

IBM SPECTRUM SCALE

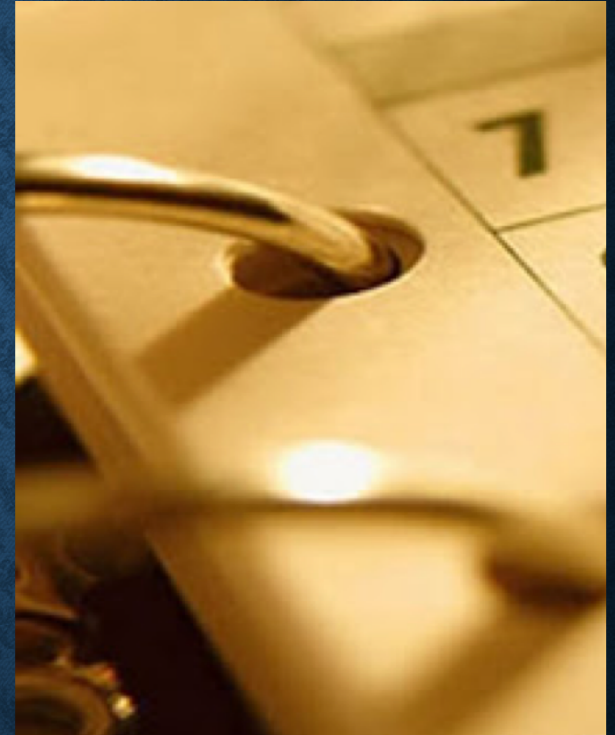
Support update, common issue and best practice

Guanglei LI
liguanglei@cn.ibm.com
March 2018



Agenda

- Better Follow-The-Sun Support Coverage
- Support Managers
- Common Issues in Field and Best Practices



Follow the sun support – Aligning support staff to customer time zone

- Spectrum Scale Support is growing to better meet customer needs.
- Beginning late 2016 we substantially grew the support team in Beijing, China, with experienced Spectrum Scale staff.
- Improved response time on severity 1 production outages; reducing customer waiting time before L2 is engaged as well as time to resolution.
- Positive impact to timely client L2 communication for severity 2, 3, and 4 PMRs within our customer time zone.
- Setup and grew EMEA support team in Germany in late 2017
- 3 major sites: North America, China, Germany
- PagerDuty was introduced this year for better PMR monitor



IBM Spectrum Scale Level 2 Support Global Time Zone Coverage



Global team locations

- North America
 - ✓ *Poughkeepsie, NY USA
 - ✓ Toronto, ON Canada
- AP
 - ✓ *Beijing, China
 - ✓ India
- Europe
 - ✓ *Germany

* Major sites



Support Delivery: Managers

1st Level: Bob Simon: ragonese@us.ibm.com; 1-845-433-7285

1st Level: Jun Hui Bu: bujunhui@cn.ibm.com; 86-10-8245-4113

1st Level: Dennis Kunkel: Dennis.Kunkel@de.ibm.com; 49-170-3387365

WW 2nd Level: Wenwei Liu: wliu@ca.ibm.com; 1-905-316-2623

Support Executive

Andrew Giblon: agiblon@ca.ibm.com; 1-905-316-2582



COMMON FIELD ISSUE AND BEST PRACTICES



DATA COLLECTION: GPFS.SNAP

1) Use the "--limit-large-files" flag to limit the amount of 'large files' collected. The 'large files' are defined to be the internal dumps, traces, and log dump files that are known to be some of the biggest consumers of space in gpfs.snap (these are files typically found in /tmp/mmfs of the form internaldump.*, trcrpt.*, logdump*.*).
Added in version 4.1.1

--limit-large-files: YYYY:MM:DD:HH:MM | Num_Days_back | 0

2) Limit the nodes on which data is collected using the '-N' flag to gpfs.snap. By default data will be collected on all nodes, with additional master data (cluster aware commands) being collected from the initiating node.

- Note: Please avoid using the -z flag on gpfs.snap unless supplementing an existing master snap or you are unable to run a master snap.

3) To clean up old data over time, it's recommended that gpfs.snap be run occasionally with the '--purge-files' flag to clean up 'large debug files' that are over the specified number of days old. added in version 4.2.0

--purge-files: KeepNumberOfDaysBack | 0



FIRST TIME DATA COLLECTION FOR PERF/HANG

1. Gather waiters and create working collective. It can be good to get multiple looks at what the waiters are and how they have changed, so doing the first `mmlsnode` command (with the `-L`) numerous times as you proceed through the steps below might be helpful (specially if issue is pure performance, no hangs).

```
mmlsnode -N waiters > /tmp/waiters.wcoll
```

```
mmdsh -N /tmp/waiters.wcoll "mkdir /tmp/mmfs 2>/dev/null"
```

```
mmlsnode -N waiters -L | sort -nk 4,4 > /tmp/mmfs/service.allwaiters.$(date +"%m%d%H%M%S")
```

2. View `allwaiters` and `waiters.wcoll` files to verify that these files are not empty. If either (or both) file(s) are empty, this indicates that the issues seen are not GPFS waiting on any of its threads. Data to be gathered in this case will vary. Do not continue with steps. Tell Service person and they will determine the best course of action and what docs will be needed.

