

Linux VFS虚拟文件系统初探

📅 2019-12-02 | 📁 技术

“关于Linux内核的VFS文件系统的一些基础”

记录一下关于VFS文件系统的一些基础学习，和一些理解，如有错误望指正！

参考资料：

《深入理解Linux内核》（第三版）

Linux 文件系统剖析

VFS中的file、dentry和inode

源码：Linux-4.19.65

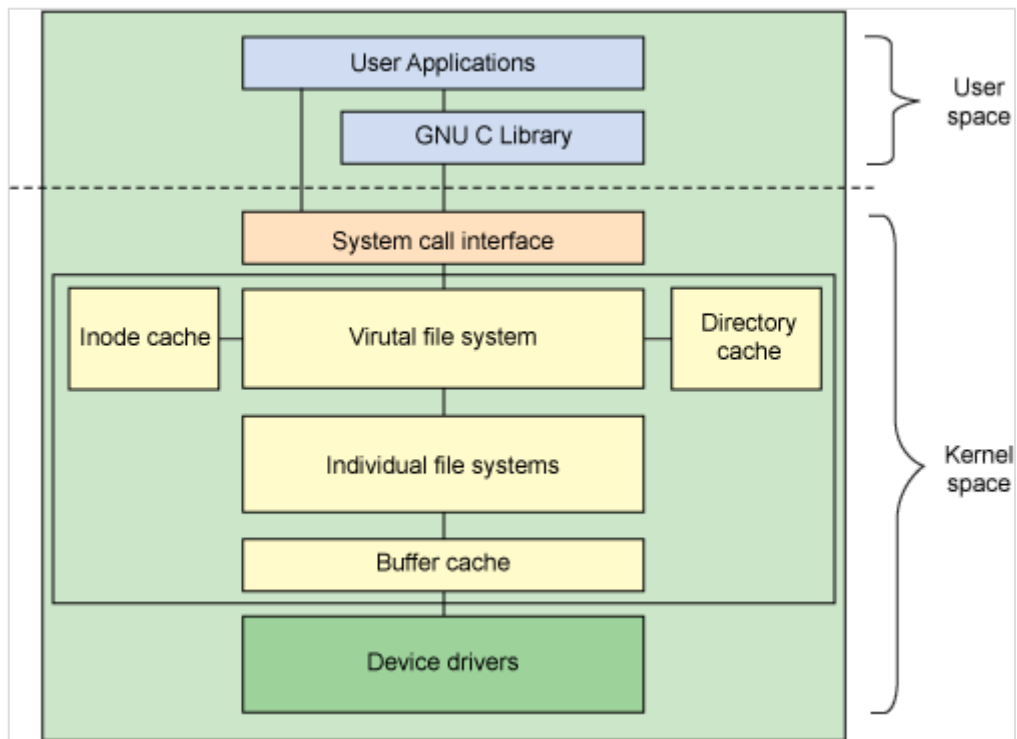
一些概念、数据结构

VFS

虚拟文件系统（Virtual Filesystem），是一个内核软件层，用来处理与Unix标准文件系统相关的所有系统调用。其健壮性表现在能为各种文件系统提供一个通用的接口。VFS是用户的应用程序与文件系统实现之间的抽象层，应用程序不需要知道操作的文件是什么文件系统类型，直接与VFS交互，VFS再与不同的文件系统交互。

VFS引入了一个通用的文件模型（common file model），这个模型能够表示所有支持的文件系统。而要实现每个具体的文件系统，必须将其物理组织结构转换为虚拟文件系统的通用文件模型。

下图所示的体系结构显示了用户空间和内核中与文件系统相关的组件之间的关系：



file_system_type

可以使用一组注册函数在 Linux 中动态地添加或删除文件系统。内核保存当前支持的文件系统的列表，可以通过 `/proc` 文件系统在用户空间中查看这个列表。这个虚拟文件还显示当前与这些文件系统相关联的设备。在 Linux 中添加新文件系统的方法是调用 `register_filesystem()`。这个函数的参数定义一个文件系统结构（`file_system_type`）的引用，这个结构定义文件系统的名称、一组属性和两个超级块函数。也可以注销文件系统。

