Ethernet Network Storage: Getting More For Less

By Patrick Khoo Program Manager Modular Connected Storage Architecture Group Network Storage Technology Division Data Storage Institute





About DSI and NST Division



DSI is a nationally funded non-profit R&D institute focusing on data storage technologies and industries



Can I do More for Less?

- University of California at Berkeley 2001
 - 12 Exabytes in mankind's history to date
 - 12 more Exabytes in next two and a half years alone!



Network Storage to the Rescue! But Does Network Storage **REALLY** Help?



Definitions

- LAN Local Area Network
- DAS Direct Attached Storage
- NAS Network Attached Storage
- SAN Storage Area Network

Components

- Servers
- Storage systems (eg. disk arrays, tape libraries, etc)
- Interconnect technologies (eg. fibre optic cables, switches etc)
- Host-bus Adapters (HBA), Network Interface Cards (NIC)
- System and Data Management Software

Why Build a SAN?

institute



The truth is, there is <u>NO</u> killer app, so stop waiting for one! But there are plenty of reasons to adopt Network Storage!

Industry Forecasts – Comparisons





Industry Forecasts – Comparisons



Asia/Pacific USD/TB Storage System Costs, 2000-2006

Source: IDC Asia/Pacific, 2001, 2002



Trends in Storage Price Erosion



Source: IDC Asia/Pacific, 2002



Industry Forecasts – The Difference

- What is the effect of the current economic climate on the storage business?
 - Storage growth has not been halted, merely delayed (by about one year)
 - At the same time, the cost of storage is dropping about 35% per year but will still stablise over the long term [IDC Asia Pacific 2002]
- Conclusion: The cost of storage <u>systems</u> has not really dropped (other than due to newer high density HDDs) – some kind of technological advancement is needed



So How?

Innovate!

. . . and think of new ways to do old things.

© 1998, Michigan Live Inc. All rights reserved BOIN When fish bungee jump.



SAN Islands are a Reality



Number of SANs Deployed in Respondents' Organizations

- Most organisations preparing or considering implementing SANs, should take into account adding new "independent" SANs in the future
- Most SANs are "Local" in nature
- Why do you need IP to go long distance? (for most people anyway)





Data Storage nstitute Why not use Ethernet?

"Ethernet the World!"



Data Storage Institute

Network and Storage Differences





You can't compare an Ethernet cable with a SCSI cable, SCSI cables transmit data in parallel! Overheard from a computer science professor

Conclusion:

Storage systems are very different from Network systems

Corollary:

Network systems providing Storage must therefore be designed differently from *normal* Network systems



Combine Two Worlds . . .



Data Storage

stitute

. . . And Think Out of the Box



2 Data Storage Institute

Yes, It Can Be Done!

HyperSCSI

The transmission of SCSI commands & data across a network and multi-technology device support



Access storage over a network

Data Storage

HyperSCSI is a new open source Network Storage Protocol

- Transmit SCSI commands and data over a network
- High performance, secure, simple, low cost solution
- Runs directly on Ethernet (No TCP/IP!)



The HyperSCSI Protocol

HyperSCSI Packet Framing / Encapsulation on Ethernet



HyperSCSI Command and Data Block

No TCP/IP!

Routeable? Secure? Reliable?





"Stealing" Components

	Network	Storage	HyperSCSI	
Flow Control	Sliding window	"Buffer Credit"-based *	Dynamically sized fixed window	
Transmission	Stream-based	Block-based	Block-based	
Data Delivery	Guaranteed Guaranteed		Guaranteed	
Channels	Single-channel (vendor specific trunking) Parallel transmission		Vendor independent multi-channel	
Addressing	"Unlimited"	Limited	Almost "Unlimited" **	
Device Discovery	Lookup-based	Broadcast-based (Bus Scan)	Broadcast-based (Local-area)	
Authentication	Multi-user challenge & response	Physical security or Zone/LUN Masking	Single-user challenge & response	
Tx Security	(Add-on) Encryption	Physical security	(Built-in) Encryption	
Scalability	"Unlimited"	Limited	Almost "Unlimited" **	
Access	Wide-area	Local-area	Local-area	



