



---

# **IBM's General Parallel File System**

## **an overview**

**Glen Corneau**  
**IBM Power Systems Advanced Technical Support**  
**[gcorneau@us.ibm.com](mailto:gcorneau@us.ibm.com)**



# Table of Contents

---

- **GPFS Features**
  - **Overview**
  - **Architectures**
  - **Performance**
  - **RAS**
  - **Supported Environments and Limits**
  - **Information Lifecycle Management**
  - **Disaster Recovery**
  - **Requirements**
  - **With Oracle RAC**
- **GPFS Terminology**
- **GPFS Administrative Tasks**
- **GPFS Debugging**
- **Additional Information**



# IBM General Parallel File System (GPFS)

---

IBM General Parallel File System (GPFS) is a scalable high-performance file management infrastructure for AIX®, Linux® and Windows systems.

**A highly available cluster architecture.**

**Concurrent shared disk access to a single global namespace.**

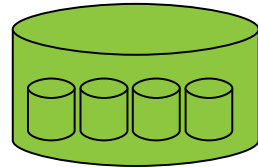
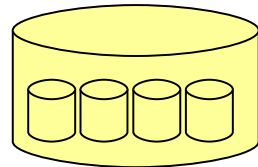
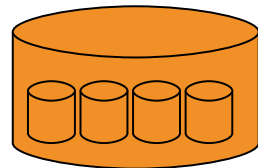
**Capabilities for high performance parallel workloads.**



# File data infrastructure optimization

GPFS enables:

- A single global namespace across platforms.
- High performance common storage.
- Eliminating copies of data.
- Improved storage use.
- Simplified file management.



**Connections**  
SAN  
TCP/IP  
InfiniBand

**Management**  
Centralized  
Monitoring  
Automated File  
Mgmt

**Availability**  
Data Migration  
Replication  
Backup

**Databases**

**File Servers**

**Backup / Archive**

**Application Servers**



## What is GPFS?

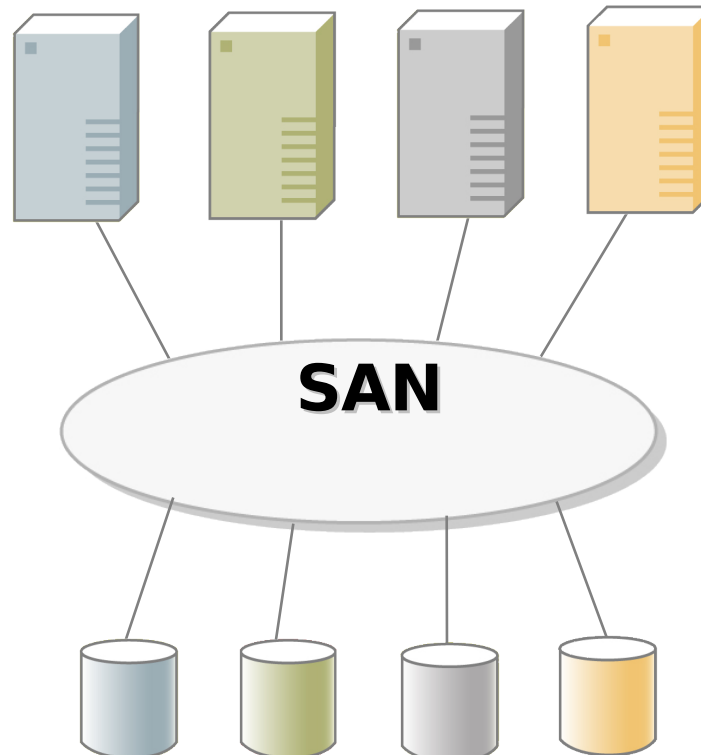
---

- General Parallel File System  
(not Global Positioning File System)
- Mature IBM product generally available since 1997
  - Current version is the 12<sup>th</sup> release (i.e. not a new product)
- ★ GPFS Version 3.3, available in Sep 2009  
[green star denotes new/updated feature]
- Adaptable to many customer environments by supporting a wide range of basic configurations and disk technologies
- Conforms to the POSIX I/O standard, but has extensions
- Provides non-POSIX advanced features (e.g., data-shipping, hints)
- Converting to GPFS does not require application code changes provided the code works in a POSIX compatible environment
- Concurrent file access requires application to be parallel file system aware



## GPFS Architectures - SAN

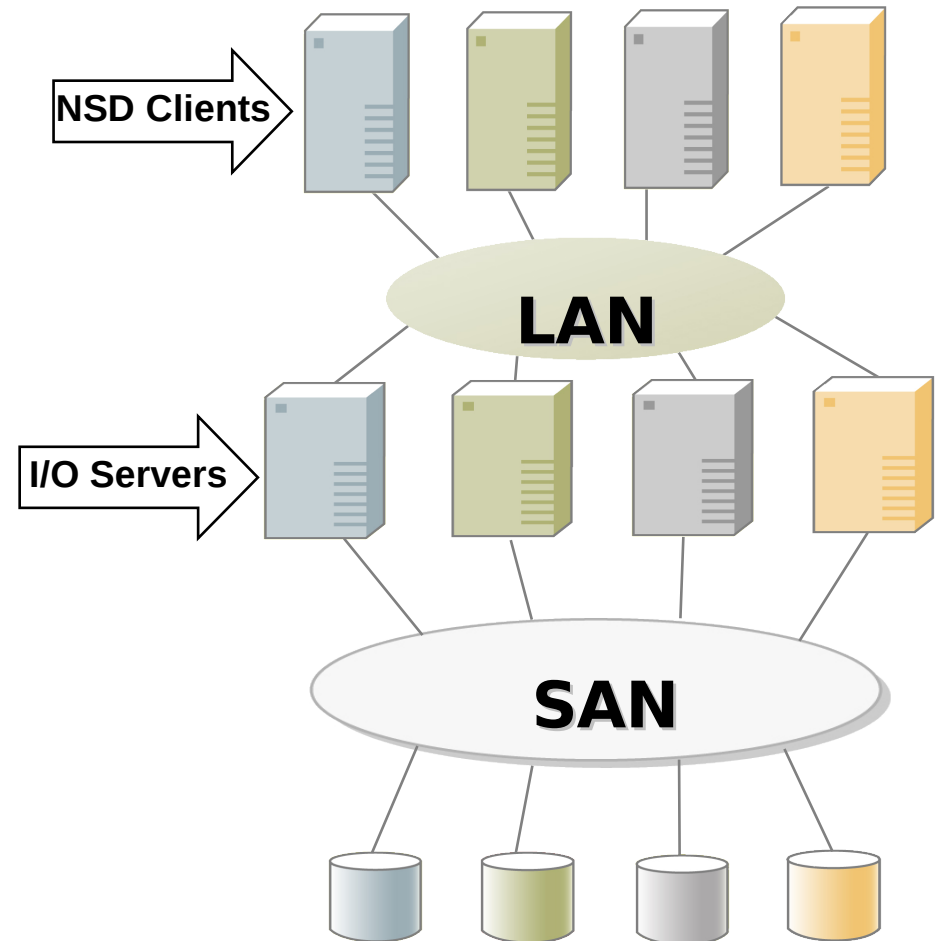
- All nodes are directly attached via the SAN
- Highest performance for all nodes





