

IBM System & Technology Group

Scale out File Services – SOFS Technology Presentation

Sven Oehme (oehmes@de.ibm.com) STG Industry Technology Leader Team

© 2007 IBM Corporation





1	NAS today
2	The Scale-out File Services (SoFS) approach
3	General Parallel File System (GPFS)
4	Inside SoFS
5	Conclusion



© 2007 IBM Corporation







Customer Pain points

- Customers have increasing demand for file services. The growth of so-called "unstructured data" outpaces that of databases by far.
 - But current NAS solutions do not scale. File systems are limited to a few Terabyte.
 - So customers have to add box by box and manage them individually
 - No way to apply Policys across this independent data islands
- Parallel access to data gets more and more common, especially in the digital media space
- We see requirements on data access rates and response times to individual files that have been unique to HPC environments before
- Making file servers clustered is easy in theory but hard to build and maintain, there
 is too much complexity
- Migration, integration or removal of storage for file services is a nightmare so a lot of legacy systems still exist and need to be maintained
- Backup windows are a big issue and get worse while the amount of data continuously increase
- Integration of ILM functions into the file system becomes more important as the amount of TB's explode





1	NAS today
2	The Scale-out File Services (SoFS) approach
3	General Parallel File System (GPFS)
4	Inside SoFS
5	Conclusion



IBM Systems & Technology Group







The SoFS Advantage





Global Namespace





Multi-node CTDB on the Example of CIFS





Does Performance matter (single node)?

[root@gpfsmum9 samba4]# bin/smbtorture //localhost/test -Uadministrator%test01 --option torture:readonly=1 --num-progs=3 --loadfile=/root/torture/read.dat BENCH-NBENCH -t30

Using seed 1166694753

Running for 30 seconds with load '/root/torture/read.dat' and warmup 1 secs

Starting 3 clients

3 clients started

3	5695	306.45 MB/sec	execute 1 sec	latency 2029.00 usec
3	9092	303.82 MB/sec	execute 2 sec	latency 1839.00 usec
3	12168	293.59 MB/sec	execute 3 sec	latency 30516.00 usec
3	17312	333.85 MB/sec	execute 4 sec	latency 33788.00 usec
3	24564	395.89 MB/sec	execute 5 sec	latency 30464.00 usec
3	31892	437.97 MB/sec	execute 6 sec	latency 791.00 usec
3	39214	468.40 MB/sec	execute 7 sec	latency 867.00 usec
3	46554	491.01 MB/sec	execute 8 sec	latency 1101.00 usec
3	53896	508.77 MB/sec	execute 9 sec	latency 641.00 usec
3	61267	523.24 MB/sec	execute 10 sec	latency 1009.00 usec
3	68578	534.58 MB/sec	execute 11 sec	latency 857.00 usec
3	76020	545.05 MB/sec	execute 12 sec	latency 934.00 usec
3	83376	553.18 MB/sec	execute 13 sec	latency 404.00 usec
3	90737	560.36 MB/sec	execute 14 sec	latency 696.00 usec
3	98060	566.19 MB/sec	execute 15 sec	latency 735.00 usec
3	105427	571.72 MB/sec	execute 16 sec	latency 996.00 usec

3	112/95	576.43 MB/Sec	execute 1/ sec	latency 904.00 usec
3	120078	580.34 MB/sec	execute 18 sec	latency 760.00 usec
3	127404	583.95 MB/sec	execute 19 sec	latency 957.00 usec
3	134717	587.14 MB/sec	execute 20 sec	latency 1118.00 usec
3	142020	590.00 MB/sec	execute 21 sec	latency 762.00 usec
3	149430	593.01 MB/sec	execute 22 sec	latency 982.00 usec
3	156699	595.27 MB/sec	execute 23 sec	latency 915.00 usec
3	163959	597.24 MB/sec	execute 24 sec	latency 928.00 usec
3	171264	599.31 MB/sec	execute 25 sec	latency 791.00 usec
3	178678	601.47 MB/sec	execute 26 sec	latency 1092.00 usec
3	186033	603.40 MB/sec	execute 27 sec	latency 894.00 usec
3	193390	605.08 MB/sec	execute 28 sec	latency 916.00 usec
3	200717	606.67 MB/sec	execute 29 sec	latency 855.00 usec
3	208014	607.97 MB/sec	cleanup 30 sec	



Saturated Windows Network link

DiskSpeed	
Disk Test File Test Data Test Settings	
File: \\192.168.11.35\ale\00-09nov1443.mxf	
Description Performs read performance tests on specified file	E Quick
Disk: 00-09nov1443.mxf 752 MB file from/to main memory async. access with command queue len 10 Block Read Lin Read Rnd Write Lin	
kB MB/sec MB/sec MB/sec +	40 30





