

Inteligent data center/next generation data center

Internet Users' Conference - CUC 2005 Dubrovnik, November 21.-23., 2005.

Wednesday, 23.11.2005. 11:00-13:00

Josip Zimet Cisco Systems

Session Number Presentation_ID

List of architectures





Session Number Presentation ID

© 2005 Cisco Systems, Inc. All rights reserved

Architectures ...

2004/2005 2001/2002 2002/2003 2005/2006 2000/2001 SDN SAFE V3PN SONA **AVVID** -Integrated security -sla - security -standardization - sla - collaborative security systems - virtualization -security -adaptiive threat defense

Next generation switches :

 $1900 \rightarrow 2900 \text{XL} \rightarrow 2950/3550 \rightarrow 500/2960/3560$

Next Generation routers :

 $2500 \rightarrow 2600 \rightarrow 2600 \text{XM} \rightarrow 2800 \text{ ISR}$

Next generation PIX :

PIX $5xx \rightarrow PIX 5xxE \rightarrow ASA 55xx$

<u>Next generation IDS, Next generation Wireless ...</u>

List of acquisitions ...



Acquisitions related to the DNA ...



DNA



Presentation ID

Data Center DNA Products



Session Number Presentation ID

© 2005 Cisco Systems, Inc. All rights reserved

Cisco DNA Infrastructure Evolution/Roadmap



The Big Picture - The Cisco Data Center



Business Ready Data Center Architecture to Topology



Presentation_ID

Enterprise Data Center Network Topology



The Evolution of the Data Center



Session Number Presentation_ID

12

Server Switching Architectural Evolution



The Evolution of the Data Center



Session Number Presentation ID

© 2005 Cisco Systems, Inc. All rights reserve

The Evolution of the Data Center

| Sorviooucere | | | |
|-------------------|-------------------|--------------------|----------------------------------|
| Layer | Cluster Mgmt | Cluster Mgmt | Virtualized Mgmt |
| | | | |
| Middleware | Data Middleware | Compute Middleware | Web Services |
| Lajor | | | |
| Resource Layer | Legacy Servers | Blade Servers | Compute-Network Farms |
| | | | |
| Network Layer | 10/100/1000 Ether | net GigE/10G | GigE/40GigE/100GigE 2G/10G FC |
| | 1990 | 2000 | 2010 |

CISCO CONFIDENTIAL

Cisco DNA Architecture



Session Number Presentation_ID

A New Category of Data Center Infrastructure-The Server Fabric Switch



Session Number Presentation ID

DNA Virtualization Vision



© 2005 Cisco Systems, Inc. All rights reser

What Makes The Server Fabric Switch Different?

High Performance Server-to-Server Interconnect

Policy-Based Dynamic Resource Mapping

Virtualization (I/O, Storage, <u>and</u> CPU)

Performance and Control

Session Number Presentation_ID

Server Fabric Switch Applications Why Performance <u>and</u> Control?



RDMA & OS Bypass, Kernel Bypass



Cluster Application Interconnect



Session Number Presentation_ID

InfiniBand Performance Measured Results



Session Number Presentation_ID

I/O Gateways for Network and Storage Eliminating Technology Islands



Programmability VFrame™

- Server Switch receives policy from VFrame[™] Director or 3rd party software.
- 2) Based on policy, Server Switch assembles the virtual server
 - Selects server(s) that meet minimum criteria (e.g. CPU, memory)
 - Boot server(s) over the network with appropriate app/os image
 - Creates virtual IPs in servers and maps to VLANs for client access.
 - Creates virtual HBAs in servers and maps to Zones, LUNs, and WWNNs for storage access



Session Number Presentation_ID

Grid or Utility Computing





Session Number Presentation_ID

VFrame



Session Number Presentation_ID

27

The Next Frontier: Datacenter Management Device and Service Provisioning + Virtualization



Session Number Presentation ID

Horizontal versus Vertical Provisioning



Session Number Presentation_ID

29

Datacenter Ecosystem



Session Number Presentation ID

What is InfiniBand?

- InfiniBand is a high speed low latency technology used to interconnect servers, storage and networks within the datacenter
- Standards Based InfiniBand Trade Association http://www.infinibandta.org
- Scalable Interconnect:
 - 1X = 2.5Gb/s (2Gb/s data)
 - 4X = 10Gb/s (8Gb/s data)
 - 12X = 30Gb/s (24Gb/s data)

High Performance Server Interconnect

- Industry Standard
- RDMA for Ultra-Low Latency
- 10Gbps Bandwidth (moving to 30Gbps)

Economics

Connection Oriented Control

Manageability

- Standard
- Connection Oriented
- Built-in Control
- Partitionable

- Boot Over IB
- Interconnect
 Agnostic
 Storage and I/O
 - -GigE, 10GigE, FC, iSCSI, etc.



