

4.6. 麻烦的扩缩容

在 glusterfs 当中，如果 glusterfs 的 volume 是使用 heketi 创建的，那么就无法正常地使用 quota 进行容量限制了，因为底层使用 lvm2 的时候，如果要扩缩容，有两种办法，一种是对底层的 vg 和 lv 进行处理，但是这种方法风险比较高，一旦操作不慎，会直接影响原来的数据，而且底层的 vg 和 lv 容量改变，信息还需要同步到上层的应用中，会比较麻烦；第二种方法，就是使用 glusterfs 的 add-brick 功能，也就是增加 brick 的方式来进行扩缩容处理，下面来了解一下。

4.6.1. add-brick 操作

```
1. root@gfs01:~# gluster volume info test-event
2.
3. Volume Name: test-event
4. Type: Distribute
5. Volume ID: 4ff63ab9-561a-4c32-a4b9-60c5bda0c6e9
6. Status: Created
7. Snapshot Count: 0
8. Number of Bricks: 1
9. Transport-type: tcp
10. Bricks:
11. Brick1: 10.0.12.9:/glusterfs/test-event
12. Options Reconfigured:
13. nfs.disable: on
14. transport.address-family: inet
15. storage.fips-mode-rchecksum: on
16.
17.
18. root@gfs01:~# gluster volume add-brick test-event 10.0.12.
    2:/glusterfs/test-event force
19. volume add-brick: success
20.
21. root@gfs01:~# gluster volume info test-event
22.
```

```
23. Volume Name: test-event
24. Type: Distribute
25. Volume ID: 4ff63ab9-561a-4c32-a4b9-60c5bda0c6e9
26. Status: Created
27. Snapshot Count: 0
28. Number of Bricks: 2
29. Transport-type: tcp
30. Bricks:
31. Brick1: 10.0.12.9:/glusterfs/test-event
32. Brick2: 10.0.12.2:/glusterfs/test-event
33. Options Reconfigured:
34. nfs.disable: on
35. transport.address-family: inet
36. storage.fips-mode-rchecksum: on
```

这里有一个单 brick 的卷 test-event,然后使用了 add-brick 进行扩容,那么这里扩容之后要注意,数据是否要进行 rebalance,因为对于数据的分布如果不进行重平衡,那么可能后续会造成数据倾斜,也就是部分 brick 的数据很多,但是新添加的 brick 可能数据很少的现象,而关于 rebalance 的问题,后面还会继续了解。

4.6.2. remove-brick 操作

有了扩容的操作,那么一般自然也会有减少 brick 的,而这个操作可能在实际的生产环境中会使用频率比较少,但是这里也有一些值得注意的问题,下面来了解一下。

```
1. root@gfs01:~# mount -t glusterfs 10.0.12.9:test-event /mnt/test-event
2.
3. root@gfs01:~# ls /mnt/test-event/
4.
5. root@gfs01:~# cd /mnt/test-event/
6.
```

```

7. root@gfs01:~# cat /mnt/test-event/smallFile.sh
8. #!/bin/bash
9.
10. for((i=1;i<=10000;i++));
11. do
12.     touch $i.txt
13.     date > $i.txt
14.     #echo $(expr $i \* 3 + 1);
15. done
16.
17. root@gfs01:~# ls /mnt/test-event/ |wc -l
18. 10001

```

为了测试缩容的效果，这里弄了一万个小文件，然后再执行 remove-brick，在执行之前，这里再核实一下不同的 brick 的文件数量分布情况。

```

1. root@gfs03:~# ls /glusterfs/test-event/ |wc -l
2. 4989
3.
4. root@gfs03:~# hostname -i
5. 10.0.12.9
6.
7. root@gfs02:~# hostname -i
8. 10.0.12.2
9.
10. root@gfs02:~# ls /glusterfs/test-event/ |wc -l
11. 5012

```

这里可以看到，在该 volume 中，文件不是和复制卷一样的，这里也没有绝对平均地分配文件数量，那么下面进行 remove-brick 操作尝试。

```

1. root@gfs01:/mnt/test-event# gluster volume remove-brick test-event 10.0.12.9:/glusterfs/test-event force
2. Remove-brick force will not migrate files from the removed bricks, so they will no longer be available on the volume.
3. Do you want to continue? (y/n) y
4. volume remove-brick commit force: success
5.

