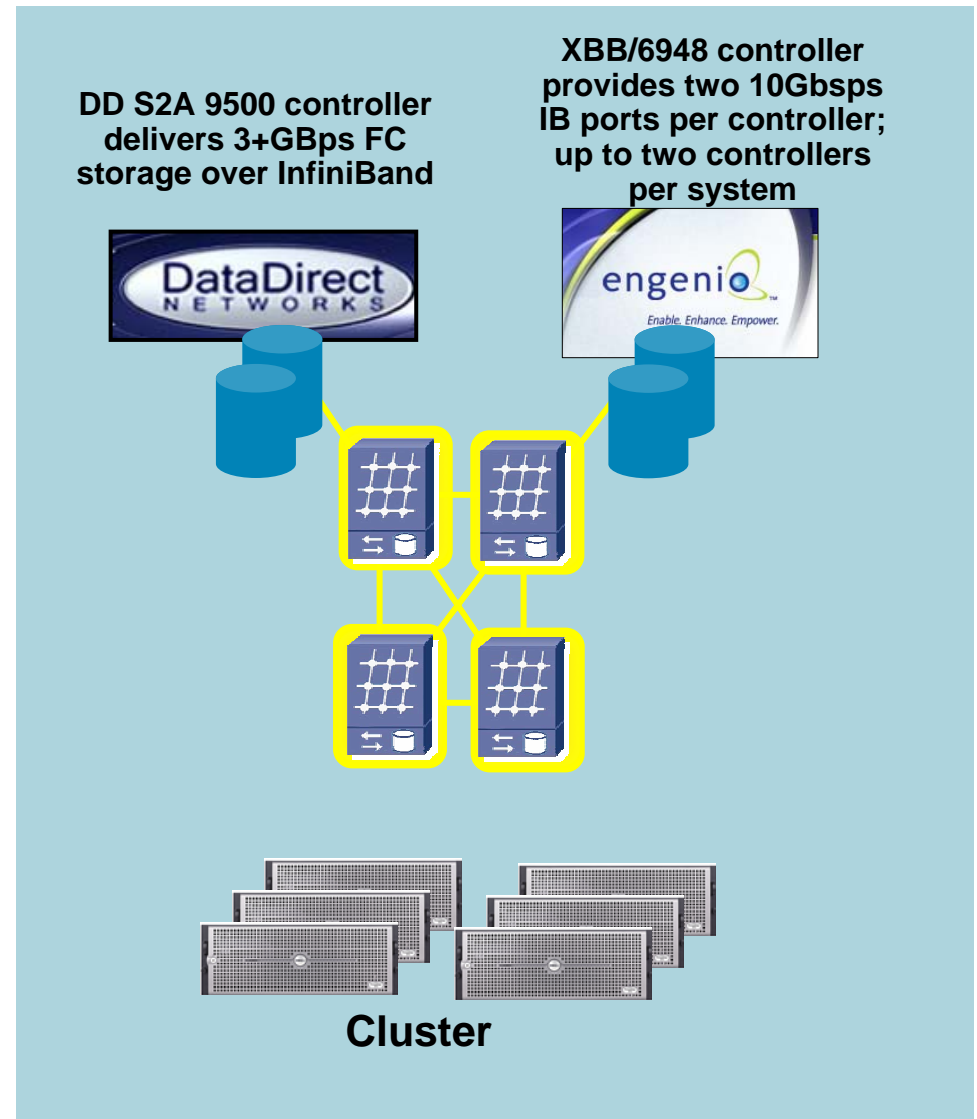
# HPC in Academia & Research Labs

- Large “Petascale” Systems  
10,000+ Node Networks
- Multi-core CPU’s  
Pushing requirements for DDR /  
QDR InfiniBand
- Higher Availability and Uptime  
Value for better engineering
- Campus-wide / WAN Grids  
Consolidation and sharing of  
Compute Infrastructure



# HPC Applications: InfiniBand Storage for Lower Cost and Higher I/O Performance

- Enable “unified fabric” with cluster and block storage over single IB fabric
- IB-attached servers get FC storage access for “free” (no FC port, HBA cost)
- SCSI RDMA protocol (SRP) moves FC block storage over IB



# Enterprise HPC switch features

- Enterprise class management features on SFS switches
  - Cisco CLI (async/telnet/ssh)
  - Access Security (TACACS or RADIUS authentication)
  - AAA - Access, Authentication, Accounting
  - SYSLOG and SNMP trap services
  - SNMP v3 monitoring and management
  - Dynamic Subnet Manager failover and database synchronization
  - CiscoWorks LMS integration (multi-box config/image management)

# Server Virtualization?



## What is VFrame™

- Cisco's data center-wide **virtualization** software suite
- **Delivers** the end-to-end manageability, control, and virtualization **benefits of the mainframe on top of** today's commodity components and **the Cisco IIN**
- Provides **virtualization, orchestration, and provisioning** for the data center resources that sit **between the "OS" and the "wire"**

