# 文件系统

by zenhumany

2012-05-19——2012-06-03

# 目录

# 1. 概述

本章中所涉及的源代码全部来至 linux 内核 2.6.34。

操作系统最重要的部分莫过于文件系统和进程管理了。一个完整的 linux 操作系统，通常包含数千到数百万个文件。

每种操作系统都至少有一种"标准文件系统"，提供一些功能，可以高效而可靠的执行所需的任务。Linux 中标准的文件系统为 ext2/ext3。

为了支持各种文件系统，且在同时允许访问其他操作系统的文件，linux 内核在用户进程和文件系统中引入了一个抽象层。该抽象层成为虚拟文件系统（VFS，Virtual File System）。

虚拟文件系统接口和数据结构构成了一个结构，各个文件系统都在该框架下运转。

本章将讲解虚拟文件系统的接口和相关数据结构，同事讨论在该框架下实现的具体文件系统 ext2。

# 2. 虚拟文件系统和相关数据结构

Linux 系统能够快速的成长，与其支持各种类型的文件系统不无关系。Linux 除了本身的文件系统 ext2 外，可以支持其他各种不同文件系统。要实现这个目的，就需要将对各种不同文件系统的操作和管理纳入到一个统一的框架中，然内核中的文件系统界面成为一条"总线"，使得用户程序可以通过同一个文件系统操作界面，也就是同一组系统调用，对各种不同的文件系统进行操作。这就是 VFS 所要达到的目的。下面看和 VFS 相关的数据结构。

## 2.1 inode文件索引节点

在 linux 系统中，磁盘上的一个文件由 indoe 索引节点唯一的标识。

include/linux/fs.h

```
724 struct inode {
725        struct hlist_node      i_hash;
726        struct list_head       i_list;        /* backing dev IO list */
727        struct list_head       i_sb_list;
728        struct list_head       i_dentry;
729        unsigned long          i_ino;
730        atomic_t               i_count;
731        unsigned int           i_nlink;
732        uid_t                  i_uid;
733        gid_t                  i_gid;
734        dev_t                  i_rdev;
735        unsigned int           i_blkbits;
736        u64            i_version;
737        loff_t             i_size;
738 #ifdef __NEED_I_SIZE_ORDERED
739        seqcount_t         i_size_seqcount;
740 #endif
741        struct timespec        i_atime;
742        struct timespec        i_mtime;
743        struct timespec        i_ctime;
744        blkcnt_t           i_blocks;
745        unsigned short             i_bytes;
746        umode_t            i_mode;
747        spinlock_t         i_lock; /* i_blocks, i_bytes, maybe i_size */
748        struct mutex           i_mutex;
749        struct rw_semaphore i_alloc_sem;
750        const struct inode_operations     *i_op;
751        const struct file_operations       *i_fop; /* former
->i_op->default_file_ops */
752        struct super_block   *i_sb;
753        struct file_lock       *i_flock;
```

```
754        struct address_space      *i_mapping;
755        struct address_space      i_data;
756 #ifdef CONFIG_QUOTA
757        struct dquot              *i_dquot[MAXQUOTAS];
758 #endif
759        struct list_head      i_devices;
760        union {
761            struct pipe_inode_info   *i_pipe;
762            struct block_device *i_bdev;
763            struct cdev        *i_cdev;
763            struct cdev        *i_cdev;
764        };
765
766        __u32                 i_generation;
767
768 #ifdef CONFIG_FSNOTIFY
769        __u32                 i_fsnotify_mask; /* all events this inode cares about */
770        struct hlist_head     i_fsnotify_mark_entries; /* fsnotify mark entries */
771 #endif
772
773 #ifdef CONFIG_INOTIFY
774        struct list_head      inotify_watches; /* watches on this inode */
775        struct mutex          inotify_mutex;   /* protects the watches list */
776 #endif
777
778        unsigned long         i_state;
779        unsigned long         dirtied_when;    /* jiffies of first dirtying */
780
781        unsigned int          i_flags;
782
783        atomic_t              i_writecount;
```

```
784 #ifdef CONFIG_SECURITY
785      void            *i_security;
786 #endif
787 #ifdef CONFIG_FS_POSIX_ACL
788      struct posix_acl   *i_acl;
789      struct posix_acl   *i_default_acl;
790 #endif
791      void            *i_private; /* fs or device private pointer */
792 };
```

- i_hash：散列链表指针
- i_list：索引节点当前转台的链表指针
- i_sb_list：超级块索引节点指针
- i_dentry：引用索引节点的目录对象链表头
- i_ino：索引节点在磁盘上的索引节点号。
- i_op：为具体文件系统和 inode 相关的操作提供一个插口
- i_fop：为具体文件系统和文件相关的操作提供一个插口
- i_sb：指向该索引节点所属的超级块对象。
- i_mapping：指向缓存 address_space 对象的指针，该结构是文件系统缓存的关键。
- i_data：嵌入在 inode 中的 address_space 对象。

每个 inode 的成员可以分为两类，一类为描述文件状态的元数据，一类为保存实际文件的数据段。

每个 inode 都有一个"i 节点号"i_ino，在同一文件系统中该节点号是唯一的。

inode 结构中有两个设备号，即 i_dev 和 i_cdev。除了特殊文件外，一个索引节点总的存储在某个设备上，这就是 i_dev。其次，如果索引节点代表的不是常规文件，而是某个设备，那就要有个设备号，这个就是 i_rdev。

就像人可以有别名一样，文件也可以有别名，也就是将一个以创建的文件"连接"到另一个文件名，可以沿着 i_dentry 找到所有这个文件的"别名"。

在应用层，总是通过文件名来标识文件，表示文件名的数据结构式
dentry。

## 2.2 目录项缓存dentry

为了加快文件的访问速度，内核需要缓存已经打开过的文件的信息，用
来缓存文件信息的数据结构就是 dentry。

include/linux/dcache.h

```
89 struct dentry {
 90      atomic_t d_count;
 91      unsigned int d_flags;           /* protected by d_lock */
 92      spinlock_t d_lock;          /* per dentry lock */
 93      int d_mounted;
 94      struct inode *d_inode;          /* Where the name belongs to - NULL is
 95                          * negative */
 96      /*
 97       * The next three fields are touched by __d_lookup.   Place them here
 98       * so they all fit in a cache line.
 99       */
100      struct hlist_node d_hash;    /* lookup hash list */
101      struct dentry *d_parent;      /* parent directory */
102      struct qstr d_name;
103
104      struct list_head d_lru;       /* LRU list */
105      /*
106       * d_child and d_rcu can share memory
107       */
108      union {
109          struct list_head d_child;    /* child of parent list */
110          struct rcu_head d_rcu;
111      } d_u;
```

```
112        struct list_head d_subdirs; /* our children */

113        struct list_head d_alias;    /* inode alias list */

114        unsigned long d_time;          /* used by d_revalidate */

115        const struct dentry_operations *d_op;

116        struct super_block *d_sb;    /* The root of the dentry tree */

117        void *d_fsdata;              /* fs-specific data */

118

119        unsigned char d_iname[DNAME_INLINE_LEN_MIN];        /* small
names */

120 };
```

- d_inode：指向目录项对应的 inode。
- d_parent：指向父目录。
- d_name：文件名。
- d_subdirs：对目录而言，子目录项的链表的头。
- d_alias：用于与同一索引节点（别名）相关的目录项链表的指针。
- d_op：为具体的文件系统提供目录项操作接口。
- d_sb：文件超级快对象。
- d_iname：短文件名存储空间

```
33 struct qstr {

 34        unsigned int hash;

 35        unsigned int len;

 36        const unsigned char *name;

 37 };
```

- hash：文件名 hash 值
- len：文件名长度
- name：指向存放文件名空间的指针。

## 2.3 文件对象

文件对象描述进程与一个打开文件的交互，是打开文件的一个上下文。
文件对象在磁盘上没有对应的映射。

```c
913 struct file {
914     /*
915      * fu_list becomes invalid after file_free is called and queued via
916      * fu_rcuhead for RCU freeing
917      */
918     union {
919         struct list_head        fu_list;
920         struct rcu_head         fu_rcuhead;
921     } f_u;
922     struct path        f_path;
923 #define f_dentry     f_path.dentry
924 #define f_vfsmnt     f_path.mnt
925     const struct file_operations     *f_op;
926     spinlock_t        f_lock;   /* f_ep_links, f_flags, no IRQ */
927     atomic_long_t         f_count;
928     unsigned int         f_flags;
929     fmode_t          f_mode;
930     loff_t               f_pos;
931     struct fown_struct   f_owner;
932     const struct cred    *f_cred;
933     struct file_ra_state     f_ra;
934
935     u64          f_version;
936 #ifdef CONFIG_SECURITY
937     void             *f_security;
938 #endif
939     /* needed for tty driver, and maybe others */
940     void             *private_data;
941
942 #ifdef CONFIG_EPOLL
943     /* Used by fs/eventpoll.c to link all the hooks to this file */
```

```
944      struct list_head    f_ep_links;
945 #endif /* #ifdef CONFIG_EPOLL */
946      struct address_space    *f_mapping;
947 #ifdef CONFIG_DEBUG_WRITECOUNT
948      unsigned long f_mnt_write_state;
949 #endif
950 };
```

include/linux/path.h

```
4 struct dentry;
5 struct vfsmount;
6
7 struct path {
8      struct vfsmount *mnt;
9      struct dentry *dentry;
10 };
```

- f_path.dentry：指向 file 对应的 dentry 结构
- f_paht.mnt：指向 file 所在文件的文件系统类型。
- f_op：为具体的文件系统提供文件操作插口。
- f_pos：文件当前的读写位置，也就是文件的上下文。
- f_mapping：指向文件地址空间对象的指针。
- private_data：指向特定文件系统或设备驱动程序所需的数据的指针

## 2.4 特定于进程的数据结构

用户进程都是通过一个文件描述符（就是整数）来在一个进程内唯一的标识打开的文件。

这就需要内核可以在用户进程中的描述符和内核之间内部使用结构建立一种关联。

每个进程的 task_struct 都包含了用于完成该工作的成员。

Include/linux/sched.h

```
1170 struct task_struct {
……
```

```
1329 /* filesystem information */
1330      struct fs_struct *fs;
1331 /* open file information */
1332      struct files_struct *files;
1333 /* namespaces */
1334      struct nsproxy *nsproxy;
1335 /* signal handlers */
……
1508 };
```

fs：文件系统信息

files：打开文件信息

## 2.4.1  fs_struct

include/linux/fs_struct.h

```
6 struct fs_struct {
 7      int users;
 8      rwlock_t lock;
 9      int umask;
10      int in_exec;
11      struct path root, pwd;
12 };
```

include/linux/path.h

```
4 struct dentry;
 5 struct vfsmount;
 6
 7 struct path {
 8      struct vfsmount *mnt;
 9      struct dentry *dentry;
10 };
```

fs_struct 中的连个类型为 path 的变量 root，pwd 为关键的结构。

path 结构中 dentry 指向一个目录项，mnt 表示 dentry 代表的目录项的安装的数据结构 vfsmount。

pwd.dentry,root.dentry 分别表示进程当前所在的目录，进程的"根目录"，也就是用户登录系统时看到的目录。这两个目录有可能并不在同一文件系统中，进程的根目录通常是安装于"/"节点上的 ext2 文件系统，而当前工作目录则可能安装于/dosc 的一个 DOS 文件系统。在文件系统中这些安装点起着非常重要的作用，所以 pwd.mnt，root.mnt 分别指向对应 dentry 的安装点。

