

2.1.2. 有趣的 VFS 和文件系统

VFS 在 linux 的出现 ,不得不惊叹是一个神来之笔 ,下面先来看一下内核代码中的 ext4 文件系统与 VFS 关系图。

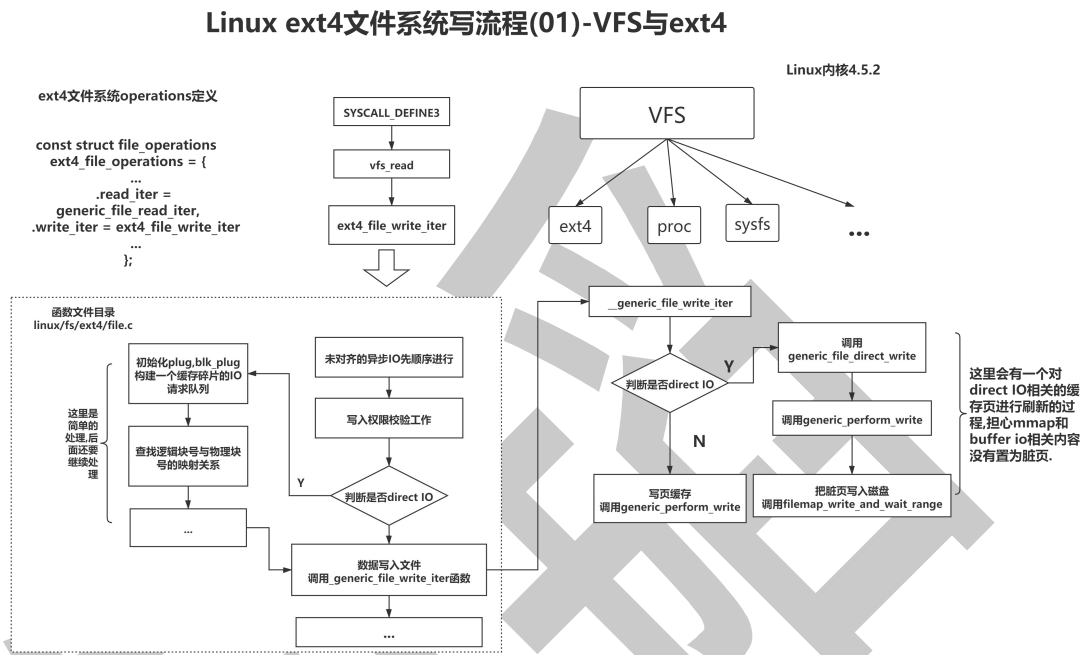


图 2.1.1-1 vfs 与 ext4 代码关系图

首先图中的一个读请求,是经过了 `vfs_read` 然后调用了 `ext4_file_write_iter`, 而后者则是 `ext4` 文件系统注册在操作系统中的 FOP,这里在调用的时候,会根据文件系统注册的函数进行回调,然后该函数就是真正地进入了文件系统的操作逻辑了,这里也就是有了两个分支,就是前面提到的是否使用操作系统内核缓存的问题了。下面为了让大家更好地理解一下 VFS 的元素,将具体介绍一下内容。

2.1.2.1. VFS inode

为了更好地理解一下 inode 信息,下面给出一个普通文件的 `stat` 命令结果,该命令可

以列出文件的一些属性信息。

```
1. $ sudo ls -li test.sh
2. 1328 test.sh
3.
4.
5. $ sudo stat test.sh
6. File: test.sh
7. Size: 614      Blocks: 19      IO Block: 1024   regular file
8. Device: 74h/116d Inode: 1328      Links: 1
9. Access: (0644/-rw-r--r--)  Uid: (  0/   root)   Gid: (  0/   root)
10. Access: 2020-11-07 03:21:30.777590296 +0000
11. Modify: 2020-11-07 03:21:26.485583231 +0000
12. Change: 2020-11-07 03:21:26.493583244 +0000
13. Birth: -
```

这里 inode 号为 1328,这里还要留意一下 change,modify 和 access 三个时间,这里也就是 ctime,mtime 和 atime,在 glusterfs 中的文件,也是利用了这三个信息进行一些文件更新校验的工作。另外这里的 size,也是在 glusterfs 中使用的一个很重要的信息。

change: 也就是 ctime,最后一次改变文件或目录(属性)的时间,例如执行 chmod, chown 等命令。

modity: 也就是 mtime,是最后一次修改文件或目录(内容)的时间。

access: 也就是 atime,是最后一次访问文件或目录的时间。

```
1. $ sudo stat dir
2. File: dir
3. Size: 2      Blocks: 1      IO Block: 131072 directory
4. Device: 74h/116d Inode: 131881    Links: 2
5. Access: (0755/drwxr-xr-x)  Uid: (  0/   root)   Gid: (  0/   root)
6. Access: 2021-06-03 03:21:17.577886842 +0000
7. Modify: 2021-06-03 03:21:17.577886842 +0000
8. Change: 2021-06-03 03:21:17.577886842 +0000
9. Birth: -
```