1	NAS today
2	The Scale-out File Services (SoFS) approach
3	General Parallel File System (GPFS)
4	Inside SoFS
5	Conclusion





GPFS: Overview

- IBM's General Parallel Filesystem available 1996
- Used in other IBM Products (BIA for SAP, VTS)
- Product available on AIX 5.1 (Power5) and Linux (x86/Power5) clusters.
- Used on many of the largest supercomputers in the world.
 - *Cluster*: 1000+ nodes, fast reliable communication, common admin domain.
 - *Shared disk*: all data and metadata on disk accessible from any node through disk I/O interface.
 - Parallel: data and metadata flows from all of the nodes to all of the disks in parallel.
- High performance
 - Multi-Terabyte files, Multi-Petabyte file systems.
 - Wide striping, large blocks, many GB/s to single file
- Highly Reliable
 - Can survive Disk and Node failures
 - Allows Split site Configurations





GPFS: ASC Purple/C Supercomputer (2005)

- 1536-node, 100 Teraflop IBM BlueGene/L cluster at Lawrence Livermore National Laboratory
- 2 PB GPFS file system (one mount point)
- 500 RAID controller pairs, 11000 disk drives
- 126 GB/s parallel I/O measured to a single file (134GB/s to multiple files)









GPFS: Information Lifecycle Management

- ILM Support introduced with GPFS 3.1
 - Storage pool group of LUNs
 - Fileset define subtrees of a file system
 - Policies for rule based management of files inside the storage pools
- What does it offer
 - One global file system name space across a pool of independent Storage
 - Files in the same directory can be in different pools
 - Files placed in storage pools at create time using policies
 - Files can be moved between pools for policy reasons
 - Can be used for hierarchical arranged storage based on files
 - Allows classification of data according to SLAs





GPFS: Storage Policies

- Rules to control the placement, migration, and retention of files
- Declarative SQL-like language
- Rule types:
 - Placement policies, evaluated at file creation, example
 rule 'hq' set pool 'gold' for fileset 'hqfs'
 rule 'otherfiles' set pool 'silver'
 - Migration policies, evaluated periodically

```
rule 'cleangold' migrate from pool 'gold'
   threshold (90,70) to pool 'silver'
```

- rule 'hsm' migrate from pool 'sata' threshold(90,85) weight(current_timestamp access_time) to pool 'hsm' where file_size > 1024kb
- rule 'cleansilver' when day_of_week()=Monday
 migrate from pool 'silver' to pool 'bronze'
 where access_age > 30 days
- Deletion policies, evaluated periodically
 - rule 'purgebronze' when day_of_month()=1 delete
 from pool 'bronze' where access_age>365 days









GPFS: Replication & Cross-cluster mounts

- Synchronous Replication
 - Synchronous intra-cluster replication on Blocklevel handled by GPFS nodes
 - Replication can be set up for subsets of the file system
- Cross-cluster mounts
 - Multiple GPFS clusters can mount remote file system owned by other clusters
 - All locking and metadata operations are routed through the owning cluster, block I/O can be done directly by the cross-mounting cluster via a SAN or fast Interconnect (IB, myrinet)
 - The cross-mounting cluster can mount the file system into one of its own file systems (single mount point).



GPFS snapshots

/fs1/file1
/fs1/file2
/fs1/subdir1/file3
/fs1/subdir1/file4
/fs1/subdir2/file5

/fs1/file1
/fs1/file2
/fs1/subdir1/file3
/fs1/subdir1/file4
/fs1/subdir2/file5
/fs1/.snapshots/snap1/file1
/fs1/.snapshots/snap1/file2
/fs1/.snapshots/snap1/subdir1/file3
/fs1/.snapshots/snap1/subdir1/file4
/fs1/.snapshots/snap1/subdir2/file5

Creating a snapshop

mmcrsnapshot fs1 snap1

Writing dirty data to disk Quiescing all file system operations Writing dirty data to disk again Creating snapshot. Resuming operations.