This is not meant to be a complete or accurate depiction of the components or mappings, but merely as an illustration of the differences between systems and components

Flow Control

A*STAR



Easy Management

A+STAR

🛃 137.132.29.10 - PuTTY		
hs-server/etc/hscsi: hs-ser	ver status	
** Displaying Status		
HYPERSCSI server module - 2	0020725	
Status: HyperSCSI server mo	dule has been initialized!	
HS DEVICE CONFIGURATION SUM	MARY:	
CONFIGURATION:		
Name Device-type Host Ch	annel ID LUN Capacity(KB) Group_Name	
hat The high o		
nda IDE-DISK O	A 137.132.29.10 - PuTTY	
CONNECTION:	hs-client/root: hs-client status	
Mac-address Name Device-t	** Displaying Status	
	HYPERSCSI client module - 20020725	
	Status: HyperSCSI client module has been initialized!	
** Displaying Recent Messag	HyperSCSI DEVICE CONFIGURATION SUMMARY:	
Network device ethO is goin		
The network MTU is 1497	CONNECTION:	
There is no SCSI Host.	Mac-address Device_ID Name Host Channel ID LUN Group_Name	
hscsi_pkt_window: 5		
	0001034476E7 O sda 1 O O MCSA-HYPERSCSI	
HyperSCSI Server module ini	그는 그는 것 같은 것 같아. 같은 것 같은 것 같아. 그는 것 같아. 같은 것은 것은 것은 것 같아.	
	** Displaying Recent Messages	
ns-server/etc/nscs1:	Network device etho is going to be used.	
	nscsi_sg_taplesize: 16	
	[ns client] ic_window : s	
	hypersest crient module inicialized:	
	hs-client/root:	
40		
Data Storage		
Institute		

Easy Management

🚰 137.132.29.10 - PuTTY	ويتابعها وتدولا وعقاليه زوار والأوسا			
hs-server/etc/hscsi: hs-server status				
** Displaying Status				
HYPERSCSI server module - 20020725				
Status: HyperSCSI server module has been initia	alized!			
HS DEVICE CONFIGURATION SUMMARY:	[HYPERSCSI-SERVER-CONFI	G-VERSION-20021024]		
	# Sample Config File for HyperSCSI Server - Modify Before Use!			
CONFIGURATION:	# Optimised for Fast Et	hernet		
Name Device-type Host Channel ID LUN Capac:	# Last Updated - 24 July 2002			
hda IDE-DISK O OOO 19	[ADD]			
CONDECTION				
Meg-address News Device-ture Hest Chennel	[MODULE_DEF]			
hac-address Name Device-cype host channel.	# For GE, try PKT_WINDC	W_SIZE 32		
0001025868F7 bda IDF_DISK 0 0	# PKT_WINDOW_SIZE:	32		
ST ST	PKT_WINDOW_SIZE:	5		
	MULTI_RCV_THREAD:	2		
** Displaying Recent Messages	MULTI_XMIT_THREAD:	3		
Network device ethO is going to be used.	REXMIT_COUNT:	2		
The network MTU is 1497	DIRECI_MC.	0		
There is no SCSI Host.	[קיס ד טען			
hscsi pkt window: 5		SDV		
	VOL_1:	SDA		
HyperSCSI Server module initialized!	[NETWORK DEF]			
	LAN 1:	ETHO		
hs-server/etc/hscsi:				
	[GROUP_DEF]			
	GROUP_NAME:	MCSA-HYPERSCSI		
	PASSWORD:	0123456789		
	NET:	LAN_1		
	IP_ON:	0		
	VOL_NAME:	VOL_1		
*	VOL_OPT:	0:0		
Data Storage				