与文件的不同,目录的 stat 信息中的 Links 注意是 2 的,其他信息则相似。当然 inode 中的信息远远不只这么少,还包括了其他很多的信息,例如包括了用户和所属组的权限,关

联的 superblock 等信息。具体想了解的可以看看 inode 的数据结构，该结构体定义在 <linux/fs.h> 中。

```
1.  struct inode
2.  {
3.      /* 哈希表 */
4.      struct hlist_node    i_hash;
5.      /* 索引节点链表 */
6.      struct list_head     i_list;
7.      /* 目录项链表 */
8.      struct list_head     i_dentry;
9.      ...
10.     uid_t                 i_uid;
11.     gid_t                 i_gid;
12.     ...
13.     struct timespec       i_atime;
14.     struct timespec       i_mtime;
15.     struct timespec       i_ctime;
16.     ...
17.     /* 索引节点操作表 */
18.     struct inode_operations *i_op;
19.     /* 默认的索引节点操作 */
20.     struct file_operations *i_fop;
21.     /* 相关的超级块 */
22.     struct super_block     *i_sb;
23.     /* 文件锁链表 */
24.     struct file_lock       *i_flock;
25.     /* 相关的地址映射 */
26.     struct address_space   *i_mapping;
27.     ...
28. }
```

2.1.2.2. VFS superblock

对于超级块来说,可以使用 df -i 命令来简单查看一下操作系统层面的 inode 使用情况。

```
1.  $ sudo df -i
2.  Filesystem      Inodes   IUsed   IFree IUse% Mounted on
3.  udev             2023350    620  2022730    1% /dev
```

4.	tmpfs	2029590	1379	2028211	1% /run
5.	/dev/nvme0n1p2	30498816	1481554	29017262	5% /
6.	tmpfs	2029590	3222	2026368	1% /dev/shm
7.	tmpfs	2029590	7	2029583	1% /run/lock
8.	tmpfs	2029590	18	2029572	1% /sys/fs/cgroup

这里使用 `df -i` 命令可以直观地看到目前系统中使用的 inode 情况 ,同时当 inode 耗尽的时候 ,也是一个非常要注意的问题。这些在日常的操作系统维护中都是需要监控的内容。

另外下面给出 superblock 的对象结构。

```

1.  struct super_block {
2.      ...
3.      // 关联设备
4.      dev_t    s_dev;
5.      // 块大小
6.      unsigned long s_blocksize;
7.      ...
8.      // 文件系统类型,常见的有 ext,xfs 等
9.      struct file_system_type *s_type;
10.     // 超级块的操作函数
11.     const struct super_operations *s_op;
12.     ...
13.     // 该 super_block 中所有 inode 的 i_sb_list 成员的双向链表
14.     struct list_head s_inodes;
15.     ...
16.
17. }
```

2.1.2.3. VFS struct file

`struct file` 结构体定义在 `include/linux/fs.h` 中 ,其中每一个文件结构体表示被进程打开的一个文件 ,而系统中每个被打开的文件 ,都可以在内核空间中找到关联的一个 `struct file` 对象。这个对象由内核在打开文件时创建 ,并且传递给在文件上进程操作的函数。在文件的所有实例都关闭后 ,内核释放这个数据结构。下面先简单感受一下在操作系统中 ,文件打开后的变化。

```

1. [root@gfsclient01 ~]# lsof |grep test
2.  COMMAND    PID  TID   USER  FD      TYPE          DEVICE  SIZE/OFF      NODE NAME
3.  ....
4.  vim          1827      root   4u     REG          253,0      12288  100663370 /root/.test.sh.
    swp
5.  /root/.test.sh.swp
6. [root@gfsclient01 ~]# ps -ef |grep 1677
7. root        1827  1529   0 10:09 pts/0    00:00:00 vim test.sh

```

在 linux 中，打开了一个文件之后，其实并不是真正意义上直接对原文件进行修改的，而是会把原文件复制一份加上.swp 之后，对应的所有的操作都会在该临时的缓存文件中，而一旦修改完成进行合并的时候，才会替换到原文件。当然这里会有常见的文件冲突，就是当文件被打开之后，如果再次打开该文件，操作系统会检测到有临时缓存文件存在，那么就证明可能是有修改没完成，或者之前的修改中断了没有合并，这样会导致冲突，但是可以进行替换操作等恢复。