Read-only copy of directory structure and files

Only changes to the original file consume disk space



GPFS snapshots on Windows

•Integrated into Windows Explorer

NameTimedata1 on 'SauFriday, October 26, 2007, 11:00 PMdata1 on 'SauYesterday, October 30, 2007, 12:00 AMdata1 on 'SauToday, October 31, 2007, 12:00 AMdata1 on 'SauToday, October 31, 2007, 12:05 AMdata1 on 'SauToday, October 31, 2007, 3:25 AMdata1 on 'SauToday, October 31, 2007, 3:27 AMdata1 on 'SauToday, October 31, 2007, 3:29 AMdata1 on 'SauToday, October 31, 2007, 3:29 AM	neral Security Pre To view a version fro You can a restore a Folder versions:	evious Versions a previous version of a folder, select the om the following list and then click View. also save a folder to a different location or previous version of a folder.
data1 on 'SauFriday, October 26, 2007, 11:00 PMdata1 on 'SauYesterday, October 30, 2007, 12:00 AMdata1 on 'SauToday, October 31, 2007, 12:00 AMdata1 on 'SauToday, October 31, 2007, 12:05 AMdata1 on 'SauToday, October 31, 2007, 3:25 AMdata1 on 'SauToday, October 31, 2007, 3:27 AMdata1 on 'SauToday, October 31, 2007, 3:27 AMdata1 on 'SauToday, October 31, 2007, 3:29 AM	Name	Time
 data1 on 'Sau'S Yesterday, October 30, 2007, 12:00 AM data1 on 'Sau Today, October 31, 2007, 12:00 AM data1 on 'Sau Today, October 31, 2007, 12:05 AM data1 on 'Sau Today, October 31, 2007, 3:25 AM data1 on 'Sau Today, October 31, 2007, 3:27 AM data1 on 'Sau Today, October 31, 2007, 3:29 AM data1 on 'Sau 	🜌 data1 on 'Sau.	Friday, October 26, 2007, 11:00 PM
Image: Second	🔀 data1 on 'Sau.'\S	Yesterday, October 30, 2007, 12:00 AM
 data1 on 'Sau Today, October 31, 2007, 12:05 AM. data1 on 'Sau Today, October 31, 2007, 3:25 AM data1 on 'Sau Today, October 31, 2007, 3:27 AM data1 on 'Sau Today, October 31, 2007, 3:29 AM 	📚 data1 on 'Sau	Today, October 31, 2007, 12:00 AM
 data1 on 'Sau Today, October 31, 2007, 3:25 AM data1 on 'Sau Today, October 31, 2007, 3:27 AM data1 on 'Sau Today, October 31, 2007, 3:29 AM 	훒 data1 on 'Sau	Today, October 31, 2007, 12:05 AM
Sata1 on 'Sau Today, October 31, 2007, 3:27 AM data1 on 'Sau Today, October 31, 2007, 3:29 AM	🧝 data1 on 'Sau	Today, October 31, 2007, 3:25 AM
🛣 data1 on 'Sau Today, October 31, 2007, 3:29 AM	훒 data1 on 'Sau	Today, October 31, 2007, 3:27 AM
	\overline 🗟 data1 on 'Sau	Today, October 31, 2007, 3:29 AM
View Copy Restore		View Copy Restore





1	NAS today
2	The Scale-out File Services (SoFS) approach
3	General Parallel File System (GPFS)
4	Inside SoFS
5	Conclusion





SoFS component view

- GPFS 3.2 IBM's High end clustered file system
- CTDB Clusterd Trivial database Daemon, Controls the cluster and the file service daemons
- Enhanced Samba Server to provide CIFS export
- RHEL5.1 Base OS, provides NFS, FTP and HTTP daemons
- SoFS Package Provides Mangement GUI, Apache file server module, acceleration tools, etc.
- IBM Hardware





SoFS Storage Hardware

- DCS9950 Support *
 - Mid-range: Equipped with up to 960TB using SATA disks or 448TB with FC, RAID-6, 8 FC host ports, 3 GB/sec read/write speed for sequential I/O
- DS 3200 *
 - Entry level: Equipped with up 14.4 TB using SAS disks, up to 1 GB cache, 6 SAS host ports
- DS 3400
 - Entry level: Equipped with up 14.4 TB using SAS disks, up to 1 GB cache, 4 FC host ports
- DS 4200
 - Mid-range: Equipped with up to 84 TB using SATA disks, 2 GB cache, 4 FC host ports.
- DS 4700
 - Mid-range: Equipped with up to 84 TB using SATA disks or 33.6 TB using FC disks, up to 8 GB cache, up to 8 FC host ports.
- DS 4800
 - Mid-range: Equipped with up to 168 TB using SATA disks or 67.2 TB using FC disks, up to 16 GB cache, 8 FC host ports
- DS 8x00

Enterprise: Equipped with up to 512 TB using FATA disks or 307.2 TB using FC disks, up to 32 GB cache, up to 128 FC host ports

*Available Q1/08





SoFS BladeCenter Hardware

- BladeCenter H chassis
 - Up to 6 Ethernet switch modules
 - Either GbE or 10 GbE uplinks
 - Two 4 Gb Fibre Channel switch modules
 - Two management modules
 - Up to 14 HS-21 blades (one dedicated for management)
 - Intel Xeon Quadcore, 8 GB RAM





SoFS Rack mounted Hardware

- X3650 *
 - 2 * Intel Xeon Quadcore
 - 2 GigE Adapter On-board
 - Up to 6* Ethernet or FC Adapter optional
 - 4 48 GB RAM



*Available Q1/08



Why to use the IBM BladeCenter ?