Secure Storage

C:\My Documents\hacking-log.txt - Viewer _ | **D |** × Files Edit Search View Options Coding Help /var/log/messages ______ Nov 18 04:02:01 nst syslogd 1.4-0: restart. ----> Hacker Start Nov 19 04:18:59 nst ftp(pam unix)[20988]: check pass; user unknown Nov 19 04:18:59 nst ftp(pam unix)[20988]: authentication failure; logname= uid=0 euid=0 tty= ruser= rhost=AMarseille-101-1-1-243.abo.wanadoo.fr Nov 19 04:19:01 nst ftpd: AMarseille-101-1-1-243.abo.wanadoo.fr: connected: IDLE [20988]: failed login from AMarseille-101-1-1-243.abo.wanadoo.fr [193.252.177.243 Nov 19 04:19:12 nst ftpd: AMarseille-101-1-1-243.abo.wanadoo.fr: connected: IDLE [20988]: lost connection to AMarseille-101-1-1-243.abo.wanadoo.fr [193.252.177.24 31 Nov 19 04:19:12 nst ftpd: AMarseille-101-1-1-243.abo.wanadoo.fr: connected: IDLE [20988]: FTP session closed ----> Hacker Stop ----> Hacker Start Nov 20 03:40:01 nst ftpd[23171]: ACCESS DENIED (not in any class) TO AMarseille-1 01-1-1-243.abo.wanadoo.fr [193.252.177.243] Nov 20 03:40:01 nst ftpd[23171]: FTP LOGIN REFUSED (access denied) FROM AMarseill e-101-1-1-243.abo.wanadoo.fr [193.252.177.243], anonymous Nov 20 03:40:02 nst ftpd[23171]: FTP session closed ----> Hacker Stop



This can't happen to your storage, IF your storage doesn't have TCP/IP

High Performance



6

20

0

Benchmark - cp

Write (MB/s)

Rewrite (MB/s)

Read (MB/s)

Reread (MB/s)

NFS iSCSI

HyperSCSI

Tests conducted on a clean Gigabit Ethernet network with Jumbo frames, single initiator and target, 8 hard disks configured in RAIDO and using only common off-the-shelf hardware and software without special tweaks or optimisations.

ata Storage istitute

Lower Embedded Overheads





Single drive on Embedded Board running Embedded Linux Slow GE Performance is noted due to poor memory speed on IQ80310 EVB (well documented HW bug)

File Systems and Multiple Clients



Multi-Channel Technology



Multi-Channel Multi-Client





RAIDO on Client, Two HDDs on Server, One GE channel on server, add one FE channel on each client in pairs

- * Weighted Fair Queuing algorithm used
- * Max GE Performance is measured from one GE channel between client and server

Key Features

- Runs on raw Ethernet does not use TCP/IP
- Supports various storage devices (including disk, optical, removable media and tape) and interface technologies (including SCSI, Fibre Channel, IDE and USB)
- Supports various network technologies (including Fast Ethernet, Gigabit Ethernet, GE + Jumbo Frames and Wireless LAN 802.11b)
- HyperSCSI on Ethernet (HS/eth) can be deployed on a multi-protocol network environment safely
- Includes built-in 128-bit Encryption
- Supports active device discovery for plug and play ad-hoc network storage
- Independent of hardware or vendor products runs with wide variety of disk/tape/optical devices and arrays, SCSI/FC HBAs, NICs, and switches
- Designed for easy deployment, and above all, to provide users with the freedom of choice, to implement network storage the way they want to, with the options they need



What Others Have to Say

"That is a very good news that you managed to get HyperSCSI running over Wireless LAN. Sustaining ~100MB/s throughput is very impressive. I am very impressed by your work. We are rebuilding our cluster nodes to use 2.4.16 kernel and will install the HyperSCSI software soon. I cannot wait to see HyperSCSI working in our test bed."