在注册新的文件系统时，会把这个文件和它的相关信息添加到 `file_systems` 列表中（见图 2 和 `linux/include/linux/mount.h`）。这个列表定义可以支持的文件系统。在命令行上输入 `cat /proc/filesystems`，就可以查看这个列表。

from <https://www.ibm.com/developerworks/cn/linux/l-linux-filesystem/index.html>

源码中对 `file_system_type` 结构的定义如下：

```
1  /* include/linux/fs.h */
2
3  struct file_system_type {
4      const char *name;
5      int fs_flags;
6  #define FS_REQUIRES_DEV      1
7  #define FS_BINARY_MOUNTDATA  2
8  #define FS_HAS_SUBTYPE      4
9  #define FS_USERNS_MOUNT      8      /* Can be mounted by userns root */
```

```

10  #define FS_RENAME_DOES_D_MOVE 32768 /* FS will handle d_move() during rename()
11      struct dentry *(*mount) (struct file_system_type *, int,
12          const char *, void *);
13      void (*kill_sb) (struct super_block *);
14      struct module *owner;
15      struct file_system_type * next; //下一个file_system_type, 可以看到所有文件系
16      struct hlist_head fs_supers;
17
18      struct lock_class_key s_lock_key;
19      struct lock_class_key s_umount_key;
20      struct lock_class_key s_vfs_rename_key;
21      struct lock_class_key s_writers_key[SB_FREEZE_LEVELS];
22
23      struct lock_class_key i_lock_key;
24      struct lock_class_key i_mutex_key;
25      struct lock_class_key i_mutex_dir_key;
26  };

```

↑ 0%

vfsmount

这个结构提供当前挂装的文件系统。

```

1  /* /include/linux/mount.h */
2  struct vfsmount {
3      struct dentry *mnt_root; /* root of the mounted tree */
4      struct super_block *mnt_sb; /* pointer to superblock */
5      int mnt_flags;
6  } __randomize_layout;

```

super_block 超级块

存放已安装文件系统的有关信息。它包含管理文件系统所需的信息，包括文件系统名称（比如 ext2）、文件系统的大小和状态、块设备的引用和元数据信息（比如空闲列表等等）。超级块通常存储在存储媒体上，但是如果超级块不存在，也可以实时创建它。可以在/include/linux/fs.h 中找到超级块结构。

```

1  /* /include/linux/fs.h */
2  struct super_block {
3      struct list_head s_list; /* Keep this first */
4      dev_t s_dev; /* search index; _not_ kdev_t */
5      unsigned char s_blocksize_bits;
6      unsigned long s_blocksize;
7      loff_t s_maxbytes; /* Max file size */
8      struct file_system_type *s_type;

```

↑ 0%

```
9      const struct super_operations  *s_op;
10     const struct dquot_operations  *dq_op;
11     const struct quotactl_ops      *s_qcop;
12     const struct export_operations *s_export_op;
13     unsigned long                   s_flags;
14     unsigned long                   s_iflags;      /* internal SB_I_* flags */
15     unsigned long                   s_magic;
16     struct dentry                   *s_root;
17     struct rw_semaphore             s_umount;
18     int                             s_count;
19     atomic_t                        s_active;
20 #ifdef CONFIG_SECURITY
21     void                            *s_security;
22 #endif
23     const struct xattr_handler **s_xattr;
24 #if IS_ENABLED(CONFIG_FS_ENCRYPTION)
25     const struct fscrypt_operations *s_cop;
26 #endif
27     struct hlist_bl_head            s_roots;        /* alternate root dentries for NFS
28     struct list_head                s_mounts;        /* list of mounts; _not_ for fs use
29     struct block_device             *s_bdev;
30     struct backing_dev_info         *s_bdi;
31     struct mtd_info                 *s_mtd;
32     struct hlist_node               s_instances;
33     unsigned int                    s_quota_types;   /* Bitmask of supported quota type
34     struct quota_info               s_dquot;        /* Diskquota specific options */
35
36     struct sb_writers               s_writers;
37
38     char                            s_id[32];        /* Informational name */
39     uuid_t                          s_uuid;          /* UUID */
40
41     void                            *s_fs_info;      /* Filesystem private info */
42     unsigned int                    s_max_links;
43     fmode_t                         s_mode;
44
45     /* Granularity of c/m/atime in ns.
46        Cannot be worse than a second */
47     u32                             s_time_gran;
48
49     /*
50      * The next field is for VFS *only*. No filesystems have any business
51      * even looking at it. You had been warned.
52      */
53     struct mutex s_vfs_rename_mutex;                /* Kludge */
54
55     /*
```

```

56     * Filesystem subtype. If non-empty the filesystem type field
57     * in /proc/mounts will be "type.subtype"
58     */
59     char *s_subtype;
60
61     const struct dentry_operations *s_d_op; /* default d_op for dentries */
62
63     /*
64      * Saved pool identifier for cleancache (-1 means none)
65      */
66     int cleancache_poolid;
67
68     struct shrinker s_shrink;          /* per-sb shrinker handle */
69
70     /* Number of inodes with nlink == 0 but still referenced */
71     atomic_long_t s_remove_count;
72
73     /* Pending fsnotify inode refs */
74     atomic_long_t s_fsnotify_inode_refs;
75
76     /* Being remounted read-only */
77     int s_readonly_remount;
78
79     /* AIO completions deferred from interrupt context */
80     struct workqueue_struct *s_dio_done_wq;
81     struct hlist_head s_pins;
82
83     /*
84      * Owning user namespace and default context in which to
85      * interpret filesystem uids, gids, quotas, device nodes,
86      * xattrs and security labels.
87      */
88     struct user_namespace *s_user_ns;
89
90     /*
91      * Keep the lru lists last in the structure so they always sit on their
92      * own individual cachelines.
93      */
94     struct list_lru          s_dentry_lru ____cacheline_aligned_in_smp;
95     struct list_lru          s_inode_lru ____cacheline_aligned_in_smp;
96     struct rcu_head          rcu;
97     struct work_struct        destroy_work;
98
99     struct mutex              s_sync_lock;    /* sync serialisation lock */
100
101     /*
102      * Indicates how deep in a filesystem stack this SB is

