

Minerva Storage and File System Upgrade

Town Hall Presentation
Scientific Computing

Patricia Kovatch

Bhupender Thakur, PhD

Francesca Tartaglione, MS

Dansha Jiang, PhD

Eugene Fluder, PhD

Hyung Min Cho, PhD

Lili Gai, PhD

Aug 10, 2017



**Mount
Sinai**

The Urgency of the Upgrade

Short of storage space:

8 PB is used (75% full) and is fully allocated.

Reaching end of disk life:

IBM GSS was purchased at 2013

and neither IBM nor Lenovo will provide any more support for it.

DDN 10K was purchased in 2012

and DDN will not provide any more support for it next year.

Out of support:

Spectrum scale (GPFS) 3.5 is out of support.



Features of Elastic Storage System (ESS)

More storage space: 6 PB usable data space

- ▶ Six Storage 4U Enclosures, each with 58 8 TB SAS drives.

Note: After the GSS and DDN 10K are removed, only 1PB will be added to the overall file system.

Higher throughput:

Three pairs of I/O nodes each with

- ▶ Two 10 core power8 processors, 128 GB DRAM,
- ▶ 6 EDR IB high speed data connections (100Gb/s),
- ▶ 3 multilane SAS controllers.

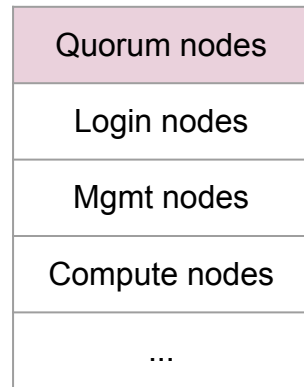
Higher code version: IBM Spectrum Scale (GPFS) 4.2

- ▶ A better way of organizing massive data.
- ▶ More functions for data management.
- ▶ Features such as local caches, NFS and samba support.
- ▶ Numerous bugs are fixed that improve reliability and performance

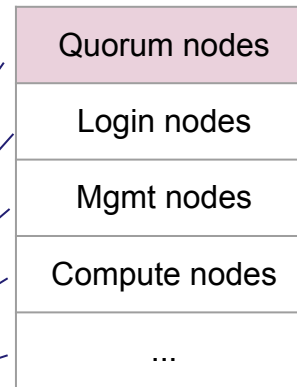
Previous Storage and File System Structure

Client side:

Cluster Manda
(GPFS 3.5)



Cluster Mothra
(GPFS 3.5)



1G Ethernet network

Infiniband (IB) network

File-system side:

/sc/orga:
(GPFS 3.5)

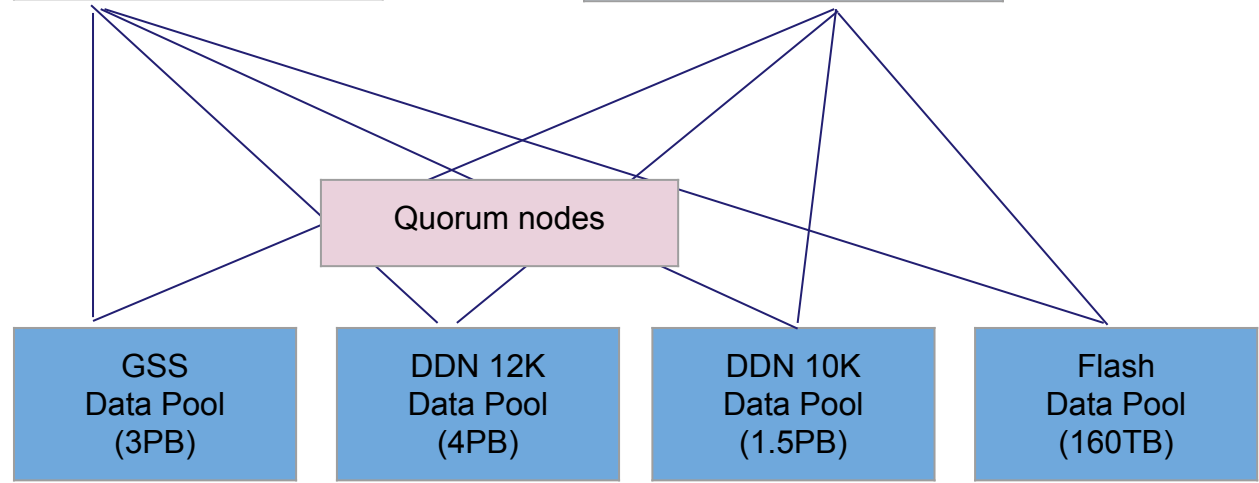
Quorum nodes

GSS
Data Pool
(3PB)