## 2.4.2  files_struct

```
40 /*
41    * Open file table structure
42    */
43 struct files_struct {
44    /*
45     * read mostly part
46     */
47      atomic_t count;
48      struct fdtable *fdt;
49      struct fdtable fdtab;
50    /*
51     * written part on a separate cache line in SMP
52     */
53      spinlock_t file_lock ____cacheline_aligned_in_smp;
54      int next_fd;
55      struct embedded_fd_set close_on_exec_init;
56      struct embedded_fd_set open_fds_init;
57      struct file * fd_array[NR_OPEN_DEFAULT];
58 };
```

fdtab 嵌入到 files_struct 中，并由 fdt 指向它。

```
31 struct fdtable {
32      unsigned int max_fds;
33      struct file ** fd;          /* current fd array */
34      fd_set *close_on_exec;
35      fd_set *open_fds;
36      struct rcu_head rcu;
37      struct fdtable *next;
38 };
```

fd 指向 file* fd_array[]的数组，当 fd_array 中没有空闲的元素时，扩充 fdtable，并由 next 连接。



## 2.5 文件系统file_system_type

内核支持的每一种文件系统都有一个 file_system_type 实例。

```
1735 struct file_system_type {
1736      const char *name;
1737      int fs_flags;
1738      int (*get_sb) (struct file_system_type *, int,
1739                     const char *, void *, struct vfsmount *);
1740      void (*kill_sb) (struct super_block *);
1741      struct module *owner;
```

```
1742        struct file_system_type * next;

1743        struct list_head fs_supers;

1744

1745        struct lock_class_key s_lock_key;

1746        struct lock_class_key s_umount_key;

1747

1748        struct lock_class_key i_lock_key;

1749        struct lock_class_key i_mutex_key;

1750        struct lock_class_key i_mutex_dir_key;

1751        struct lock_class_key i_alloc_sem_key;

1752 };
```

name：文件系统的名称

get_sb：获取超级快的函数。

## 2.6 安装点vfsmount

把一个设备安装到一个目录节点时需要用到一个 vfsmount 数据结构作为"连接件"。

```
50 struct vfsmount {

51      struct list_head mnt_hash;

52      struct vfsmount *mnt_parent;      /* fs we are mounted on */

53      struct dentry *mnt_mountpoint;   /* dentry of mountpoint */

54      struct dentry *mnt_root;      /* root of the mounted tree */

55      struct super_block *mnt_sb; /* pointer to superblock */

56      struct list_head mnt_mounts;      /* list of children, anchored here */

57      struct list_head mnt_child; /* and going through their mnt_child */

58      int mnt_flags;

59      /* 4 bytes hole on 64bits arches */

60      const char *mnt_devname;      /* Name of device e.g. /dev/dsk/hda1 */

61      struct list_head mnt_list;
```

```
62      struct list_head mnt_expire;      /* link in fs-specific expiry list */
63      struct list_head mnt_share; /* circular list of shared mounts */
64      struct list_head mnt_slave_list;/* list of slave mounts */
65      struct list_head mnt_slave; /* slave list entry */
66      struct vfsmount *mnt_master;      /* slave is on master->mnt_slave_list */
67      struct mnt_namespace *mnt_ns;     /* containing namespace */
68      int mnt_id;             /* mount identifier */
69      int mnt_group_id;           /* peer group identifier */
70      /*
71       * We put mnt_count & mnt_expiry_mark at the end of struct vfsmount
72       * to let these frequently modified fields in a separate cache line
73       * (so that reads of mnt_flags wont ping-pong on SMP machines)
74       */
75      atomic_t mnt_count;
76      int mnt_expiry_mark;          /* true if marked for expiry */
77      int mnt_pinned;
78      int mnt_ghosts;
79 #ifdef CONFIG_SMP
80      int __percpu *mnt_writers;
81 #else
82      int mnt_writers;
83 #endif
84 };
```

- mnt_parent：指向安装点所在设备当初安装时的 vfsmount，也就是上一层的 vfsmount。
- mnt_mountpoint，mnt_root：安装点的 dentry 数据结构，mount_root 则指向所安装设备上根目录的 dentry 结构，这连个目录在目录树中处于同一层次。
- mnt_sb：所安装设备的超级块。

- mnt_mounts，mnt_child：当上一层的 vfsmount 存在时，通过 mnt_child 链入到 mnt_mounts 中。
- mnt_list：系统中有个总的 vfsmount 结构队列 vfsmntlist，所有以安装的 vfsmount 结构都通过 mnt_list 链入到 vfsmntlist 中。

## 2.7 super_block超级块

文件系统的超级块记录了文件系统的大部分信息。

```
1319 struct super_block {
1320     struct list_head     s_list;      /* Keep this first */
1321     dev_t                s_dev;       /* search index; _not_ kdev_t */
1322     unsigned char        s_dirt;
1323     unsigned char        s_blocksize_bits;
1324     unsigned long        s_blocksize;
1325     loff_t               s_maxbytes; /* Max file size */
1326     struct file_system_type *s_type;
1327     const struct super_operations     *s_op;
1328     const struct dquot_operations     *dq_op;
1329     const struct quotactl_ops     *s_qcop;
1330     const struct export_operations *s_export_op;
1331     unsigned long        s_flags;
1332     unsigned long        s_magic;
1333     struct dentry        *s_root;
1334     struct rw_semaphore s_umount;
1335     struct mutex         s_lock;
1336     int          s_count;
1337     int          s_need_sync;
1338     atomic_t         s_active;
1339 #ifdef CONFIG_SECURITY
1340     void                     *s_security;
1341 #endif
1342     struct xattr_handler     **s_xattr;
```

```
1343
1344    struct list_head    s_inodes;    /* all inodes */
1345    struct hlist_head   s_anon;       /* anonymous dentries for (nfs)
exporting */
1346    struct list_head    s_files;
1347    /* s_dentry_lru and s_nr_dentry_unused are protected by
dcache_lock */
1348    struct list_head    s_dentry_lru;    /* unused dentry lru */
1349    int          s_nr_dentry_unused; /* # of dentry on lru */
1350
1351    struct block_device *s_bdev;
1352    struct backing_dev_info *s_bdi;
1353    struct mtd_info     *s_mtd;
1354    struct list_head    s_instances;
1355    struct quota_info    s_dquot;    /* Diskquota specific options */
1356
1357    int          s_frozen;
1358    wait_queue_head_t    s_wait_unfrozen;
1359
1360    char s_id[32];              /* Informational name */
1361
1362    void            *s_fs_info; /* Filesystem private info */
1363    fmode_t          s_mode;
1364
1365    /* Granularity of c/m/atime in ns.
1366      Cannot be worse than a second */
1367    u32          s_time_gran;
1368
1369    /*
1370     * The next field is for VFS *only*. No filesystems have any business
1371     * even looking at it. You had been warned.
1372     */
```

```
1373      struct mutex s_vfs_rename_mutex;      /* Kludge */
1374
1375      /*
1376       * Filesystem subtype.   If non-empty the filesystem type field
1377       * in /proc/mounts will be "type.subtype"
1378       */
1379      char *s_subtype;
1380
1381      /*
1382       * Saved mount options for lazy filesystems using
1383       * generic_show_options()
1384       */
1385      char *s_options;
1386 };
```

- s_list：指向超级块链表的指针。

- s_dev：设备标识符

- s_blocksize：以字节为单位的块大小

- s_blocksize_bits：以位为单位的块大小

- s_dirt：是否为脏标识符

- s_maxbytes：文件的最长长度

- s_op：具体文件系统超级块插口

- type：文件系统类型

- s_root：文件系统根目录的目录项对象

- s_inodes：所有索引节点的链表头

- s_anon：用于处理远程网络文件系统的匿名目录项的链表。

- s_files：文件对象的链表。

- s_fs_info：指向具体的文件系统的超级块信息。

# 3. ext2 文件系统

虚拟文件系统（VFS）接口和数据结构构成了一个框架，需要具体的文件系统来填充框架中的内容。本节以 linux 文件系统 ext2 为例讲解。

## 3.1 Ext2 磁盘数据结构

要学习 EXT2 文件系统，首先需要学习 EXT2 的磁盘设计。也就是 EXT2 将一个磁盘分区格式化成一个什么样子。

现代的磁盘驱动器都是多片的，所以不同的盘面上的相同磁道合在一起就形成了"柱面"（cylinder）的概念。从磁盘读出多个记录块时，如果从同一柱面中读取就比较快，因为这种情况下不需要移动磁头（磁头组）。互相连续的记录块实际上分布在同一柱面的各个盘面上，只有将一个柱面用完



后才会进入下一个柱面。所以，在许多文件系统中都把整个设备划分为若干柱面组，将反映着盘面存储空间的组织与管理的信息分散后就近存储于各个柱面中。相比之下，早期的文件系统往往将这些信息集中存储在一起，使得磁头在文件访问时来回"疲于奔命"而降低了效率。

ext2 文件系统在硬盘上的布局如下：

| 启动块 | 块组 0 | 块组 1 | …… | 块组 n |
|---|---|---|---|---|

每一个块组的划分如下：

| 超级块 | 组描述符 | 数据位图 | inode 位图 | inode 表 | 数据块 |
|---|---|---|---|---|---|

在 linux 下面，目录也是作为一种特定的文件来实现的。

- 图中组描述块记录着全部的组描述符结构，具体占用的块数取决于设备的大小。

- 数据块位图是本组的记录块位图，每一位对应一个记录块，1 表示已分配，0 表示空闲。
- Inode 位图是本块组中 inode 分配的记录位图，一般 inode 位图小于一个记录块。
- Inode 表：本块组中用来存放 inode 的记录块集合
- 数据库：本块组中用来存放数据块的记录块集合。

