

Disaster Recovery using GPFS

Vladimir Sapunenko (CNAF)

Antonio Budano (Roma3)

for

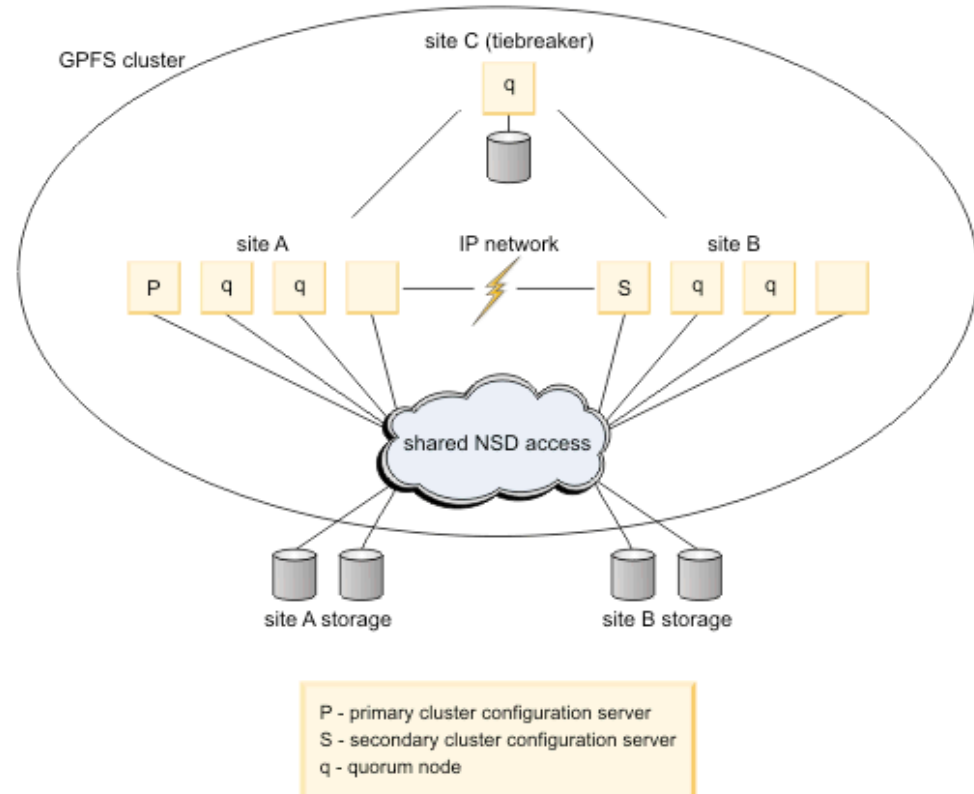
CCR Disaster Recovery Working Group

HA features of GPFS

- There is a number of features inside GPFS to facilitate implementation of HA environments against catastrophic HW failures
 - **Replication** of the file system's data at a geographically-separated site ensures the data availability in the event of a total failure of the primary (production) site
 - **Snapshot** function allows a backup process to run concurrently with user updates and assures consistency of the data used for backup
 - **AFM** enables sharing data across unreliable or high latency networks. Location and flow of file data between GPFS clusters can be automated