```

```

6. root@gfs01:/mnt/test-event# gluster volume info test-event
7.
8. Volume Name: test-event
9. Type: Distribute
10. Volume ID: 4ff63ab9-561a-4c32-a4b9-60c5bda0c6e9
11. Status: Started
12. Snapshot Count: 0
13. Number of Bricks: 1
14. Transport-type: tcp
15. Bricks:
16. Brick1: 10.0.12.2:/glusterfs/test-event
17. Options Reconfigured:
18. performance.client-io-threads: on
19. nfs.disable: on
20. transport.address-family: inet
21. storage.fips-mode-rchecksum: on
22.
23. root@gfs01:/mnt/test-event# gluster volume rebalance test-event status
24. volume rebalance: test-event: failed: Volume test-event is not a distribute volume or contains only 1 brick.
25. Not performing rebalance

```

这里可以看到，因为只有一个 brick，那么是无法进行 rebalance 的，同时这里如果只剩下了最后一个 brick，会有数据丢失的风险的，因此要特别注意。

4.6.3. replace-brick 操作

这里的除了增加和删除之外，其实还有一个叫 replace-brick 的操作，这里在什么时候需要使用该操作呢？主要是在当前节点有问题的时候，想转移这个 brick 的数据到另外一个节点上，那么下面可以了解一下。

```

1. root@gfs01:~# gluster volume info test-replica
2.
3. Volume Name: test-replica
4. Type: Replicate
5. Volume ID: d2614e89-9aba-46f6-bf04-984782ac6d6f

```

```

6. Status: Started
7. Snapshot Count: 0
8. Number of Bricks: 1 x 3 = 3
9. Transport-type: tcp
10. Bricks:
11. Brick1: 10.0.12.2:/glusterfs/test-replica
12. Brick2: 10.0.12.9:/glusterfs/test-replica
13. Brick3: 10.0.12.12:/glusterfs/test-replica
14. Options Reconfigured:
15. features.quota-deem-statfs: on
16. features.inode-quota: on
17. features.quota: on
18. cluster.granular-entry-heal: on
19. storage.fips-mode-rchecksum: on
20. transport.address-family: inet
21. nfs.disable: on
22. performance.client-io-threads: off
23.
24.
25.
26. root@gfs01:~# gluster volume replace-brick test-replica 1
    0.0.12.2:/glusterfs/test-replica 10.0.12.2:/glusterfs/test
    -replica-new commit force
27. volume replace-brick: success: replace-brick commit force
    operation successful
28.
29.
30. root@gfs01:~# gluster volume info test-replica
31.
32. Volume Name: test-replica
33. Type: Replicate
34. Volume ID: d2614e89-9aba-46f6-bf04-984782ac6d6f
35. Status: Started
36. Snapshot Count: 0
37. Number of Bricks: 1 x 3 = 3
38. Transport-type: tcp
39. Bricks:
40. Brick1: 10.0.12.2:/glusterfs/test-replica-new
41. Brick2: 10.0.12.9:/glusterfs/test-replica
42. Brick3: 10.0.12.12:/glusterfs/test-replica
43. Options Reconfigured:
44. features.quota-deem-statfs: on
45. features.inode-quota: on
46. features.quota: on

```

```
47. cluster.granular-entry-heal: on
48. storage.fips-mode-rchecksum: on
49. transport.address-family: inet
50. nfs.disable: on
51. performance.client-io-threads: off
```

那么注意，这里因为使用 replace-brick 之后，因为不是分布式卷，因此是无法进行 rebalance 的，下面可以看到内容。

```
1. root@gfs01:~# gluster volume rebalance test-replica status
2. volume rebalance: test-replica: failed: Volume test-replica is not a distribute volume or contains only 1 brick.
3. Not performing rebalance
```

那么这里还可以留意一下是否数据已经转移了。

```
1. root@gfs02:~# ls -l /glusterfs/test-replica
2. total 1822728
3. -rw-r--r-- 2 root root 1866465280 Jun 23 20:51 CentOS-8.3.2011-x86_64-minimal.iso
4.
5. root@gfs02:~# ls -l /glusterfs/test-replica-new/
6. total 1822728
7. -rw-r--r-- 2 root root 1866465280 Jun 23 20:51 CentOS-8.3.2011-x86_64-minimal.iso
8.
9. root@gfs02:~# md5sum /glusterfs/test-replica/CentOS-8.3.2011-x86_64-minimal.iso
10. 8934d42a86d8589342ac9bdfec82d6b4 /glusterfs/test-replica/CentOS-8.3.2011-x86_64-minimal.iso
11.
12. root@gfs02:~# md5sum /glusterfs/test-replica-new/CentOS-8.3.2011-x86_64-minimal.iso
13. 8934d42a86d8589342ac9bdfec82d6b4 /glusterfs/test-replica-new/CentOS-8.3.2011-x86_64-minimal.iso
14.
15. root@gfs02:~# hostname -i
16. 10.0.12.2
```