### 3.1.1 磁盘超级块

磁盘超级块描叙了磁盘格式化后各方面的信息，存储了用于管理磁盘的元信息。

EXT2 在磁盘上的超级块的数据结构为 ext2_super_block。

include/linux/ext2_fs.h

```
373 struct ext2_super_block {
374     __le32  s_inodes_count;      /* Inodes count */
375     __le32  s_blocks_count;      /* Blocks count */
376     __le32  s_r_blocks_count;    /* Reserved blocks count */
377     __le32  s_free_blocks_count;     /* Free blocks count */
378     __le32  s_free_inodes_count;     /* Free inodes count */
379     __le32  s_first_data_block; /* First Data Block */
380     __le32  s_log_block_size;    /* Block size */
381     __le32  s_log_frag_size;     /* Fragment size */
382     __le32  s_blocks_per_group; /* # Blocks per group */
383     __le32  s_frags_per_group;   /* # Fragments per group */
384     __le32  s_inodes_per_group; /* # Inodes per group */
385     __le32  s_mtime;             /* Mount time */
386     __le32  s_wtime;             /* Write time */
387     __le16  s_mnt_count;             /* Mount count */
388     __le16  s_max_mnt_count;     /* Maximal mount count */
389     __le16  s_magic;             /* Magic signature */
390     __le16  s_state;             /* File system state */
391     __le16  s_errors;            /* Behaviour when detecting errors */
```

```
392      __le16   s_minor_rev_level;   /* minor revision level */
393      __le32   s_lastcheck;              /* time of last check */
394      __le32   s_checkinterval;       /* max. time between checks */
395      __le32   s_creator_os;           /* OS */
396      __le32   s_rev_level;             /* Revision level */
397      __le16   s_def_resuid;           /* Default uid for reserved blocks */
398      __le16   s_def_resgid;            /* Default gid for reserved blocks */
399      /*
400       * These fields are for EXT2_DYNAMIC_REV superblocks only.
401       *
402       * Note: the difference between the compatible feature set and
403       * the incompatible feature set is that if there is a bit set
404       * in the incompatible feature set that the kernel doesn't
405       * know about, it should refuse to mount the filesystem.
406       *
407       * e2fsck's requirements are more strict; if it doesn't know
408       * about a feature in either the compatible or incompatible
409       * feature set, it must abort and not try to meddle with
410       * things it doesn't understand...
412      __le32   s_first_ino;           /* First non-reserved inode */
413      __le16    s_inode_size;         /* size of inode structure */
414      __le16   s_block_group_nr;    /* block group # of this superblock */
415      __le32   s_feature_compat;    /* compatible feature set */
416      __le32   s_feature_incompat;      /* incompatible feature set */
417      __le32   s_feature_ro_compat;      /* readonly-compatible feature set */
418      __u8      s_uuid[16];          /* 128-bit uuid for volume */
419      char      s_volume_name[16];   /* volume name */
420      char      s_last_mounted[64];       /* directory where last mounted */
421      __le32   s_algorithm_usage_bitmap; /* For compression */
422      /*
423       * Performance hints.   Directory preallocation should only
```

```
424        * happen if the EXT2_COMPAT_PREALLOC flag is on.
425        */
426       __u8    s_prealloc_blocks;   /* Nr of blocks to try to preallocate*/
427       __u8    s_prealloc_dir_blocks;   /* Nr to preallocate for dirs */
428       __u16   s_padding1;
429       /*
430        * Journaling support valid if EXT3_FEATURE_COMPAT_HAS_JOURNAL set.
431        */
432       __u8    s_journal_uuid[16]; /* uuid of journal superblock */
433       __u32   s_journal_inum;        /* inode number of journal file */
434       __u32   s_journal_dev;         /* device number of journal file */
435       __u32   s_last_orphan;         /* start of list of inodes to delete */
436       __u32   s_hash_seed[4];        /* HTREE hash seed */
437       __u8    s_def_hash_version; /* Default hash version to use */
438       __u8    s_reserved_char_pad;
439       __u16   s_reserved_word_pad;
440       __le32  s_default_mount_opts;
441       __le32  s_first_meta_bg;      /* First metablock block group */
442       __u32   s_reserved[190];      /* Padding to the end of the block */
443 };
```

- s_inodes_count：索引节点的总数
- s_blocks_count：块总数
- s_r_blocks_count：保留的块数
- s_free_blocks_count：空闲的块数
- s_free_indoes_count：空闲的索引节点数
- s_first_data_block：第一个使用的块号
- s_log_block_size：块的大小
- s_log_frag_size：片的大小
- s_blocks_per_group：每组中的块数
- s_frags_per_group：没组中的片数
- s_inodes_per_group：每组中的索引节点数

- s_mtime：最后一次安装操作的时间
- s_wtime：最后一次写操作的时间
- s_mnt_count：被执行安装操作的次数
- s_max_mnt_count：检查之前安装操作的次数
- s_magic：文件系统魔术签名
- s_state：文件系统的状态标志
- s_first_ino：第一个非保留的索引节点号
- s_inode_size：磁盘上索引节点结构的大小
- s_block_group_nr：这个超级快的块组号
- s_prealloc_blocks：预分配的块数
- s_prealloc_dir_blocks：为目录预分配的块数。

## 3.1.2 组描述符

每个块组都有一个组描述符，用于描述本块组的元信息。

```
130 /*
131    * Structure of a blocks group descriptor
132    */
133 struct ext2_group_desc
134 {
135      __le32   bg_block_bitmap;        /* Blocks bitmap block */
136      __le32   bg_inode_bitmap;         /* Inodes bitmap block */
137      __le32   bg_inode_table;       /* Inodes table block */
138      __le16   bg_free_blocks_count;    /* Free blocks count */
139      __le16   bg_free_inodes_count;    /* Free inodes count */
140      __le16   bg_used_dirs_count; /* Directories count */
141      __le16   bg_pad;
142      __le32   bg_reserved[3];
143 };
```

- bg_block_bitmap：块位图的块号
- bg_inode_bitmap：索引节点位图的块号

- bg_inode_table：索引节点表块的起始块号
- bg_free_blocks_count：组中空闲块的个数
- bg_free_inodes_count：组中空闲索引节点的个数
- bg_used_dirs_count：组中目录的个数

## 3.1.3 索引节点

索引节点记录了磁盘上一个文件的元信息。

```
239 /*
240   * Structure of an inode on the disk
241   */
242 struct ext2_inode {
243     __le16  i_mode;       /* File mode */
244     __le16  i_uid;        /* Low 16 bits of Owner Uid */
245     __le32  i_size;       /* Size in bytes */
246     __le32  i_atime;      /* Access time */
247     __le32  i_ctime;      /* Creation time */
248     __le32  i_mtime;       /* Modification time */
249     __le32  i_dtime;       /* Deletion Time */
250     __le16  i_gid;         /* Low 16 bits of Group Id */
251     __le16  i_links_count;   /* Links count */
252     __le32  i_blocks;     /* Blocks count */
253     __le32  i_flags;      /* File flags */
254     union {
255         struct {
256             __le32   l_i_reserved1;
257         } linux1;
258         struct {
259             __le32   h_i_translator;
260         } hurd1;
261         struct {
```

```c
262            __le32   m_i_reserved1;
263        } masix1;
264    } osd1;                /* OS dependent 1 */
265    __le32   i_block[EXT2_N_BLOCKS];/* Pointers to blocks */
266    __le32   i_generation;    /* File version (for NFS) */
267    __le32   i_file_acl; /* File ACL */
268    __le32   i_dir_acl;   /* Directory ACL */
269    __le32   i_faddr;      /* Fragment address */
270    union {
271        struct {
272            __u8     l_i_frag;    /* Fragment number */
273            __u8     l_i_fsize;   /* Fragment size */
274            __u16    i_pad1;
275            __le16   l_i_uid_high;    /* these 2 fields      */
276            __le16   l_i_gid_high;    /* were reserved2[0] */
277            __u32    l_i_reserved2;
278        } linux2;
279        struct {
280            __u8     h_i_frag;    /* Fragment number */
281            __u8     h_i_fsize;   /* Fragment size */
282            __le16   h_i_mode_high;
283            __le16   h_i_uid_high;
284            __le16   h_i_gid_high;
285            __le32   h_i_author;
286        } hurd2;
287        struct {
288            __u8     m_i_frag;    /* Fragment number */
289            __u8     m_i_fsize;   /* Fragment size */
290            __u16    m_pad1;
291            __u32    m_i_reserved2[2];
292        } masix2;
```

```
293        } osd2;                  /* OS dependent 2 */
294 };
```

- i_mode：文件类型和访问权限
- i_uid：拥有者标示符
- i_size：以字节为单位的文件长度
- i_atime：最后一次访问文件的时间
- i_ctime：索引节点最后改变的时间
- i_mtime：文件内容最后改变的时间
- i_dtime：文件删去的时间
- i_gid：用户组标识符低 16 位
- i_links_count：硬链接计数器
- i_blocks：文件的数据库数
- i_flags：文件标志
- i_block：指向数据块的指针

i_size 存放的是以字节为单位的文件的有效长度，而 i_blocks 字段存放的是已分配给文件的数据块数（以 512 字节为单位）。

i_size 和 i_blocks 的值没有必然的联系。因为一个文件总是存放在整数块中，一个非空文件至少接受一个数据块且 i_size 可能小于 512*i_blocks。另一方面，一个文件可能包含有空洞，i_size 可能大于 512*i_blocks。

i_blocks 数 inode 结构中最重要的一个字段，因为它记录了内存中文件的块号（i_blocks 的下标）与磁盘逻辑块号（i_blocks 数组中的元素）之间的对应关系。

## 3.1.4  目录项

目录项记录了文件名称与索引节点之间的关系。在 linux 下，目录项也是做为一个文件来实现的，该文件的内容是一个个 ext2_dir_entry_2 的结构。

```
549 /*
550    * The new version of the directory entry.    Since EXT2 structures are
551    * stored in intel byte order, and the name_len field could never be
552    * bigger than 255 chars, it's safe to reclaim the extra byte for the
553    * file_type field.
```

```
554   */
555 struct ext2_dir_entry_2 {
556      __le32   inode;              /* Inode number */
557      __le16   rec_len;            /* Directory entry length */
558      __u8     name_len;            /* Name length */
559      __u8     file_type;
560      char     name[EXT2_NAME_LEN];      /* File name */
561 };
```

- inode:索引节点号
- rec_len：本目录项的长度。
- name_len：文件名长度
- file_type：文件类型
- name：保存文件名的空间

因为文件名的长度不一致，所以 ext2_dir_entry_2 在硬盘上的实际大小根据目录的长度而变，rec_len 记录了本目录项的实际大小。

## 3.2 VFS接口数据结构

3.1 节中分析了 EXT2 文件系统在硬盘上的相关结构，内核如何与这些数据结构建立联系呢，这些是通过 VFS 接口数据来实现的，VFSZ 中的接口数据只存在于内存中，为了与磁盘上的数据结构区分，称磁盘上对应的数据结构在内存中的实现为内存数据结构。

## 3.2.1 EXT2 超级块对象（内存）

VFS 超级块 super_block 有一个字段 s_fs_info 指向一个文件系统信息的数据结构。对于 EXT2 文件系统，该字段指向 ext2_sb_info 类型的结构：

include/linux/ext2_fs_sb.h

```
68 /*
69   * second extended-fs super-block data in memory
70   */
71 struct ext2_sb_info {
72      unsigned long s_frag_size;   /* Size of a fragment in bytes */
```

```
73      unsigned long s_frags_per_block;/* Number of fragments per block */
74      unsigned long s_inodes_per_block;/* Number of inodes per block */
75      unsigned long s_frags_per_group;/* Number of fragments in a group */
76      unsigned long s_blocks_per_group;/* Number of blocks in a group */
77      unsigned long s_inodes_per_group;/* Number of inodes in a group */
78      unsigned long s_itb_per_group;   /* Number of inode table blocks per group */
79      unsigned long s_gdb_count;   /* Number of group descriptor blocks */
80      unsigned long s_desc_per_block; /* Number of group descriptors per block */
81      unsigned long s_groups_count;    /* Number of groups in the fs */
82      unsigned long s_overhead_last;   /* Last calculated overhead */
83      unsigned long s_blocks_last;      /* Last seen block count */
84      struct buffer_head * s_sbh; /* Buffer containing the super block */
85      struct ext2_super_block * s_es; /* Pointer to the super block in the buffer */
86      struct buffer_head ** s_group_desc;
87      unsigned long   s_mount_opt;
88      unsigned long s_sb_block;
89      uid_t s_resuid;
90      gid_t s_resgid;
91      unsigned short s_mount_state;
92      unsigned short s_pad;
93      int s_addr_per_block_bits;
94      int s_desc_per_block_bits;
95      int s_inode_size;
96      int s_first_ino;
97      spinlock_t s_next_gen_lock;
98      u32 s_next_generation;
99      unsigned long s_dir_count;
100     u8 *s_debts;
101     struct percpu_counter s_freeblocks_counter;
```

```
102        struct percpu_counter s_freeinodes_counter;

103        struct percpu_counter s_dirs_counter;

104        struct blockgroup_lock *s_blockgroup_lock;

105        /* root of the per fs reservation window tree */

106        spinlock_t s_rsv_window_lock;

107        struct rb_root s_rsv_window_root;

108        struct ext2_reserve_window_node s_rsv_window_head;

109 };
```