- Management efficiency:
 - All network interconnect (Ethernet & Fibre Channel) is hard-wired on the midplanes. No chance for wrong or missing cabling.
 - (In a SoFS cluster with 14 blades the midplanes replace 112 network cables)
 - Blades can be discovered, monitored, started, stopped etc. through BC management modules.
 - Switches are an integrated part of the cluster and can be monitored through BC management modules.
- Space efficiency: Fits 14 servers in an 9U package
- Power efficiency: The BC power modules up to 50% more efficient then smaller power supplies found in rack-mounted servers.
- Scalability: Add up to 14 blades, the infrastructure (power, cooling, network, management) is already there.
- Price: Initial investment looks expensive, but equipped with approximately >=5 blades the TCO is lower than with rack mounted servers.





SOFS CIFS Enhancements

- Clustering
 - Multiple exports of the same file system over multiple nodes including distributed lock, share and lease support
 - Failover capabilities on the server no Client side changes needed
 - Integration with NFS, FTP, HTTP daemons in regard of locking, failover and authorization
- Performance optimization for GPFS backend
- NTFS ACL Support in Samba using the native GPFS NFSv4 ACL Support
- HSM support within Samba to allow destaging of files to tape and user transparent recall.
- Simple install and configuration tools
- Snapshot Support Integrated into the Windows Explorer







Simple SoFS single-site setup





SoFS with synchronous replication on-site





SoFS with synchronous replication across sites





SoFS with cross-cluster mount





Availability Management

🌒 Integrated Solutions Console - Mozilla F	irefox						_ = ×
<u>File E</u> dit <u>V</u> iew Hi <u>s</u> tory <u>B</u> ookmarks <u>T</u> ools	<u>H</u> elp					💷 oehmes@goo	glemail.com - 🔿
 	ost/ibm/console/login	.do?action=secure			🚳 🔹 🕨 🕞	Google	Q)
Integrated Solutions Console welcome root					Help Logout		IBM.
Welcome	Node Availability						
Console Settings	Node List						- 0
Scale-Out File Services (SOFS) Customer Information Node Availability System Utilization GPFS Management	Node availab	ility report.					
GPFS Cluster Selection Electron Cluster Management							
Create new GPFS Cluster	overall state ♀	Host name≎	Description 🗘	IP Address 🗘	Conn. state 🗘	GPFS state 🛇	Choose
Self-Service Space Creation	Уок	node1		9.155.61.20	ок	active	
	₩ ^{OK}	node2		9.155.61.21	ок	active	
	У ОК	node3		9.155.61.22	ок	active	
	WOK.	node4		9.155.61.23	ок	active	
	Total: 4					Fri Jun 22 07:37	7:00 EDT 2007
	0						
	Node History Selected node: 0/22/07 4:58:01 AM or 0/22/07 4:57:51 AM or	nn.state GPFS state	Select a chart:	Overall Availability I	II availabil	ity	ок
Done				localh	iost 🚘 🝓 1.160s	🔓 M 0 oehme	s@googlemail.com



Node Management

🥘 Integrated Solutions Console - Mozilla	Firefox					_ = ×
<u>F</u> ile <u>E</u> dit <u>V</u> iew Hi <u>s</u> tory <u>B</u> ookmarks <u>T</u> ools	<u>H</u> elp			¦i oe	hmes@googler	nail.com 🔹 🔿
 	host/ibm/console/login.do?action=se	cure		🚔 🔹 🕨 💽 🕞 Goog	gle	Q)
Integrated Solutions Console welcome root				Help Logout		IBM.
Welcome	GPPS Node Page					
Console Settings	GPPS Nodes					2 - 🗆
Scale-Out File Services (SOFS)						
Customer Information						
Node Availability System Utilization	Active Cluster: apfs1.fscc.m	nainz.de.lbm.com Active	e node: node1			
GPF5 Management						
GPFS Cluster Selection	Start Stop Roman					
Cluster Management						
Cluster Configuration						
Disks & NSDs	Select Name Description	OS Environment	IP Address	GPFS Status	Status	Choose
IGPFS Nodes	node1	SLES 9 1386	9.155.61.20	active	ок	
GPFS Quotas Delice based data management	node2	SLES 9 1386	9.155.61.21	active	ок	
Create new GPFS Cluster	node3	SLES 9 1386	9.155.61.22	active	ок	
Gelf-Service Space Creation	node4	SLES 9 1386	9.155.61.23	active	ок	ā
	Total: 4				Fri Jun 22 07:39:00 ED	OT 2007
	Node Details GPFS Node Settings HostName: Description: Description:	node1		Status History 	jsta act	7 – D
	Product Version: Client: NodeNumber: OsName: DaemonAdress: HostAdress: Manager: Quorum: UserName:	3.1.0.7 false 1 5LE9 91306 9.155.61.20 9.155.61.20 • root				
https://localhost/ibm/console/navigation.do?pa	Password: Apply Settings geID=com.ibm.fscc.gpfsgui.nav-node	 es&moduleRef=com.ibm.fscc.	gpfsgui localhost	🖻 🖏 1.237s 🛛 🖓	M 0 oehmes@g	ooglemail.com



