PRODUCT > MAPR CORE (/SUPPORT/S/TOPIC/0TO0L000000... (/SUPPORT/S

(/SUPPORT/S/TOPIC/0TO0L0000...

> <u>CLDB/MAPR-FS</u> <u>(/SUPPORT/S/TOPIC/0TO0L0000</u>...

(/support/s/login)

Notify Me

How to recover from disk failures caused by slow MapR-FS / MFS disks

In the faileddisk.log file, you will see information on the cause of failure. Failures can be due to CRC error, I/O error / time-out, tolerance conf limits settings for disk speed, disk lost form OS etc. Notice that the log file also provides instructions for removing disks and adding them back to MapR-FS. This article is compilation of relevant information for disk slowness that are distributed across maprdocs.mapr.com. This article is merely in the interest of making disk troubleshooting easier.

() Feb 1, 2018 · How do I

Task or use-case

When a disk fails, MapR raises the disk failure Alarm (NODE_ALARM_DISK_FAILURE) on the node where the disk fails and takes whole storage pool offline. This article is compilation of relevant information for disk slowness that are distributed across maprdocs.mapr.com. (http://maprdocs.mapr.com.) This article is merely in the interest of making disk troubleshooting easier.

Environment

· Any MapR core version

Solution

When a disk fails, MapR raises the disk failure Alarm (NODE_ALARM_DISK_FAILURE) on the node where the disk fails and takes whole storage pool offline.

1. Identify the failure reason

When a disk fails, you can check /opt/mapr/logs/faileddisk.log and check the reason of failure.

Disk : sde

Failure Reason : I/O error

Time of Failure : Sat Nov 11 03:51:42 PST 2017

Lost Volumes : mapr.node1.local.mapred

Resolution

Please refer to MapR's online documentation at http://mapr.com/doc/display/MapR/Disks (http://mapr.com/doc/display/MapR/Disks) on how to handle disk failures.

If you have further questions, please either post on answers.mapr.com (http://answers.mapr.com) or contact MapR technical support.

Disk : sdy

Failure Reason : Slow disk

Time of Failure : Tue Feb 28 14:36:39 GMT 2017

Lost Volumes : mapr.node1.local.mapred

Resolution:

Please refer to MapR's online documentation at http://mapr.com/doc/display/MapR/Disks (http://mapr.com/doc/display/MapR/Disks (http://mapr.com/doc/display/MapR/Disks) on how to handle disk failures.

If you have further questions, please either post on answers.mapr.com (http://answers.mapr.com) or contact MapR technical support.

Disk : sdz

Failure Reason : I/O time out

Time of Failure : Thu Feb 23 12:52:19 GMT 2017

Lost Volumes : mapr.node1.local.mapred

Resolution:

Please refer to MapR's online documentation at http://mapr.com/doc/display/MapR/Disks (http://mapr.com/doc/display/MapR/Disks (http://mapr.com/doc/display/MapR/Disks) on how to handle disk failures.

If you have further questions, please either post on answers.mapr.com (http://answers.mapr.com) or contact MapR technical support.

The above disk failures are are due to I/O errors (sde) and slow I/O of the disks (sdy and sdz). For the latter, the default criteria for slow MFS disk IO is 60 seconds. This is defined by the property "mfs.disk.io (http://mfs.disk.io).timeout" in /opt/mapr/conf/mfs.conf.

#mfs.num.compress.threads=1

#mfs.max.aio.events=5000

#mfs.disable.periodic.flush=0

#mfs.io.disk.timeout (http://mfs.io.disk.timeout)=60

There could be many reasons for the disk to be slow. For example:.

- 1. Raid controller is failing
- 2. Disk is slow or failing.
- 3. MFS disks are under heavy processing in the disks which possibly leads to IO timeout
- 4. If disks are exported from SAN environment and attached to each node, any kind of network or SAN issue can cause the I/O or slow disk

2. Recover the failed SP using fsck if the hardware is not failed

From the /opt/mapr/logs/mfs.log-3 logs you can check below message, which means that in order to bring this SP(SP2) online you need to run FSCK with repair option.

