# Experiences using a multi-tiered GPFS file system at Mount Sinai

Bhupender Thakur

Patricia Kovatch

Francesca Tartagliogne

Dansha Jiang
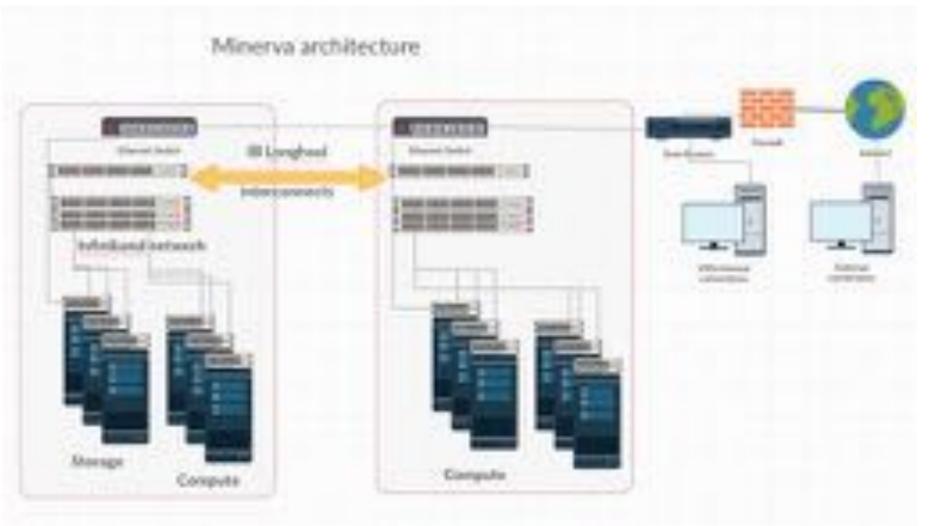
Mount Sinai
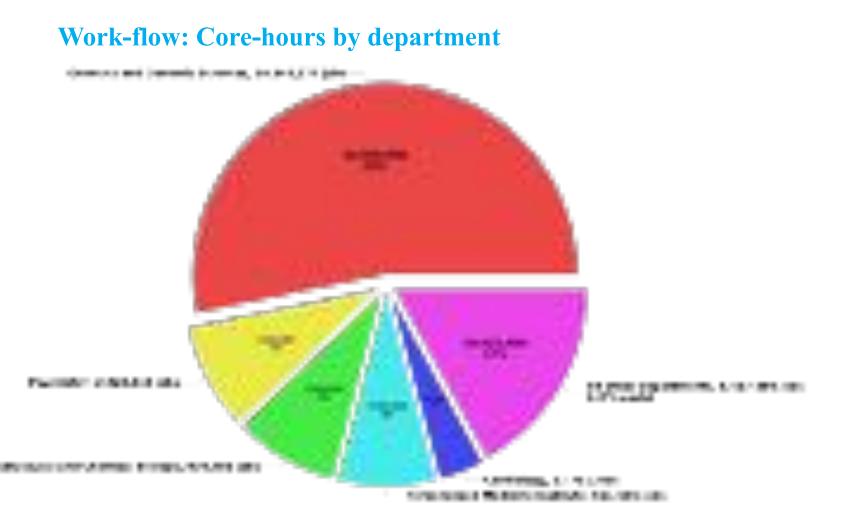
# Outline