Filesystem Management

🔋 Integrated Solutions Console - Mozilla Firefox 📃 🗆 🗙											
<u>F</u> ile <u>E</u> dit <u>V</u> iew Hi <u>s</u> tory <u>B</u> ookmarks <u>T</u> ools	Help						<mark>¦i=</mark> oehme	s@google	mail.com 🔹 🔿		
🐗 • 🗼 - 💽 💿 🏠 🗋 https://locall	host/ibm/console/login.do?	action=secure				<u>-</u>	G• Google		Q)		
Integrated Solutions Console welcome root						Help Log	jout		IBM.		
Welcome	GPPS Filesystem Page								^		
Console Settings	File Systems								2 - 0		
Scale-Out File Services (SOFS)											
Customer Information											
Node Availability System Utilization Active Cluster: gpfs1.fscc.mainz.de.lbm.com Active filesystem: gpfs0											
System Utilization Active Cluster: gpfs1.fscc.mainz.de.lbm.com Active filesystem: gpfs0											
GPES Cluster Selection	I										
Cluster Management	Mount Unmount Remo	we									
Cluster Configuration File Systems											
Disks & NSDs	Select Name	Mountpoint	Size	Usage	Inode usage	Pool usage	Pool usage ma	* Mounted	Choose		
GPFS Nodes	apho	/anfs	2.10 TB	258.28 GB	14	12	12	4 Nodes			
Policy-based data management		· = = -				-		st of Filesyst	ems		
Create new GPFS Cluster	<u>qps</u>	167.77 GB	140.99 MD	2	0	0	4 nodes				
Self-Service Space Creation	gpfs3	/gpfs3	68.08 GB	165.12 MB	6	0	0	4 Nodes			
		/gpfs4	115.34 GB	1.21 GB	7	0	1	3 Nodes			
	Total: 4						Fri Jun	22 07:39:02 EDT	2007 🗘		
	1. <u></u>										
	Add Filesystem										
	-										
	File System configuration				2 - 11 - 1	le System disks			2 - 11		
	Filesystem Settin	nas	Mount & Perfor	mance Setti	nas			0.00			
	i iicoyoteini oettii	.90	riount di l'unon	nunco occu	iigo	DISKS IN T	ne Filesyst	em			
						1					
	Device:	g pfsO	Mountpoint:*	/opfs		Remove					
	RemoteDeviceName:	gpts0	AutoMount:	Yes 💌			444				
	BlockAllocationType:	cluster	DriveLetter:								
	ReadWribe:	true	OtherMountOptions:	none		Select Name		Usage Type			
	Туре:	local	Inode5ize:	512		NSD1	301	dataAndMetad dataOnIv	ata		
	AcIType:	NFSv4	NumInodes:	1,953,027		NSD10		dataAndMetad	iata		
	LockingType:	NFSv4	NumNodes:	30		NSD2		dataOnly			
	StripeMethod:	roundRobin	MinimumFragmentSize:	16384		NSD3		dataOnly			
	Dmap(Enabled:	false	Blocksize:	262144		NSD8		dataAndMetad	lata		
	LargeLunSupport:	true	IndirectBlocksize:	16384		Total: 7					
	SuppressAtime:	1									
	ExactMtime:	~				Add disk to file	esystem				
	Version:	9.03	Mount Informat	ion							
	Replication & Qu	iota Settings				lesystem Usage			7 - 0		
	2		MountedOn								
https://localbost/ibm/console/pavigation.do?pac	alD=com ibm fscc anfsa	ii nav-filesystem	SmoduleRef=com ib	m fscc anfe	localbost	a 🖉 2 0/34) oobmor@c			
https:///ocalifostibiti/console/flavigatioff.do/pag	gene-contributiti sec.gprsgc	all a venico y sterri	amodulerter=com.ib	macc.gpis	locariost	🔤 👿 2.9435		o venimes@g	oogiernali.com		