# Three Categories of Server Virtualization

- Virtual Machine: **Splits a servers into independent virtual servers.**

*VMWare, XEN, MSFT*

**Main value is higher server utilization.**

- Virtual SMP: **Combines servers together into a single managed powered entity.**

*Virtual Iron, Qlusters*

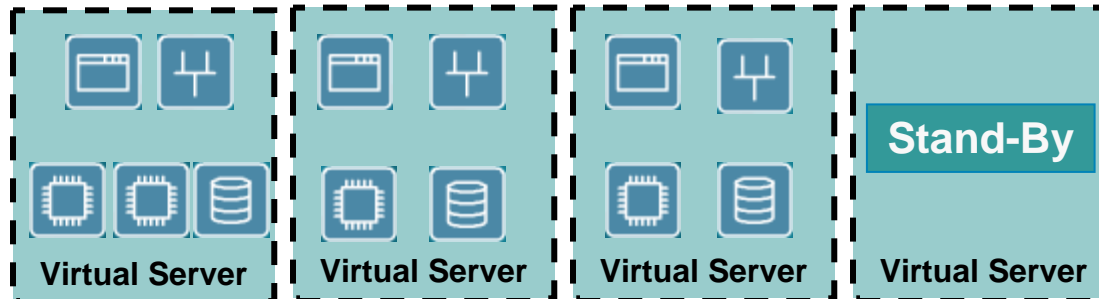
**Main value is scaling mission critical apps on commodity HW.**

- Physical Server Virtualization: **Makes servers stateless by moving server identity into the network, including storage and I/O subsystem.**

*Cisco VFrame™, Egenera*

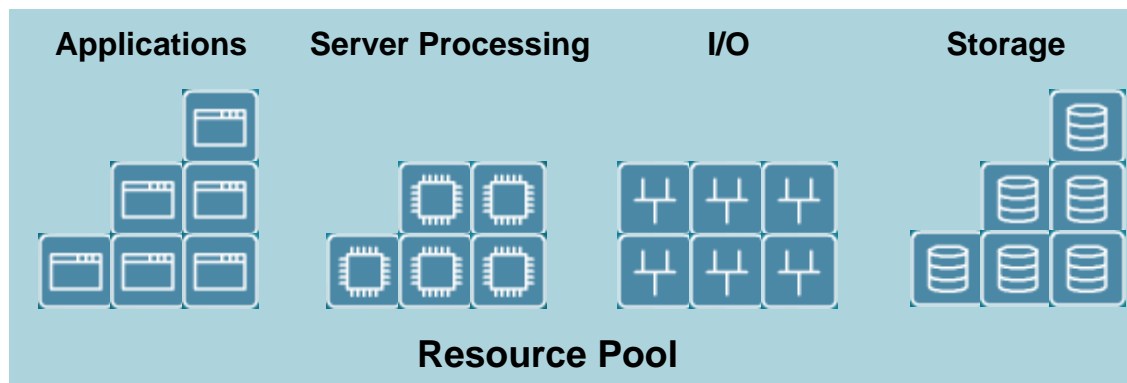
**Main value is making infrastructure change easier in heterogeneous environment.**