3. Gather `internaldump` from all nodes in the working collective

```
mmdsh -N /tmp/waiters.wcoll "/usr/lpp/mmfs/bin/mmfsadm dump all > /tmp/mmfs/service.\$(hostname -s).dumpall.\$(date +"%m%d%H%M%S")"
```



FIRST TIME DATA COLLECTION FOR PERF/HANG CONT.

4. Gather kthreads from all nodes in the working collective

```
mmdsh -N /tmp/waiters.wcoll "/usr/lpp/mmfs/bin/mmfsadm dump kthreads > /tmp/mmfs/service.\$(hostname  
-s).kthreads.\$(date +"%m%d%H%M%S")"
```

*note:

If running Linux OS on SpectrumScale (formerly GPFS) 4.1 or higher - this step could be skipped.

5. If this is a performance problem, get 60 seconds mmfs trace from the nodes in the working collective.

If AIX ...

```
mmtracectl --start --aix-trace-buffer-size=256M --trace-file-size=512M -N /tmp/waiters.wcoll ; sleep 60;
```

```
mmtracectl --stop -N /tmp/waiters.wcoll
```

If Linux ..

```
mmtracectl --start --trace-file-size=512M -N /tmp/waiters.wcoll ; sleep 60; mmtracectl --stop -N /tmp/waiters.wcoll
```

6. Run gpfs.snap to collect all the data generated

```
gpfs.snap -N /tmp/waiters.wcoll
```



PERFORMANCE TUNING

- 1) **pagepool** - cache user file data and file system metadata
Needs to understand the IO pattern on client nodes when tuning pagepool:
Sequential IO, Random IO, Direct IO
- 2) **maxFilesToCache** - controls how many file descriptors each node can cache.
 - Needs large value if there will be many files opened concurrently, e.g., 1M for NFS & Samba service. Large value can improve the performance of user interactive operations like running "ls"
 - Small value with many files being accessed will cause high CPU usage
 - Increasing maxFilesToCache in a large cluster with hundreds of nodes increases the number of tokens a token manager needs to store. Ensure that the manager node has enough memory and tokenMemLimit is increased when running GPFS version 4.1.1 and earlier.
- 3) **workerThreads** - controls an integrated group of variables that tune the file system performance
 - New in GPFS 4.2.0.3 to simplify tuning. Some variables are auto-calculated when WorkerThreads is enabled. e.g, worker1Threads, worker3Threads
 - You can manually adjust external variables to avoid auto-tuned by workerThreads when Spectrum Scale computed from WorkerThreads are not suitable for your workload
 - Default 48. Increase to 512 or 1024 if there will be many threads access GPFS file system on that node. e.g., running NFS and Samba service on that node



PERFORMANCE TUNING CONT.

1. defaultHelperNodes – Specify the nodes to be used for distributed commands
 - Command list: mmaddddisk, mmapplypolicy, mmbackup, mmchdisk, mmcheckquota, mmdefragfs, mmdeldisk, mmdelsnapshot, mmfileid, mmfsck, mmimgbackup, mmimgrestore, mmrestorefs, mmrestrieps, mmrpldisk
 - Example: runningmmrestrieps on limited nodes including NSD servers
2. maxMBps - indicates the maximum throughput in megabytes per second that GPFS can submit into or out of a single node
 - It's a hint GPFS uses to calculate how many prefetch/writebehind threads should be scheduled
 - Set client nodes maxMBpS based on IO throughput. 2x of total IO throughput divided by # of client nodes



FS CORRUPTION

1) MMFS_FSSTRUCT error

- It will be printed into system log if GPFS detect FS corruption when access the file system.
- `fsstructlx.awk(Linux)` `fsstruct.awk(AIX)` under `/lpp/mmfs/samples/debugtools/` to decode the MMFS_FSSTRUCT message in system log:
`fsstructlx.awk /var/log/messages > fsstruct.message`
- `mmhealth` will report FS corruptions

2) Offline mmfsck to check file system and generate report

- GPFS file system needs to be unmounted from all nodes.
- Use patch file option (from 4.1.1) to avoid two rounds of long running `mmfsck`:

```
mmfsck -nV --patch-file /tmp/fsck.patch
```

- **Online mmfsck**

- run `mmfsck` with `-o` option while FS is mounted
- Can only fix the lost blocks – data block marked as used but not referenced¹² by any file/dir



FS CORRUPTION CONT.

3. Upload mmfsck output and patch file for IBM to review. Additional output may be required:

- tsfindinode to identify the pathname for corrupted inodes. Needs to mount FS
- tsdbfs output for inode dumps

4. Run offline mmfsck fix under guidance of IBM support

- If patch is used, run it with:

```
mmfsck <fs> -V --patch-file /tmp/mmfs/fsck.patch --patch
```

5. Log recovery failure

- mmfsck <fs > -xk
 - Needs to unmount FS
 - Supported in ver >=4.2
 - Run it after confirmed with IBM support.