2017-11-13 13:24:34,5236 ERROR IOMgr spinit.cc (http://spinit.cc):296 SP SP2:/dev/sde online failed, it was previously marked with disk ERROR: I/O error, -5. To brir

Before running FSCK please make sure to check the following

- 1. Respective storage pool is offline
- 2. There should not be any volume unavailable Alarm.

To run fsck with the repair option the recommended syntax is:

\$ /opt/mapr/server/fsck -n <SP name> -r -m <allocate memory in MB>

In this case, it will be:

\$ /opt/mapr/server/fsck -n SP2 -r -m 10000.

Ex:

```
Using logfile /opt/mapr/logs/fsck.log.fsck.log.2017-11-20.18:07:45.2781
tcmalloc: large alloc 107374182400 bytes == 0x205a000 @ 0xa0c413 0xa28cb1 0x8cacb6
tcmalloc: large alloc 1572864000 bytes == 0x2290000 @ 0xa0c413 0xa28d3e 0x88884c
tcmalloc: large alloc 107374190592 bytes == 0x19080e2000 @ 0xa0c413 0xa28d3e 0x888ac4
FSCK Repair start (initialize storage pool and replay log) ...
2017-11-13 13:34:32,6272 Disk /dev/sde GUID 939915D9-0645-B4F5-7CA7-04AFE8E05900 hit IO Read error Input/output error -5 at block 942978 count 1
Allocator init: 14897g (1952645120 blocks) in 29795 groups
1: SG: f 97%: 0 [n 0 0%, r 0] --> 1619 [n 65536 100%, r 0]
FSCK phase 1 (initialize cache and verify log) ...
fs/server/fsck/phase1.cc (http://phase1.cc):39:FSINF Superblock is marked with error -5
fs/server/fsck/phase1.cc (http://phase1.cc):141:FSINF Init Log
FSCK phase 2 and 3 (verify all containers and inodes) ...
  working on container 2930 of 2973 ....
fs/server/fsck/phase2.cc (http://phase2.cc):2652:FSINF bad magic on cid 3931831 inode 894044 block 0x398de2 : exp 0x4d415049 found 0x0
fs/server/fsck/phase2.cc (http://phase2.cc):2501:FSINF bad crc or itype 0 on cid 3931831 inode 894044 block 0x398de2 : exp 0 found 0
fs/server/fsck/phase2.cc (http://phase2.cc):2652:FSINF bad magic on cid 3931831 inode 894045 block 0x398de2 : exp 0x4d415049 found 0x0
fs/server/fsck/phase2.cc (http://phase2.cc):2501:FSINF bad crc or itype 0 on cid 3931831 inode 894045 block 0x398de2 : exp 0 found 0
fs/server/fsck/phase2.cc (http://phase2.cc):2652:FSINF bad magic on cid 3931831 inode 894046 block 0x398de2 : exp 0x4d415049 found 0x0
fs/server/fsck/phase2.cc (http://phase2.cc):2501:FSINF bad crc or itype 0 on cid 3931831 inode 894046 block 0x398de2 : exp 0 found 0
fs/server/fsck/phase2.cc (http://phase2.cc):2652:FSINF bad magic on cid 3931831 inode 894047 block 0x398de2 : exp 0x4d415049 found 0x0
fs/server/fsck/<u>phase2.cc (http://phase2.cc)</u>:2501:FSINF bad crc or itype 0 on cid 3931831 inode 894047 block 0x398de2 : exp 0 found 0
fs/server/fsck/<u>btreewalk.cc (http://btreewalk.cc)</u>:1063:FSINF itype 2 stype 0
fs/server/fsck/btreewalk.cc (http://btreewalk.cc):1068:FSINF alias: 0x1cc681 block: 0x6187db6 idx: 0 shared: 0 visited: 1 verified: 1