## Network Block I/O – Using NSDs

- I/O Server is directly attached to the disks
- NSD clients see local block device same as direct attached
- I/O is sent over the LAN: Ethernet, Infiniband, Myrinet
- Mixed Environments: Nodes with highest bandwidth needs direct attached to SAN. Other nodes attached using NSD





## Performance Features

---

- File systems are striped across the disks of which it is comprised
- Large blocks (with support for sub-blocks) support (up to 4MB)
- Byte range locking (rather than file or extent locking) available with the proper application file open flags.
- Access pattern optimizations (e.g., sequential, strided, backward)
- File caching (called pagepool) separate from local OS file cache
  - Default size is 64MB, can be up to 256GB
  - pinned kernel memory
  - vmo / ioo does not directly affect pagepool (not like JFS/JFS2 tunables)
- Multi-threaded 64-bit daemon that supports pre-fetch and write-behind
- Distributed management functions (e.g., metadata, tokens)
- Multi-pathing (i.e., multiple, independent paths to the same file data from anywhere in the cluster)





## RAS Features

---

- If a node providing GPFS management functions fails, an alternative node automatically assumes responsibility preventing loss of file system access.
- Disks may be served to the cluster via a primary and secondary server, or be seen by a group (or all) of nodes. Therefore losing a server node or an adapter does not result in loss of access to data.
  - Up to 8 NSD servers for each NSD
- Data and metadata may also be separately replicated (similar to mirroring) to remove a storage subsystem as a SPOF.
- GPFS utilizing storage subsystems with built-in RAID protects against loss of data and loss of data access.
- Online and dynamic system management
  - Add/remove disk, restriping, replication, filesystem modification
- In file system snapshots (limited by amount of change + free space)
  - Up to 256 in GPFS V3.3



## RAS Features

---

- ★ Restripping of a file system is parallelized for improved performance (v3.3)
- Rolling upgrades supported in Version 3
- SCSI-3 Persistent Reserve can be used for fast failover of 2-node clusters
- SNMP monitoring with GPFS for Linux
- Clustered NFS solution with GPFS for Linux
- GPFS daemon will attempt to re-connect when socket connections are broken
- GPFS can be used in Disaster Recovery configurations
- NSD access can fail over to network from SAN and back, if so desired
- ★ Parallel backup integration with TSM now uses policy engine and supports incremental forever



# Supported Environments

---

<b>GPFS Version:</b>	<b>3.2</b>	<b>3.2.1</b>	<b>3.3</b>
<b>AIX 5.2</b>	Y	Y	n
<b>AIX 5.3</b>	Y	Y	Y
<b>AIX 6</b>	Y	Y	Y
<b>Linux on POWER</b>	RHEL4, 5	RHEL4, 5	RHEL4, 5
	SLES9, 10	SLES9, 10,11	SLES9, 10, 11
<b>Multiplatform Linux</b>	RHEL4, 5	RHEL4, 5	RHEL4, 5
	SLES9, 10	SLES9, 10,11	SLES9, 10, 11
<b>Windows Server 2003 R2 64-bit</b>	n	V3.2.1.5+	n
<b>Windows Server 2008 SP2 64-bit</b>	n	n	Y

- Check the FAQ for details on kernel levels, patches, etc



## Supported Environments

---

- Hardware Environments:
  - POWER-based systems: AIX & Linux
  - AMD-based IBM systems: Linux
  - Intel-based System x systems: Linux or Windows (AMD x64, EMT64T)
- The GPFS cluster interconnect must be invariant
  - i.e. cannot be managed by PowerHA, etc
  - For availability, look to EtherChannel or 802.3ad Aggregates!
- For storage supported and/or tested, look at the GPFS FAQ.
- GPFS supports the use of PowerVM features (i.e. Virtual Ethernet, Virtual SCSI, Shared Ethernet Adapter).
- GPFS has not been officially tested with Live Partition Mobility or NPIV



## Supported Environments

---

- GPFS for Windows limitations (i.e. “no”):
  - Storage devices directly attached to Windows systems (i.e. network-based NSD access only)
  - Remote mounting of file systems across GPFS clusters
  - "manager" and "quorum" node roles
  - File systems that are DMAPI-enabled
  - File systems created with prior GPFS releases (3.2.1.5+)
  - GPFS Application Programming Interfaces
  - Only 64 Windows nodes in the cluster only (>32 requires GPFS Development's assessment)
- Role Based Access Control (RBAC) is not supported by GPFS and is disabled by default.
- Workload Partitions (WPARs) or storage protection keys are not exploited by GPFS.



## Supported Interconnects

---

- Linux, AIX and/or Windows (mixed cluster):
  - Ethernet: 100Mb, 1Gb and 10Gb
- Linux only cluster
  - Myrinet (IP only)
  - Infiniband: IP or VERBS RDMA (only Multiplatform, not Power)
- AIX only cluster
  - Myrinet (IP only)
  - eServer HPS (homogenous AIX clusters)
  - Infiniband (IP only)
  - Virtual Ethernet/Shared Ethernet Adapter
- GPFS supports both inter-cluster and intra-cluster network definitions
  - i.e. local cluster can use local, high-speed interconnect



## GPFS Limits

---

- Very Large file system design:  $2^{99}$  bytes or around 633,825 Yottabytes ( $2^{80}$ ) !! (that's a lotta bytes, ha!)
  - Tested file system limit: 4 Petabytes (4000 TB)
- Maximum number of files in filesystem (FS size dependent):
  - 2,147,483,648
- Supported node scaling support:
  - AIX: 1530 nodes (>128 nodes requires IBM review)
  - Multiplatform for Linux: 3794 nodes (>512 nodes requires IBM review)
  - Windows: 64 nodes (no homogenous Windows clusters)
  - AIX + Multiplatform for Linux: 3906 (3794 Linux, 112 AIX nodes)
- Largest disk size:
  - 64-bit AIX/Linux: limited by the device driver & OS (>2TB)
  - 32-bit AIX: 1TB
  - 32-bit Linux: 2TB.
- Maximum number of file systems is 256



## Other GPFS Features

---

- Journaling (logging) File System - logs information about operations performed on the file system metadata as atomic transactions that can be replayed
- Data Management API (DMAPI) - Industry-standard interface allows third-party applications (e.g. TSM) to implement hierarchical storage management
- NFSv4 ACL support
- External storage pools allow automated migration of files to/from HSM systems (i.e. tape).
- Backward compatibility for multi-cluster environments
- ★ GPFS can be configured with a smaller number of “trusted” nodes that have remote shell capabilities to all other nodes (admin nodes)
  - Remote shell can be prompted or non-prompted
  - Can utilize ssh-agent in this configuration
- ★ GPFS allows defined user scripts to be executed after specific GPFS events