Usage/Quota Management

🥹 Integrated Solutions Console - Mozilla Firefox 📃 🗆 🗙													
<u>F</u> ile <u>E</u> dit <u>V</u> iew Hi <u>s</u> tory <u>B</u> ookmarks <u>T</u> ools	<u>H</u> elp									30	oehmes@g	jooglema	l.com - 🔿
💠 - 💽 💿 🏠 🗋 https://local	lhost/ibm	n/console/login	.do?a	ction=secu	ire				<u> </u>	• 🕨 G•	Google		٩
Integrated Solutions Console welcome root										Help Logout			IBM.
Welcome													
Console Settings	GPPS Qu	uotas											2 - 🗆
Scale-Out File Services (SOFS)													
Customer Information													
Node Availability	0.0	ctive Cluster	apfe 1	fore mail	nz de ihm i	Activ	e filosy	stamia	ofc0				
System Utilization	A	Active_cluster:gptst.tscc.mdllZ.de.jbm.com											
GPFS Management													
Cluster Management	C	2 🖸 🖬 🐺	*0										
Cluster Configuration		-		-									
File Systems				Filter by:	Name 🗾	Filter Term:		1.00	Go				
Disks & NSDs GPFS Nodes	Sele	ict Name≎	Device	Soft limit (kl)Hard limit (kb)≎	Grace time 🛇	Used (kb) ≎	Usage %	Soft limit (inodes)	Hard limit (inodes)≎	Grace time 🗘	Used (inodes)	Usage % ≎
Policy-based data management] 🚨 aa	g pfs0	61.44 MB	71.68 MB	N/A	20.22 MB	28 %	1800	2000	N/A	2	0 %
Create new GPES Cluster	0.0	a											
Self-Service Space Creation		gu000000	g pfs0	30.72 MB	40.96 MB	N/A	20.49 MB	50 %	N/A	N/A	N/A	21	N/A
] 🚨 gu000007	g pfs0	7.17 MB	7.17 MB	66 days, 0:32:57 ago	7.18 MB	100 %	N/A	N/A	N/A	8	N/A
] 🚨 gu007393	g pfs0	15.36 MB	20.48 MB	30 days, 0:36:46 ago	18.44 MB	90 %	N/A	N/A	N/A	19	N/A
] 🚨 gu019980	g pfs0	N/A	N/A	N/A	6.15 MB	N/A	7	8	N/A	7	66 %
] <u>3</u> gu019981	g pfs0	N/A	N/A	N/A	8.20 MB	N/A	7	8	60 days, 1:58:39 ago	9	113 %
] 🚨 gu019996	g pfs0	17.15 MB	17.15 MB	79 days, 20:32:30 ago	17.42 MB	102 %	N/A	N/A	N/A	18	N/A
] 🚨 gu019997	g pfs0	18.43 MB	20.48 MB	N/A	15.11 MB	74 %	N/A	N/A	N/A	12	N/A
] 🤱 gu019996	g pfs0	18.43 MB	20.48 MB	N/A	10.25 MB	50 %	N/A	N/A	N/A	11	N/A
] 🚨 gu019999	g pfs0	18.43 MB	20.48 MB	N/A	5.13 MB	25 %	N/A	N/A	N/A	6	N/A
] 💰 quotatst	g pfs0	28.67 MB	30.72 MB	58 days, 22:41:45 ago	30.49 MB	99 %	N/A	N/A	N/A	29	N/A
] 🤷 aa	g pfs4	40.96 MB	51.20 MB	N/A	N/A	N/A	800	1000	N/A	N/A	N/A
	Tob	tal: 12	Pag	ge 1 of 1						[Fri Jun 22	05:01:12 EDT 20	o7 €0
	Exp	port to CSV										No	tifications
https://localhost/ibm/console/navigation.do?pa	geID=col	m.ibm.fscc.gr	fsgui.	nav-quota	&moduleRet	f=com.ibm.fsc	c.gpfsqu	i I	ocalhost 🚗 🛙	2.303s	🔒 M 0 oeh	mes@good	lemail.com



System Utilization Reports







1	NAS today
2	The Scale-out File Services (SoFS) approach
3	General Parallel File System (GPFS)
4	Inside SoFS
5	Conclusion





SoFS vs Classic NAS

- SoFS combines a Clustered Filesystem with a Global Namespace. It's deployed centrally, managed centrally, backed up centrally and grown centrally.
- SoFS scales horizontally not vertically. When you need more capacity, you just add more disks. When you need more performance you add nodes and/or disks. There are almost no architectural limits for future growth.
- NAS boxes growth vertically. You can add controllers and storage up to the maximum configuration of the box. Beyond that you have to buy a new one.
- SoFS provides integrated Information Lifecycle Management
 - Different tiers of storage, e.g. FC, SATA and tape
 - Policy driven placement and migration of files over there entire Lifetime
- SoFS offers synchronous Block Level replication for data and Metadata
- SoFS offers multi-site configurations via the cross-cluster mount functionality.
- SoFS allows to build Highly Available Systems including Disaster resistant Systems across multiple sites