1. Storage summary

2. Planning and Migration

3. Challenges

4. Conclusions/Feedback

# Storage

# Storage
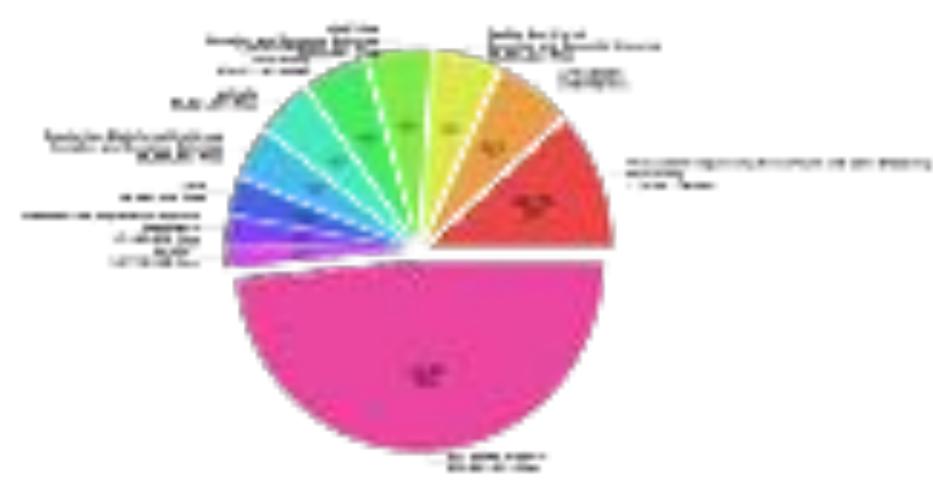


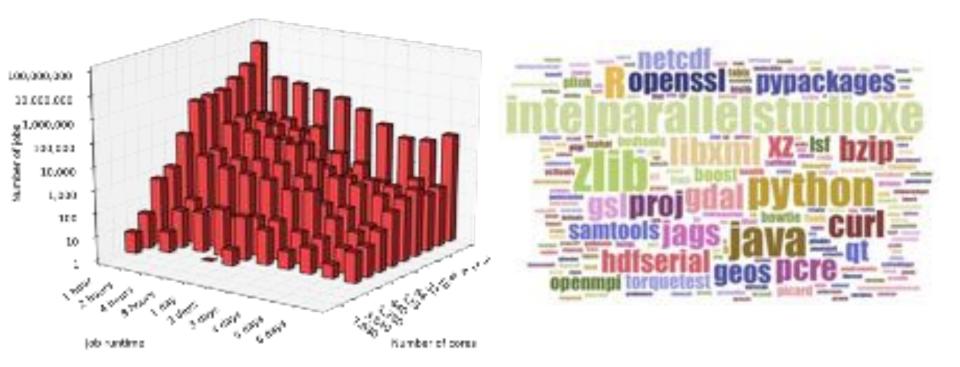Minerva architecture

# Work-flow: Core-hours by department



61,206,676 core-hours used by 28,672,552 jobs from 89 departments

# Work-flow: GPFS Project storage (/sc/orga/)



5.6 PB used by from 222 projects

# Job mix



- Mostly serial, pipeline based workflows
- Few MPI and multi-node jobs

# Storage structure

- Single filesystem available as "project" directories

- Single IB fabric across two datacenters ~0.5 km apart

- Flash tier for metadata

- Storage charge per Sinai policy. Charge based on usage rather than quota.

GSS:
4xGL6
Building
blocks
3 PB

FLASH
8xFS820
150 TB

DDN:
SFA10k
1.4 PB

DDN:
SFA12k
3.9 PB

ESS-BE
3xGL6
Building
blocks
5.6 PB

ESS-LE
1xGL6S
Building
block
4PB

FLASH
2xGSS2S
260 TB

2014

2018

2020

# Storage structure

The GSS and Flash sub-systems have been our primary workhorses.

| Storage | Avg. Read Ops/Day | Avg. Write Ops/Day | Avg. Reads /day | Writes/day |
|---|---|---|---|---|
| IBM GSS: Data Subsystem (3.5PB) | ~650 million | ~385 million | ~300 TB | ~100TB |
| IBM FLASH: Small File/ Metatadata Subsystem (150TB) | ~1050 million | ~600 million | - | - |
| Total(Including DDN Subsystems) | ~1.8 billion | ~1 billion | ~500TB | ~175TB |

# Storage structure

- The GSS and Flash sub-systems have been our primary workhorses.

- Organic growth despite recent charge-by-use model.