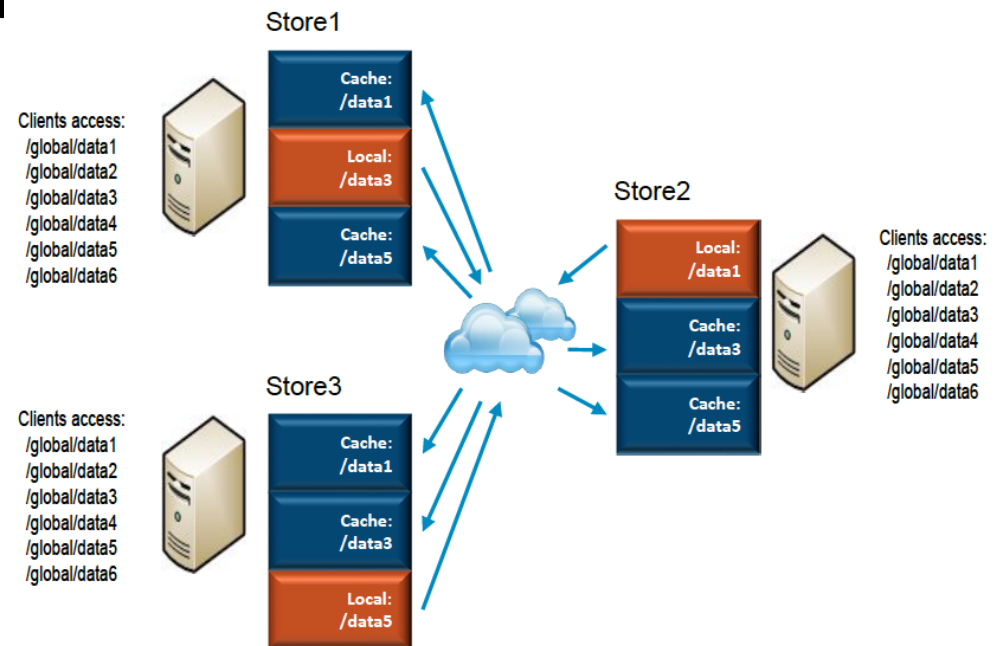
Synchronous mirroring using GPFS Replication

- Data and metadata replication of GPFS can be used to implement synchronous mirroring between a pair of geographically separate sites



Active File Management

- can be used to create a global namespace within a data center or between data centers located around the world.
- AFM is designed to enable efficient data transfers over wide area network (WAN) connections.
- Transfer home -> cache can happen in parallel within a node called a *gateway* or *across multiple gateway nodes*.



Replication + Snapshot + AFM = Complete Solution

- 3 or 4 geo separated sites

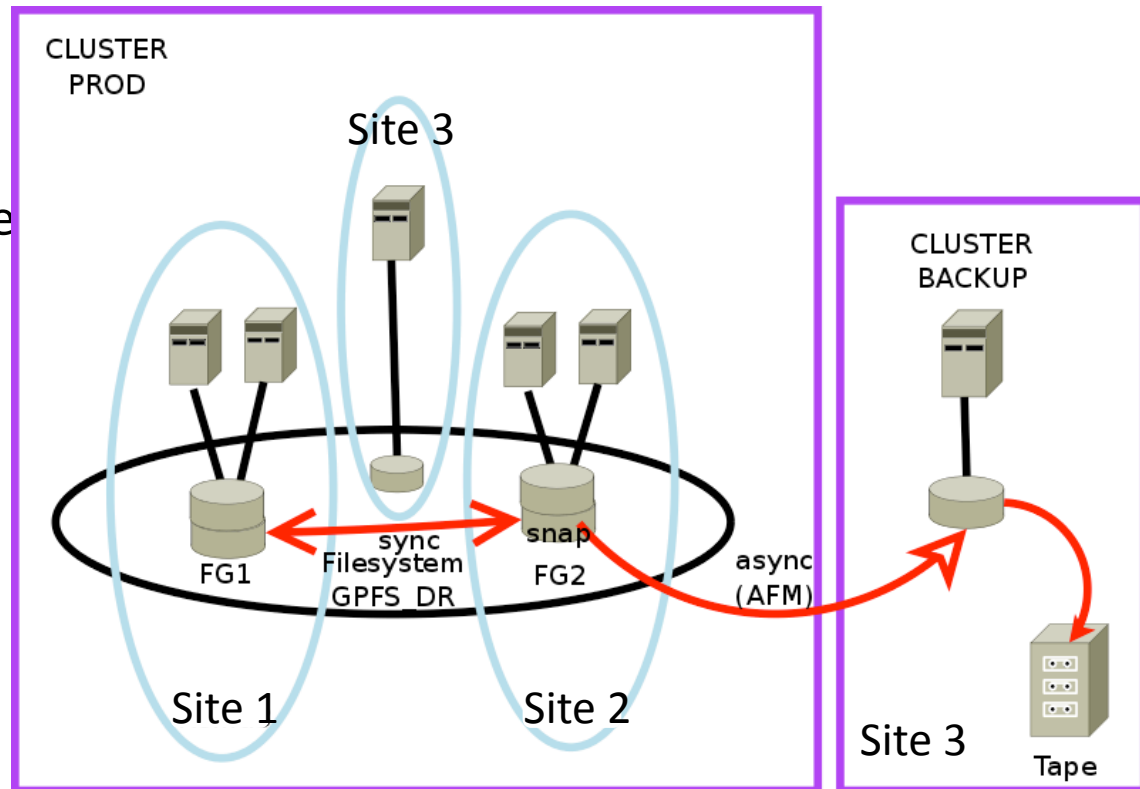
- 2 sites in close vicinity to compose HA cluster