# Licensing

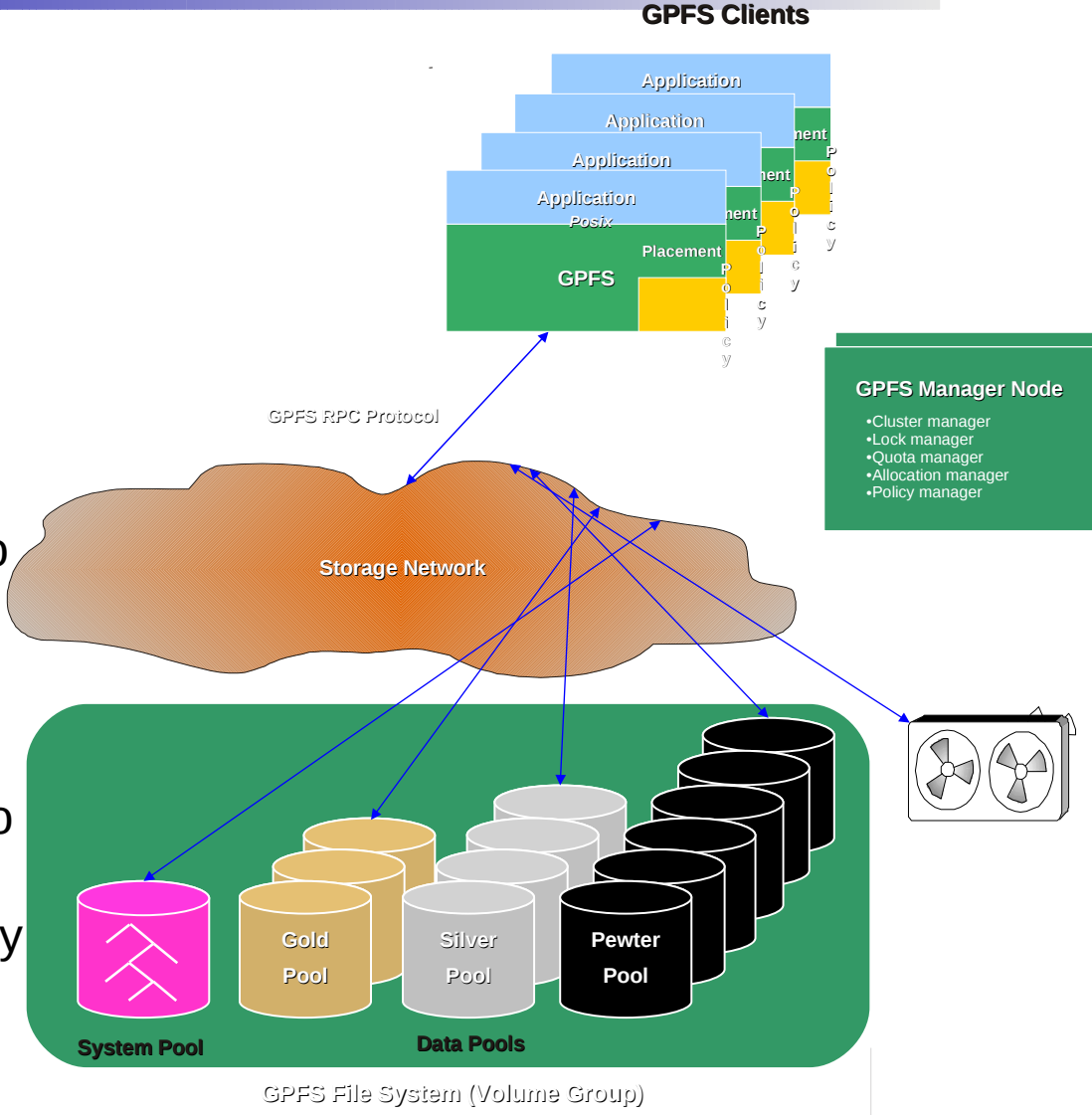
---

- GPFS Server license:
  - Mount file systems from NSD servers or local block devices
  - Perform GPFS management: quorum node, file system manager, cluster configuration manager, NSD server
  - Share data through any application such as NFS, CIFS, FTP or HTTP
- GPFS client license
  - Mount file systems from NSD servers or local block devices
  - Exchanging data between nodes that locally mount the same GPFS file system
- Per PVU on System x, Intel or AMD. Ordered via Passport Advantage
- Per core on Power Systems. Ordered via eConfig (AAS).
- See GPFS FAQ for examples



# ILM features in GPFS

- GPFS Version 3 adds support for Information Lifecycle Management (ILM) abstractions: storage pools, filesets, policies
  - Storage pool – group of LUNs
  - Fileset: named subtree of a file system
  - Policy – rules for placing files into storage pools
- Examples of policy rules
  - Place new files on fast, reliable storage, move files as they age to slower storage, then to tape
  - Place media files on video-friendly storage (fast, smooth), other files on cheaper storage
  - Place related files together, e.g. for failure containment





# Disaster Recovery

## Active/Active Dispersed Cluster

- Figure at left is one geographically dispersed cluster
  - All nodes in either site have SAN/NSD access to the disk
  - Site A storage is duplicated in site B with GPFS replication
- Simple recovery actions in case of site failure (more involved if you lose tiebreaker site as well)

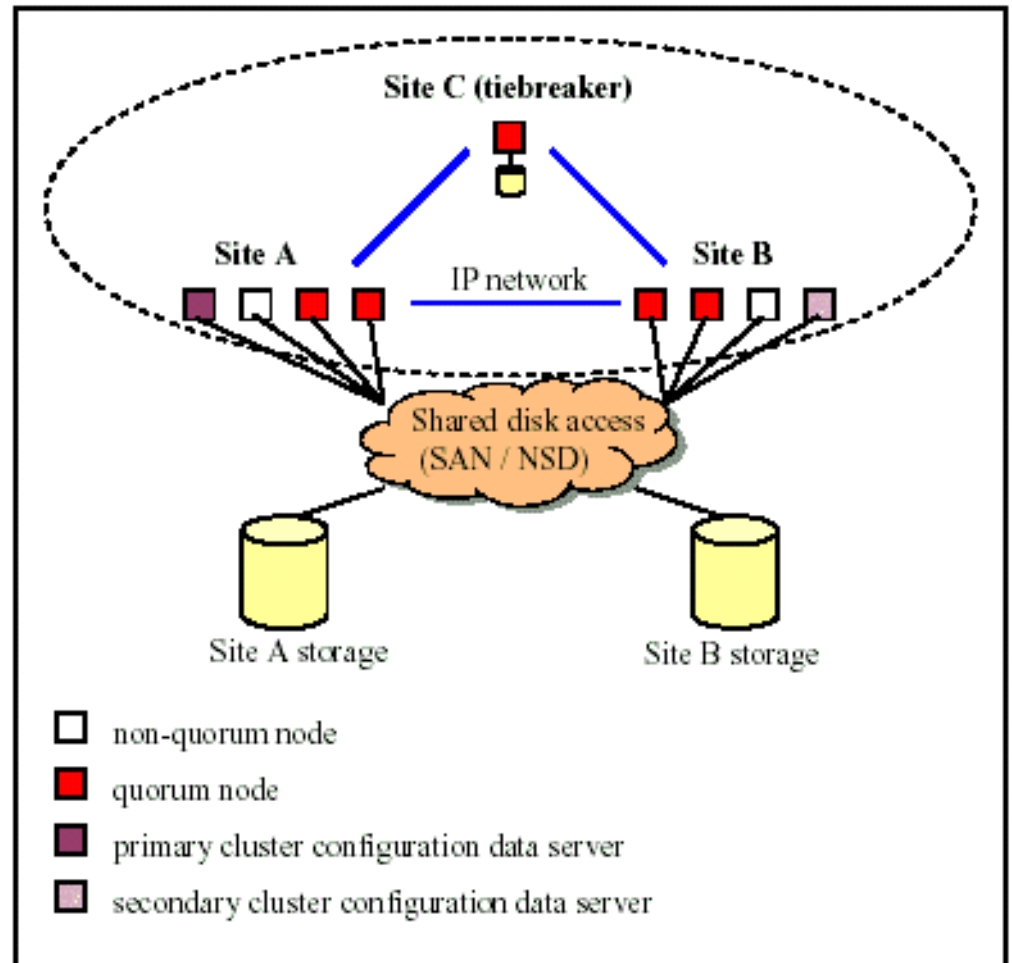


Figure 1. High-level organization of a replicated geographically dispersed GPFS cluster