fs/server/fsck/keyvisit.cc (http://keyvisit.cc):452:FSINF dangling dirent output.04761 Regular 1 child 4294967295.894030.47090922 in fid 3931831.854883.47071696
fs/server/fsck/btreewalk.cc (http://btreewalk.cc):1062:FSINF btree node error 53 on fid: 3931831.854883.47071696 nlevels: 2 level: 1
fs/server/fsck/<u>btreewalk.cc (http://btreewalk.cc)</u>:1063:FSINF itype 2 stype 0
fs/server/fsck/<u>btreewalk.cc (http://btreewalk.cc)</u>:1068:FSINF alias: 0x1cc670 block: 0x6187db7 idx: 0 shared: 0 visited: 1 verified: 1
fs/server/fsck/keyvisit.cc (http://keyvisit.cc):452:FSINF dangling dirent output.04748 Regular 1 child 4294967295.894041.47090870 in fid 3931831.854883.47071696
fs/server/fsck/keyvisit.cc (http://keyvisit.cc):452:FSINF dangling dirent output.04765 Regular 1 child 4294967295.894038.47090938 in fid 3931831.854883.47071696
fs/server/fsck/btreewalk.cc (http://btreewalk.cc):1062:FSINF btree node error 53 on fid: 3931831.854883.47071696 nlevels: 2 level: 1
fs/server/fsck/btreewalk.cc (http://btreewalk.cc):1063:FSINF itype 2 stype 0
fs/server/fsck/btreewalk.cc (http://btreewalk.cc):1068:FSINF alias: 0x1cc662 block: 0x6187db8 idx: 0 shared: 0 visited: 1 verified: 1
fs/server/fsck/<u>kevvisit.cc (http://kevvisit.cc)</u>:452:FSINF dangling dirent output.04751 Regular 1 child 4294967295.894047.47090882 in fid 3931831.854883.47071696
fs/server/fsck/keyvisit.cc (http://keyvisit.cc):452:FSINF dangling dirent output.04769 Regular 1 child 4294967295.894046.47090954 in fid 3931831.854883.47071696
fs/server/fsck/btreewalk.cc (http://btreewalk.cc):1062:FSINF btree node error 53 on fid: 3931831.854883.47071696 nlevels: 2 level: 1
fs/server/fsck/btreewalk.cc (http://btreewalk.cc):1063:FSINF itype 2 stype 0
fs/server/fsck/btreewalk.cc (http://btreewalk.cc):1068:FSINF alias: 0x1cc65c block: 0x6187dba idx: 0 shared: 0 visited: 1 verified: 1
fs/server/fsck/keyvisit.cc (http://keyvisit.cc):452:FSINF dangling dirent output.04750 Regular 1 child
```

```
fs/server/fsck/btreewalk.cc (http://btreewalk.cc):1063:FSINF itype 2 stype 0
fs/server/fsck/btreewalk.cc (http://btreewalk.cc):1068:FSINF alias: 0x1cc67c block: 0x6187dbd idx: 0 shared: 0 visited: 1 verified: 1
fs/server/fsck/btreewalk.cc (http://btreewalk.cc):1522:FSINF btree on fid: 39343774.88858.47077755 nkeys expected: 4900 found: 4871
fs/server/fsck/phase2.cc (http://phase2.cc):1623:FSINF container cid 398786 oblocks expected: 1583491 found: 1583462
fs/server/fsck/phase2.cc (http://phase2.cc):1636:FSINF container cid 39768766 lblocks expected: 1583491 found: 1583462

FSCK phase 4 (verify namespace and orphanage) ...

FSCK phase 5 (verify allocation bitmap) ...
fs/server/fsck/phase5.cc (http://phase5.cc):146:FSINF Mismatch in allocation bitmap(1551)

FSCK completed without errors.
```

It is recommended to always run fsck in the background so that it is not interrupted in the event the console session is terminated early.