BEST PRACTICE: NSD MISSING

1) Disk Missing

- Use “mmlsnsd -X” to check if any disk reported as “(not found)”
- Use “tspreparedisk -s” on each node to check if a NSD could be identified.
- `mmnsddiscover -a -N all`
- User exit of `/var/mmfs/etc/nsddevices` could affect NSD discovery
- Disk type mismatch: `mmchconfig updateNsdType=<nsd_type_file>`

2) Disk Header Missing

- There are 3 parts in NSD header: NSD desc, Disk desc, FS desc.
- “`mmfsadm test readdescraw /dev/dev_name`” could be used to show headers.
- Use `tspreparedisk` & `dd` command to restore NSD header. Do this under guidance of IBM support, and not able to restore in some cases.
- A common cause for header missing: disk header erased by UEFI driver update [link](#)



BEST PRACTICE: EXPEL

1) Network

- GPFS will send out pings before expel a node:
... is being expelled because of an expired lease. Pings sent: 60. Replies received: 0
- Common causes
 - Mis-matched MTU size: Jumbo Frames enabled on some or all nodes but not on the network switch.
 - Old adapter firmware levels and/or incorrect OFED software are utilized
 - OS specific (TCP/IP, Memory) tuning has not been re-applied.
 - verbsRdmaSend is enabled for SS ver < 5.0. It has scaling issue in GPFS 3.x and 4.x [link1](#) [link2](#)
 - Node A can't talk with Node B. Node A will ask Cluster Manager to expel Node B. Node A or Node B will be expelled.

2) Node load

- GPFS cluster manager is too busy to handle incoming lease request. Avoid overloading cluster manager on large scale cluster
- GPFS >= 4.2.3 support Prioritization of critical RPCs including lease request
- Increase failure detection time for node expel:
mmchconfig minMissedPingTimeout=120 (default is 3)
mmchconfig maxMissedPingTimeout=120 (default is 60)
mmchconfig leaseRecoveryWait=120 (default is 35)



BEST PRACTICE: EXPEL CONT.

1) Expel auto data collection from 4.1.1

- **When a node is about to be expelled for unknown reasons, debug data is collected automatically to help find the root cause**
- **Controlled by config parameter: `expelDataCollectionDailyLimit`, `expelDataCollectionMinInterval`**
- **Expel debug data will be collected on cluster manager and involved nodes.**

2) Auto data collection for unhealthy TCP connections from 4.2.3.

- **GPFS log(`var/adm/ras/mmfs.log.laest`):**
The TCP connection to IP address 192.168.38.52 c38f2bc1n02 <c0n4> (socket 45) state is unexpected:
ca_state=0 unacked=46 rto=25856000
- **Controlled by expel Data collection parameters.**



SPECTRUM SCALE ANNOUNCE FORUMS

Monitor the Announce forums for news on the latest problems fixed, technotes, security bulletins and Flash advisories.

<https://www.ibm.com/developerworks/community/forums/html/forum?id=11111111-0000-0000-0000-000000001606&ps=25>

Subscribe to IBM notifications (for PTF availability, Flashes/Alerts):

<https://www-947.ibm.com/systems/support/myview/subscription/css.wss/subscriptions>



ADDITIONAL RESOURCES

Tuning parameters change history:

https://www.ibm.com/support/knowledgecenter/STXKQY_4.2.2/com.ibm.spectrum.scale.v4r22.doc/blladm_changehistory.htm?cp=STXKQY

ESS best practices:

https://www.ibm.com/support/knowledgecenter/en/SSYSP8_3.5.0/com.ibm.spectrum.scale.raid.v4r11.adm.doc/blladv_planning.htm

Tuning Parameters:

[https://www.ibm.com/developerworks/community/wikis/home?lang=en#!/wiki/General%20Parallel%20File%20System%20\(GPFS\)/page/Tuning%20Parameters](https://www.ibm.com/developerworks/community/wikis/home?lang=en#!/wiki/General%20Parallel%20File%20System%20(GPFS)/page/Tuning%20Parameters)

Share Nothing Environment Tuning Parameters:

<https://www.ibm.com/developerworks/community/wikis/home?lang=en#!/wiki/General%20Parallel%20File%20System%20%28GPFS%29/page/IBM%20Spectrum%20Scale%20Tuning%20Recommendations%20for%20Shared%20Nothing%20Environments>

Further Linux System Tuning:

[https://www.ibm.com/developerworks/community/wikis/home?lang=en#!/wiki/Welcome%20to%20High%20Performance%20Computing%20\(HPC\)%20Central/page/Linux%20System%20Tuning%20Recommendations](https://www.ibm.com/developerworks/community/wikis/home?lang=en#!/wiki/Welcome%20to%20High%20Performance%20Computing%20(HPC)%20Central/page/Linux%20System%20Tuning%20Recommendations)



THANK YOU!