# Migration plan

| Pool | Total(PB) | Used(PB) | % Used |
|------|-----------|----------|--------|
| DDN10k | 1.4 | 0.8 | 61.02% |
| DDN12k | 3.9 | 1.8 | 47.78% |
| GSS | 3.1 | 2.7 | 85.83% |
| ESS - BE | 4.1 | | |

*Migration (Plan A)*

1. Install ESS at lower ESS/GPFS(4.1) code level to be compatible with existing GSS (max v.3.5.25 allowed for Lenovo GSS)
2. Split ESS – BE recovery group disks between two pools: new "ess" pool and old "gss".
3. Migrate GSS data within gss pool
   - Suspend GSS disks
   - Restripe and rebalance
4. Migrate data from gss -> ess pool
5. Repeat for ddn10k and add remaining ess disks to new "ess" pool
6. Upgrade GPFS code to 4.1 -> 4.2, update release version, fs version

# Migration plan

| Pool | Total(PB) | Used(PB) | % Used |
|------|-----------|----------|--------|
| DDN10k | 1.4 | 0.8 | 61.02% |
| DDN12k | 3.9 | 1.8 | 47.78% |
| GSS | 3.1 | 2.7 | 85.83% |
| ESS - BE | 4.1 | | |

*Migration (Plan A: Take 2)*

1. Remove ESS – BE Rebuild test cluster with ESS v5.2.0, downgrade GPFS to 4.1.1.15
2. Test/Debug new filesystem.
   - Debug IB issues (To bond or not to bond)
   - Test compatible client OFED levels
   - Upgrade FLASH/DDN GPFS levels
3. Remove disks from test cluster and add them to a new pool in production cluster.

# Migration plan

| Pool | Total(PB) | Used(PB) | % Used |
|------|-----------|----------|--------|
| DDN10k | 1.4 | 0.8 | **61.02%** |
| DDN12k | 3.9 | 1.8 | 47.78% |
| GSS | 3.1 | 2.7 | **85.83%** |
| ESS - BE | 4.1 | | |

*Migration (Plan B)*

1. Squeeze allocations. Migrate projects data from 52 known projects.
2. Migrate as much s possible to ddn12k
   - from ddn10k -> ddn12k
   - from gss -> ddn12k
3. Migrate (*)
   - from ddn10k -> ess
   - from gss -> ess
4. Cleanup
   - Change default pool from gss to ess
   - Remove ddn10k, gss disks

# Policy moves ( pool: ddn10k)

| Pool | Total (PB) | Before | After | Before(%) | After(%) |
|---|---|---|---|---|---|
| ddn10k1 | 1.4 | 0.8 | 0 | 57.13% | 3.24% |
| ddn12k1 | 3.9 | 2.3 | 2.3 | **59.22%** | **59.15%** |
| ess | 4.1 | 0.2 | 0.9 | **4.21%** | **22.41%** |
| gss | 3.1 | 2.3 | 2.3 | 73.61% | 73.51% |

Evaluating policy rules with CURRENT_TIMESTAMP =
   **2018-01-08@15:39**:23 UTC
[I] **2018-01-08@16:49:52.650** Directory entries scanned: 1888732606.
…
[I] Directories scan: 1457460636 files, 271560756 directories, 159711214 other objects,
[I] 2018-01-08@17:05:23.719 Statistical candidate file choosing. 498694530 records processed.

[I] Summary of Rule Applicability and File Choices:
Rule#    Hit_Cnt         KB_Hit          Chosen       KB_Chosen         KB_Ill  Rule
0          498694530      741777261952.  498694530.  741777261952      0          RULE 'migrate10ktoess.1'
MIGRATE FROM POOL 'ddn10k1' TO POOL 'ess'
[I] A total of 498694530 files have been migrated, deleted or processed by an EXTERNAL EXEC/script;
       8336739 'skipped' files and/or errors.
…
[I]**2018-01-21@00:16:50.361** Exiting. RC=0.