# Compute Networking and Virtualization —How Does It Work?

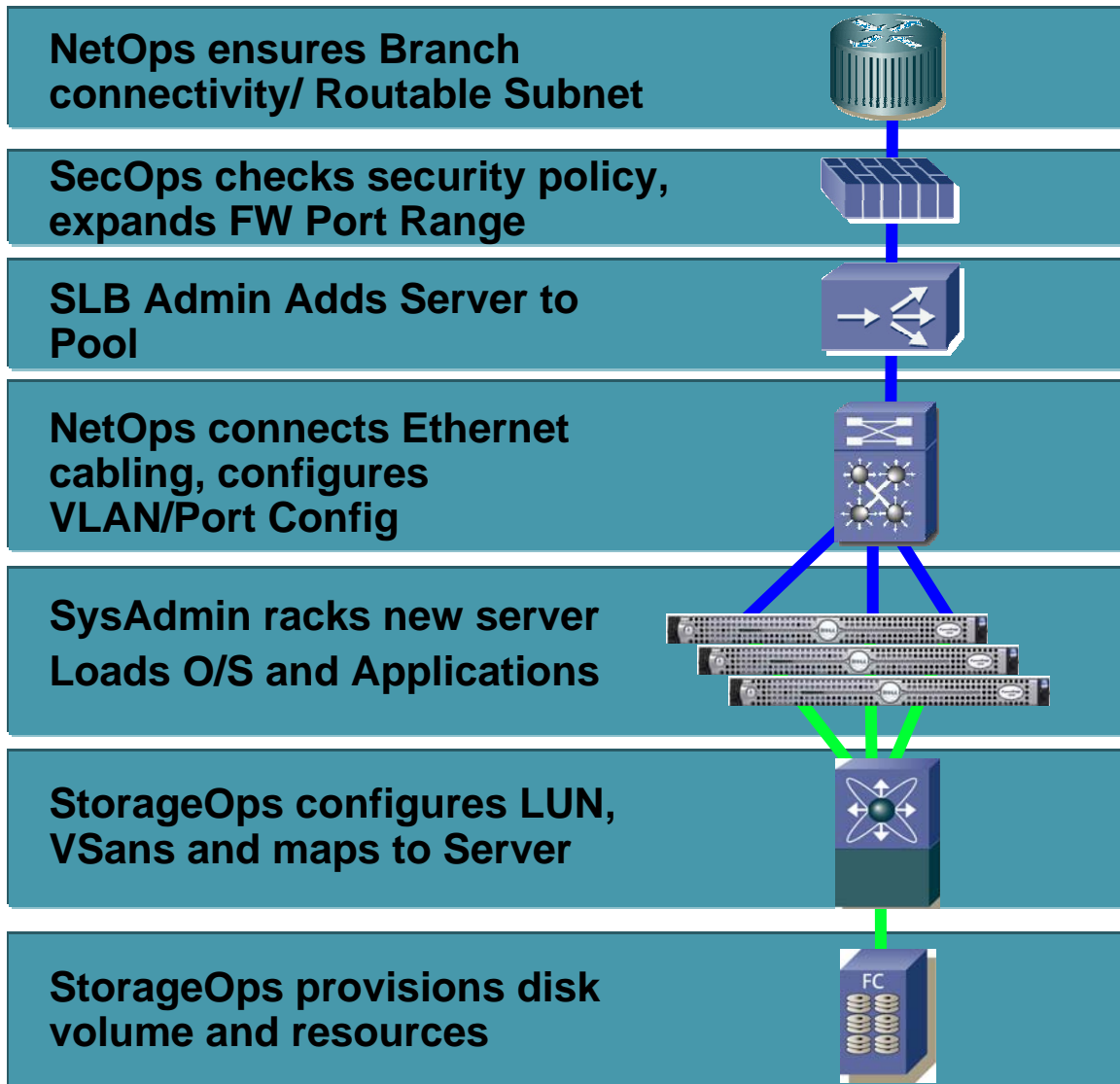


- Server is “taken apart” into its basic components—I/O, applications, compute power and storage
- Fabric re-assembles pools on demand to create “Virtual Servers” out of components
- Unified over an Intelligent Interconnect Fabric