```

```

103         */
104         int s_stack_depth;
105
106         /* s_inode_list_lock protects s_inodes */
107         spinlock_t s_inode_list_lock ____cacheline_aligned_in_smp;
108         struct list_head s_inodes; /* all inodes */
109
110         spinlock_t s_inode_wblist_lock;
111         struct list_head s_inodes_wb; /* writeback inodes */
112     } __randomize_layout;

```

super_operations

super_block中的一个重要字段是 `const struct super_operations *s_op`，这个结构定义一组用来管理这个文件系统中inode的函数，其中实现了一些高级操作，比如删除文件或安装磁盘。

每个具体的文件系统都可以定义自己的超级块操作。当VFS需要调用其中一个操作时，比如 `write_inode()`，它执行下列操作：

```
1 sb->s_op->write_inode();
```

super_operations 结构的定义如下：

```

1  /* /include/linux/fs.h */
2
3  struct super_operations {
4      struct inode *(*alloc_inode)(struct super_block *sb);
5      void (*destroy_inode)(struct inode *);
6
7      void (*dirty_inode) (struct inode *, int flags);
8      int (*write_inode) (struct inode *, struct writeback_control *wbc);
9      int (*drop_inode) (struct inode *);
10     void (*evict_inode) (struct inode *);
11     void (*put_super) (struct super_block *);
12     int (*sync_fs)(struct super_block *sb, int wait);
13     int (*freeze_super) (struct super_block *);
14     int (*freeze_fs) (struct super_block *);
15     int (*thaw_super) (struct super_block *);
16     int (*unfreeze_fs) (struct super_block *);
17     int (*statfs) (struct dentry *, struct kstatfs *);
18     int (*remount_fs) (struct super_block *, int *, char *);
19     void (*umount_begin) (struct super_block *);
20

```

```

21     int (*show_options)(struct seq_file *, struct dentry *);
22     int (*show_devname)(struct seq_file *, struct dentry *);
23     int (*show_path)(struct seq_file *, struct dentry *);
24     int (*show_stats)(struct seq_file *, struct dentry *);
25 #ifdef CONFIG_QUOTA
26     ssize_t (*quota_read)(struct super_block *, int, char *, size_t, loff_t);
27     ssize_t (*quota_write)(struct super_block *, int, const char *, size_t, loff_t);
28     struct dquot **(*get_dquots)(struct inode *);
29 #endif
30     int (*bdev_try_to_free_page)(struct super_block *, struct page *, gfp_t);
31     long (*nr_cached_objects)(struct super_block *,
32                               struct shrink_control *);
33     long (*free_cached_objects)(struct super_block *,
34                                struct shrink_control *);
35 };

```

↑ 0%

inode 索引节点

inode（索引节点对象）表示文件系统中的对象，它具有唯一标识符。各个文件系统提供将文件名映射为唯一 inode 标识符和 inode 引用的方法，文件名可以随时更改，但是索引节点对文件是唯一的，并且随文件的存在而存在。

inode 结构的定义如下：

```

1  /* /include/linux/fs.h */
2  /*
3   * Keep mostly read-only and often accessed (especially for
4   * the RCU path lookup and 'stat' data) fields at the beginning
5   * of the 'struct inode'
6   */
7  struct inode {
8      umode_t                i_mode;
9      unsigned short         i_opflags;
10     kuid_t                  i_uid;
11     kgid_t                  i_gid;
12     unsigned int            i_flags;
13
14 #ifdef CONFIG_FS_POSIX_ACL
15     struct posix_acl         *i_acl;
16     struct posix_acl         *i_default_acl;
17 #endif
18
19     const struct inode_operations *i_op;
20     struct super_block        *i_sb; //对应的超级块对象