   *700T @ 50T/day in 2 weeks*

# Policy moves ( pool: gss)

| Pool_Name | Total | Used PB (Before) | Used PB (After) | Before(%) | After(%) |
|---|---|---|---|---|---|
| ddn10k1 | 1.38 | 0 | 0 | 0.02% | 0.02% |
| ddn12k1 | 3.87 | 2.75 | 2.75 | 71.20% | 71.20% |
| ess | 4.07 | 0.9 | 2.73 | **22.17%** | **66.90%** |
| gss | 3.12 | 1.84 | 0.01 | 58.90% | 0.45% |

Evaluating policy rules with CURRENT_TIMESTAMP = 2018-02-05@16:35:39 UTC

…
[I] 2018-02-05@18:05:04.305 Statistical candidate file choosing. 2813187 records processed.
[I] Summary of Rule Applicability and File Choices:
Rule#    Hit_Cnt    KB_Hit    Chosen    KB_Chosen    KB_Ill    Rule    0    2813187    1822704400896    2813187    1822704400896    0    RULE 'migrategssktoess.1' MIGRATE FROM POOL 'gss' TO POOL 'ess' WHERE(.)
…
[I] 2018-02-16@14:12:49.891 Executing file list:
[E] Error on writeEOR to work file localhost:42624: Broken pipe

[I] A total of 2503335 files have been migrated, deleted or processed by an EXTERNAL EXEC/script;    0 'skipped' files and/or errors.

[I] 2018-03-20@15:44:28.062 Policy execution. 2503335 files dispatched.
*2PB @ 35TB/day in 7 weeks*

# Policy moves: Things to consider

Things for consideration:

- Running multiple policies(such as purging) can be problematic.
- On a single fabric, HCA/Switch firmware upgrade can cause several hard IB sweeps.
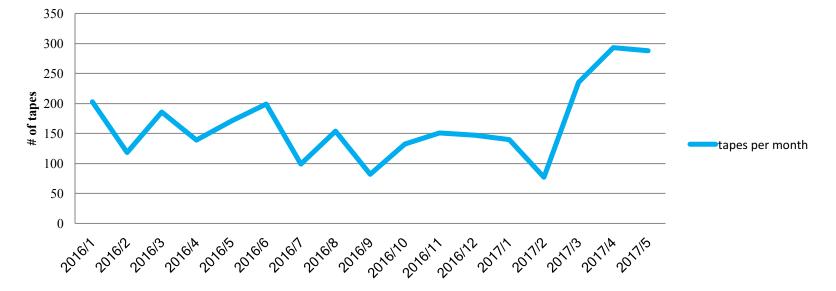- GNR/OFED parameters needed to be carefully checked.

# Archival storage

# Archival storage history

| Tape usage | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Month** | | | | | | | | | | | | | |
| **Year** | **1** | **2** | **3** | **4** | **5** | **6** | **7** | **8** | **9** | **10** | **11** | **12** | **Grand Total** |
| 2016 | 203 | 118 | 186 | 139 | 171 | 199 | 99 | 154 | 82 | 132 | 151 | 147 | 1,781 |
| 2017 | 140 | 77 | 235 | 293 | 288 | | | | | | | | 1,033 |
| **Grand Total** | **343** | **195** | **421** | **432** | **459** | **199** | **99** | **154** | **82** | **132** | **151** | **147** | **2,814** |

## # of Tapes per month

# Archival storage history

**Statistics for the last 365 days**

Number of newly registered users= 340
Number of users that did archives = 78
Number of archives = 23,771
Number of retrieves = 7,065

**Current storage usage**

Archived data in TB = 1,495
Primary storage pool in TB = 4,495
Number of tapes used = 5,750
Number of tapes onsite = 2,968
Number of tapes offsite = 2,782

# Challenges

- Sinai growth: Ten-fold increase in users/storage since our first GPFS install in 2013.
- Sinai future directions: Fully HIPAA compliant cluster.

- Every new ESS install is starting to require a lot more thought and planning
  - xCAT support outside of regular
  - Ethernet network integration: VLAN and switch considerations
  - High-speed network integration
  - Data migration and planning
- GNR configurations present an added complexity for upgrades
- ESS recovery group stability

# Imminent challenges

- Integrate ESS BE and ESS LE clusters
  - Lack of "official" support in merging xCAT networks
    GS8052 and Mellanox 8831-S52 trunk
  - xCAT table updates and configuration changes (name servers)
  - HMC

- Migrating the metadata servers

# Conclusions/Feedback

- Its more challenging for admins to maintain a unified filesystem, but its easier for customers.

- Long time from "Deploy" status(code 20) to actual production use.

- Would like to thank IBM CE and Lab Services and dev teams who have been very patient in handling our issues.