同磁盘超级块的字段

s_indoes_per_block：每块可以包含的 inode 数目

s_itb_per_group：块组中 indoe 表占用的块数。

s_gdb_count：块组中组描述符占用的块数。

s_desc_per_block：每块中可以包含的组描述符数目

s_groups_count：文件系统中的块组数

s_sbh：指向包含超级块的块缓存

s_es：指向磁盘上的超级块结构

s_group_desc：指向包含组描述符的缓存

## 3.2.2 EXT2 索引节点对象（内存）

```
13 /*
14    * second extended file system inode data in memory
15    */
16 struct ext2_inode_info {
17        __le32   i_data[15];
18        __u32    i_flags;
19        __u32    i_faddr;
20        __u8     i_frag_no;
21        __u8     i_frag_size;
22        __u16    i_state;
23        __u32    i_file_acl;
24        __u32    i_dir_acl;
25        __u32    i_dtime;
26
27        /*
28         * i_block_group is the number of the block group which contains
29         * this file's inode.   Constant across the lifetime of the inode,
30         * it is used for making block allocation decisions - we try to
31         * place a file's data blocks near its inode block, and new inodes
32         * near to their parent directory's inode.
33         */
34        __u32    i_block_group;
35
36        /* block reservation info */
37        struct ext2_block_alloc_info *i_block_alloc_info;
38
39        __u32    i_dir_start_lookup;
40 #ifdef CONFIG_EXT2_FS_XATTR
41        /*
```

```
42        * Extended attributes can be read independently of the main file
43        * data. Taking i_mutex even when reading would cause contention
44        * between readers of EAs and writers of regular file data, so
45        * instead we synchronize on xattr_sem when reading or changing
46        * EAs.
47        */
48       struct rw_semaphore xattr_sem;
49 #endif
50       rwlock_t i_meta_lock;
51
52       /*
53        * truncate_mutex is for serialising ext2_truncate() against
54        * ext2_getblock().   It also protects the internals of the inode's
55        * reservation data structures: ext2_reserve_window and
56        * ext2_reserve_window_node.
57        */
58       struct mutex truncate_mutex;
59       struct inode      vfs_inode;
60       struct list_head i_orphan;   /* unlinked but open inodes */
61 };
```

- ■ i_data：与磁盘索引节点中 i_block 对应。
- ■ i_block_group：包含该索引节点的块组号

# 4. ext2 辅助操作函数

## 4.1 ext2_fill_super函数

ext2_fill_super 函数分配必要的数据结构，并从磁盘中读入磁盘超级块对象用来初始化 super_block。

```
739 static int ext2_fill_super(struct super_block *sb, void *data, int silent)
```

```
740 {
741     struct buffer_head * bh;
742     struct ext2_sb_info * sbi;
743     struct ext2_super_block * es;
744     struct inode *root;
745     unsigned long block;
746     unsigned long sb_block = get_sb_block(&data);
```

sb_block 超级块的起始块号

```
747     unsigned long logic_sb_block;
748     unsigned long offset = 0;
749     unsigned long def_mount_opts;
750     long ret = -EINVAL;
751     int blocksize = BLOCK_SIZE;
752     int db_count;
753     int i, j;
754     __le32 features;
755     int err;
756
757     sbi = kzalloc(sizeof(*sbi), GFP_KERNEL);
```

分配 ext2 超级块内存数据结构

```
758     if (!sbi)
759         return -ENOMEM;
760
761     sbi->s_blockgroup_lock =
762         kzalloc(sizeof(struct blockgroup_lock), GFP_KERNEL);
763     if (!sbi->s_blockgroup_lock) {
764         kfree(sbi);
765         return -ENOMEM;
766     }
767     sb->s_fs_info = sbi;
```

初始化 sb 的 s_fs_info 指针

| | |
|---|---|
| 768 | sbi->s_sb_block = sb_block; |
| 769 | |
| 770 | /* |
| 771 | * See what the current blocksize for the device is, and |
| 772 | * use that as the blocksize.   Otherwise (or if the blocksize |
| 773 | * is smaller than the default) use the default. |
| 774 | * This is important for devices that have a hardware |
| 775 | * sectorsize that is larger than the default. |
| 776 | */ |
| 777 | blocksize = sb_min_blocksize(sb, BLOCK_SIZE); |
| 778 | if (!blocksize) { |
| 779 | ext2_msg(sb, KERN_ERR, "error: unable to set blocksize"); |
| 780 | goto failed_sbi; |
| 781 | } |
| 782 | |
| 783 | /* |
| 784 | * If the superblock doesn't start on a hardware sector boundary, |
| 785 | * calculate the offset. |
| 786 | */ |
| 787 | if (blocksize != BLOCK_SIZE) { |
| 788 | logic_sb_block = (sb_block*BLOCK_SIZE) / blocksize; |
| 789 | offset = (sb_block*BLOCK_SIZE) % blocksize; |
| 790 | } else { |
| 791 | logic_sb_block = sb_block; |
| 792 | } |

　　首先确定设备上记录块的大小。EXT2 文件系统的记录块一般大小是 1K 字节，但是为提高读写效率也可以采用 2K 字节或者 4K 字节。内核中有一个以主设备号为下标的指针数组 hardsect_size。如果这个数组中设置了值，且大于 BLOCK_SIZE，则以该值为准。

　　超级块通常是设备上的 1 号记录块（即第二个记录块），所以 sb_block 设置为 1，在记录块大小为 BLOCK_SIZE 的设备上其逻辑块号

logic_sb_block 也为 1。但在记录块大于 BLOCK_SIZE 的设备上，由于超级块的大小仍为 BLOCK_SIZE，就要通过计算确定其所在的记录块，以及在块内的偏移。此时超级块虽然称为"块"，但实际上自身记录块中的一部分。确定了这两个参数后，就可以通过 sb_bread 将超级块读入内存了。

| 793 | |
|---|---|
| 794 | if (!(bh = sb_bread(sb, logic_sb_block))) { |
| 795 | ext2_msg(sb, KERN_ERR, "error: unable to read superblock"); |
| 796 | goto failed_sbi; |
| 797 | } |

调用 sb_bread 在缓冲区中分配一个缓冲区和缓冲区首部。然后从磁盘读入超级块在缓冲区中。

| 802 | es = (struct ext2_super_block *) (((char *)bh->b_data) + offset); |
|---|---|

es 指向磁盘超级块在内存中的地址。

| 803 | sbi->s_es = es; |
|---|---|

设置 sbi 的 s_es 字段。

| 804 | sb->s_magic = le16_to_cpu(es->s_magic); |
|---|---|
| 805 | |
| 806 | if (sb->s_magic != EXT2_SUPER_MAGIC) |
| 807 | goto cantfind_ext2; |
| 808 | |
| 809 | /* Set defaults before we parse the mount options */ |
| 810 | def_mount_opts = le32_to_cpu(es->s_default_mount_opts); |
| 811 | if (def_mount_opts & EXT2_DEFM_DEBUG) |
| 812 | set_opt(sbi->s_mount_opt, DEBUG); |
| 813 | if (def_mount_opts & EXT2_DEFM_BSDGROUPS) |
| 814 | set_opt(sbi->s_mount_opt, GRPID); |
| 815 | if (def_mount_opts & EXT2_DEFM_UID16) |
| 816 | set_opt(sbi->s_mount_opt, NO_UID32); |
| 817 | #ifdef CONFIG_EXT2_FS_XATTR |
| 818 | if (def_mount_opts & EXT2_DEFM_XATTR_USER) |
| 819 | set_opt(sbi->s_mount_opt, XATTR_USER); |

```
820 #endif
821 #ifdef CONFIG_EXT2_FS_POSIX_ACL
822     if (def_mount_opts & EXT2_DEFM_ACL)
823         set_opt(sbi->s_mount_opt, POSIX_ACL);
824 #endif
825
826     if (le16_to_cpu(sbi->s_es->s_errors) == EXT2_ERRORS_PANIC)
827         set_opt(sbi->s_mount_opt, ERRORS_PANIC);
828     else if (le16_to_cpu(sbi->s_es->s_errors) ==
EXT2_ERRORS_CONTINUE)
829         set_opt(sbi->s_mount_opt, ERRORS_CONT);
830     else
831         set_opt(sbi->s_mount_opt, ERRORS_RO);
833     sbi->s_resuid = le16_to_cpu(es->s_def_resuid);
834     sbi->s_resgid = le16_to_cpu(es->s_def_resgid);
835
836     set_opt(sbi->s_mount_opt, RESERVATION);
837
838     if (!parse_options((char *) data, sb))
839         goto failed_mount;
840
841     sb->s_flags = (sb->s_flags & ~MS_POSIXACL) |
842         ((EXT2_SB(sb)->s_mount_opt & EXT2_MOUNT_POSIX_ACL) ?
843          MS_POSIXACL : 0);
844
845     ext2_xip_verify_sb(sb); /* see if bdev supports xip, unset
846                     EXT2_MOUNT_XIP if not */
847
848     if (le32_to_cpu(es->s_rev_level) == EXT2_GOOD_OLD_REV &&
849         (EXT2_HAS_COMPAT_FEATURE(sb, ~0U) ||
850          EXT2_HAS_RO_COMPAT_FEATURE(sb, ~0U) ||
851          EXT2_HAS_INCOMPAT_FEATURE(sb, ~0U)))
```

```
852            ext2_msg(sb, KERN_WARNING,
853                    "warning: feature flags set on rev 0 fs, "
854                    "running e2fsck is recommended");
855        /*
856         * Check feature flags regardless of the revision level, since we
857         * previously didn't change the revision level when setting the flags,
858         * so there is a chance incompat flags are set on a rev 0 filesystem.
859         */
860        features = EXT2_HAS_INCOMPAT_FEATURE(sb, ~EXT2_FEATURE_INCOMPAT_SUPP);
861        if (features) {
862            ext2_msg(sb, KERN_ERR,    "error: couldn't mount because of "
863                        "unsupported optional features (%x)",
864                le32_to_cpu(features));
865            goto failed_mount;
866        }
867        if (!(sb->s_flags & MS_RDONLY) &&
868            (features = EXT2_HAS_RO_COMPAT_FEATURE(sb, ~EXT2_FEATURE_RO_COMPAT_SUPP))){
869            ext2_msg(sb, KERN_ERR, "error: couldn't mount RDWR because of "
870                        "unsupported optional features (%x)",
871                    le32_to_cpu(features));
872            goto failed_mount;
873        }
875        blocksize = BLOCK_SIZE << le32_to_cpu(sbi->s_es->s_log_block_size);
```

计算真正的块大小。

```
876
877        if (ext2_use_xip(sb) && blocksize != PAGE_SIZE) {
878            if (!silent)
879                ext2_msg(sb, KERN_ERR,
```

| | |
|---|---|
| 880 | "error: unsupported blocksize for xip"); |
| 881 | goto failed_mount; |
| 882 | } |
| 883 | |
| 884 | /* If the blocksize doesn't match, re-read the thing.. */ |
| 885 | if (sb->s_blocksize != blocksize) { |
| 886 | brelse(bh); |
| 887 | |
| 888 | if (!sb_set_blocksize(sb, blocksize)) { |
| 889 | ext2_msg(sb, KERN_ERR, "error: blocksize is too small"); |
| 890 | goto failed_sbi; |
| 891 | } |