Limits

	SoFS	IBM System Storage N	Windows 2k3 Server
Number of nodes per cluster	13 (limit will be lifted with upcoming releases) Multiple clusters can export the same file system	2	n/a
Number of CPUs per cluster	52 (limit will be lifted with upcoming releases)	16	n/a
Max. capacity	33554432 Yobibytes (2 ¹⁰⁵ Bytes)	504 Terabyte (~2 ⁴⁹ Bytes)	n/a
Max. size of single file system	524288 Yobibytes (2 ⁹⁹ Bytes)	16 Tebibyte (2 ⁴⁴ Bytes)	256 Tebibyte (2 ⁴⁸ Bytes)
Max. number of file systems	256	200	n/a
Max. size of single file	16 Exibytes (2 ⁶⁴ Bytes)	16 Tebibyte (2 ⁴⁴ Bytes)	16 Tebibyte (2 ⁴⁴ Bytes)
Max. number of files per file system	2 billion (2 ³¹)	??	4 billion (2 ³²)
Max. number of snapshots per fs	31	256	n/a

IBM Systems & Technology Group



What's next?





Realtime Performance Analysis

- IO response time seen at NSD client (green curve) and NSD server (blue curve)
- Top chart: Maximum IO response time (noisy)
- Bottom chart: Average IO response time (better problem indicator)
- Region 1: Disk controller bottleneck (slowness seen both at NSD client and server)
- Region 2: Network/node bottleneck (slowness seen only at NSD client)



Maximum IO Response Time



Average IO Response Time



SOFS Online Transparent Data Migration





SOFS Global Namespace consolidation





SoFS with asynchronous replication across sites











SoFS Roadmap





Upcoming features 2008

January

- HTTP(S) access to file system with full authentication and ACL enforcement
- GUI enhancements regarding ease of operation. E.g. consolidated logs, monitoring panels, automated alerts, …
- Integration of TSM 5.5 release and Director 5.20.2 release
- Support of DCS9550 storage system
- March
 - Support of IBM x3650 rack mounted SoFS nodes with SAS attached DS3200 storage systems.
 - Support of common LDAP user directories for authentication (in addition to MS Active Directory).



Upcoming features 2008

May

- Lift cluster size limit to 25 nodes
- Support for single external IP address of cluster using internal load-balancing instead of DNS round-robin
- Asynchronous replication of file system
- Customer managed operation
- 10 Gb Ethernet support
- SAN Volume Controller support
- November
 - Lift cluster size limit to 40 nodes
 - IPv6 support
 - InfiniBand support for intra-cluster and client-cluster communication
 - GUI enhancements for monitoring
 - Simplification of installation & upgrade process



SOFS Architecture





Special Notices

Copyright IBM Corporation, 2006

- This presentation was produced in the United States. IBM may not offer the products, programs, services or features discussed herein in other countries, and the information may be subject to change without notice. Consult your local IBM business contact for information on the products, programs, services, and features available in your area. Any reference to an IBM product, program, service or feature is not intended to state or imply that only IBM's product, program, service or feature may be used. Any functionally equivalent product, program, service or feature that does not infringe on any of IBM's intellectual property rights may be used instead of the IBM product, program, service or feature.
- The [e(logo)server] brand consists of the established IBM e-business logo followed by the descriptive term "server".
- Information in this presentation concerning non-IBM products was obtained from the suppliers of these products, published announcement material or other publicly available sources. Sources for non-IBM list prices and performance numbers are taken from publicly available information including D.H. Brown, vendor announcements, vendor WWW Home Pages, SPEC Home Page, GPC (Graphics Processing Council) Home Page and TPC (Transaction Processing Performance Council) Home Page. IBM has not tested these products and cannot confirm the accuracy of performance, compatibility or any other claims related to non-IBM products. Questions on the capabilities of non-IBM products should be addressed to the suppliers of those products.
- IBM may have patents or pending patent applications covering subject matter in this presentation. The furnishing of this
 presentation does not give you any license to these patents. Send license inquires, in writing, to IBM Director of Licensing, IBM
 Corporation, New Castle Drive, Armonk, NY 10504-1785 USA.
- All statements regarding IBM's future direction and intent are subject to change or withdrawal without notice, and represent goals and objectives only. Contact your local IBM office or IBM authorized reseller for the full text of a specific Statement of General Direction.
- The information contained in this presentation has not been submitted to any formal IBM test and is distributed "AS IS". While each item may have been reviewed by IBM for accuracy in a specific situation, there is no guarantee that the same or similar results will be obtained elsewhere. The use of this information or the implementation of any techniques described herein is a customer responsibility and depends on the customer's ability to evaluate and integrate them into the customer's operational environment. Customers attempting to adapt these techniques to their own environments do so at their own risk.
- IBM is not responsible for printing errors in this presentation that result in pricing or information inaccuracies.
- The information contained in this presentation represents the current views of IBM on the issues discussed as of the date of
 publication. IBM cannot guarantee the accuracy of any information presented after the date of publication.
- All prices shown are IBM's suggested list prices; dealer prices may vary.
- IBM products are manufactured from new parts, or new and serviceable used parts. Regardless, our warranty terms apply.
- Information about non-IBM products was obtained from suppliers of those products. IBM makes no representations or warranties regarding these products. Non-IBM products are offered and warranted by third-parties, not IBM.