- **Performance implication:** GPFS has no knowledge of a replica's physical locality. There is no way to specify disk access priority (i.e. Local storage first)



# Disaster Recovery

## Active/Passive with Asynchronous Storage Replication

- Uses “mmfsctl syncFSconfig” to set up the filesystem definitions at recovery site
- Storage subsystem replication keeps the LUNs in sync
- Failover requires production site to be down (i.e. must shutdown GPFS daemons if not a total site failure)
- More involved configuration and failover

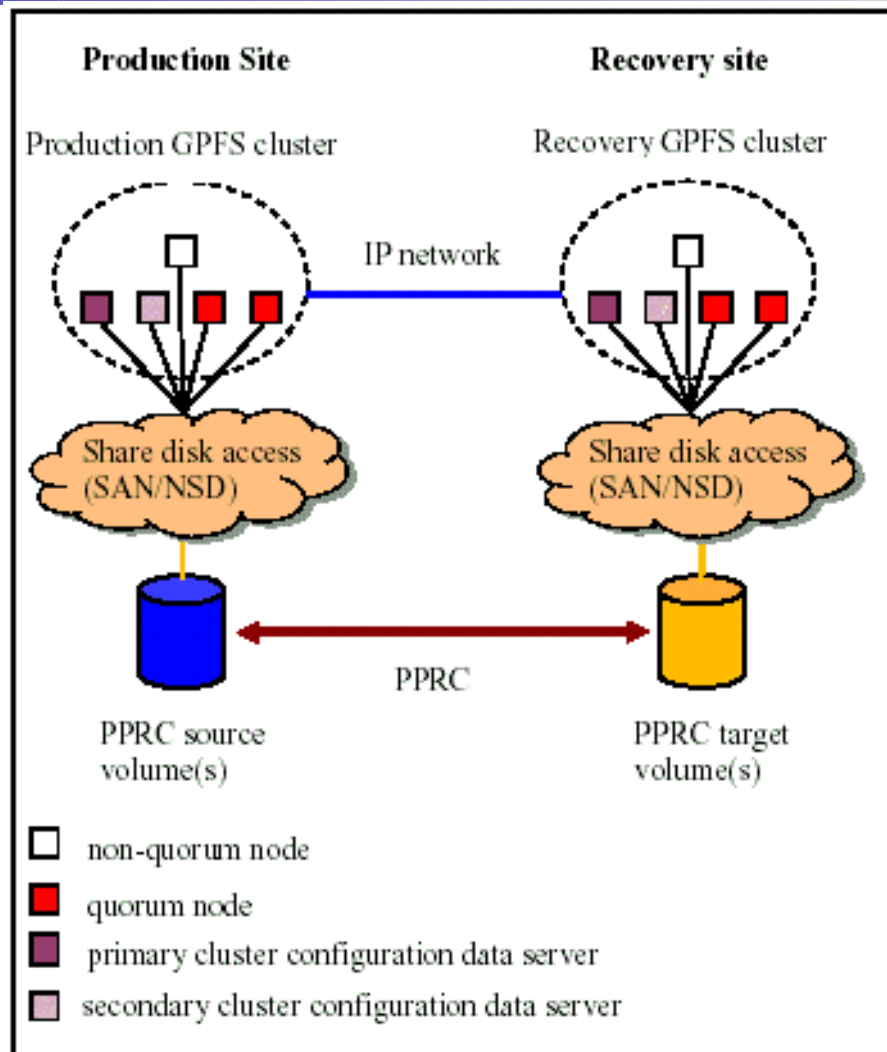


Figure 2. An active/passive disaster-tolerant GPFS environment using PPRC mirroring



# Disaster Recovery

## Active/Active with Asynchronous Storage Replication

---

- Same node layout as active/active using GPFS replication:
  - Nodes in three sites, primary A, secondary B and tiebreaker C.
- Same disk layout as Active/Passive with PPRC
  - DiskA is SAN-attached and accessible from sites A & B
  - DiskB is SAN-attached and accessible from site B only
  - PPRC relationship from diskA to diskB
  - Consistency groups should be defined over all logical subsystems at the primary site.
- Failover involves disengaging the diskA access from siteB (via SAN or nsddevices user exit script)



# Disaster Recovery

## Online Backup with FlashCopy

---

- Can be used to make off-site point-in-time copy of the LUNs that comprise a GPFS filesystem
- Not to be confused with GPFS snapshots which make a in-filesystem point-in-time copy for backup purposes.
  - i.e. the GPFS mmbackup command that works with TSM utilizes this under the covers.
- Requires temporary suspension of primary GPFS volumes when initiating FlashCopy commands (flushes all buffers/cache to disk for a consistent file system image)
  - “mmfsctl <device> suspend | resume”
- Can be used both for availability (DR backup) or other purposes (spin off copies for slow “to tape” backups, additional data analysis, etc).
- Can have a pseudo-Active/Active configuration with second site live at the same time as primary site



## GPFS with Oracle RAC

---

- GPFS Versions 3.1 and 3.2 are certified on AIX 5.3 and AIX 6.1 with
  - Oracle 10 (version 10.2.0.3 or later)
  - Oracle 11.1
- Obtain a copy of Oracle Metalink Article 302806.1, entitled “IBM General Parallel File System (GPFS) and Oracle RAC on AIX 5L and IBM eServer pSeries” and Article 282036.1
- Oracle RAC detects the usage of GPFS for its database files and will open them in Direct I/O mode. This bypasses GPFS's pagepool for DB files, but it is still used for other files.
- HACMP is not required for Oracle RAC 10g (or later) implementations.



# GPFS with Oracle RAC

## Tuning Recommendations

---

- When running Oracle RAC 10g, it is suggested you increase the value for `OPROCD_DEFAULT_MARGIN` to at least 500 to avoid possible random reboots of nodes.
- Read the section “GPFS use with Oracle” in the GPFS Planning and Installation Guide for details on threads and AIO.
- Suggested that Voting and OCR not be in GPFS file systems, but rather in shared raw devices (hdisk)





# GPFS with Oracle RAC

## Tuning Recommendations

---

- For file systems holding large Oracle databases, set the GPFS file system block size to a large value:
  - 512 KB is generally suggested.
  - 256 KB is suggested if there is activity other than Oracle using the file system and many small files exist which are not in the database.
  - 1 MB is suggested for file systems 100 TB or larger.
- ★ The large block size makes the allocation of space for the databases manageable and has no affect on performance when Oracle is using the Asynchronous I/O (AIO) and Direct I/O (DIO) features of AIX.