如果之前设置的块大小与真实的块大小不等，则重新设置
sb->s_blocksize 字段。

| | |
|---|---|
| 892 | |
| 893 | logic_sb_block = (sb_block*BLOCK_SIZE) / blocksize; |
| 894 | offset = (sb_block*BLOCK_SIZE) % blocksize; |
| 895 | bh = sb_bread(sb, logic_sb_block); |

重新读入块。

| | |
|---|---|
| 896 | if(!bh) { |
| 897 | ext2_msg(sb, KERN_ERR, "error: couldn't read" |
| 898 | "superblock on 2nd try"); |
| 899 | goto failed_sbi; |
| 900 | } |
| 901 | es = (struct ext2_super_block *) (((char *)bh->b_data) + offset); |
| 902 | sbi->s_es = es; |
| 903 | if (es->s_magic != cpu_to_le16(EXT2_SUPER_MAGIC)) { |
| 904 | ext2_msg(sb, KERN_ERR, "error: magic mismatch"); |
| 905 | goto failed_mount; |
| 906 | } |

设置相应字段。

| 907 | } |

885-907 行是对记录块大小的修正。前面已经确定了记录块的大小，但那未必是第一手资料。现在有了来至设备超级块的的信息，该信息更为准确，应该以此为准。

现在有一个问题了，既然原来的参数不对，那么根据不正确的参数读入的超级块为何是正确的呢？既然读入长度超级块时正确的，何必有重读一次呢？

上述问题的原因在于，在 sb_block 等于 1 的前提下计算出来的 logic_sb_block 和 offset 只有两种结果。当 sb->s_blocksize 大于 BLOCK_SIZE 时，logic_sb_block 总是 0 而 offset 为 BLOCK_SIZE,而与 sb->s_blocksize 的具体数值无关。当 sb->s_blocksize 等于 BLOCK_SIZE 时，logic_sb_block 为 1，offset 为 0。所以，只要记录块大小等于 BLOCK_SIZE 时超级块存放在第二块上，而在记录块大小大于 BLOCK_SIZE 时，则除将超级块的放在第二块的开头处外再在第一块中偏移为 BLOCK_SIZE 处放上一个副本，就不会报错了。这里重新读一遍，只不过是让缓冲区中含有整个记录块，而不只是超级块而已。

| 908 | |
| 909 | sb->s_maxbytes = ext2_max_size(sb->s_blocksize_bits); |
| 910 | |
| 911 | if (le32_to_cpu(es->s_rev_level) == EXT2_GOOD_OLD_REV) { |
| 912 | sbi->s_inode_size = EXT2_GOOD_OLD_INODE_SIZE; |
| 913 | sbi->s_first_ino = EXT2_GOOD_OLD_FIRST_INO; |
| 914 | } else { |
| 915 | sbi->s_inode_size = le16_to_cpu(es->s_inode_size); |
| 916 | sbi->s_first_ino = le32_to_cpu(es->s_first_ino); |
| 917 | if ((sbi->s_inode_size < EXT2_GOOD_OLD_INODE_SIZE) || |
| 918 | !is_power_of_2(sbi->s_inode_size) || |
| 919 | (sbi->s_inode_size > blocksize)) { |
| 920 | ext2_msg(sb, KERN_ERR, |
| 921 | "error: unsupported inode size: %d", |
| 922 | sbi->s_inode_size); |
| 923 | goto failed_mount; |
| 924 | } |

```
925        }
926
927        sbi->s_frag_size = EXT2_MIN_FRAG_SIZE <<
928                        le32_to_cpu(es->s_log_frag_size);
929        if (sbi->s_frag_size == 0)
930            goto cantfind_ext2;
931        sbi->s_frags_per_block = sb->s_blocksize / sbi->s_frag_size;
932
933        sbi->s_blocks_per_group = le32_to_cpu(es->s_blocks_per_group);
934        sbi->s_frags_per_group = le32_to_cpu(es->s_frags_per_group);
935        sbi->s_inodes_per_group = le32_to_cpu(es->s_inodes_per_group);
936
937        if (EXT2_INODE_SIZE(sb) == 0)
938            goto cantfind_ext2;
939        sbi->s_inodes_per_block = sb->s_blocksize /
EXT2_INODE_SIZE(sb);
940        if (sbi->s_inodes_per_block == 0 || sbi->s_inodes_per_group == 0)
941            goto cantfind_ext2;
942        sbi->s_itb_per_group = sbi->s_inodes_per_group /
943                        sbi->s_inodes_per_block;
944        sbi->s_desc_per_block = sb->s_blocksize /
945                        sizeof (struct ext2_group_desc);
946        sbi->s_sbh = bh;
```

设置 s_sbh 指针指向超级块的缓存。

```
947        sbi->s_mount_state = le16_to_cpu(es->s_state);
948        sbi->s_addr_per_block_bits =
949            ilog2 (EXT2_ADDR_PER_BLOCK(sb));
950        sbi->s_desc_per_block_bits =
951            ilog2 (EXT2_DESC_PER_BLOCK(sb));
952
```

```c
953        if (sb->s_magic != EXT2_SUPER_MAGIC)
954            goto cantfind_ext2;
955
956        if (sb->s_blocksize != bh->b_size) {
957            if (!silent)
958                ext2_msg(sb, KERN_ERR, "error: unsupported blocksize");
959        goto failed_mount;
960    }
961
962        if (sb->s_blocksize != sbi->s_frag_size) {
963            ext2_msg(sb, KERN_ERR,
964                "error: fragsize %lu != blocksize %lu"
965                "(not supported yet)",
966                sbi->s_frag_size, sb->s_blocksize);
967        goto failed_mount;
968    }
969
970        if (sbi->s_blocks_per_group > sb->s_blocksize * 8) {
971            ext2_msg(sb, KERN_ERR,
972                "error: #blocks per group too big: %lu",
973                sbi->s_blocks_per_group);
974        goto failed_mount;
975    }
976        if (sbi->s_frags_per_group > sb->s_blocksize * 8) {
977            ext2_msg(sb, KERN_ERR,
978                "error: #fragments per group too big: %lu",
979                sbi->s_frags_per_group);
980        goto failed_mount;
981    }
982        if (sbi->s_inodes_per_group > sb->s_blocksize * 8) {
983            ext2_msg(sb, KERN_ERR,
```

| 984 | "error: #inodes per group too big: %lu", |
| 985 | sbi->s_inodes_per_group); |
| 986 | goto failed_mount; |
| 987 | } |
| 988 | |
| 989 | if (EXT2_BLOCKS_PER_GROUP(sb) == 0) |
| 990 | goto cantfind_ext2; |
| 991 | sbi->s_groups_count = ((le32_to_cpu(es->s_blocks_count) - |
| 992 | le32_to_cpu(es->s_first_data_block) - 1) |
| 993 | / EXT2_BLOCKS_PER_GROUP(sb)) + 1; |

计算所有文件系统中块组的总数目。计算公式：

（总的块数—起始的数据块块号-1）/块组中包含的块数 + 1

| 994 | db_count = (sbi->s_groups_count + EXT2_DESC_PER_BLOCK(sb) - 1) / |
| 995 | EXT2_DESC_PER_BLOCK(sb); |

db_count 表示存放 s_groups_count 个块组描述符所需要的记录块数。

| 996 | sbi->s_group_desc = kmalloc (db_count * sizeof (struct buffer_head *), GFP_KERNEL); |

初始化 s_group_sesc 函数指针。指向一个数组，数组中的每个元素都是一个类型为 struct buffer_head*的指针，数组的元素个数为 db_count。

| 997 | if (sbi->s_group_desc == NULL) { |
| 998 | ext2_msg(sb, KERN_ERR, "error: not enough memory"); |
| 999 | goto failed_mount; |
| 1000 | } |
| 1001 | bgl_lock_init(sbi->s_blockgroup_lock); |
| 1002 | sbi->s_debts = kcalloc(sbi->s_groups_count, sizeof(*sbi->s_debts), GFP_KERNEL); |
| 1003 | if (!sbi->s_debts) { |
| 1004 | ext2_msg(sb, KERN_ERR, "error: not enough memory"); |
| 1005 | goto failed_mount_group_desc; |
| 1006 | } |
| 1007 | for (i = 0; i < db_count; i++) { |

```
1008          block = descriptor_loc(sb, logic_sb_block, i);
```

得到块组的起始块号。

```
1009          sbi->s_group_desc[i] = sb_bread(sb, block);
1010          if (!sbi->s_group_desc[i]) {
1011              for (j = 0; j < i; j++)
1012                  brelse (sbi->s_group_desc[j]);
1013              ext2_msg(sb, KERN_ERR,
1014                  "error: unable to read group descriptors");
1015              goto failed_mount_group_desc;
1016          }
1017      }
```

叫块组描述符缓存指针保存在 s_group_desc 数组中。

```
1018      if (!ext2_check_descriptors (sb)) {
1019          ext2_msg(sb, KERN_ERR, "group descriptors corrupted");
1020          goto failed_mount2;
1021      }
1022      sbi->s_gdb_count = db_count;
1023      get_random_bytes(&sbi->s_next_generation, sizeof(u32));
1024      spin_lock_init(&sbi->s_next_gen_lock);
1025
1026      /* per fileystem reservation list head & lock */
1027      spin_lock_init(&sbi->s_rsv_window_lock);
1028      sbi->s_rsv_window_root = RB_ROOT;
1029      /*
1030       * Add a single, static dummy reservation to the start of the
1031       * reservation window list --- it gives us a placeholder for
1032       * append-at-start-of-list which makes the allocation logic
1033       * _much_ simpler.
1034       */
1035      sbi->s_rsv_window_head.rsv_start =
EXT2_RESERVE_WINDOW_NOT_ALLOCATED;
```

```
1036        sbi->s_rsv_window_head.rsv_end =
EXT2_RESERVE_WINDOW_NOT_ALLOCATED;
1037        sbi->s_rsv_window_head.rsv_alloc_hit = 0;
1038        sbi->s_rsv_window_head.rsv_goal_size = 0;
1039        ext2_rsv_window_add(sb, &sbi->s_rsv_window_head);
1040
1041        err = percpu_counter_init(&sbi->s_freeblocks_counter,
1042                    ext2_count_free_blocks(sb));
1043        if (!err) {
1044            err = percpu_counter_init(&sbi->s_freeinodes_counter,
1045                    ext2_count_free_inodes(sb));
1046        }
1047        if (!err) {
1048            err = percpu_counter_init(&sbi->s_dirs_counter,
1049                    ext2_count_dirs(sb));
1050        }
1051        if (err) {
1052            ext2_msg(sb, KERN_ERR, "error: insufficient memory");
1053            goto failed_mount3;
1054        }
1055        /*
1056         * set up enough so that it can read an inode
1057         */
```

超级块已安装了足够的信息，可以读取节点了。

```
1058        sb->s_op = &ext2_sops;
1059        sb->s_export_op = &ext2_export_ops;
1060        sb->s_xattr = ext2_xattr_handlers;
1061        root = ext2_iget(sb, EXT2_ROOT_INO);
```

读取根节点。根节点的索引节点号是固定的，在 EXT2 中为
EXT2_ROOT_INO（2）。

```
1062        if (IS_ERR(root)) {
```

```
1063            ret = PTR_ERR(root);
1064            goto failed_mount3;
1065        }
1066        if (!S_ISDIR(root->i_mode) || !root->i_blocks || !root->i_size) {
1067            iput(root);
1068            ext2_msg(sb, KERN_ERR, "error: corrupt root inode, run
e2fsck");
1069            goto failed_mount3;
1070        }
```