Special Notices (contd.)

- Information provided in this presentation and information contained on IBM's past and present Year 2000 Internet Web site pages regarding
 products and services offered by IBM and its subsidiaries are "Year 2000 Readiness Disclosures" under the Year 2000 Information and
 Readiness Disclosure Act of 1998, a U.S statute enacted on October 19, 1998. IBM's Year 2000 Internet Web site pages have been and will
 continue to be our primary mechanism for communicating year 2000 information. Please see the "legal" icon on IBM's Year 2000 Web site
 (www.ibm.com/year2000) for further information regarding this statute and its applicability to IBM.
- Any performance data contained in this presentation was determined in a controlled environment. Therefore, the results obtained in other operating environments may vary significantly. Some measurements quoted in this presentation may have been made on development-level systems. There is no guarantee these measurements will be the same on generally-available systems. Some measurements quoted in this presentation may have been estimated through extrapolation. Actual results may vary. Users of this presentation should verify the applicable data for their specific environment.
- The following terms are registered trademarks of International Business Machines Corporation in the United States and/or other countries: AIX, AIXwindows, AS/400, C Set++, CICS, CICS/6000, DataHub, DataJoiner, DB2, DEEP BLUE, DYNIX, DYNIX/ptx, e(logo), ESCON, IBM, IBM(logo), Information Warehouse Intellistation, IQ-Link, LANStreamer, LoadLeveler, Magstar, MediaStreamer, Micro Channel, MQSeries, Net.Data, Netfinity, NUMA-Q, OS/2, OS/390, OS/400, Parallel Sysplex, PartnerLink, PartnerWorld, POWERparallel, PowerPC, PowerPC(logo), ptx/ADMIN, RISC System/6000, RS/6000, S/390, Scalable POWERparallel Systems, SecureWay, Sequent, SP2, System/390, The Engines of e-business, ThinkPad, Tivoli(logo), TURBOWAYS, VisualAge, WebSphere. The following terms are trademarks of International Business Machines Corporation in the United States and/or other countries: AIX/L, AIX/L(logo), AIX, SL, AIX PVMe, Application Region Manager, AS/400e, Blue Gene, Chipkill, ClusterProven, DB2 OLAP Server, DB2 Universal Database, e(logo)business, GigaProcessor, HACMP/6000, Intelligent Miner, iSeries, Network Station, NUMACenter, PowerPC Architecture, PowerPC 604, POWER2 Architecture, pSeries, Sequent (logo), SequentLINK, Service Director, Shark, SmoothStart, SP, Tivoli Enterprise, TME 10, Videocharger, Visualization Data Explorer, xSeries, zSeries. A full list of U.S. trademarks owned by IBM may be found at http://jlswww.nas.ibm.com/wpts/trademarks/trademar.htm.
- Lotus and Lotus Notes are registered trademarks and Domino and Notes are trademarks of Lotus Development Corporation in the United States and/or other countries.
- NetView, Tivoli and TME are registered trademarks and TME Enterprise is a trademark of Tivoli Systems, Inc. in the United States and/or other countries.
- Microsoft, Windows, Windows NT and the Windows logo are registered trademarks of Microsoft Corporation in the United States and/or other countries.
- UNIX is a registered trademark in the United States and other countries licensed exclusively through The Open Group.
- LINUX is a registered trademark of Linus Torvalds.
- Intel and Pentium are registered trademarks and MMX, Itanium, Pentium II Xeon and Pentium III Xeon are trademarks of Intel Corporation in the United States and/or other countries.
- Java and all Java-based trademarks and logos are trademarks of Sun Microsystems, Inc. in the United States and/or other countries.
- Other company, product and service names may be trademarks or service marks of others.