## Intelligent Interconnect Fabric



# Current Enterprise Server Provisioning

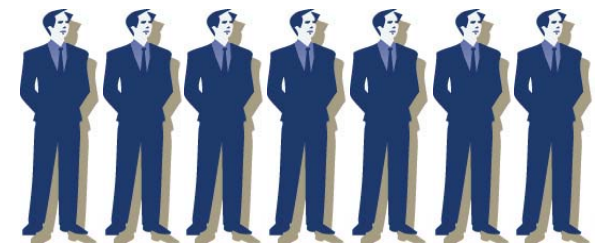


Assume you just want to add one server to a web-farm...

The challenge is one of 'coordination delays'. This type of simple scale-out of an existing serve often takes enterprises 90-days.

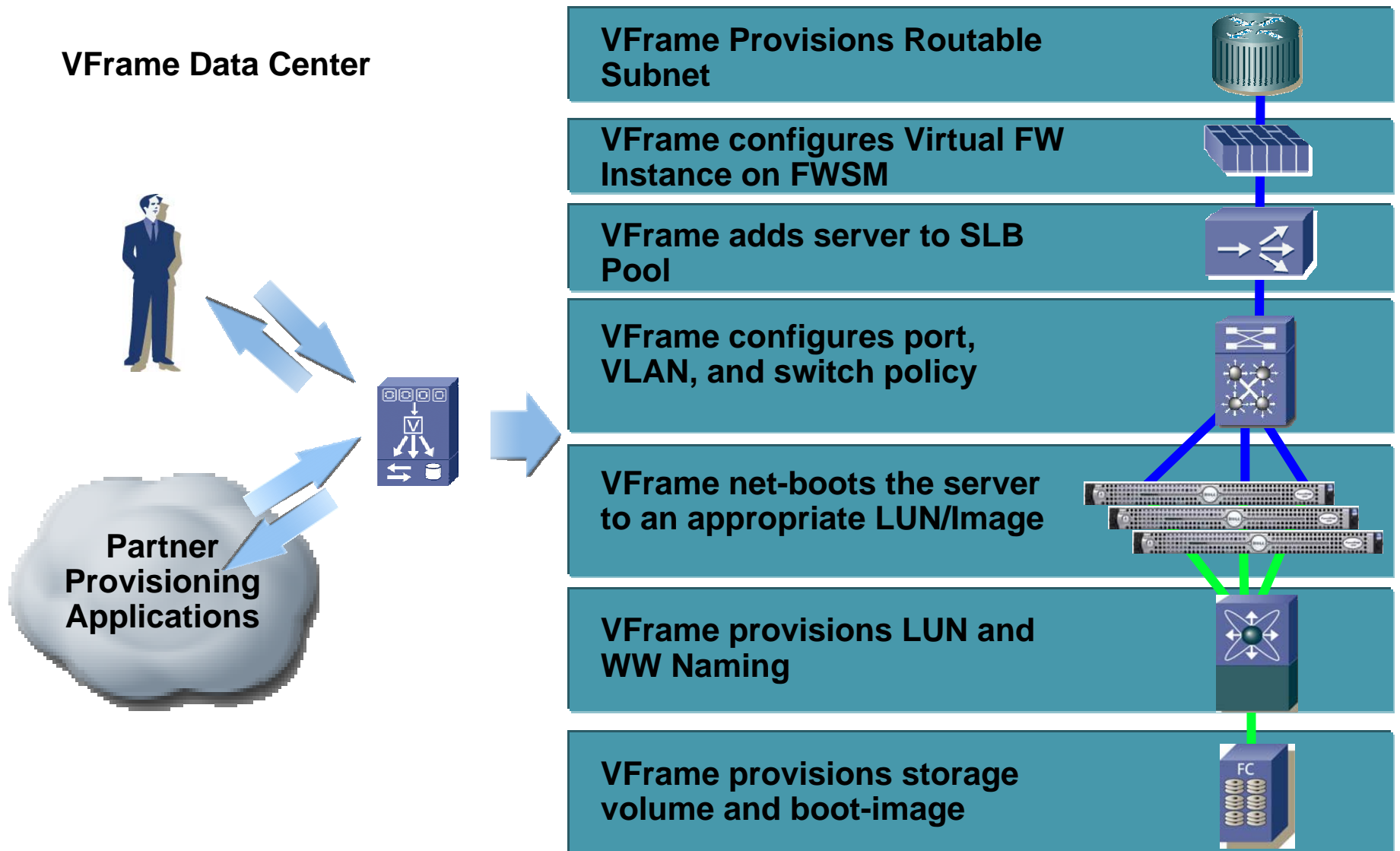
New service turn-ups, after the application has been developed, often take months of planning