```

```

21     struct address_space    *i_mapping;
22
23 #ifdef CONFIG_SECURITY
24     void                    *i_security;
25 #endif
26
27     /* Stat data, not accessed from path walking */
28     unsigned long          i_ino; //索引节点号
29     /*
30      * Filesystems may only read i_nlink directly. They shall use the
31      * following functions for modification:
32      *
33      * (set|clear|inc|drop)_nlink
34      * inode_(inc|dec)_link_count
35      */
36     union {
37         const unsigned int i_nlink;
38         unsigned int __i_nlink;
39     };
40     dev_t                  i_rdev;
41     loff_t                 i_size;
42     struct timespec64      i_atime;
43     struct timespec64      i_mtime;
44     struct timespec64      i_ctime;
45     spinlock_t             i_lock; /* i_blocks, i_bytes, maybe i_size */
46     unsigned short         i_bytes;
47     u8                    i_blkbits;
48     u8                    i_write_hint;
49     blkcnt_t              i_blocks;
50
51 #ifdef __NEED_I_SIZE_ORDERED
52     seqcount_t             i_size_seqcount;
53 #endif
54
55     /* Misc */
56     unsigned long          i_state;
57     struct rw_semaphore    i_rwsem;
58
59     unsigned long          dirtied_when; /* jiffies of first dirtying */
60     unsigned long          dirtied_time_when;
61
62     struct hlist_node      i_hash;
63     struct list_head       i_io_list; /* backing dev IO list */
64 #ifdef CONFIG_CGROUP_WRITEBACK
65     struct bdi_writeback   *i_wb; /* the associated cgroup wb */
66
67     /* foreign inode detection, see wbc_detach_inode() */

```


↑ 0%

```
68         int                i_wb_frn_winner;
69         u16                 i_wb_frn_avg_time;
70         u16                 i_wb_frn_history;
71     #endif
72         struct list_head     i_lru;           /* inode LRU list */
73         struct list_head     i_sb_list;
74         struct list_head     i_wb_list;       /* backing dev writeback list */
75         union {
76             struct hlist_head i_dentry;
77             struct rcu_head    i_rcu;
78         };
79         atomic64_t            i_version;
80         atomic_t              i_count;        //引用计数器
81         atomic_t              i_dio_count;
82         atomic_t              i_writecount;
83     #ifdef CONFIG_IMA
84         atomic_t              i_readcount;    /* struct files open RO */
85     #endif
86         const struct file_operations *i_fop; /* former ->i_op->default_file_op;
87         struct file_lock_context *i_flctx;
88         struct address_space i_data;
89         struct list_head     i_devices;
90         union {
91             struct pipe_inode_info *i_pipe;
92             struct block_device *i_bdev;
93             struct cdev *i_cdev;
94             char *i_link;
95             unsigned i_dir_seq;
96         };
97
98         __u32                 i_generation;
99
100    #ifdef CONFIG_FSNOTIFY
101         __u32                 i_fsnotify_mask; /* all events this inode cares at
102         struct fsnotify_mark_connector __rcu *i_fsnotify_marks;
103    #endif
104
105    #if IS_ENABLED(CONFIG_FS_ENCRYPTION)
106         struct fscrypt_info *i_crypt_info;
107    #endif
108
109         void                  *i_private; /* fs or device private pointer */
110    } __randomize_layout;
```