"I had been researching a solution for about a month before I stumbled across HyperSCSI. I started out looking at FiberChannel but the protocol is cryptic, the Linux support is poor, and I object to the cost being 6-20x that of Gig-Ethernet when its basically the same technology. I looked at iSCSI but again the protocols are over engineered, the Linux support is poor, the cost of equipment is robbery, and there are a lot of research papers that suggest block protocols over TCP/IP are a poor choice. The Linux network block driver is a hack and not a very good one. I looked into doing sharing with a SCSI bus, but the lack of target mode support in Linux killed that idea quickly. So I started researching how to do raw Ethernet access in Linux with the intent of writing a Ethernet based protocol to allow machines concurrent access to raw disk. And then I found HyperSCSI. Its simple. Its elegant. You can implement it using commodity hardware. It works today. End of search."



Porting to Windows 2000 / XP

👪 MS-DOS Prom	pt	
∄ r 6 x 10 .		
Packet length, c: 00000000 : 31 31 00000010 : 35 36 00000020 : d8 d8 00000030 : d8 d8	aptured portion: 60, 60 31 31 31 31 00 04 76 de 84 4a 30 37 38 39 30 31 32 33 34 35 36 3 d8 d8 d8 d8 d8	▲ 1 32 33 34 111111v .JJ1234 7 38 39 30 5678901234567890 3 d8 d8 d8
Packet length, c: 00000000 : 31 31 00000010 : 35 36 00000020 : 00 00 00000030 : 00 00	aptured portion: 60, 60 31 31 31 31 00 04 76 de 84 4a 33 37 38 39 30 31 32 33 34 35 36 33 00 00 00 00 00 00 00 00 00 00 00 00 00	1 32 33 34 111111v┃.J1234 7 38 39 30 5678901234567890 0 00 00 00
Packet length, c: 000000000 : 31 31 00000010 : 35 36 00000020 : 28 28 00000030 : 28 28	aptured portion: 60, 60 31 31 31 31 00 04 76 de 84 4a 33 37 38 39 30 31 32 33 34 35 36 33 28	1 32 33 34 111111v .J1234 7 38 39 30 5678901234567890 3 28 28 28 ((((((((((()(())
4		▼ ▶

- Virtual SCSI Card
- HyperSCSI client on Windows





Is This for Real? <u>Yes</u> for Consumers!

-	2000	PTT ERICSSON			
(tref, labbr (they)	: 192. 568. 11.64)	4.3	NOW	ERL.	*
e cat. An rither an Wester and Table Jar Hagar (11)	ner/meni/meni metropoli 10 Onorrel i 00 Idi Dilandi funeli 0 Dilandi funeli 0 Dilandi funeli 0 Dilandi funeli 0 Dilandi Control Marti	an Lant (0) C 1945-CD2, M M ² 2y DH SIMSHI 1	1 Hill	n let m	
The - 400 constitute a function without # []	ers Device 10 for		1		



- Wireless Network Storage
- Personal and Home Networks
- Consumer Electronics
- Entertainment and Content Distribution





Is This for Real? Yes for Corporates!





- Storage Area Networks (SAN) and Network Attached Storage (NAS)
- Information Continuance, Performance, Security and Reliability



Is This for Real? Yes for Clusters!

- Storage for HPC / Grid
- Remote Boot, Concurrent Access, Cluster File Systems







Is This for Real? Yes for Convergence!



Data Storage

Is This for Real? Yes for Production!





The Network Storage Vision



Institute

Conclusion

- Some key points:
 - New Technologies like HyperSCSI, iSCSI, TOE, iSCSI HBA, blah, blah, blah
 - New applications for Network Storage like Wireless LAN, Home Networks, Personal Networks
 - Lower costs, ease of use, network storage for everyone!

Ethernet Storage is **NOT** coming, It's already here!

And yes, you WILL get more for less



Thank You

http://nst.dsi.a-star.edu.sg/mcsa/