根索引节点必须表示满足如下条件：

- 必须是一个目录索引
- 文件的块数不能为 0
- 文件大小不能为 0

```
1071
1072        sb->s_root = d_alloc_root(root);
```

初始化根目录项。调用 d_alloc_root 函数，完成根目录项的设置。

```
1096 /**
1097  * d_alloc_root - allocate root dentry
1098  * @root_inode: inode to allocate the root for
1099  *
1100  * Allocate a root ("/") dentry for the inode given. The inode is
1101  * instantiated and returned. %NULL is returned if there is insufficient
1102  * memory or the inode passed is %NULL.
1103  */
1104
1105 struct dentry * d_alloc_root(struct inode * root_inode)
1106 {
1107        struct dentry *res = NULL;
1108
1109        if (root_inode) {
1110            static const struct qstr name = { .name = "/", .len = 1 };
```

设置根目录的名称为"/"，名称长度为 1。

```
1111
1112         res = d_alloc(NULL, &name);
1113         if (res) {
1114             res->d_sb = root_inode->i_sb;
1115             res->d_parent = res;
```

根目录的父目录指向自己。

```
1116             d_instantiate(res, root_inode);
```

安装根目录。

```
1117         }
1118     }
1119     return res;
1120 }
1121 EXPORT_SYMBOL(d_alloc_root);
```

```
1073     if (!sb->s_root) {
1074         iput(root);
1075         ext2_msg(sb, KERN_ERR, "error: get root inode failed");
1076         ret = -ENOMEM;
1077         goto failed_mount3;
1078     }
1079     if (EXT2_HAS_COMPAT_FEATURE(sb,
EXT3_FEATURE_COMPAT_HAS_JOURNAL))
1080         ext2_msg(sb, KERN_WARNING,
1081             "warning: mounting ext3 filesystem as ext2");
1082     ext2_setup_super (sb, es, sb->s_flags & MS_RDONLY);
1083     return 0;
1084
1085 cantfind_ext2:
1086     if (!silent)
```

```
1087        ext2_msg(sb, KERN_ERR,
1088            "error: can't find an ext2 filesystem on dev %s.",
1089            sb->s_id);
1090    goto failed_mount;
1091 failed_mount3:
1092    percpu_counter_destroy(&sbi->s_freeblocks_counter);
1093    percpu_counter_destroy(&sbi->s_freeinodes_counter);
1094    percpu_counter_destroy(&sbi->s_dirs_counter);
1095 failed_mount2:
1096    for (i = 0; i < db_count; i++)
1097        brelse(sbi->s_group_desc[i]);
1098 failed_mount_group_desc:
1099    kfree(sbi->s_group_desc);
1100    kfree(sbi->s_debts);
1101 failed_mount:
1102    brelse(bh);
1103 failed_sbi:
1104    sb->s_fs_info = NULL;
1105    kfree(sbi->s_blockgroup_lock);
1106    kfree(sbi);
1107    return ret;
1108 }
```

总结一下 ext2_fill_super 所作的工作如下：

- 分配一个 ext2_sb_info（ext2 超级块内存对象），将其地址保存到 sb 的 s_fs_info 字段中。

- 调用 sb_bread 在缓冲区页中分配一个缓冲区和缓冲区首部，然后从磁盘读入超级块存放在缓存中。

- 根据超级块的内容，调整之前设置的块记录大小，有必要的话重新读入超级块所在的记录块。

- 分配一个数组用来存放缓冲区首部指针，这些指针用来指向内存中的组描述符。该数组存放在 s_group_desc 中。

- 重复调用 sb_bread 分配缓冲区，从磁盘读入包含 EXT2 组描述符的块，把缓冲区首部地址存放在上一部得到的 s_group_sesc 数组中。
- 安装 sb 中的 s_op，s_export_op 指针，准备好根目录，分配一个索引节点和目录项对象。

## 4.2 inode的分配 ext2_new_inode

ext2 函数创建目标文件在存储设备上的索引节点和在内存中的 inode 结构。

/fs/ext2/ialloc.c

```
438 struct inode *ext2_new_inode(struct inode *dir, int mode)
439 {
440     struct super_block *sb;
441     struct buffer_head *bitmap_bh = NULL;
442     struct buffer_head *bh2;
443     int group, i;
444     ino_t ino = 0;
445     struct inode * inode;
446     struct ext2_group_desc *gdp;
447     struct ext2_super_block *es;
448     struct ext2_inode_info *ei;
449     struct ext2_sb_info *sbi;
450     int err;
451
452     sb = dir->i_sb;
453     inode = new_inode(sb);
```

调用如下的 new_inode 函数，分配一个 inode 结构。

```
647 struct inode *new_inode(struct super_block *sb)
648 {
649     /*
650      * On a 32bit, non LFS stat() call, glibc will generate an EOVERFLOW
```

| | |
|---|---|
| 651 |    * error if st_ino won't fit in target struct field. Use 32bit counter |
| 652 |    * here to attempt to avoid that. |
| 653 |    */ |
| 654 | static unsigned int last_ino; |
| 655 | struct inode *inode; |
| 656 | |
| 657 | spin_lock_prefetch(&inode_lock); |
| 658 | |
| 659 | inode = alloc_inode(sb); |
| 660 | if (inode) { |
| 661 |    spin_lock(&inode_lock); |
| 662 |    __inode_add_to_lists(sb, NULL, inode); |
| 663 |    inode->i_ino = ++last_ino; |

设置 inode 的索引节点号

| | |
|---|---|
| 664 |    inode->i_state = 0; |
| 665 |    spin_unlock(&inode_lock); |
| 666 | } |
| 667 | return inode; |
| 668 | } |

现代的块设备通常都是很大的。为了提高访问效率，就把存储介质划分为许多"块组"。一般来说，文件就应该与其所在目录存储在同一块组中，这样可以提高访问效率。另一方面，文件的内容和文件的索引节点也应该存储在同一块组中，所以在创建文件系统（格式化）时已经注意到了每个块组在索引节点和记录数块数量之间的比例。这个比例来自于统计信息，取决于平均的文件大小。此外，每一个块组中平均有多少个目录，也就是所每个目录平均有多少文件，大致也有个比例。所以，如果要创建的是文件，就应该首先考虑将它的索引节点分配在与其所在目录所处的块组中。如果要创建的是目录，则要考虑将来是否能够将其属下的文件都容纳在同一块组中，所以应该找一个器空闲节点的数量超过真个设备上的平均值这么一个块组，而不惜离开器父节点所在的块组。下面看这是如何实现的。

| | |
|---|---|
| 454 | if (!inode) |

```
455          return ERR_PTR(-ENOMEM);
456
457     ei = EXT2_I(inode);
458     sbi = EXT2_SB(sb);
459     es = sbi->s_es;
```

ei 指向 ext2 在内存中的 inode。

```
460     if (S_ISDIR(mode)) {
461          if (test_opt(sb, OLDALLOC))
462               group = find_group_dir(sb, dir);
463          else
464               group = find_group_orlov(sb, dir);
465     } else
466          group = find_group_other(sb, dir);
```

根据所要创建的是目录还是文件，调用不同的函数计算 inode 应该在的块组。具体的函数下面介绍。

```
467
468     if (group == -1) {
469          err = -ENOSPC;
470          goto fail;
471     }
472
473     for (i = 0; i < sbi->s_groups_count; i++) {
474          gdp = ext2_get_group_desc(sb, group, &bh2);
```

根据块组号，得到组描述符 gdp。

```
475          brelse(bitmap_bh);
476          bitmap_bh = read_inode_bitmap(sb, group);
```

得到本块组的索引位图记录块。

```
477          if (!bitmap_bh) {
478               err = -EIO;
479               goto fail;
480          }
```

```
481              ino = 0;
482
483 repeat_in_this_group:
484              ino = ext2_find_next_zero_bit((unsigned long *)bitmap_bh->b_data,
485                                      EXT2_INODES_PER_GROUP(sb), ino);
```

查找位图中第一个下一个为 0 的 bit 位，并记录其与索引位图起始位的距离 ino。

```
486              if (ino >= EXT2_INODES_PER_GROUP(sb)) {
```

如果 ino 大于库组中位图的数目

```
487                      /*
488                       * Rare race: find_group_xx() decided that there were
489                       * free inodes in this group, but by the time we tried
490                       * to allocate one, they're all gone.   This can also
491                       * occur because the counters which find_group_orlov()
492                       * uses are approximate.   So just go and search the
493                       * next block group.
494                       */
495                      if (++group == sbi->s_groups_count)
496                          group = 0;
497                      continue;
```

本组中找不到合适的，则循环前进到下一个块组。

```
498              }
499              if (ext2_set_bit_atomic(sb_bgl_lock(sbi, group),
500                                      ino, bitmap_bh->b_data)) {
```

Inode 被别人抢占，则继续到下一个块组。

```
501                      /* we lost this inode */
502                      if (++ino >= EXT2_INODES_PER_GROUP(sb)) {
503                          /* this group is exhausted, try next group */
504                          if (++group == sbi->s_groups_count)
505                              group = 0;
```

```
506                continue;
507            }
508            /* try to find free inode in the same group */
509            goto repeat_in_this_group;
510        }
511        goto got;
```

运行导次，则找到了 inode 索引。

```
512    }
513
/*
515     * Scanned all blockgroups.
516     */
517    err = -ENOSPC;
518    goto fail;
519 got:
520    mark_buffer_dirty(bitmap_bh);
```

由于修改了位图索引，必须将其缓存设置为 dirty，以便回写到磁盘。

```
521    if (sb->s_flags & MS_SYNCHRONOUS)
522        sync_dirty_buffer(bitmap_bh);
523    brelse(bitmap_bh);
524
525    ino += group * EXT2_INODES_PER_GROUP(sb) + 1;
```

计算 inode 索引的在文件系统中的全局编号。

```
526    if (ino < EXT2_FIRST_INO(sb) || ino >
le32_to_cpu(es->s_inodes_count)) {
527        ext2_error (sb, "ext2_new_inode",
528                "reserved inode or inode > inodes count - "
529                "block_group = %d,inode=%lu", group,
530                (unsigned long) ino);
531        err = -EIO;
532        goto fail;
```

| 533 | } |
| --- | --- |
| 534 | |
| 535 | percpu_counter_add(&sbi->s_freeinodes_counter, -1); |
| 536 | if (S_ISDIR(mode)) |
| 537 | percpu_counter_inc(&sbi->s_dirs_counter); |
| 538 | |
| 539 | spin_lock(sb_bgl_lock(sbi, group)); |
| 540 | le16_add_cpu(&gdp->bg_free_inodes_count, -1); |

减少块组中空闲节点的计数。

| 541 | if (S_ISDIR(mode)) { |
| --- | --- |
| 542 | if (sbi->s_debts[group] < 255) |
| 543 | sbi->s_debts[group]++; |
| 544 | le16_add_cpu(&gdp->bg_used_dirs_count, 1); |

增加块组中目录项的计数。

| 545 | } else { |
| --- | --- |
| 546 | if (sbi->s_debts[group]) |
| 547 | sbi->s_debts[group]--; |
| 548 | } |
| 549 | spin_unlock(sb_bgl_lock(sbi, group)); |
| 550 | |
| 551 | sb->s_dirt = 1; |
| 552 | mark_buffer_dirty(bh2); |