inode_operations

inode_operations 定义直接在 inode 上执行的操作

```
1  /* /include/linux/fs.h */
2  struct inode_operations {
3      struct dentry * (*lookup) (struct inode *,struct dentry *, unsigned int);
4      const char * (*get_link) (struct dentry *, struct inode *, struct delayed_c
5      int (*permission) (struct inode *, int);
6      struct posix_acl * (*get_acl)(struct inode *, int);
7
8      int (*readlink) (struct dentry *, char __user *,int);
9
10     int (*create) (struct inode *,struct dentry *, umode_t, bool);
11     int (*link) (struct dentry *,struct inode *,struct dentry *);
12     int (*unlink) (struct inode *,struct dentry *);
13     int (*symlink) (struct inode *,struct dentry *,const char *);
14     int (*mkdir) (struct inode *,struct dentry *,umode_t);
15     int (*rmdir) (struct inode *,struct dentry *);
16     int (*mknod) (struct inode *,struct dentry *,umode_t,dev_t);
17     int (*rename) (struct inode *, struct dentry *,
18                     struct inode *, struct dentry *, unsigned int);
19     int (*setattr) (struct dentry *, struct iattr *);
20     int (*getattr) (const struct path *, struct kstat *, u32, unsigned int);
21     ssize_t (*listxattr) (struct dentry *, char *, size_t);
22     int (*fiemap)(struct inode *, struct fiemap_extent_info *, u64 start,
23                   u64 len);
24     int (*update_time)(struct inode *, struct timespec64 *, int);
25     int (*atomic_open)(struct inode *, struct dentry *,
26                         struct file *, unsigned open_flag,
27                         umode_t create_mode);
28     int (*tmpfile) (struct inode *, struct dentry *, umode_t);
29     int (*set_acl)(struct inode *, struct posix_acl *, int);
30 } ____cacheline_aligned;
```

file_operations

file_operations 定义与文件和目录相关的方法（标准系统调用）。

```
1  /* /include/linux.h */
2  struct file_operations {
3      struct module *owner;
4      loff_t (*llseek) (struct file *, loff_t, int);
5      ssize_t (*read) (struct file *, char __user *, size_t, loff_t *);
6      ssize_t (*write) (struct file *, const char __user *, size_t, loff_t *);
7      ssize_t (*read_iter) (struct kiocb *, struct iov_iter *);
```

```

8      ssize_t (*write_iter) (struct kiocb *, struct iov_iter *);
9      int (*iterate) (struct file *, struct dir_context *);
10     int (*iterate_shared) (struct file *, struct dir_context *);
11     __poll_t (*poll) (struct file *, struct poll_table_struct *);
12     long (*unlocked_ioctl) (struct file *, unsigned int, unsigned long);
13     long (*compat_ioctl) (struct file *, unsigned int, unsigned long);
14     int (*mmap) (struct file *, struct vm_area_struct *);
15     unsigned long mmap_supported_flags;
16     int (*open) (struct inode *, struct file *);
17     int (*flush) (struct file *, fl_owner_t id);
18     int (*release) (struct inode *, struct file *);
19     int (*fsync) (struct file *, loff_t, loff_t, int datasync);
20     int (*fasync) (int, struct file *, int);
21     int (*lock) (struct file *, int, struct file_lock *);
22     ssize_t (*sendpage) (struct file *, struct page *, int, size_t, loff_t *, int);
23     unsigned long (*get_unmapped_area)(struct file *, unsigned long, unsigned long, int, int);
24     int (*check_flags)(int);
25     int (*flock) (struct file *, int, struct file_lock *);
26     ssize_t (*splice_write)(struct pipe_inode_info *, struct file *, loff_t *, loff_t *, int);
27     ssize_t (*splice_read)(struct file *, loff_t *, struct pipe_inode_info *, int, int);
28     int (*setlease)(struct file *, long, struct file_lock **, void **);
29     long (*fallocate)(struct file *file, int mode, loff_t offset,
30                       loff_t len);
31     void (*show_fdinfo)(struct seq_file *m, struct file *f);
32 #ifndef CONFIG_MMU
33     unsigned (*mmap_capabilities)(struct file *);
34 #endif
35     ssize_t (*copy_file_range)(struct file *, loff_t, struct file *,
36                               loff_t, size_t, unsigned int);
37     int (*clone_file_range)(struct file *, loff_t, struct file *, loff_t,
38                             u64);
39     int (*dedupe_file_range)(struct file *, loff_t, struct file *, loff_t,
40                             u64);
41     int (*fadvise)(struct file *, loff_t, loff_t, int);
42 } __randomize_layout;