- 1 tie breaker site

- Keeping also FS descriptor and cluster configuration

- 1 backup site

- Can coincide con tie breaker



- Backup can be done from a Snapshot copied to backup site via AFM
 - Backup window = time to stop/sync/start application
 - All data transferred to backup site in background (asynchronously)
 - Backup will be kept in 4 copies (2 on disk in prod cluster, 1 on disk and 1 on tape in backup cluster)

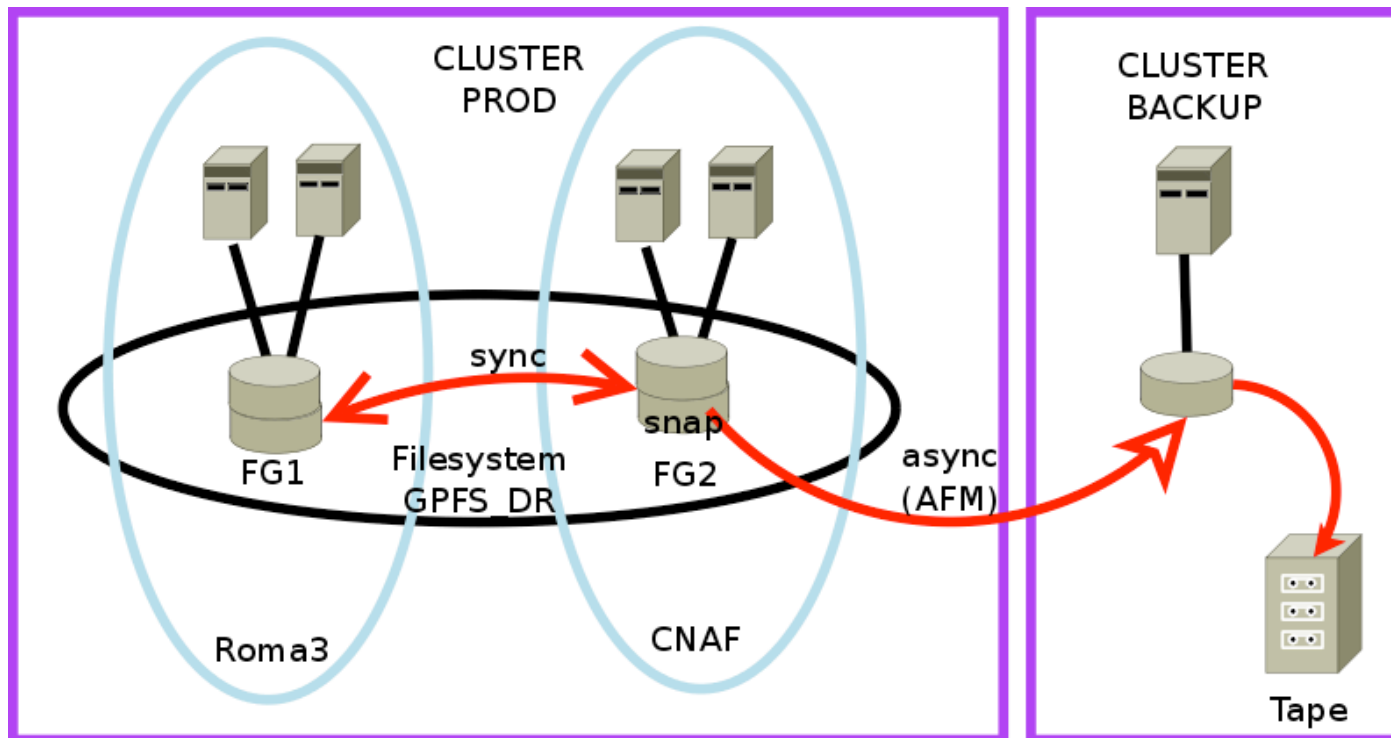
Failure scenarios

failures	effects	actions	down time
Disk on site1	Switching to access disk remotely from site 2	non	0
WAN network connection to site 1	no access to data, application crashes or hangs	ensure that application is not running on site 1, restart application on site2	t1
Site1 failure		restart application on site2	t1
Site3 (tiebreaker) failure	non	non	0
site2 and site3 failure	No access to data, file system down application crashes or hangs	reconfigure quorum nodes, restart application	t1+ 1min

t1 = time to restart application

Testing Sync Replication between CNAF and Roma3

- Bandwidth 1 Gbit/s, RTT=6ms
- 2 servers + 2TB of disk from each side
- Sync writes (RTT penalty in case of random or IO intensive operations)
- All reads are local (no any influence from the remote site)



Tests performed

- Sequential I/O using “*dd*”
- Sequential and Random I/O using “*IOzone*”
- Failure of remote site (disabling Eth port on remote servers) while writing
- Failure of local disk (disabling FC ports on storage system) while writing
- Running *IOzone* in *VMware* based VM
- Running SI applications on VM (KVM) resided on the geo-replicated file system

Some numbers from IOzone

	Sequential write, 1MB blocks, MB/s	Sequential read, 1MB blocks, MB/s	Random write, 1MB blocks, MB/s	Random Read, 1MB blocks, MB/s
Local disk, Roma3	168	160	160	84
Replicated FS From Roma3	51	168	74	55