Orchestration is designed to eliminate these delays and automate the provisioning of services

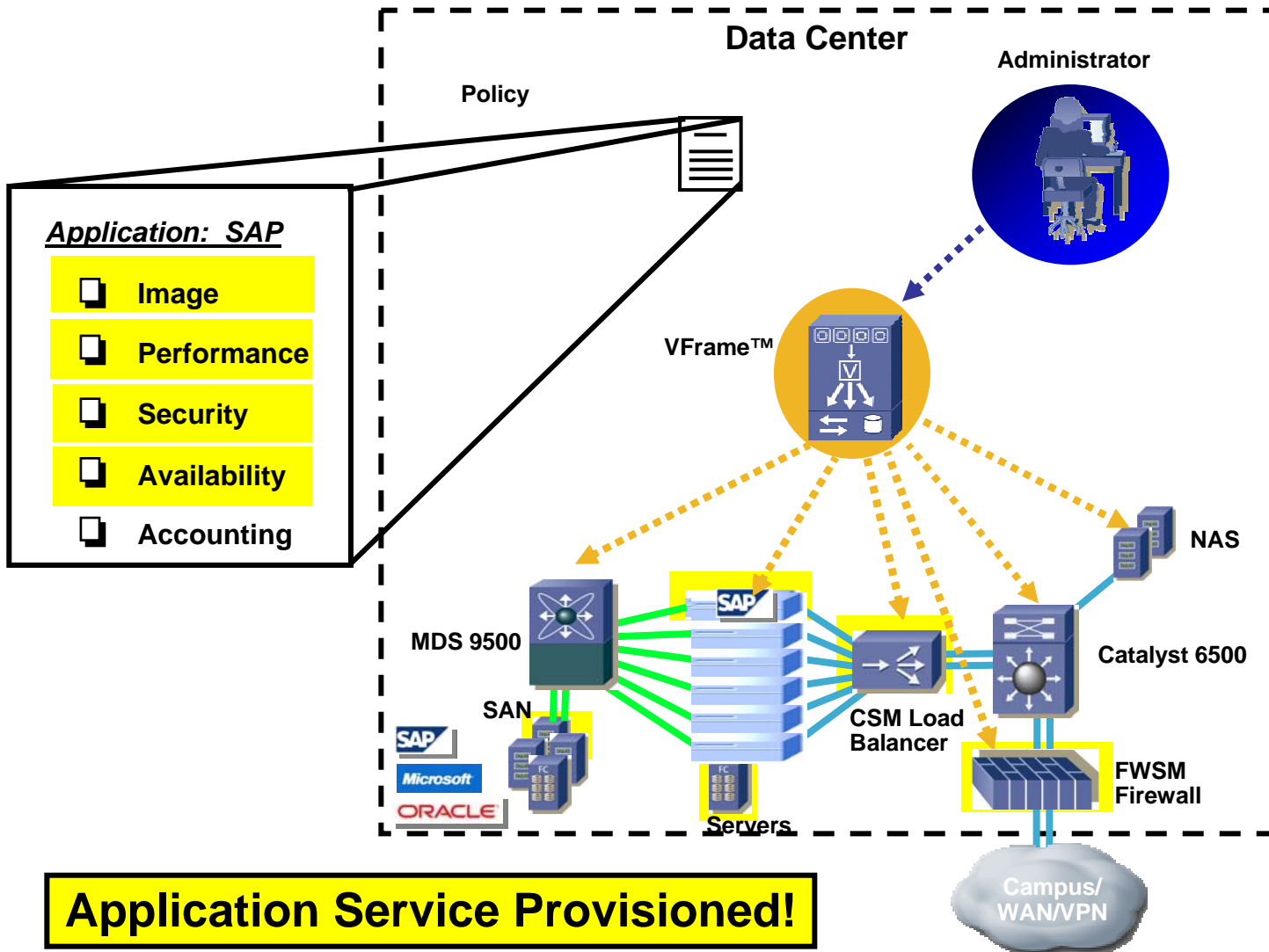




# VFrame Enterprise Service Provisioning



# Vframe Data Center 1.1 Creating a Virtual Fabric



Define application services and pass policy to VFrame

VFrame translates policies to actions and passes to infrastructure

VFrame identifies right App / OS Image From storage

VFrame picks server with right criteria to run application and boots server

VFrame gives new server right VLAN and LUN info so it can find/be found by right clients and storage

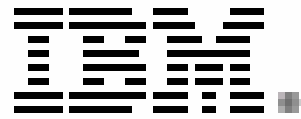
VFrame provisions security policies to FWSM

VFrame provisions CSM to add new server to load balancing pool

# VFrame Benefits

- Manage the data center from a **service-oriented, application-centric** perspective
- **Eliminate** number of **layers/devices** required to be touched to provision or modify
- **Ensure** security policies are enforced for compliancy and regulation.
- Treat the entire data center infrastructure (from the “OS” to the “wire”) as one manageable entity of **shared virtualized resources** (Virtual Mainframe)
- Expose a **single orchestration and provisioning interface** for all data center infrastructure
- Dramatically **reduce TCO**

# Cisco HPC and Top Tier Server Vendors



- Cisco has built relationships with server vendors to deliver integrated HPC and cluster solutions, jointly testing for solution delivery

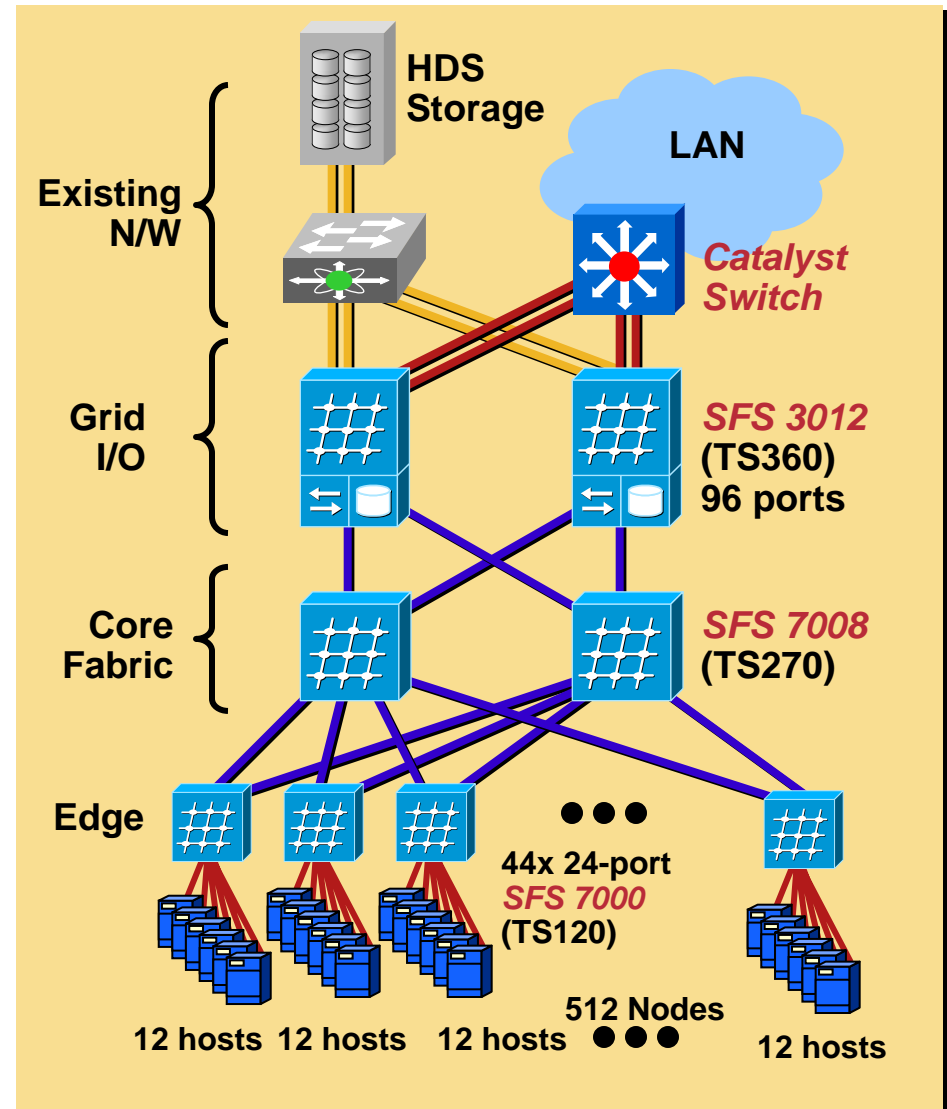
# InfiniBand – Who uses it today?