设置超级块、组描述符所在缓存为脏。

| 553 | inode->i_uid = current_fsuid(); |
| --- | --- |
| 554 | if (test_opt (sb, GRPID)) |
| 555 | inode->i_gid = dir->i_gid; |
| 556 | else if (dir->i_mode & S_ISGID) { |
| 557 | inode->i_gid = dir->i_gid; |
| 558 | if (S_ISDIR(mode)) |
| 559 | mode |= S_ISGID; |
| 560 | } else |

| | |
|---|---|
| 561 | inode->i_gid = current_fsgid(); |
| 562 | inode->i_mode = mode; |
| 563 | |
| 564 | inode->i_ino = ino; |
| 565 | inode->i_blocks = 0; |

设置 inode 索引节点号，文件的块数。

| | |
|---|---|
| 566 | inode->i_mtime = inode->i_atime = inode->i_ctime = CURRENT_TIME_SEC; |
| 567 | memset(ei->i_data, 0, sizeof(ei->i_data)); |

设置 i_data（记录文件块号的数组）为空。

| | |
|---|---|
| 568 | ei->i_flags = |
| 569 | ext2_mask_flags(mode, EXT2_I(dir)->i_flags & EXT2_FL_INHERITED); |
| 570 | ei->i_faddr = 0; |
| 571 | ei->i_frag_no = 0; |
| 572 | ei->i_frag_size = 0; |
| 573 | ei->i_file_acl = 0; |
| 574 | ei->i_dir_acl = 0; |
| 575 | ei->i_dtime = 0; |
| 576 | ei->i_block_alloc_info = NULL; |
| 577 | ei->i_block_group = group; |
| 578 | ei->i_dir_start_lookup = 0; |
| 579 | ei->i_state = EXT2_STATE_NEW; |
| 580 | ext2_set_inode_flags(inode); |
| 581 | spin_lock(&sbi->s_next_gen_lock); |
| 582 | inode->i_generation = sbi->s_next_generation++; |
| 583 | spin_unlock(&sbi->s_next_gen_lock); |
| 584 | if (insert_inode_locked(inode) < 0) { |
| 585 | err = -EINVAL; |
| 586 | goto fail_drop; |
| 587 | } |

```
588
589     dquot_initialize(inode);
590     err = dquot_alloc_inode(inode);
591     if (err)
592         goto fail_drop;
593
594     err = ext2_init_acl(inode, dir);
595     if (err)
596         goto fail_free_drop;
597
598     err = ext2_init_security(inode,dir);
599     if (err)
600         goto fail_free_drop;
601
602     mark_inode_dirty(inode);
603     ext2_debug("allocating inode %lu\n", inode->i_ino);
604     ext2_preread_inode(inode);
605     return inode;
```

返回 inode。

```
606
607 fail_free_drop:
608     dquot_free_inode(inode);
609
610 fail_drop:
611     dquot_drop(inode);
612     inode->i_flags |= S_NOQUOTA;
613     inode->i_nlink = 0;
614     unlock_new_inode(inode);
615     iput(inode);
616     return ERR_PTR(err);
617
```

```
618 fail:

619     make_bad_inode(inode);

620     iput(inode);

621     return ERR_PTR(err);

622 }
```

## 4.2.1 ext2_get_group_desc

根据块组号的到块组描述符结构。

```
39 struct ext2_group_desc * ext2_get_group_desc(struct super_block * sb,

40                          unsigned int block_group,

41                          struct buffer_head ** bh)
```

参数

- sb：超级块

- block_group：块组号

- buffer_head：bh（out 参数），记录 block_group 对应的组描述符
  的缓存指针。

```
42 {

43     unsigned long group_desc;

44     unsigned long offset;

45     struct ext2_group_desc * desc;

46     struct ext2_sb_info *sbi = EXT2_SB(sb);

47

48     if (block_group >= sbi->s_groups_count) {

49         ext2_error (sb, "ext2_get_group_desc",

50                 "block_group >= groups_count - "

51                 "block_group = %d, groups_count = %lu",

52                 block_group, sbi->s_groups_count);

53

54         return NULL;

55     }

56
```

```
57        group_desc = block_group >> EXT2_DESC_PER_BLOCK_BITS(sb);
```

记录块的索引

```
58        offset = block_group & (EXT2_DESC_PER_BLOCK(sb) - 1);
```

记录块中的偏移

```
59        if (!sbi->s_group_desc[group_desc]) {
60            ext2_error (sb, "ext2_get_group_desc",
61                        "Group descriptor not loaded - "
62                        "block_group = %d, group_desc = %lu, desc = %lu",
63                         block_group, group_desc, offset);
64            return NULL;
65        }
66
67        desc = (struct ext2_group_desc *)
sbi->s_group_desc[group_desc]->b_data;
68        if (bh)
69            *bh = sbi->s_group_desc[group_desc];
70        return desc + offset;
71 }
```

## 4.2.2  read_inode_bitmap

根据块组号读取该块组中索引节点位图所在的记录块。

```
45 static struct buffer_head *
46 read_inode_bitmap(struct super_block * sb, unsigned long block_group)
47 {
48      struct ext2_group_desc *desc;
49      struct buffer_head *bh = NULL;
50
51      desc = ext2_get_group_desc(sb, block_group, NULL);
```

得到组描述符

```
52      if (!desc)
```

| | |
|---|---|
| 53 | goto error_out; |
| 54 | |
| 55 | bh = sb_bread(sb, le32_to_cpu(desc->bg_inode_bitmap)); |

根据组描述符中 bg_inode_bitmap 字段记录的块号，读取索引节点位图记录块。

| | |
|---|---|
| 56 | if (!bh) |
| 57 | ext2_error(sb, "read_inode_bitmap", |
| 58 | "Cannot read inode bitmap - " |
| 59 | "block_group = %lu, inode_bitmap = %u", |
| 60 | block_group, le32_to_cpu(desc->bg_inode_bitmap)); |
| 61 error_out: | |
| 62 | return bh; |
| 63 } | |

## 4.2.3  find_group_orlov

如果创建的是目录，则通过 find_group_orlov 函数查找分配索引节点的块组号。

| | |
|---|---|
| 266 static int find_group_orlov(struct super_block *sb, struct inode *parent) | |
| 267 { | |
| 268 | int parent_group = EXT2_I(parent)->i_block_group; |
| 269 | struct ext2_sb_info *sbi = EXT2_SB(sb); |
| 270 | struct ext2_super_block *es = sbi->s_es; |
| 271 | int ngroups = sbi->s_groups_count; |
| 272 | int inodes_per_group = EXT2_INODES_PER_GROUP(sb); |
| 273 | int freei; |
| 274 | int avefreei; |
| 275 | int free_blocks; |
| 276 | int avefreeb; |
| 277 | int blocks_per_dir; |
| 278 | int ndirs; |
| 279 | int max_debt, max_dirs, min_blocks, min_inodes; |

```
280        int group = -1, i;
281        struct ext2_group_desc *desc;
282
283        freei =
percpu_counter_read_positive(&sbi->s_freeinodes_counter);
284        avefreei = freei / ngroups;
```

avefreei：每组的平均的空闲 inode

```
285        free_blocks =
percpu_counter_read_positive(&sbi->s_freeblocks_counter);
286        avefreeb = free_blocks / ngroups;
```

avefreeb：每组平均空闲块数

```
287        ndirs = percpu_counter_read_positive(&sbi->s_dirs_counter);
```

```
288
289        if ((parent == sb->s_root->d_inode) ||
290            (EXT2_I(parent)->i_flags & EXT2_TOPDIR_FL)) {
```

父目录为根节点：以文件系统根 root 为父目录的目录项应该分散在各个块组中。

```
291            struct ext2_group_desc *best_desc = NULL;
292            int best_ndir = inodes_per_group;
293            int best_group = -1;
294
295            get_random_bytes(&group, sizeof(group));
296            parent_group = (unsigned)group % ngroups;
297            for (i = 0; i < ngroups; i++) {
298                group = (parent_group + i) % ngroups;
299                desc = ext2_get_group_desc (sb, group, NULL);
300                if (!desc || !desc->bg_free_inodes_count)
301                    continue;
302                if (le16_to_cpu(desc->bg_used_dirs_count) >= best_ndir)
303                    continue;
```

```
304            if (le16_to_cpu(desc->bg_free_inodes_count) < avefreei)
305                continue;
306            if (le16_to_cpu(desc->bg_free_blocks_count) < avefreeb)
307                continue;
308            best_group = group;
309            best_ndir = le16_to_cpu(desc->bg_used_dirs_count);
310            best_desc = desc;
311        }
```

遍历块组，在其中查找一个：

- 空闲的索引节点大于平均索引节点
- 空闲的块数比平均的块数大。
- 已经使用的目录数比平均的目录树小。

```
312        if (best_group >= 0) {
313            desc = best_desc;
314            group = best_group;
315            goto found;
316        }
317        goto fallback;
318    }
319
320    if (ndirs == 0)
321        ndirs = 1;   /* percpu_counters are approximate... */
322
323    blocks_per_dir = (le32_to_cpu(es->s_blocks_count)-free_blocks) /
ndirs;
324
325    max_dirs = ndirs / ngroups + inodes_per_group / 16;
326    min_inodes = avefreei - inodes_per_group / 4;
327    min_blocks = avefreeb - EXT2_BLOCKS_PER_GROUP(sb) / 4;
328
329    max_debt = EXT2_BLOCKS_PER_GROUP(sb) / max(blocks_per_dir,
BLOCK_COST);
```

```
330      if (max_debt * INODE_COST > inodes_per_group)
331          max_debt = inodes_per_group / INODE_COST;
332      if (max_debt > 255)
333          max_debt = 255;
334      if (max_debt == 0)
335          max_debt = 1;
336
337      for (i = 0; i < ngroups; i++) {
338          group = (parent_group + i) % ngroups;
339          desc = ext2_get_group_desc (sb, group, NULL);
340          if (!desc || !desc->bg_free_inodes_count)
341              continue;
342          if (sbi->s_debts[group] >= max_debt)
343              continue;
344          if (le16_to_cpu(desc->bg_used_dirs_count) >= max_dirs)
345              continue;
346          if (le16_to_cpu(desc->bg_free_inodes_count) < min_inodes)
347              continue;
348          if (le16_to_cpu(desc->bg_free_blocks_count) < min_blocks)
349              continue;
350          goto found;
351      }
```

遍历块组，找一个：

- 负债小
- 使用的目录少
- 空闲节点大于 min_inodes
- 空闲块数大于 min_blocks

```
352
353 fallback:
354      for (i = 0; i < ngroups; i++) {
```

```
355          group = (parent_group + i) % ngroups;

356          desc = ext2_get_group_desc (sb, group, NULL);

357          if (!desc || !desc->bg_free_inodes_count)

358              continue;

359          if (le16_to_cpu(desc->bg_free_inodes_count) >= avefreei)

360              goto found;

361      }

362
```

退一步计划，从包含父目录中的块组开始，选择第一个满足：空闲索引节点数比平均空闲索引节点数大的块组。

```
363      if (avefreei) {

364          /*

365           * The free-inodes counter is approximate, and for really small

366           * filesystems the above test can fail to find any blockgroups

367           */

368          avefreei = 0;

369          goto fallback;

370      }
```

如果退一步计划失败，且 avefreei 不为 0，则设置 avefreei 为 0，跳转到退一步计划中。

```
371

372      return -1;

373

374 found:

375      return group;

376 }
```