Local disk CNAF	164	176	160	64
Replicated FS, from CNAF	108	196	90	90

Observations

- Sequential writes are limited by network bandwidth between sites
- Sequential reads are limited by local disk performance
- Random IO are mostly affected by latency (RTT) between sites
- Recovery in case of secondary site failure depends on configuration parameter (default 60 sec)
- Recovery in case of local disk failure almost instant
- VM performance depends on partitioning and format of the VM disks
- SI applications and services (including Oracle) are starting correctly (no performance measurements done so far)

Conclusions

- GPFS provides all necessary means to build robust Disaster Recovery solution
- Verified and widely used in industry
- Such a solution can guarantee continuity of operations by
 - Instant Failover to the secondary site if the primary goes down
 - Failover to backup site with data recovery from backup if both Production sites become inaccessible
 - There is also a possibility to recover all data from the previous backup locally from the snapshot (for example when some data were deleted from disk because of human error)

Gruppo

Coordinatore : Stefano Zani

Componenti : Sandro Angius
Massimo Donatelli
Claudio Galli
Guido Guizzunti
Dael Maselli
Massimo Pistoni
Claudio Soprano
Riccardo Veraldi

+ Collaborazione :

Nunzio Amanzi
Alessandro De Salvo
(indicazioni su distributed File System)
(GPFS) Vladimir Sapunenko
Antonio Budano

Aree di lavoro



DNS {distribuito + geo-replica}

MAILING {distribuito + mail relay }

SISTEMA INFORMATIVO :

Contabilità (CNAF) (→R12)

Portale Utente (CNAF)

Gestione Presenze (CNAF)

Documentale Alfresco (CNAF) [new]

Business Intelligence BI (CNAF)

Protocollo (CNAF) [new]

AAI + GODIVA (LNF)

Stipendiale + Sxgest2 (LNF)

Stipendiale + Cezanne (LNF) [new]

Protocollo (LNF)

Documentale (LNF)

Portale Unico (LNF) [new]

...