```

↑ 0%

file 文件对象

存放打开文件与进程之间进行交互的有关信息。当进程打开一个文件时，会创建一个 file 文件对象，它描述进程怎样与一个打开的文件进行交互。

严格来讲，这个结构体不仅仅和文件系统相关，它也和进程相关联。对于虚拟文件系统而言，它的重要性一般，而相反，这个结构体对进程管理更重要，它记录了进程打开文件的上下文。该结构体，浮于表面，更贴近用户，更贴近进程。

```
1  /* /include/linux/fs.h */
2  struct file {
3      union {
4          struct llist_node    fu_llist;
5          struct rcu_head      fu_rcuhead;
6      } f_u;
7      struct path              f_path;
8      struct inode             *f_inode;      /* cached value */
9      const struct file_operations *f_op;
10
11     /*
12      * Protects f_ep_links, f_flags.
13      * Must not be taken from IRQ context.
14      */
15     spinlock_t                f_lock;
16     enum rw_hint              f_write_hint;
17     atomic_long_t             f_count;
18     unsigned int              f_flags;
19     fmode_t                   f_mode;
20     struct mutex               f_pos_lock;
21     loff_t                    f_pos;
22     struct fown_struct         f_owner;
23     const struct cred          *f_cred;
24     struct file_ra_state      f_ra;
25
26     u64                        f_version;
27 #ifdef CONFIG_SECURITY
28     void                       *f_security;
29 #endif
30     /* needed for tty driver, and maybe others */
31     void                       *private_data;
32
33 #ifdef CONFIG_EPOLL
34     /* Used by fs/eventpoll.c to link all the hooks to this file */
35     struct list_head          f_ep_links;
36     struct list_head          f_tfile_llink;
37 #endif /* #ifdef CONFIG_EPOLL */
38     struct address_space      *f_mapping;
39     errseq_t                  f_wb_err;
40 } __randomize_layout
```

有一个重要信息是文件指针 `loff_t f_pos`，即文件中当前的位置，下一个操作将在该位置发生。由于几个进程可能同时访问同一文件，因此文件指针必须存放在文件对象而不是索引节点对象中。

↑ 0%

file结构体中的file_operations

每个文件系统都有其自己的文件操作集合，执行诸如读写文件这样的操作。当进程打开一个文件时，VFS会用对应的 `inode->i_fop` 中的 `file_operations` 结构体 来初始化新文件对象的 `file->fop`，使得对文件操作的后续调用能够使用这些函数。如果需要VFS也可以通过修改这个字段而修改文件操作的集合。

[illegible]

```

37         int (*clone_file_range)(struct file *, loff_t, struct file *, loff_t,
38                                 u64);
39         int (*dedupe_file_range)(struct file *, loff_t, struct file *, loff_t,
40                                 u64);
41         int (*fadvise)(struct file *, loff_t, loff_t, int);
42     } __randomize_layout;

```

↑ 0%

path

file 结构体中的 struct path f_path 字段，记录了这个文件对象对应的目录项（dentry）和文件系统（vfsmount）

```

1  /* /include/linux/path.h */
2  struct path {
3      struct vfsmount *mnt; //含有该文件的已安装文件系统
4      struct dentry *dentry; //与文件相关的目录项对象
5  } __randomize_layout;

```

dentry 目录项对象

dentry是目录项缓存，是一个存放在内存里的缩略版的磁盘文件系统目录树结构，他是directory entry的缩写。它存放目录项（也就是文件的特定名称）与对应文件进行链接的有关信息。一旦目录项被读入内存，VFS就把它转换成基于 dentry 结构的一个目录项对象。对于进程查找的路径名中的每个分量，内核都为其创建一个目录项对象；目录对象将每个分量与其对应的索引节点相联系。例如在查找路径名 /tmp/test 时，内核为根目录 / 创建一个目录对象，为根目录下的 tmp 项创建一个第二级目录项对象，对 /tmp 目录下的 test 项创建一个第三级目录项对象。

目录项对象存放在名为 dentry_cache 的slab分配器高速缓存中。因此，目录项对象的创建和删除是通过 kmem_cache_alloc() 和 kmem_cache_free() 实现的。

```

1  /* /include/linux/dcache.h */
2  struct dentry {
3      /* RCU lookup touched fields */
4      unsigned int d_flags;           /* protected by d_lock */
5      seqcount_t d_seq;              /* per dentry seqlock */
6      struct hlist_bl_node d_hash;    /* lookup hash list */
7      struct dentry *d_parent;        /* parent directory */
8      struct qstr d_name;
9      struct inode *d_inode;          /* Where the name belongs to - NULL is
10                                     * negative */
11      unsigned char d_iname[DNAME_INLINE_LEN]; /* small names */