Industry	Organizations	Description
Financial Services	<ol style="list-style-type: none"> <li>1. Bank 1</li> <li>2. Bank 2</li> <li>3. Fitch Ratings</li> </ol>	<ol style="list-style-type: none"> <li>1. Risk calculations applications, MonteCarlo simulation</li> <li>2. Market data, back-end trading, hedge-fund pricing</li> <li>3. Oracle database</li> </ol>
Manufacturing	<ol style="list-style-type: none"> <li>1. Auto 1 2 3</li> <li>2. Airline Mfg</li> </ol>	<ol style="list-style-type: none"> <li>1. 600 - 800+ servers, Finite Element Analysis, CFD; ISV apps (Fluent, LS0-Dyna, PAM-crash)</li> <li>2. Computational fluid dynamics</li> </ol>
BIO / Pharma	<ol style="list-style-type: none"> <li>1. Private</li> </ol>	<ol style="list-style-type: none"> <li>1. 1000+ server, non-blocking cluster; used for protein folding research</li> </ol>
Service Providers	<ol style="list-style-type: none"> <li>1. Telstra, Australia</li> <li>2. EDS</li> <li>3. Sun</li> </ol>	All providers are using the InfiniBand interconnect to provide a cost-effective, flexible grid that can be used for a variety of applications and customers; 1000+ servers
Telcos	<ol style="list-style-type: none"> <li>1. British Telecom</li> </ol>	<ol style="list-style-type: none"> <li>1. Oracle database used for billing system</li> </ol>
Research Labs	<ol style="list-style-type: none"> <li>1. Sandia National Lab</li> <li>2. SARA, Netherlands</li> </ol>	<ol style="list-style-type: none"> <li>1. 4700 servers (Largest IB cluster in production)</li> <li>2. 512 nodes; Used</li> </ol>
Academia	<ol style="list-style-type: none"> <li>1. NCSA</li> <li>2. Universities multiple</li> </ol>	<ol style="list-style-type: none"> <li>1. Provides computation for Oil &amp; Gas clients; 500+ servers</li> <li>2. Academic research</li> </ol>