# GPFS Terminology

---

The following pages provide quick definitions for the following GPFS terms and concepts:

- Network Shared Disks
- File Systems
- Failure Groups & Replication
- Node Roles:
  - Cluster Data Server
  - Configuration Manager
  - File System Manager
- Node Quorum



# NSD

## NSD: Network Shared Disk

- This comes from the Linux port of GPFS. An NSD is typically built upon a raw (i.e. not a volume group/logical volume) disk device that is made available to remote clients as part of a GPFS file system across an IP network.
- Can be used over standard interconnects (specialized networks not required) and utilizes TCP/1191 (IANA registered port)
  - Gigabit Ethernet minimum recommended
  - EtherChannel supported
- What about LVs on AIX?
  - Supported if VG/LV managed manually, or if migrated from previous versions of GPFS (i.e. V2.2 or earlier)
  - Also typically used in DR configurations for file system descriptors on 3<sup>rd</sup> site node.



# Disks

- Disks are the smallest units that comprise a GPFS file system

**Disks defined to GPFS cannot belong to more than one file system (not like VGs with separate LVs) !**

- For storage connectivity where you have multiple (typically 2) paths to the disk (i.e. two Fibre Channel adapters) and using a device driver that is multi-path aware you will typically have a virtual device that represents both paths. For example:
  - MPIO on AIX = hdisk
  - SDD (non-MPIO) on AIX = vpath
  - PowerPath by EMC = hdiskpower
  - HDLM by Hitachi = dlmfdrv



# Disks

---

- It's these virtual devices that are used as input to GPFS to define GPFS “disks”.
- GPFS uses the concept of failure groups for availability of the disk data (data and metadata) that comprises file systems
- Disks that have the same “path” (or rely on the same physical device) to the host are in the same failure group. Examples:
  - all disks defined in a single storage subsystem
  - all disks accessed via the same virtual path
- A failure group number can be automatically or manually assigned to a disk



## File Systems

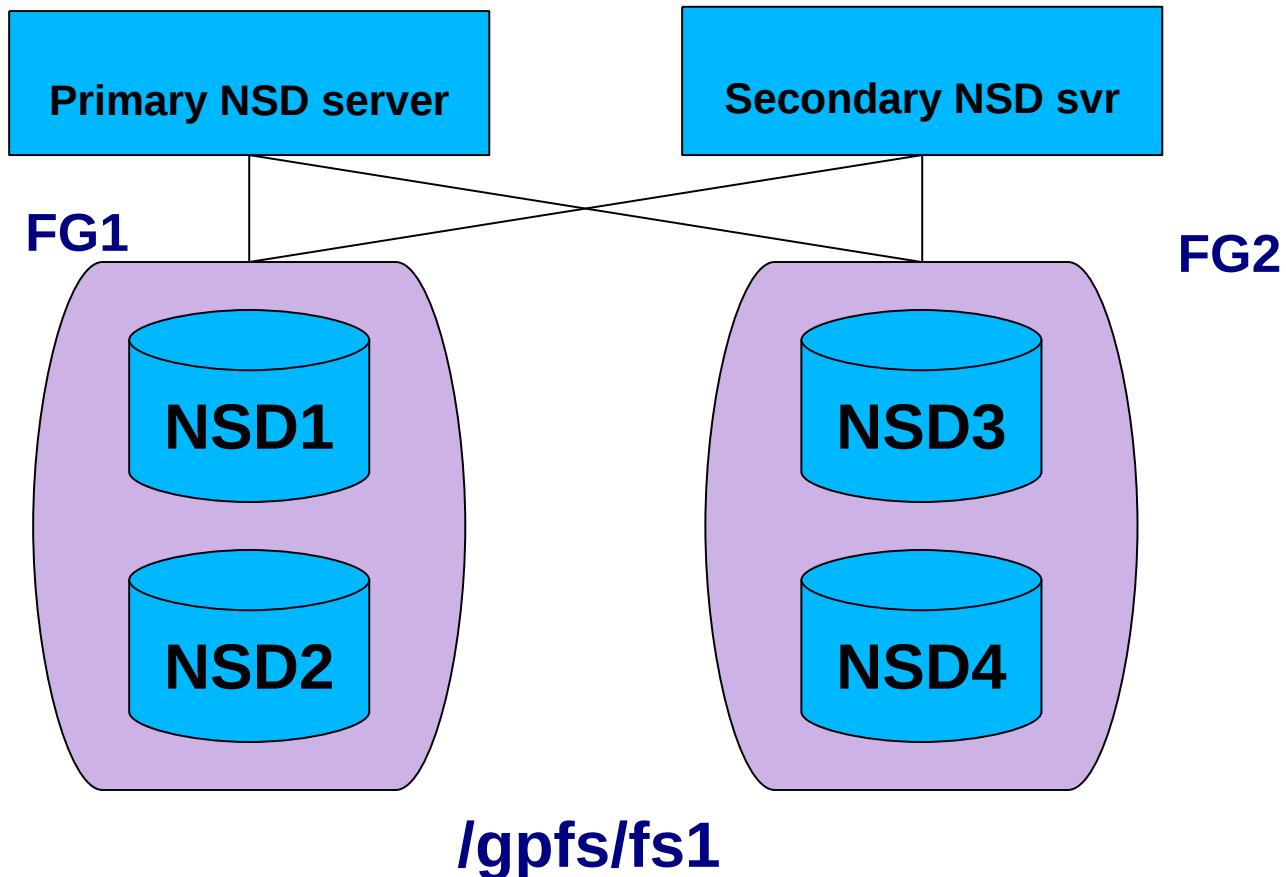
- File system availability is provided (at the FS level) through the use of data and metadata replication.
- Replication is the duplication of either data and/or metadata (typically both) on GPFS disks in different failure groups
- Requires 2x the storage
- Only way to remove a single storage subsystem as Single Point of Failure (SPOF) for a file system.
  - Just using more than one storage subsystem without replication will not provide GPFS file system availability

**GPFS disks cannot be mirrored at the logical volume level!**

- Disk subsystems that use various RAID levels internally are okay!



# Failure Groups & Replication



- With replication enabled, two copies of data and/or metadata are kept, each on NSDs in separate failgroup groups.
- Failure of a disk will cause GPFS to mark the disk down and continue to utilize the other copy exclusively until a repair or admin action is taken