DDN 12K
Data Pool
(4PB)

DDN 10K
Data Pool
(1.5PB)

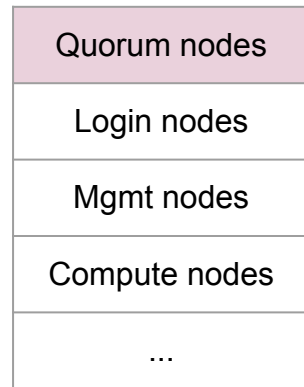
Flash
Data Pool
(160TB)



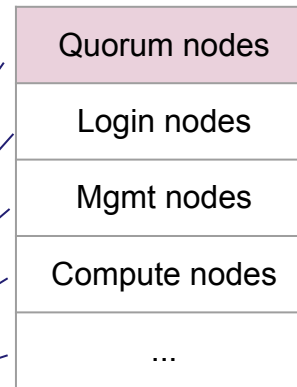
Storage and File System Upgrade Plan

Client side:

Cluster Manda
(GPFS 4.1=>4.2)



Cluster Mothra
(GPFS 4.1=>4.2)



1G Ethernet network

Infiniband (IB) network

File-system side:

/sc/orga:
(GPFS 4.1=>4.2)

Quorum nodes

ESS
Data Pool
(6PB)

~~GSS
Data Pool
(3PB)~~

DDN 12K
Data Pool
(4PB)

~~DDN 10K
Data Pool
(1.5PB)~~

Flash
Data Pool
(160TB)



Data Transfer

ESS Upgrade PLAN Overview

By IBM, DDN, Mellanox and us

- 1. Install ESS with GPFS downgrade to 4.1 and test ESS as stand alone cluster.
 - ➔ 2. Migrate data off GSS to ESS.*
 - ➔ 3. Upgrade client cluster to GPFS 4.1.
 - 4. Upgrade DDN 12K/10K and Flash to GPFS 4.1.
 - 5. Remove GSS from ORGA.
 - 6. === Switch whole cluster to GPFS 4.1 =====
-
- 1. Upgrade client cluster and file system cluster to GPFS 4.2.
 - 2. May also retire DDN 10K.
 - 3. === Switch whole cluster to GPFS 4.2 =====

Goal: Minimal number of outages



Technical Difficulties Experienced Over the Last Week

ESS not setup correctly:

ESS passed tests but is behaving unexpectedly, taking itself out of the cluster seemingly randomly.

Compatibility of code version:

ESS comes with GPFS 4.2, which is not compatible with GPFS 3.5 in our current system.

We downgraded the ESS GPFS to 4.1 which is supposed to be compatible with GPFS 3.5.

Mixed vendor/GPFS version file system pools:

- ▶ IBM (ESS, GSS, Flash)
- ▶ DDN (DDN10k, DDN12k)

Data migration:

- Policy engine: move data across different data pools but does not work reliably on our cluster though IBM cannot explain why.
- Restripe: data transfer within a data-pool to improve layout and performance.

Hidden IB-net issues:

IB-switches are behind in firmware, and ofed versions.

Nodes and switches are running with mixed code versions. Upgrading means more outages.

Mixed EDR, FDR, QDR cables and connections.

Fixes Applied During This Week's Outages

- ▶ Upgraded GPFS from 4.1.1.8 to 4.1.1.15 on the ESS I/O nodes to fix the deadman lock bug and resolved the kernel panic.
- ▶ Changed parameters on ESS to increase the cache.
- ▶ Updated firmware and upgraded OFED on ESS I/O nodes and ESS switches.
- ▶ Changed IB parameters and settings to resolve IB issue.

Current status:

Data migration stopped after the first outage.

Resumed GSS disks and suspended all ESS disks.

Plan For Future:

- ▶ Crit-Sit (Critical Situation) case registered with IBM.
- ▶ Revisit our entire ESS upgrade plan with IBM.
- ▶ Actively engage IBM/DDN/Mellanox for inspection of parameters and settings cluster-wide.
- ▶ Test ESS heavily as a stand-alone file system.
- ▶ **Schedule weekly Preventative Maintenance periods to test changes to the system:**
 - Tuesdays from 8 AM-8 PM weekly (will cancel if not needed).**

We are open for your suggestions!

Thank you for your understanding and patience.

Thank you!



Minerva supercomputer @Mount Sinai