3. Online the recovered SP

Once fsck is completed refresh the loaded SPs to bring the SP online.

```
$ /opt/mapr/server/mrconfig sp refresh
```

Verify the SP is online.

```
$ /opt/mapr/server/mrconfig sp list -v
```

It is possible that MFS will mark the disk(s) as failed again after some time after running fsck with repair option. This is typically because of persistent slow disks and if so check if the disks or other underlying hardware is still reliable. If everything is healthy increase the value of "mfs.io (http://mfs.io).disk.timeout parameter" to 120 or 240 in /opt/mapr/conf/mfs.conf

```
#mfs.num.compress.threads=1
#mfs.max.aio.events=5000
#mfs.disable.periodic.flush=0
mfs.io.disk.timeout (http://mfs.io.disk.timeout)=120
```

MFS must be restarted via a restart of warden after this change in /opt/mapr/conf/mfs.conf for the change in the disk IO timeout to take effect. This change will help to mitigate repeated disk failures due to IO timeout caused by slow disks,

If the disks are bad then replace the disks, please follow below steps to replace the disks.

- 1. Verify no 'data unavailable' alarms are raised for any critical volumes in the cluster.
- 2. Run the command "maprcli disk remove" to remove the slow disk(s) and the remaining disks in the storage pool.

```
# maprcli disk remove -disks /dev/sde -host `hostname -f`
message host disk
removed. node1 /dev/sde
removed. node1 /dev/sdy
removed. node1 /dev/sdz
```

2. Once old disk are removed and new disks are added from OS side, use the "maprcli disk add" command to add disks into MapR FS.

```
# maprcli disk add -disks /dev/sdd,/dev/sdh,/dev/sdi -host `hostname -f`
message host disk
added. node1 /dev/sdd
added. node1 /dev/sdh
added. node1 /dev/sdi
```

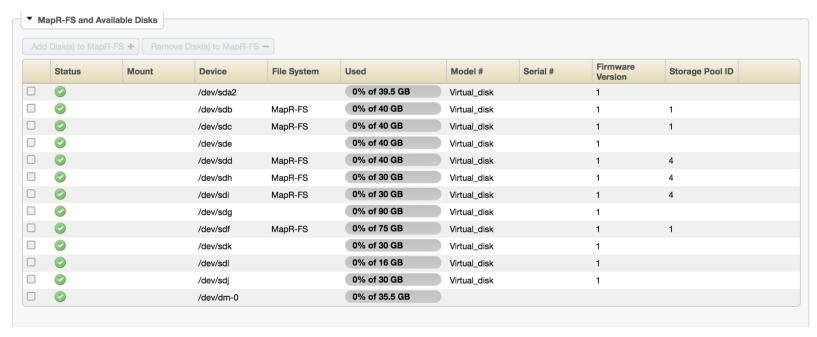
3. Once disks are added you can verify it by running this command "/opt/mapr/server/mrconfig sp list -v".

```
[root@node1 ~]# /opt/mapr/server/mrconfig sp list -v
ListSPs resp: status 0:2
No. of SPs (2), totalsize 244560 MB, totalfree 243665 MB

SP 0: name SP1, Online, size 147265 MB, free 146777 MB, path /dev/sdf, log 200 MB, port 5660, guid 30162ea8001ed3810059e0e8af0bad84, clusterUuid -8531477943840073658

SP 1: name SP4, Online, size 97294 MB, free 96887 MB, path /dev/sdd, log 200 MB, port 5660, guid 2f30bd91c30db9fc005a1374550d7e53, clusterUuid -8531477943840073658--
```

Alternatively, you can login to MCS and click on the respective node on the dashboard tab and check on MapR-FS and Available Disks.