```

```

12
13     /* Ref lookup also touches following */
14     struct lockref d_lockref;          /* per-dentry lock and refcount */
15     const struct dentry_operations *d_op;
16     struct super_block *d_sb;          /* The root of the dentry tree */
17     unsigned long d_time;              /* used by d_revalidate */
18     void *d_fsdata;                   /* fs-specific data */
19
20     union {
21         struct list_head d_lru;        /* LRU list */
22         wait_queue_head_t *d_wait;    /* in-lookup ones only */
23     };
24     struct list_head d_child;          /* child of parent list */
25     struct list_head d_subdirs;        /* our children */
26     /*
27      * d_alias and d_rcu can share memory
28      */
29     union {
30         struct hlist_node d_alias;     /* inode alias list */
31         struct hlist_bl_node d_in_lookup_hash; /* only for in-lookup ones */
32         struct rcu_head d_rcu;
33     } d_u;
34 } __randomize_layout;

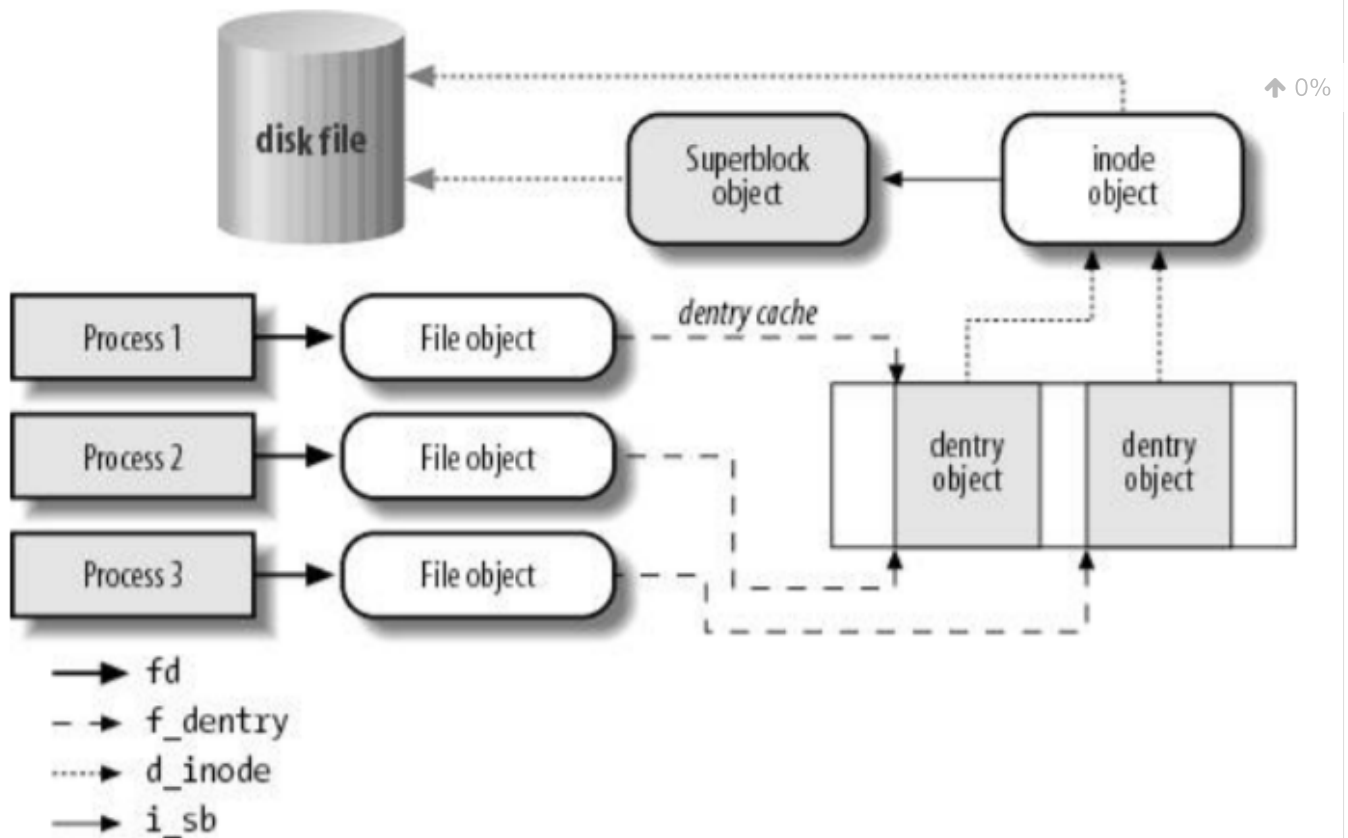
```

↑ 0%

inode dentry file 对象之间的关系

下图是一个简单的实例。三个不同的进程打开同一个文件，其中两个进程使用同一个硬链接，还有一个进程使用另一个硬链接，但是都是指向同一个文件的。在这种情况下，其中的每个进程都使用自己的文件对象，但只需要两个目录项对象，每个硬链接对应一个目录项对象。这两个目录项对象都指向同一个索引节点对象，该索引节点对象标识超级块对象。