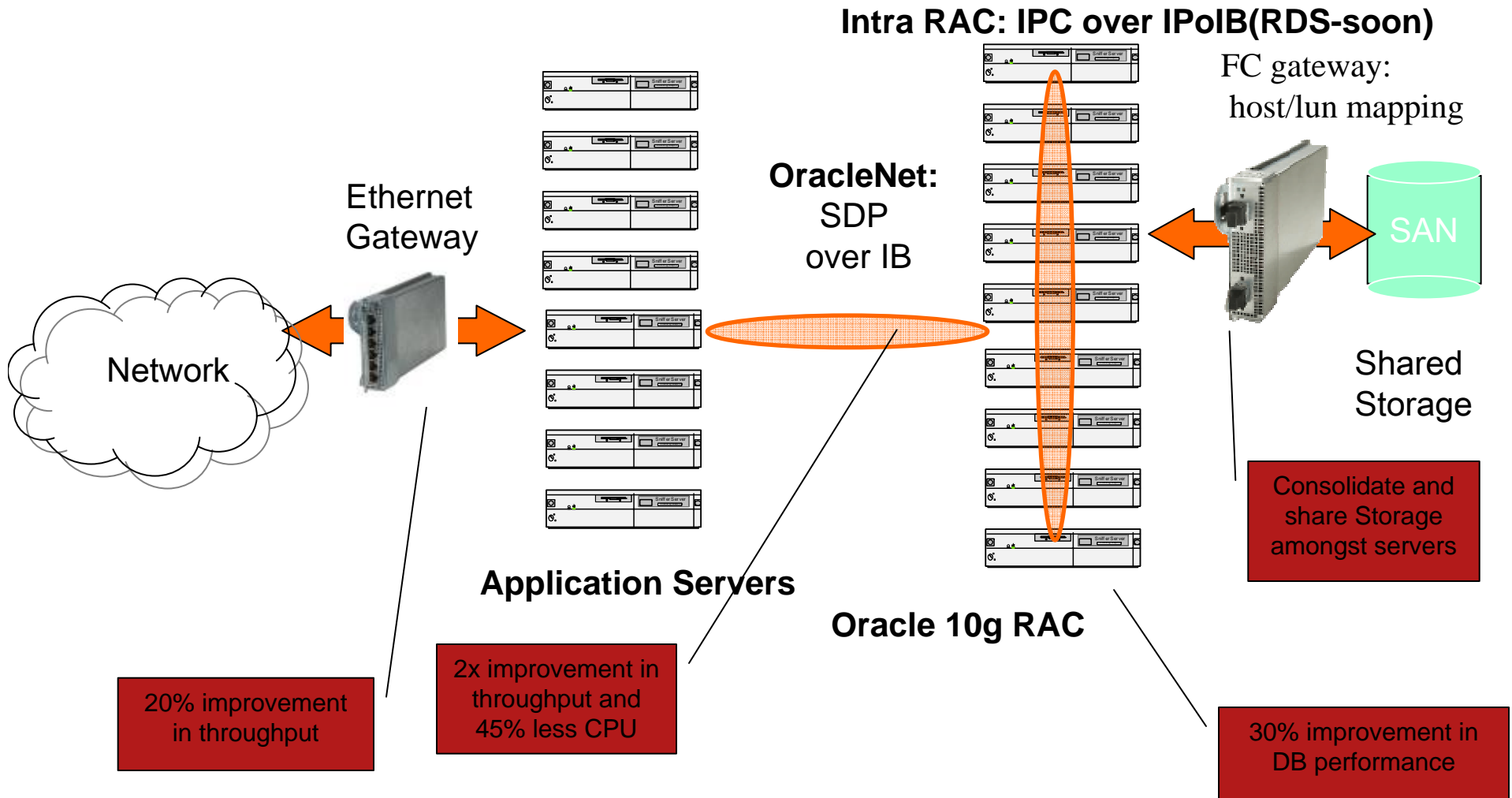
# Case Study: Large Wall Street Bank

## Enterprise Grid Computing

- Application:
  - Replace proprietary platforms with standards-based components
  - Build scalable “on-demand” compute grid for financial applications
- Environment:
  - 500+ Intel Servers per slice
  - Cisco Server Switch with Ethernet and Fibre Channel Gateways
  - Hitachi RAID Storage
  - SAN Switches
  - Ethernet Switches
- Benefits:
  - 20X Price/Performance Improvement over four years
  - 30% Application Performance Improvement
  - Standards-based solution for on-demand computing
  - Environment that scales using 500-node building blocks

