

3. 第三章 glusterfs 的核心概念

从本章开始，有了前面的关于 linux 的一些文件系统的简单认识，那么从这里开始真正去认识一下 **glusterfs** 这种无中心架构的特点，因为无中心架构，那么必然在元数据存储与节点之间通信，数据恢复等方面会与常见的架构有着非常大的不同，而这也是 **glusterfs** 的魅力之处，下面开始将一点点地深入去理解 **glusterfs** 是如何运作的。另外这里先做个简单的约定，下面的代码内容如无特殊说明，则基本是在虚拟机环境中进行的实验操作，对于集群节点则是以 **gfs01,gfs02** 这样来命令规范的，而客户端机器则是 **gfsclient01** 这样的，因此阅读代码的时候，可以进行区别开来。

关于实验的虚拟机的 ip 与 hostname 对应关系如下所示..

```
1. # cat /etc/hosts
2. ...
3. 192.168.0.110 gfs01
4. 192.168.0.111 gfs02
5. 192.168.0.112 gfs03
```

3.1. 有趣的扩展属性 gfid

3.1.1. 文件的 gfid 属性

因为 **glusterfs** 是无中心架构的，这点在前面的内容中也多次提到过，当然可能会有困惑，那么到底这个文件是存放到哪里的，由什么来决定呐？对于 **glusterfs** 来说，这里巧妙地利用了 linux 文件的扩展属性，而一个文件除了有基本的元数据以外，还支持扩展属性的，而其中有一个叫做 **gfid** 的扩展属性是关键的作用，下面来感受一下。

```
1. # gluster volume info test-replica
2.
3. Volume Name: test-replica
4. Type: Replicate
5. Volume ID: 4568c063-5b75-4304-98f6-21f3955cc138
6. Status: Started
7. Snapshot Count: 0
8. Number of Bricks: 1 x 3 = 3
9. Transport-type: tcp
10. Bricks:
11. Brick1: 192.168.0.110:/glusterfs/test-replica
12. Brick2: 192.168.0.111:/glusterfs/test-replica
13. Brick3: 192.168.0.112:/glusterfs/test-replica
14. Options Reconfigured:
```

15. cluster.granular-entry-heal: on
16. storage.fips-mode-rchecksum: on
17. transport.address-family: inet
18. nfs.disable: on
19. performance.client-io-threads: off

首先这里有一个 3 副本的复制卷，然后在一个客户端虚拟机上面进行挂载，接着查看一下扩展属性。

1. [root@gfsclient01 ~]# mount -t glusterfs -o aux-gfid-mount 192.168.0.110:test-replica /mnt/test-replica
2. [root@gfsclient01 ~]# cd /mnt/test-replica/
3. [root@gfsclient01 test-replica]# getfattr -n glusterfs.gfid.string a.txt
4. # file: a.txt
5. glusterfs.gfid.string="71b9fb49-53ae-42f8-bb1b-b53af17521fc"

从这里可以看到 **glusterfs** 的文件是有一个叫做 **gfid** 的扩展属性的，那么这里还可以到 **brick** 所在的目录下面进行查看一下信息。

1. [root@gfs01 test-replica]# getfattr -d -m . -e hex a.txt
2. # file: a.txt
3. trusted.gfid=0x71b9fb4953ae42f8bb1bb53af17521fc
4. trusted.gfid2path.640ca896bde02acc=0x30303030303030302d303030302d303030302d303030302d303030303030303030303030312f612e747874
5. trusted.glusterfs.mdata=0x0100000000000000000000000060ba3cdf000000001ca8c3c90000000060ba3cdf000000001ca8c3c90000000060ba3cdf000000001ca8c3c9
- 6.
7. [root@gfs02 ~]# cd /glusterfs/test-replica
8. [root@gfs02 test-replica]# getfattr -d -m . -e hex a.txt
9. # file: a.txt
10. trusted.gfid=0x71b9fb4953ae42f8bb1bb53af17521fc
11. trusted.gfid2path.640ca896bde02acc=0x30303030303030302d303030302d303030302d303030302d303030302d30303030303030303030303030303030312f612e747874
12. trusted.glusterfs.mdata=0x0100000000000000000000000060ba3cdf000000001ca8c3c90000000060ba3cdf000000001ca8c3c90000000060ba3cdf000000001ca8c3c9

那么这里可以看到 **gfid** 是相同的，在复制卷里面，因为都是相同的文件，因此这里的 **gfid** 显示一致，那么这个 **gfid** 是不是 **inode** 呢？其实不是的，**gfid** 是客户端计算出来的一个独一无二的类似 **uuid** 这样的字符串，另外下面还可以查看 **inode** 和其他信息。

1. [root@gfsclient01 test-replica]# ls -i a.txt
2. 13482569174227427836 a.txt
3. [root@gfsclient01 test-replica]# ls -l a.txt
4. -rw-r--r--. 1 root root 0 Jun 4 10:46 a.txt