Figure 12-2. Interaction between processes and VFS objects



与进程相关的文件

fs_struct

文件描述符 `task_struct` 中有一个字段 `struct fs_struct *fs`，它记录了进程当前的工作目录和根目录等信息，这是内核用来表示进程与文件系统相互作用所必须维护的数据。

```
1  /* /include/linux/fs_struct.h */
2  struct fs_struct {
3      int users;
4      spinlock_t lock;
5      seqcount_t seq;
6      int umask;
7      int in_exec;
8      struct path root, pwd; //root为根目录, pwd为当前工作的目录
9  } __randomize_layout;
10
11 /* /include/linux/path.h */
12 struct path {
13     struct vfsmount *mnt;
14     struct dentry *dentry;
15 } __randomize_layout;
```


files_struct

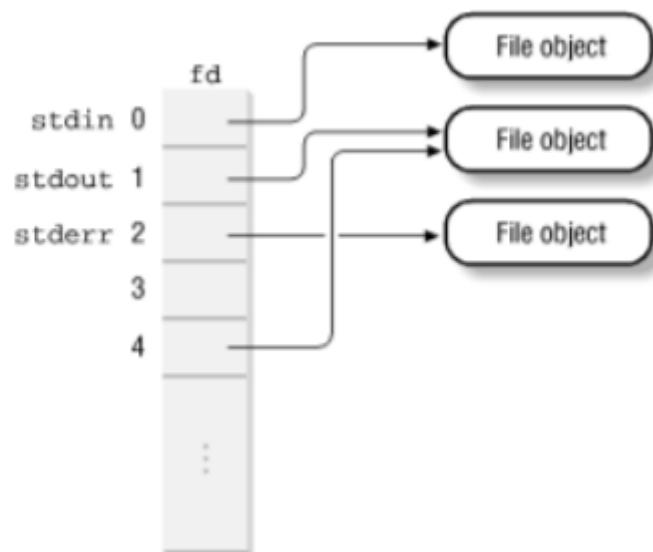
文件描述符 `task_struct` 中有一个字段 `struct files_struct *files`，它记录了进程当前打开的文件的信息。

↑ 0%

```
1  /* /include/linux/fdtable.h */
2  /*
3   * Open file table structure
4   */
5  struct files_struct {
6      /*
7       * read mostly part
8       */
9          atomic_t count;
10         bool resize_in_progress;
11         wait_queue_head_t resize_wait;
12
13         struct fdtable __rcu *fdt;
14         struct fdtable fdtab;
15     /*
16      * written part on a separate cache line in SMP
17      */
18         spinlock_t file_lock ____cacheline_aligned_in_smp;
19         unsigned int next_fd;
20         unsigned long close_on_exec_init[1];
21         unsigned long open_fds_init[1];
22         unsigned long full_fds_bits_init[1];
23         struct file __rcu * fd_array[NR_OPEN_DEFAULT];
24 };
25
26 struct fdtable {
27     unsigned int max_fds;
28     struct file __rcu **fd;      /* current fd array */
29     unsigned long *close_on_exec;
30     unsigned long *open_fds;
31     unsigned long *full_fds_bits;
32     struct rcu_head rcu;
33 };
```

`fdtable` 中的 `fd` 是一个指针，指向一个数组，数组中的每一个元素都是一个 `struct file` 类型的指针，对于在 `fd` 数组中有相应对象指针的文件来说，数组的索引下标就是文件描述符(file descriptor)。通常，数组的第一个元素（索引为0）是进程的标准输入文件，数组的第二个元素（索引为1）是进程的标准输出文件，数组的第三个元素（索引为2）是进程的标准错误文件。

The fd array



↑ 0%

请注意，借助于dup()、dup2()和fcntl()系统调用，两个文件描述符可以指向同一个打开的文件，也就是说，数组的两个元素可以指向同一个文件对象。例如，当用户将标准错误文件重定向到标准输出文件上时。

相关文章

- [Linux mmap与DMA 初探 【未完成】](#)
- [Linux内核学习参考资料汇总](#)
- [Linux内核漏洞利用 Bypass SMEP](#)
- [QEMU+busybox 搭建Linux内核运行环境](#)
- [Linux内核-进程、线程初探](#)
- [初探Linux内核UAF](#)
- [初探Linux内核栈溢出](#)

本文作者： 孙小空

本文链接： <http://www.sunxiaokong.xyz/2019-12-02/lzx-01-babyVFS/>

版权声明： 本博客所有文章除特别声明外，均采用 [CC BY-NC-SA](#) 许可协议。转载请注明出处！

[Linux内核](#)