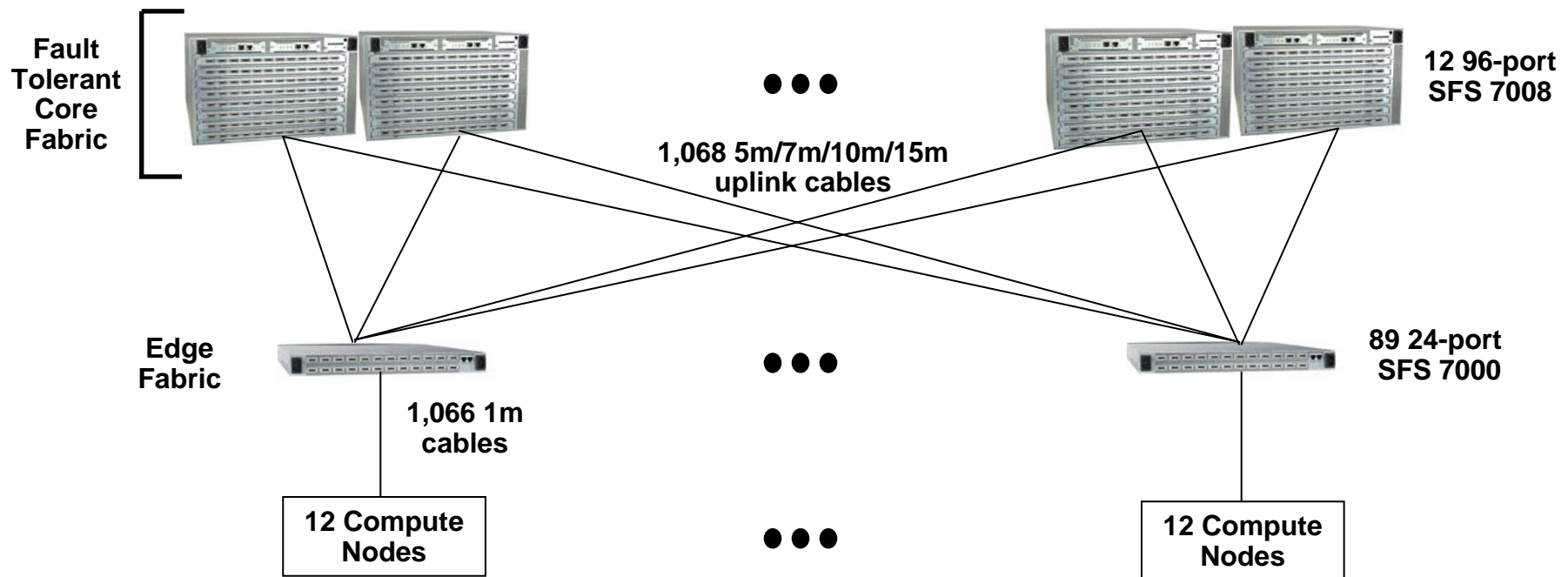


# Oracle RAC 10g: Scope of IB Benefits



# Bio-Informatics Cluster: 1,066 Node Supercomputer

## 1,066 Fully Non-Blocking Fault Tolerant IB Cluster



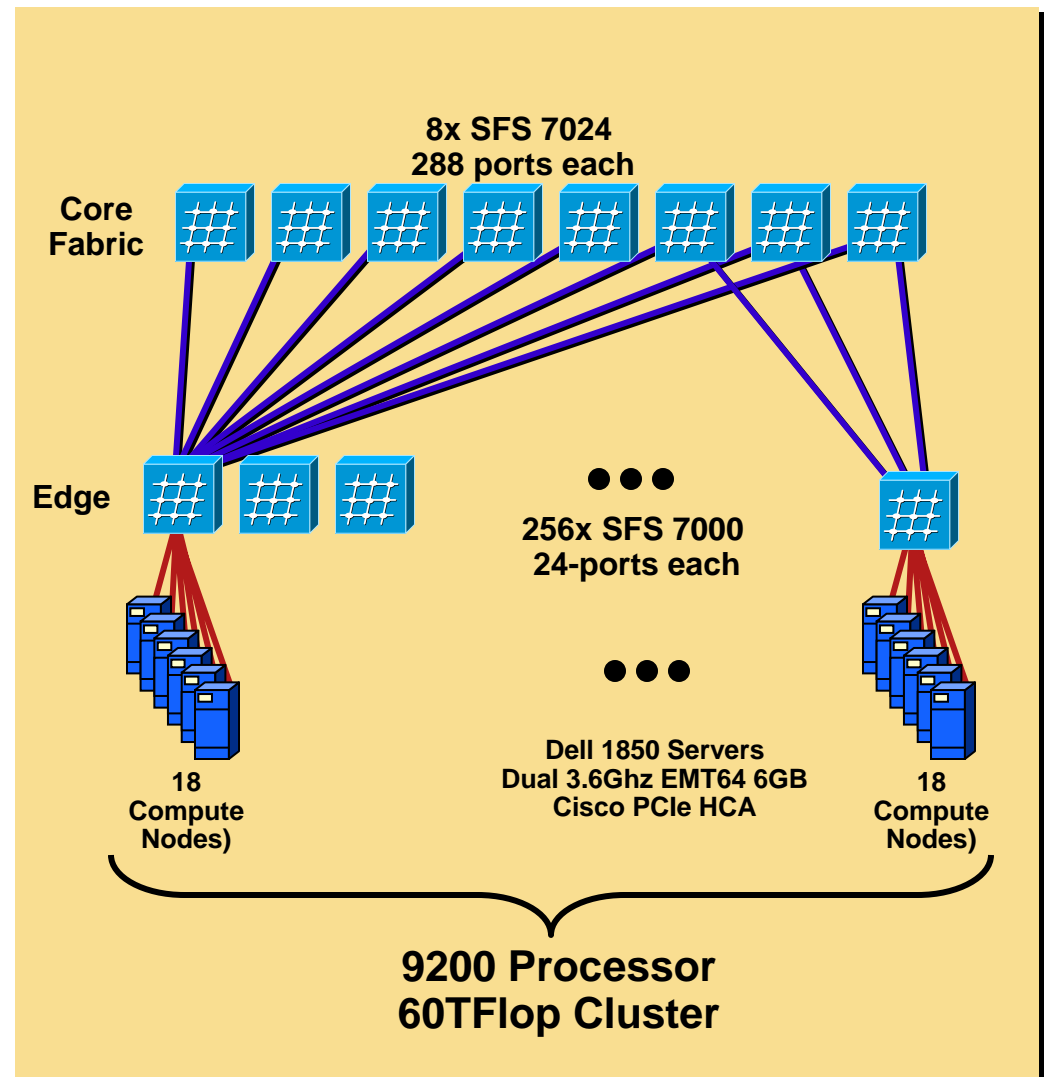
### Key decision factors:

- Cisco benchmarked and tuned customer MPI application
- Best operational experience with large clusters – best references
- “Rapid Service” architecture proved 2-min vs. 2-day MTTR.



# Sandia National Labs – 4600 Nodes Cluster

- Application:  
High Performance SuperComputing Cluster
- Environment:  
4600 Dell Servers  
50% Blocking Ratio  
8 SFS 7024  
256 SFS 7000's
- Benefits:  
Compelling Price/Performance  
Largest IB Cluster ever built  
3<sup>rd</sup> Largest Supercomputer in the world (Top500)



## Key Takeaways

- Cisco provides most complete HPC solution encompassing InfiniBand, Ethernet switching and storage.
- Cisco is a leading manufacturer of InfiniBand, Ethernet and Storage networking switches that are deployed in some of the Worlds largest clusters.
- Only network manufacturer with global channels, expertise and support for Complete HPC networking solutions.