# File Systems

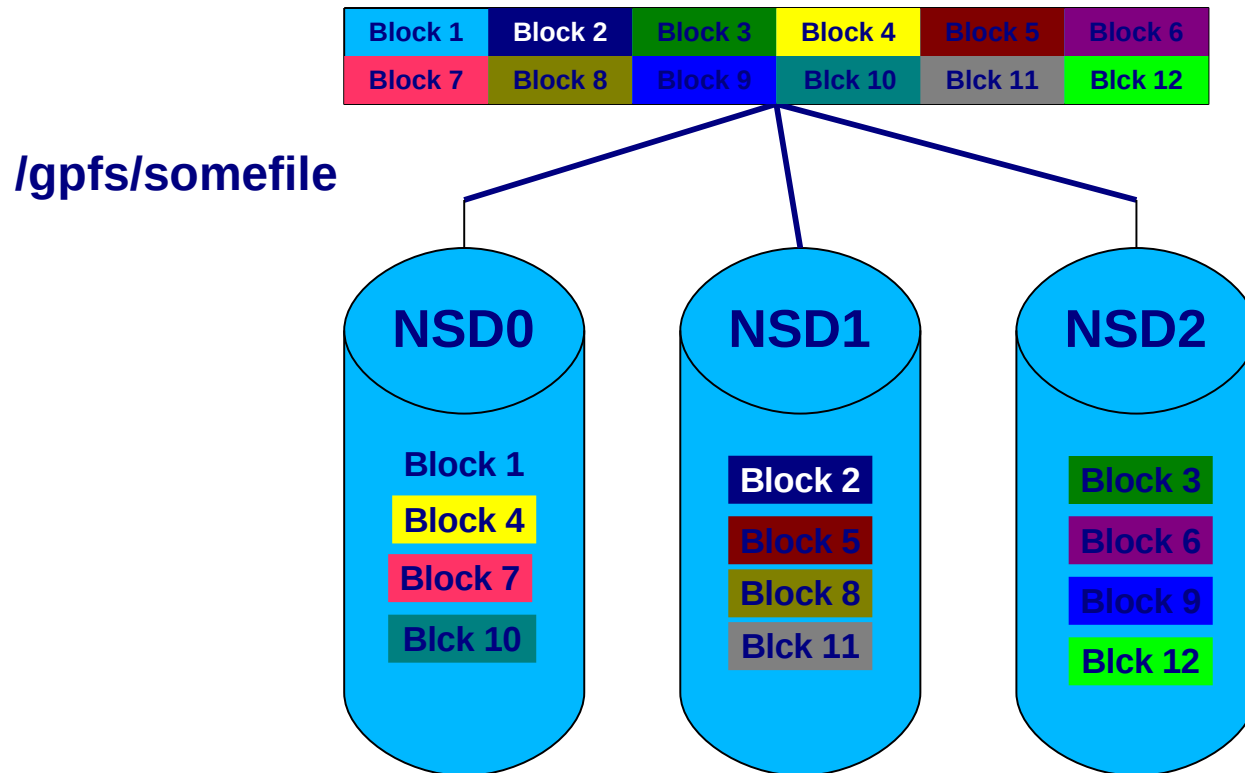
---

- GPFS file systems are created and striped across the GPFS disks that comprise the file system.
  - Filesystems created on one disk are not striped!
  - There is no other way to disable striping other than to use only a single disk!
- File systems are striped across disks using the block size specified during FS creation. Block size cannot be changed after file system creation!
  - 16K, 64K, 128K, 256K, 512K, 1024K (1M), 2M, 4M
- GPFS file systems are synonymous with the device they're created upon.
  - for example: GPFS file system /bigfs is mounted on the device /dev/bigdev
- Some dynamic operations on file systems include:
  - adding disks, deleting disks, restriping, increasing i-nodes





# File System Striping



- GPFS manages the stripe distribution internally.
- Addition or removal of NSDs from an existing file system can optionally restripe the data.
  - Restriping can be parallelized by specifying multiple nodes to participate in the restriping
  - Restriping can be very I/O intensive



## Node Roles: Cluster Data Server

---

- At minimum a primary cluster data server must be defined to act as the primary repository of the GPFS cluster configuration information file (/var/mmfs/etc/mmsdrfs).
  - A secondary GPFS cluster configuration server is highly recommended
- If your primary server fails and you have not designated a secondary server, the GPFS cluster configuration data files are inaccessible and any GPFS administrative commands that need access to the configuration file fail. Similarly, when the GPFS daemon starts up, at least one of the two GPFS cluster configuration servers must be accessible.



## Node Roles: Configuration Manager

---

- The "oldest continuously operating node" is automatically selected as the configuration manager. (Should it fail, another node is automatically selected)
  - Can be chosen via "mmchmgr -c" command
- The Configuration Manager:
  - Selects the file system manager for each file system from available manager nodes
  - Determines whether a cluster quorum exists
- Quorum – two algorithms

Quorum is the minimum number of nodes in a cluster that can be running for the GPFS daemon (i.e., mmfsd) to operate. For most clusters:

- Standard, multi-node:  $\text{quorum} = 1 + \text{sizeof}(\text{quorum nodes})/2$
- Alternative, tie breaker disks: more on this later...



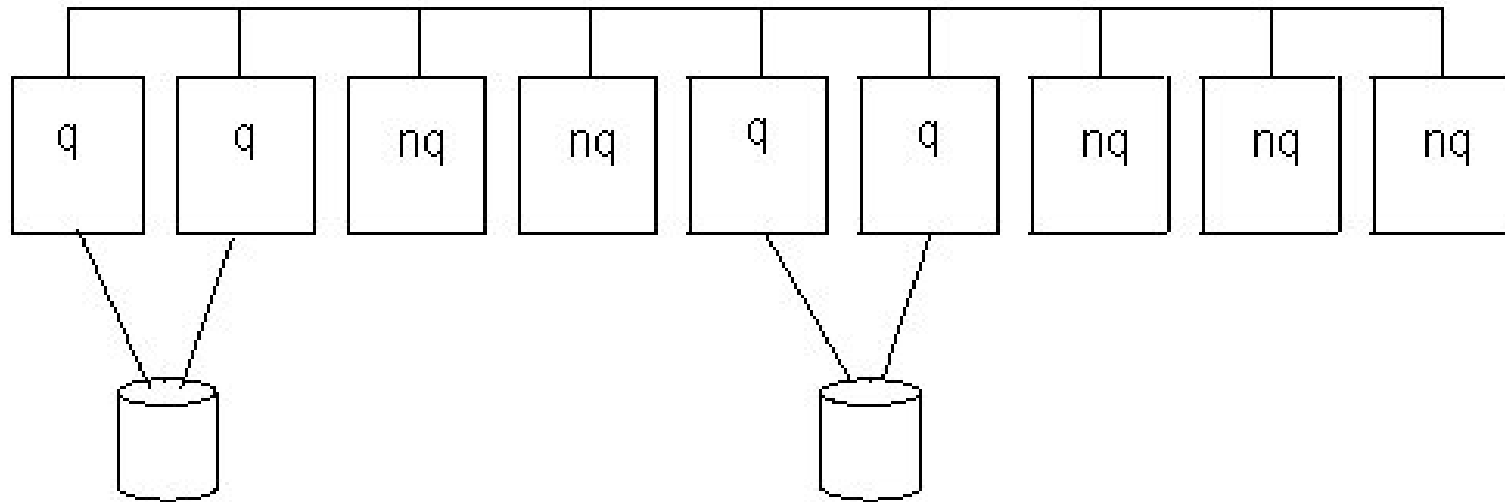
## Node Roles: File System Manager

---

- Each file system is assigned a file system manager.
- The FS manager is responsible for:
  - File system configuration (adding disks, changing disk availability, repairing the file system)
  - Mount and unmount processing is performed on both the file system manager and the node requesting the service.
  - Management of disk space allocation (Controls which regions of disks are allocated to each node, allowing effective parallel allocation of space.)
  - Token management (this can be distributed in Version 3)
  - Quota management
- Failure of the file system manager node will cause another node to automatically be assigned the role