下面给出一种使用文件句柄来恢复数据的方式。当该文件有被进程引用 fd，那么就有机会恢复原来的数据。

```

1. [root@gfsclient01 ~]# lsof |grep test_file.txt
2.  tail        2165      root   3r     REG          253,0        10  100730896 /root/test_file.
    txt
3. [root@gfsclient01 ~]# ls -l /proc/2165/fd/3
4. lr-x-----. 1 root root 64 Jun  3 10:20 /proc/2165/fd/3 -> /root/test_file.txt
5. [root@gfsclient01 ~]# cat /proc/2165/fd/3
6. #123
7.
8. 321
9. [root@gfsclient01 ~]# rm -fr test_file.txt
10. [root@gfsclient01 ~]# cat /proc/2165/fd/3 > abc.txt
11. [root@gfsclient01 ~]# cat abc.txt
12. #123
13.
14. 321

```

这里首先创建了一个名字为 test_file.txt 的文件，然后在一个终端使用 tail -f 来阅读该文件，这样可以保证文件一直在被打开的状态，同时 fd 中也可以看到是 read 状态。接着通过 lsof 找到该文件对应的引用进程，在 proc 下可以找到该文件的缓存内容，这样就可以进行数据恢复了。

同样的，为了进一步了解 struct file 的数据结构，该结构定义在内核代码中的 include/linux/fs.h 中，下面给出部分代码内容。

```
1.  struct file {
2.      ...
3.      struct path      f_path;
4.      struct inode      *f_inode;
5.      //和文件关联的操作.
6.      const struct file_operations  *f_op;
7.      spinlock_t      f_lock;
8.      atomic_long_t      f_count;
9.      //文件标志,有 O_RDONLY 和 O_SYNC 等.
10.     unsigned int      f_flags;
11.     //文件模式
12.     fmode_t      f_mode;
13.     struct mutex      f_pos_lock;
14.     //当前读写位置
15.     loff_t      f_pos;
16.     struct fown_struct f_owner;
17.     const struct cred      *f_cred;
18.     struct file_ra_state  f_ra;
19.     ...
20. }
```

2.1.2.4. VFS dentry

dentry，也就是目录项，是多个文件或者目录的链接，通过这个链接可以找寻到目录之下的文件或者是目录项。dentry 在文件系统里是极其重要的一个概念，dentry 结构体

在 linux 内核里也是用处广泛，这个结构体定义在 include/linux/dcache.h 里。

```
1.  struct dentry {
2.      //目录项标志
3.      unsigned int d_flags;
4.      seqcount_t d_seq;
5.      //散列表表项的指针
6.      struct hlist_bl_node d_hash;
7.      //父目录的目录项
8.      struct dentry *d_parent;
9.      struct qstr d_name;
10.     // 与文件名关联的索引节点
11.     struct inode *d_inode;
12.     unsigned char d_iname[DNAME_INLINE_LEN];
13.
14.
15.     struct lockref d_lockref;
16.     //dentry 的操作函数
17.     const struct dentry_operations *d_op;
18.     //文件的超级块对象
19.     struct super_block *d_sb;
20.     unsigned long d_time;
21.     void *d_fsdata;
22.
23.     struct list_head d_lru;
24.     //父列表的子级
25.     struct list_head d_child;
26.     struct list_head d_subdirs;
27.
28.     union {
29.         struct hlist_node d_alias;
30.         struct rcu_head d_rcu;
31.     } d_u;
32. };
```

为了观察操作系统中目录的关系结构，可以使用 tree 命令进行查看。

```
1.  [root@gfsclient01 ~]# tree /tmp/
2.  /tmp/
3.  ├── ks-script-YbuVG1
4.  └── systemd-private-4556ac58b212443eb846ec7ab9a7f059-chronyd.service-soC2Da
```

```

5. |   └─ tmp
6. └─ yum.log
7.   └─ yum_save_tx.2021-06-03.09-45.cn30by.yumtx
8.
9. 2 directories, 3 files

```

2.1.2.5. 文件对象

有了前面的四个核心基本元素，那么如果要打开一个文件，就必须要使用文件对象了，它的作用是描述进程和文件交互的关系，说白了就是要知道哪个进程在打开这个文件，并且打开读取的文件位置等信息。文件的结构定义如下所示。

```

1. struct file {
2.     //指向文件对应的 dentry 结构
3.     struct dentry                *f_dentry;
4.     //指向文件所属于的文件系统的 vfsmount 对象,和挂载 有关
5.     struct vfsmount              *f_vfsmnt;
6.     //指向文件的 file operation
7.     const struct file_operations *f_op;
8.     ...
9.     //记录进程对文件操作的位置。对文件读取前 n 字节，那么其指向 n+1 字节位置
10.    loff_t                        f_pos;
11.
12.    struct fown_struct             f_owner;
13.    //表示文件的用户的 uid 和 gid
14.    unsigned int                   f_uid;
15.    unsigned int                   f_gid;
16.    //用于记录文件预读的设置
17.    struct file_ra_state           f_ra;
18.
19.    //这个结构是指向一个封装文件的读写缓存页面的结构
20.    struct address_space           *f_mapping;
21.
22. }

```