Cluster Application Interconnect



Session Number Presentation_ID

Price / Performance Comparative

InfiniBand Offers the Best Price / Performance for HPC

| | InfiniBand PCI-Express | Myrinet D | Myrinet E | 10GbE | GbE | GbE/RNIC |
|------------------------------------|---------------------------|-----------|-----------|-----------|-------------|-------------|
| Data Bandwidth (Large Messages) | 950MB/s | 245MB/s | 495MB/s | 900MB/s | 100MB/s | 100MB/s |
| MPI Latency (Small Messages) | 5us | 6.5us | 5.7us | 50us | 50us | 18us |
| HCA Cost (Street Price) | \$550 | \$535 | \$880 | \$2K-\$5K | Free | \$500 |
| Switch Port | \$250 | \$400 | \$400 | \$2K-\$6K | \$100-\$300 | \$100-\$300 |
| Cable Cost (3m Street Price) | \$100 | \$175 | \$175 | \$50 | \$25 | \$25 |

Note: MPI "User Space" to "User Space" latency – switch latency is less

* Myrinet pricing data from Myricom Web Site (Dec 2004) utilizing Myrinet's latest switches ** InfiniBand pricing data based on Topspin avg. sales price (Dec 2004) *** Myrinet, GigE, and IB performance data from independent June 2004 OSU study **** 10GigE and GigE Cost and Performance data from Cisco Internal document

InfiniBand Performance Measured Results



InfiniBand Protocol Summary

| Protocol / Application | Summary | Application Example |
|--|--|---|
| IPoIB (IP over InfiniBand) | Enables IP-based applications to run over InfiniBand transport. | Standard IP-based applications. When used in conjunction with Ethernet Gateway, allows connectivity between IB network |
| SDP (Sockets Direct Protocol) | Accelerates sockets-based applications using RDMA. | and LAN 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. |
| uDAPL (Direct Access Programming Library) | Enables maximum advantage of RDMA flexible programming API. | Used for IPC communication between cluster nodes for Oracle 10G RAC. |
| MPI (Message Passing Interface) | Low latency protocol used widely in HPC environments. | HPC applications. |

The InfiniBand Driver Architecture



IP over InfiniBand

- Transmission of IP over Infiniband
 - Use IB as a link layer for IP
 - Define data link and link layer address
 - Encapsulation for ARP, IPv4 and IPv6
 - Address resolution
 - Transport IP multicast over IB
- Provides highest level of application compatibility.
- Applications do not need to be re-written or re-compiled
- Standard IP utilities and applications work as usual:
 - Ifconfig, ping, telnet, File sharing (NFS, CIFS); Login access (ssh, telent, etc); Cluster heartbeat
 - DHCP over IB
 - IP over InfiniBand MIB

How IP over InfiniBand works



* Notes: Uses standard Berkeley TCP/IP libraries

Sockets Direct Protocol

- Sockets Direct Protocol
- Runs socket based TCP/IP traffic with TCP and copy offload
- Highly configurable:
 - By process
 - By port
 - By destination
 - By environment variable
- No application recompile or rework necessary
- Zero copy capability using Asynchronous I/O (AIO)

Tangible Benefits of Server Fabric Switching

- Purchase 20X more compute power for same dollars (pay as you grow, moore's law, expense vs capitalization)
- 50% cost savings from resource consolidation delivers instantaneous ROI (Single Server Fabric- eliminates adapter, cables, ports)
- Dramatically Reduce TCO Manage enterprise-wide Server GRID centrally (wire once, *control* servers over the network)
- **Provision** New Servers in seconds, not days (or weeks)
- Help Eliminate Server Downtime (Failover provision, add/remove I/O or storage bandwidth on the fly)
- Control Ballooning Investments in Real Estate, Power & Cooling

(capitalize on dense server packaging and Blade architectures)

• Political Power and "Self-Rule" for the Server Team (Eliminate dependence on other teams to get apps provisioned quickly)

Cisco DNA Impact : Improved Server Utilization



Session Number Presentation_ID

Dramatic Cisco DNA Impact: Application Acceleration Improvement on Response Times

| Application | Software | Before | After | Transaction Time Reduction |
|----------------------|------------------|---------|---------|-------------------------------|
| Call center | PeopleSoft | 63 sec | 23 sec | □ 63% (□270%) |
| "JIT" manufacturing | SAP | 76 sec | 22 sec | □ 71% (□350%) |
| Store management | IBM WebSphere | 46 sec | 16 sec | □ 66% (□290%) |
| Claims management | IBM WebSphere | 42 sec | 19 sec | □ 55% (□220%) |
| Collaboration | Lotus iNotes | 90 sec | 28 sec | □ 68% (□320%) |
| Employee portal | Plumtree | 204 sec | 59 sec | □ 71% (□350%) |
| Portal consolidation | SunOne, Vignette | 43 sec | 6 sec | □ 85% (□670%) |
| Employee portal | PeopleSoft | 103 sec | 32 sec | □ 69% (□320%) |
| CRM | Siebel | 389 sec | 133 sec | □ 66% (□290%) |