# Node Quorum



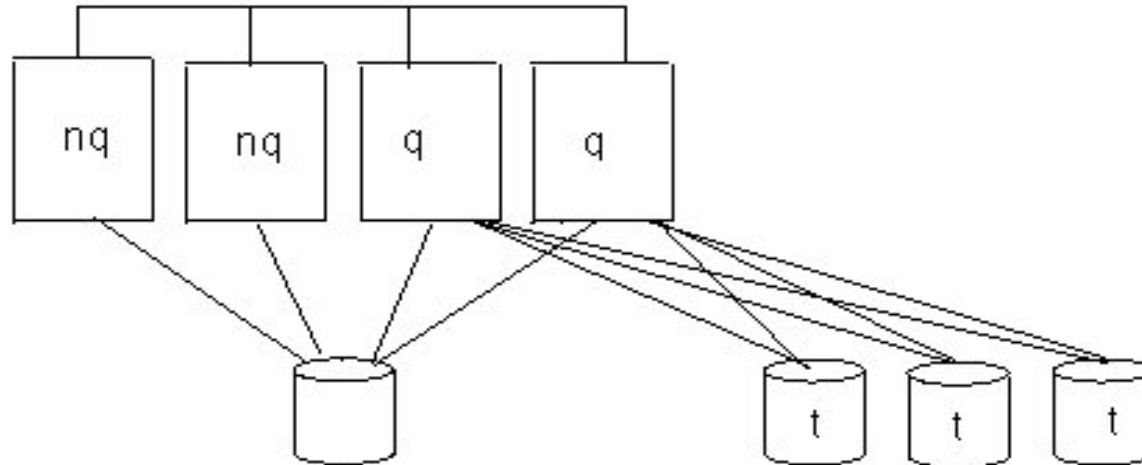
**q - quorum node**

**nq - non quorum node**

- GPFS Node Quorum allows some subset of the total node population to be assigned as explicit quorum nodes.
- Large clusters achieve quorum faster and can be hardened against failure more readily with fewer quorum nodes
- Typically 7 nodes or less... odd numbers are good.



## Node Quorum with Tiebreaker Disks



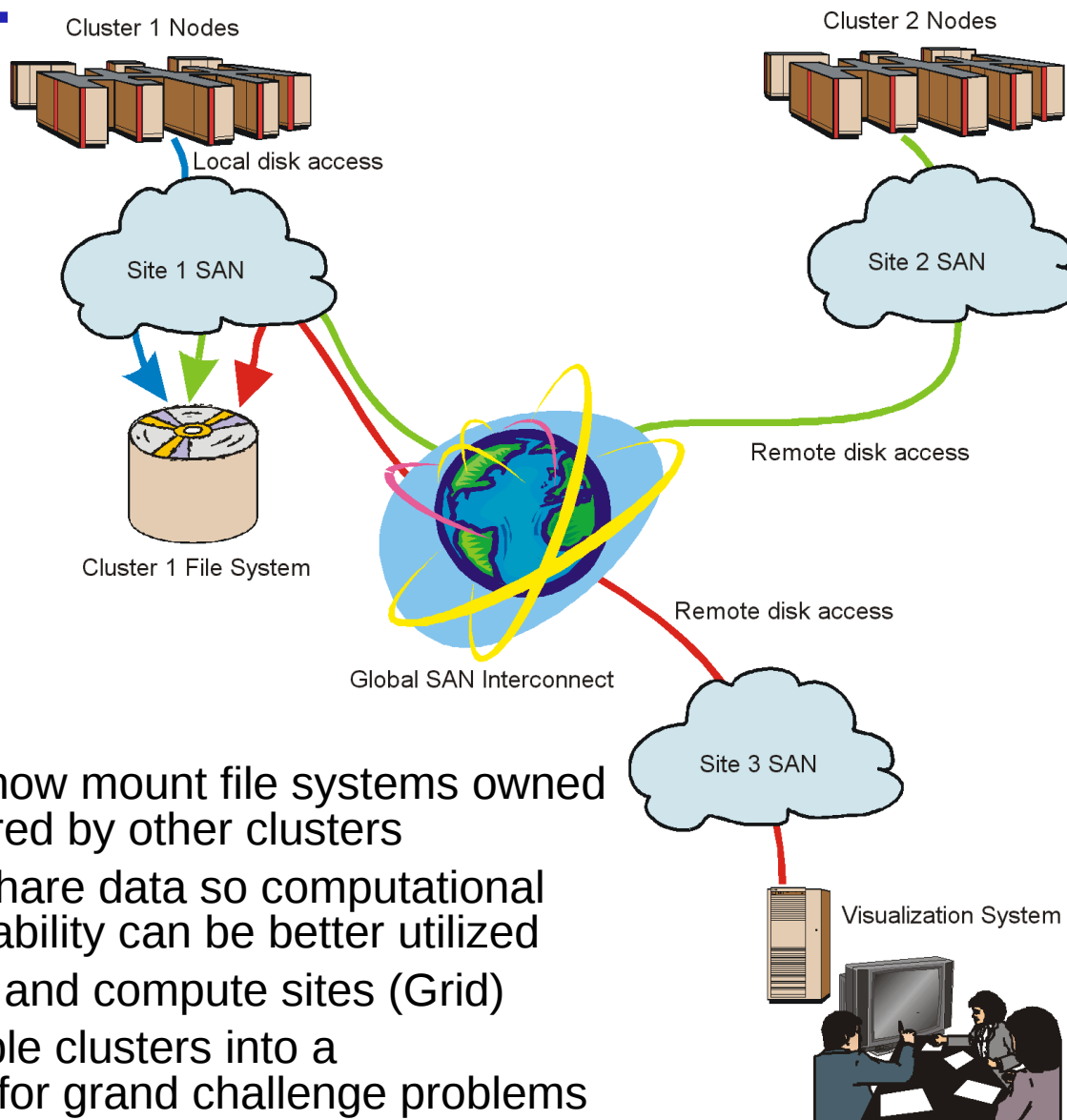
q - quorum node  
t - tiebreaker disk

nq - non quorum node

- No disk hardware SCSI-3 persistent reserve dependency.
- Clusters at Version 3 may contain 8 quorum nodes
  - Clusters may also contain any number of non quorum nodes.
- From one to three disks may be used as tiebreakers, odd is good.



# Cross Cluster Mounts



- Clusters may now mount file systems owned and administered by other clusters
- Clusters can share data so computational resource availability can be better utilized
- Separate data and compute sites (Grid)
- Forming multiple clusters into a “supercluster” for grand challenge problems



# Administrative Tasks

---

- Expanding File Systems
- Shrinking File Systems
- Graphical User Interface
- Common Commands
- Other Tasks
- Debugging





## Expanding File Systems

---

### Command: mmaddisk

- A dynamic operation to increase the size of the file system.
- Requires a modified disk descriptor file as input (i.e. NSDs)
- Can restripe (mmrestripefs) file system to include the new disk if needed (-r flag). This operation can possibly take a long time and be I/O intensive.
  - Can restripe asynchronously (i.e. In the background) with the “-a” flag.
- You can also elect which nodes take part in the restriping. The more nodes that participate, the shorter the time it takes.
- Use the mmlsnsd command to list unused NSDs



## Shrinking File Systems

### Command: `mmdeldisk`

- A dynamic operation to decrease the size of the file system.
- Requires the NSD name(s) and the FS device as input.
- Can restripe (`mmrestripefs`) the file system in addition to just moving data from the deleted disk if needed (`-r` flag). This operation can possibly take a long time and be I/O intensive.
  - Can restripe asynchronously (i.e. In the background) with the “`-a`” flag.
- You can also elect which nodes take part in the restriping. The more nodes that participate, the shorter the time it takes.
- Things to watch:
  - Replication settings (i.e. Don't want to leave a replication state unbalanced)
  - Available space to remove disk!