4.6.4. rebalance 很重要

那么提到了扩扩容的问题，就不得不说一下 rebalance 重平衡了，所谓的重平衡，这里就是对数据的分布进行重新的均衡，这样的话可以保证数据不会出现明显的数据倾斜情况，下面可以进行测试一下。

```
1. root@gfs01:~# gluster volume create rebalance-test replica 3
   10.0.12.{2,9,12}:/glusterfs/rebalance-01 force
2. volume create: rebalance-test: success: please start the volume
   to access data
3.
4. root@gfs01:~# gluster volume start rebalance-test
5. volume start: rebalance-test: success
6.
7. root@gfs01:~# mkdir -p /mnt/rebalance-test
8.
9. root@gfs01:~# mount -t glusterfs 10.0.12.2:rebalance-test /mnt/
   rebalance-test
10.
11.root@gfs01:~# cp /mnt/test-rebalance/test.sh /mnt/rebalance-test/
12.
13.root@gfs01:~# cd /mnt/rebalance-test/
14.
15.root@gfs01:/mnt/rebalance-test# ls
16.test.sh
17.
18.root@gfs01:/mnt/rebalance-test# cat test.sh
19.#!/bin/bash
20.for ((i=1; i<=10000; i++))
21.do
22.  sudo touch $i.txt
23.  sudo date > $i.txt
24.done
25.
26.root@gfs01:/mnt/rebalance-test# bash test.sh
27.
28.root@gfs01:/mnt/rebalance-test# ls -l |wc -l
29.10002
```

这里首先创建一个 3 副本的复制卷，然后使用脚本创建了 1 万个小文件，接着下面就是使用 add-brick 进行增加 brick，然后进行 rebalance 操作。

```
1. root@gfs01:/mnt/rebalance-test# gluster volume add-brick re
   balance-test 10.0.12.{2,9,12}:/glusterfs/rebalance-02 force
2. volume add-brick: success
3.
4. root@gfs01:/mnt/rebalance-test# gluster volume rebalance re
   balance-test start
5. volume rebalance: rebalance-test: success: Rebalance on reba
   lance-test has been started successfully. Use rebalance status
   command to check status of the rebalance process.
6. ID: 39101691-e7d7-4e51-847a-4fc2845eea3a
7.
8. root@gfs01:/mnt/rebalance-test# gluster volume rebalance re
   balance-test status
9.  Node Rebalanced-files size  scanned failures skipped  status
   run time in h:m:s
10. -----
11. 10.0.12.9 21 672Bytes 506 0 0 in progress 0:00:0
   4
12. gfs02 22 704Bytes 504 0 0 in progress 0:00:04
13. localhost 8 576Bytes 503 0 0 in progress 0:00:04
14. The estimated time for rebalance to complete will be unavaila
   ble for the first 10 minutes.
15. volume rebalance: rebalance-test: success
```

从这里可以看到,重平衡已经启动了，那么这里数据最后的分布会变成怎样呢？

```
1. root@gfs02:~# ls -l /glusterfs/rebalance-01/ | wc -l
2. 4990
3. root@gfs02:~# ls -l /glusterfs/rebalance-02/ | wc -l
4. 5013
```

另外对于 volume 的 rebalance 的操作，这里是有两种情况的，一种是像上面的把数据全部重新平衡，包括元数据，还有一种是单纯平衡元数据的，这里

可以根据业务场景的不同进行选择。

对于 rebalance 这里，还有一点需要注意的是，force 选项并不是强制的，加上 force 的话会考虑 brick 的数据分布情况，也就是说，如果不加上 force 的话，那么 add-brick 之后有可能数据并不会一定进行迁移和平衡的，这一点可以查看官方 issue 编号为 2571 的内容，本人在做测试的时候曾经遇到过。