Note: Sometime even after replacing faulty disk it gets marked as 'Failed' or 'dead' if there is a <diskname>failed.info (http://failed.info) file under /opt/mapr/logs/. Ex:

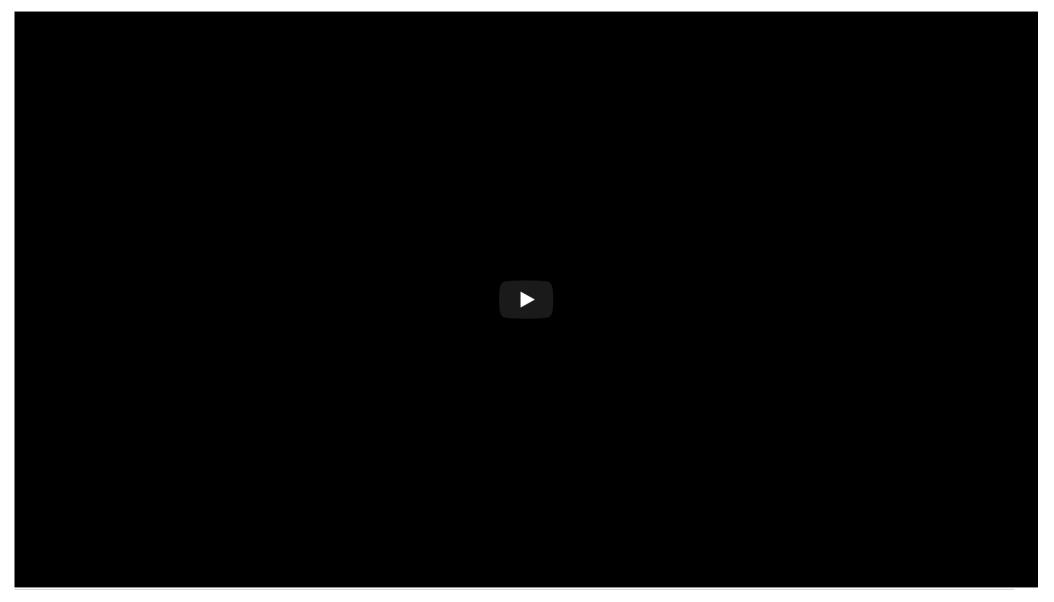
```
# ls -l | grep -i failed
-rw-rw-rw- 1 mapr mapr 515 Nov 11 05:51 faileddisk.log
-rw-rw-rw- 1 mapr mapr 0 Nov 11 05:51 sde.failed.info (http://sde.failed.info)
```

To fix this, please remove <diskname>.failed.info (http://failed.info) file and run /opt/mapr/server/disklist.sh (http://disklist.sh) on respective node.

Alternatives

General additional Info

Abizer Adenwala, Technical Support Engineer at MapR, walks you through what a storage pool is, why disks are striped, reasons disk would be marked as failed, what happens when a disk is marked failed, what to watch out for before reformatting/re-adding disk back, and what is the best path to recover from disk failure.



Links to Reference Docs

https://maprdocs.mapr.com/home/AdministratorGuide/ManagingDisks-SlowDisks.html (https://maprdocs.mapr.com/home/AdministratorGuide/ManagingDisks-SlowDisks.html)
https://maprdocs.mapr.com/home/AdministratorGuide/c-manage-disk-failures.html (https://maprdocs.mapr.com/home/AdministratorGuide/c-manage-disk-failures.html)
https://maprdocs.mapr.com/52/AdministratorGuide/c-recover-disk-failure.html (https://maprdocs.mapr.com/52/AdministratorGuide/c-recover-disk-failure.html)

Technology Group

MapR Core

Article Type

How_do_l__kav

Article Number

000002988

Article Total View Count

513

Last Published Date

2/1/2018 8:41 PM

MAPR USE-CASES

(/support/s/topic/0TO0L000004agcj...

MapR FS

(/support/s/topic/0TO0L000004ag...

CLDB/MapR-FS

(/support/s/topic/0TO0L000004agU...