# Graphical User Interface

- Available with GPFS V3.2.1.1 and later

The screenshot shows the Integrated Solutions Console in Mozilla Firefox. The main content area displays a table with one cluster entry:

Select	Name	Primary server	Secondary server
<input type="checkbox"/>	<a href="#">LAB1.dfw.ibm.com</a>	cler01.dfw.ibm.com	cler02.dfw.ibm.com

Below the table, there is a section for "Cluster settings" with the following details:

Name: LAB1.dfw.ibm.com  
ID: 654005687887710744  
Primary server: cler01.dfw.ibm.com  
Secondary server: cler02.dfw.ibm.com  
TCP port: 1191  
Remote copy path (RCP): /usr/bin/scp  
Remote shell (RSH) path: /usr/bin/ssh  
User Interface Domain (UID): dfw.ibm.com

There is also a "Cluster nodes" section with a table:

Host name	Is manager node	Is quorum node
cler01.dfw.ibm.com	✓	✓
cler02.dfw.ibm.com	✓	✓
cler03.dfw.ibm.com	—	—

Additional sections include "Cluster file systems" and "Cluster disks".

**Cluster file systems**

Name	Mount point
gpfs1v	/gpfs/with/a/really/deep/filename
gpfs1v	/gpfs

**Cluster disks**

Name	Usage	File system
descnsd	descOnly	gpfs1v
gpfs1013nsd	dataAndMetadata	gpfs1v
gpfs1014nsd	dataAndMetadata	gpfs1v
gpfs1018nsd	metadataOnly	gpfs1v
gpfs1nsd	dataOnly	gpfs1v
gpfs2nsd	dataAndMetadata	gpfs1v



## Common Commands

---

- View / Change GPFS cluster configuration
  - mmlscluster / mmchcluster
- View / Change GPFS configuration details
  - mmlsconfig / mmchconfig
- View / Change file system configuration
  - mmlsfs / mmchfs <GPFS device>
- View file system / storage pool usage
  - mmdf <required GPFS device>



## Common Commands

---

- View current GPFS threads on one node
  - mmfsadm dump waiters
- View / Change the state of a GPFS disk
  - mmlsdisk / mmchdisk
- View / Change GPFS NSD information
  - mmlsnsd / mmchnsd
- Add / remove nodes
  - mmaddnode / mmdelnode



## Other Tasks

---

**D = dynamic, ND= not dynamic**

- Modification of NSD servers ND
  - mmchnsd command  
file system must be unmounted
- Adding inodes to a FS D
  - mmchfs -F #
- Create / delete snapshots D
  - mmcrsnapshot (create, in the <fs>/snapshots directory), mmdelsnapshot (delete), mmlssnapshot (list)
- Parallel backup via TSM D
  - mmbackup (see docs for configuration)  
Utilizes snapshots under the covers.



## Other Tasks

---

**D = dynamic, ND= not dynamic**

- Start/stop GPFS daemon D
  - mmstartup / mmshutdown (-a for all)
- Configure Tiebreaker disks ND
  - mmchconfig tiebreakerDisks="nsd1;nsd2;nsd3"
    - NSDs can be part of file systems, don't have to be dedicated.
  - mmchconfig tiebreakerDisks=no
- Change GPFS cache size (pagepool) (requires stop/start of daemon) ND
  - mmchconfig pagepool=500M



## Other Tasks

---

**D = dynamic, ND= not dynamic**

- Assign the file system manager D
  - mmchmgr
- Determine which nodes have a file system mounted D
  - mmlsmount
- Mount / unmount file systems on (all) nodes D
  - mmmount / mmumount (-a)
- Restripe (rebalance) a file system D
  - mmrestripefs





## Other Tasks

---

**D = dynamic, ND= not dynamic**

- Enable SCSI Persistent Reserve ND
  - mmchconfig usePersistentReserve=yes
- ★ Assign the cluster configuration manager D
  - mmchmgr -c
- ★ Change node attributes D
  - mmchnode
- Modify actual replication state of data D
  - mmrestripefs -R



# Debugging GPFS

---

Check from the top down:

- Is it my membership in the cluster?
  - Can we communicate?
  - Is the daemon running?
- Is the file system mounted?
  - Are they all mounted?
- Is there a problem with the disks?
  - From the Operating System's point of view?
  - From a GPFS point of view?
- Performance issue?
  - Check out standard AIX performance commands
  - Examine the mmpmon command



## Debugging Tasks

---

- First, document the cluster:
  - `mmlscluster` Lists the nodes, pri/sec data servers
  - `mmlsconfig` Lists the cluster config values, FS devices
  - `mmlsnsd -M/-X` Lists all the NSDs and how they're viewed from the nodes
  - `mmlsmgr` Lists the file system managers
  - `mmlsfs <fs_dev>` Lists the attributes of the FS
- Determine the state of things
  - `mmgetstate -a`  
Lists the state of each quorum node in the GPFS cluster
  - `mmlsdisk <GPFS device> -L`  
Lists the state of each disk in the FS, as well as quorum info
  - `mmlsnsd <GPFS device> -M`



# Debugging Tasks

---

- Then read the logs:
  - `/var/adm/ras/mmfs.log.latest`      # Symbolic link to latest log
  - `/var/adm/ras/mmfs.log.previous`
  - `/var/adm/ras/mmfs.log.<timestamp>`
- Read the docs!
  - GPFS Problem Determination Guide
- Random tips:
  - If node was cloned from another GPFS node, then delete contents of `/var/mmfs/gen` and `/etc/cluster.nodes` before adding to GPFS cluster
  - If `mmgetstate` command fails for quorum nodes, check out the remote shell commands (must work any-to-any without prompting)



## Information Sources

---

The following URIs provide more information about GPFS.

- Cluster Resource Center:

- <http://publib.boulder.ibm.com/infocenter/clresctr/>

- GPFS product documentation

- [... /topic/com.ibm.cluster.gpfs.doc/gpfsbooks.html](.../topic/com.ibm.cluster.gpfs.doc/gpfsbooks.html)

- GPFS Version FAQ (very important for the latest updates):

- [... /index.jsp?topic=/com.ibm.cluster.gpfs.doc/gpfs\\_faqs/gpfsclustersfaq.html](.../index.jsp?topic=/com.ibm.cluster.gpfs.doc/gpfs_faqs/gpfsclustersfaq.html)

- GPFS Forum and mailing list:

- [http://www-128.ibm.com/developerworks/forums/dw\\_forum.jsp?forum=479&cat=13](http://www-128.ibm.com/developerworks/forums/dw_forum.jsp?forum=479&cat=13)

- <http://lists.sdsc.edu/mailman/listinfo.cgi/gpfs-general>

- GPFS Wiki

- <http://www.ibm.com/developerworks/wikis/display/hpccentral/General+Parallel+File+System+%28GPFS%29>



## Information Sources

---

The following URIs provide more information about GPFS.

- Main marketing web site, has links to Whitepapers:  
<http://www.ibm.com/systems/clusters/software/gpfs/index.html>
- Redbooks ([www.redbooks.ibm.com](http://www.redbooks.ibm.com))
  - Digital Media (SG24-6700)
  - Oracle RAC (SG24-7541)
- Fix download site:  
<http://www14.software.ibm.com/webapp/set2/sas/f/gpfs/home.html>