Notes: Bandwidth reduction averages 80-90%

All timings are customer-verified using either LoadRunner, FineGround AppScope, or in-house customer tools

Topspin Building Blocks



Gateway Modules

- InfiniBand to Ethernet
- InfiniBand to Fibre Channel



Host Channel Adapter (HCA) With upper layer protocols



- SRPSDP
- uDAPLMPI
- IPolB

Linux and Windows driver support

Integrated System and Fabric management



The Cisco SFS Product Line



Cisco InfiniBand Blade Switch Modules

For IBM eServer BladeCenter



- Plug one card into each server blade
- Plug one or two switch modules into chassis
- Each server blade gets one or two 1x IB (2.5Gbps) connections
- Target markets: HPC, Multifabric I/O (MFIO), On-Demand data centers
- See IBM Redbook for more details:

http://www.redbooks.ibm.com/redpieces/abstracts/REDP39 49.html?Open

 Plug one daughter card into each server blade

For Dell PowerEdge 1855 Blades

- Plug one or two pass through modules into chassis
- Each server blade gets one or two 4x (10Gbps) IB connections
- Target markets: HPC, Multifabric I/O (MFIO), Scalable Enterprise centers

Case Study: Leading Research Facility High Performance Computing Cluster

- Application:
 - High Performance
 Computing Cluster
 - Compute time outsourced to Commercial Enterprises (major oil & gas)
- Environment:
 - 520 Dell Servers
 - 3:1 Blocking ratio
 - 6x SFS 7008 (TS270)
 - 29x SFS 7000 (TS120)
- Benefits
 - Compelling Price & Performance
 - Measured MPI latency 5.2µs

Core Fabric: 6x SFS 7008 (TS270)



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

Topspin Server Switch with Ethernet and Fibre Channel Gateways

Hitachi RAID Storage

SAN Switches

Ethernet Switches

Benefits:

20X Price/Performance Improvement over four years

30-50% Application Performance Improvement

Standards-based solution for on-demand computing

Environment that scales using 500-node building blocks



Case Study: Major System Vendor Utility Computing Service

• Application:

Build scalable "on-demand" compute service for enterprise customers (license \$/CPU)

Key initiatives around Financial Services and Energy verticals

Environment

1024x Sun V20z Nodes

34 TS270 Server Fabric Switches

Non-Blocking CLOS Network

8 TS360's with Gateways

Sun Storage

Enterprise-Class Reliability

• Benefits:

Ability to outsource computing services to many customers with common infrastructure



Case Study- Large Government Lab Worlds 2nd Largest Super Computer

- Application:
 - High Performance SuperComputing Cluster
- Environment:
 - 4096 Dell Servers
 - 50% Blocking Ratio
 - 8 TS 740s
 - 256 TS120s
- Benefits:
 - Compelling Price/Performance
 - Largest Cluster Ever Built (by approx. 2X)
 - Expected to be 2nd Largest Supercomputer in the world



Leading UK Telecom Provider Oracle 10g Deployment (BT)

- Broadband billing application
- 12 identical regional deployments
- Running Oracle 10g
- No single-point of failure
- Each server has a single HCA, with ports dual connected to two TS90s



Sandia National Labs – 4096 Nodes Cluster

- Application:
 - High Performance SuperComputing Cluster
- Environment:

4096 Dell Servers 50% Blocking Ratio 8 SFS 7048 256 SFS 7000's

• Benefits:

Compelling Price/Performance

Largest IB Cluster ever built

Expected to be 3rd Largest Supercomputer in the world



Oracle 10g: Broad Scope of IB Benefits



Bio-Informatics Cluster: 1,066 Node Supercomputer

1,066 Fully Non-Blocking Fault Tolerant IB Cluster



Key decision factors:

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

Cisco InfiniBand Landscape Vendors working with Cisco Server Switches



Session Number Presentation ID

Topspin and Top Tier Server Vendors



Session Number Presentation_ID

Who Owns the Datacenter?



Presentation ID

© 2005 Cisco Systems, Inc. All rights reserved.



Thank You!

Session Number Presentation_ID Cisco Public 58