## 4.2.4 find_group_other

如果创建的是文件，则调用 find_group_other 来计算块组号。

```
378 static int find_group_other(struct super_block *sb, struct inode *parent)

379 {
```

```
380        int parent_group = EXT2_I(parent)->i_block_group;
381        int ngroups = EXT2_SB(sb)->s_groups_count;
382        struct ext2_group_desc *desc;
383        int group, i;
384
385        /*
386         * Try to place the inode in its parent directory
387         */
388        group = parent_group;
389        desc = ext2_get_group_desc (sb, group, NULL);
390        if (desc && le16_to_cpu(desc->bg_free_inodes_count) &&
391                le16_to_cpu(desc->bg_free_blocks_count))
392            goto found;
```

父目录中查找。

```
393
394        /*
395         * We're going to place this inode in a different blockgroup from its
396         * parent.   We want to cause files in a common directory to all land in
397         * the same blockgroup.   But we want files which are in a different
398         * directory which shares a blockgroup with our parent to land in a
399         * different blockgroup.
400         *
401         * So add our directory's i_ino into the starting point for the hash.
402         */
403        group = (group + parent->i_ino) % ngroups;
404
405        /*
406         * Use a quadratic hash to find a group with a free inode and some
407         * free blocks.
408         */
409        for (i = 1; i < ngroups; i <<= 1) {
```

```
410            group += i;
411            if (group >= ngroups)
412                group -= ngroups;
413            desc = ext2_get_group_desc (sb, group, NULL);
414            if (desc && le16_to_cpu(desc->bg_free_inodes_count) &&
415                    le16_to_cpu(desc->bg_free_blocks_count))
416                goto found;
417        }
```

如果父目录中没有找到，则使用指数跳跃方式查找。

```
418
419     /*
420      * That failed: try linear search for a free inode, even if that group
421      * has no free blocks.
422      */
423     group = parent_group;
424     for (i = 0; i < ngroups; i++) {
425            if (++group >= ngroups)
426                group = 0;
427            desc = ext2_get_group_desc (sb, group, NULL);
428            if (desc && le16_to_cpu(desc->bg_free_inodes_count))
429                goto found;
430        }
```

线性查找。

```
431
432     return -1;
433
434 found:
435     return group;
436 }
```

## 4.3 ext2_get_block 数据块分配

当内核要分配一个数据块来保存 EXT2 普通文件的数据时，就调用 ext2_get_block 函数。

### 4.3.1 文件内块号到磁盘设备上块号的映射

从文件内块号到设备上块号的映射，最简单最迅速的方法是使用一个以文件内块号为下标的线性数组，并且将数组置与索引节点 inode 结构中。这样就需要很大的数组，从而是索引节点和 inode 的结构变大。

另一种方法是使用间接寻址，也就是将上述的数组分块放在设备上本来可用来存储数据的若干记录块中，而将这些记录块的块号放在索引节点 inode 结构中。这些记录块虽然在设备上的数据区中，却并不构成文件的本身内容，而只是一些管理信息。由于索引节点和（inode 结构）应该是固定大小的，所以当文件较大时还要将这种间接寻址的结构框架做成树状或链状，这样才能随着文件本身的大小而扩展器容量，但是这种方式解决了容量的问题，但是降低了运行效率。

基于上述考虑，UNIX 早期的文件系统采用了一种折衷的方法，使用直接与间接相结合。基本思想是将文件的记录块分为几个部分来实现。

第一部分是个以文件内块号为下标的数组，这是采用直接映射的部分，对于较小的文件这一部分就够用了。由于根据文件内块号就可以在 inode 结构的数组中直接找到相应的设备上的块号，所以效率高。

至于较大的文件，开头的部分可以直接映射，但是当文件超出这部分的容量时，超出的那部分就必须采用间接寻址了。

EXT2 文件系统的直接映射的这部分为 12 个记录块。前面看到 ext2_inode_info 结构中，有个大小为 15 的整型数组 i_data，其开头 12 个元素即用于此目的。当文件大小超过这部分内容时，数组中第 13 个元素指向一个记录块，这个记录块的内容也是一个整型数组，其中的每一个元素都指向一个设备上的记录块。第 14 个元素用于二次间接寻址，第 15 个元素用于三次间接寻址。

多重间接示意图


## 4.3.2  ext2_get_block

```
720 int ext2_get_block(struct inode *inode, sector_t iblock, struct buffer_head
*bh_result, int create)
```

参数

lblock：所处理的记录块在文件中的逻辑块号

create：是否需要创建的标志

bh_result：目标块的缓存区头。

```
721 {
722        unsigned max_blocks = bh_result->b_size >> inode->i_blkbits;
```

计算最大的块数

```
723        int ret = ext2_get_blocks(inode, iblock, max_blocks,
724                        bh_result, create);
```

调用 ext2_get_blocks 实现具体的功能。

```
725     if (ret > 0) {
726         bh_result->b_size = (ret << inode->i_blkbits);
727         ret = 0;
728     }
729     return ret;
730
731 }
```

### ext2_get_blocks

```
558 /*
559  * Allocation strategy is simple: if we have to allocate something, we will
560  * have to go the whole way to leaf. So let's do it before attaching anything
561  * to tree, set linkage between the newborn blocks, write them if sync is
562  * required, recheck the path, free and repeat if check fails, otherwise
563  * set the last missing link (that will protect us from any truncate-generated
564  * removals - all blocks on the path are immune now) and possibly force the
565  * write on the parent block.
566  * That has a nice additional property: no special recovery from the failed
567  * allocations is needed - we simply release blocks and do not touch anything
568  * reachable from inode.
569  *
570  * `handle' can be NULL if create == 0.
571  *
572  * return > 0, # of blocks mapped or allocated.
573  * return = 0, if plain lookup failed.
574  * return < 0, error case.
```

```
575    */
576 static int ext2_get_blocks(struct inode *inode,
577                    sector_t iblock, unsigned long maxblocks,
578                    struct buffer_head *bh_result,
579                    int create)
580 {
581        int err = -EIO;
582        int offsets[4];
583        Indirect chain[4];
584        Indirect *partial;
585        ext2_fsblk_t goal;
586        int indirect_blks;
587        int blocks_to_boundary = 0;
588        int depth;
589        struct ext2_inode_info *ei = EXT2_I(inode);
590        int count = 0;
591        ext2_fsblk_t first_block = 0;
592
593        depth =
ext2_block_to_path(inode,iblock,offsets,&blocks_to_boundary);
```

ext2_block_to_path 根据文件中的逻辑块号，完成两项任务：

- 计算映射的深度。
- 计算每一层映射中使用的位移量，即数组的下标，保存在 offsets 数组中。

```
594
595        if (depth == 0)
596            return (err);
```

如果 ext2_block_to_path 返回 0，表示出错。

```
597
598        partial = ext2_get_branch(inode, depth, offsets, chain, &err);
```

ext2_get_branch 从磁盘上逐层读入用于间接映射的记录块。

ext2_get_branch 返回值有两种情况：

- 如果顺利的完成了映射则返回值为 NULL。

- 如果在某一层上发现映射表内的相应项为 0，则说明这个表项（记录块）原来不存在，现在因为写操作而需要扩充文件的大小。此时返回的 Indirect 结构的指针，表示映射在此处"断裂"了。

| | |
|---|---|
| 599 | /* Simplest case - block found, no allocation needed */ |
| 600 | if (!partial) { |

顺利完成映射

| | |
|---|---|
| 601 | first_block = le32_to_cpu(chain[depth - 1].key); |

硬盘上起始块的块号

| | |
|---|---|
| 602 | clear_buffer_new(bh_result); /* What's this do? */ |
| 603 | count++; |
| 604 | /*map more blocks*/ |
| 605 | while (count < maxblocks && count <= blocks_to_boundary) { |
| 606 | ext2_fsblk_t blk; |
| 607 | |
| 608 | if (!verify_chain(chain, chain + depth - 1)) { |

检测多层映射之间的映射关系是否正确

| | |
|---|---|
| 609 | /* |
| 610 | * Indirect block might be removed by |
| 611 | * truncate while we were reading it. |
| 612 | * Handling of that case: forget what we've |
| 613 | * got now, go to reread. |
| 614 | */ |
| 615 | err = -EAGAIN; |
| 616 | count = 0; |
| 617 | break; |
| 618 | } |
| 619 | blk = le32_to_cpu(*(chain[depth-1].p + count)); |

获得与起始索引出相差 count 个元素的索引号。

| | |
|---|---|
| 620 | if (blk == first_block + count) |

| 621 | count++; |
|---|---|

如果数组的两个下标之差与对应的数组元素之差相等，则 count++。

| 622 | else |
|---|---|
| 623 | break; |

否则跳出循环。

| 624 | } |
|---|---|

605 到 624 层的循环，用来检测在文件中连续的块在磁盘上的逻辑块也连续的情况。

| 625 | if (err != -EAGAIN) |
|---|---|
| 626 | goto got_it; |
| 627 | } |
| 628 | |
| 629 | /* Next simple case - plain lookup or failed read of indirect block */ |
| 630 | if (!create || err == -EIO) |
| 631 | goto cleanup; |
| 632 | |
| 633 | mutex_lock(&ei->truncate_mutex); |
| 634 | /* |
| 635 | * If the indirect block is missing while we are reading |
| 636 | * the chain(ext3_get_branch() returns -EAGAIN err), or |
| 637 | * if the chain has been changed after we grab the semaphore, |
| 638 | * (either because another process truncated this branch, or |
| 639 | * another get_block allocated this branch) re-grab the chain to see if |
| 640 | * the request block has been allocated or not. |
| 641 | * |
| 642 | * Since we already block the truncate/other get_block |
| 643 | * at this point, we will have the current copy of the chain when we |
| 644 | * splice the branch into the tree. |
| 645 | */ |
| 646 | if (err == -EAGAIN || !verify_chain(chain, partial)) { |

```
647            while (partial > chain) {
648                    brelse(partial->bh);
649                    partial--;
650            }
651            partial = ext2_get_branch(inode, depth, offsets, chain, &err);
652            if (!partial) {
653                    count++;
654                    mutex_unlock(&ei->truncate_mutex);
655                    if (err)
656                        goto cleanup;
657                    clear_buffer_new(bh_result);
658                    goto got_it;
659            }
660    }
```

646-647 行：如果 chain 已经改变，需要重新调用 ext2_get_branch 来填充 chain。

```
661
662    /*
663     * Okay, we need to do block allocation.   Lazily initialize the block
664     * allocation info here if necessary
665    */
666    if (S_ISREG(inode->i_mode) && (!ei->i_block_alloc_info))
667        ext2_init_block_alloc_info(inode);
668
669    goal = ext2_find_goal(inode, iblock, partial);
```

为目标记录块分配一个建议的块号 goal

```
670
671    /* the number of blocks need to allocate for [d,t]indirect blocks */
672    indirect_blks = (chain + depth) - partial - 1;
```

中间路径需要分配的块的数目

```
673    /*
```

| 674 | * Next look up the indirect map to count the totoal number of |
| 675 | * direct blocks to allocate for this branch. |
| 676 | */ |
| 677 | count = ext2_blks_to_allocate(partial, indirect_blks, |
| 678 | maxblocks, blocks_to_boundary); |

计算直接数据块的数目

| 679 | /* |
| 680 | * XXX ???? Block out ext2_truncate while we alter the tree |
| 681 | */ |
| 682 | err = ext2_alloc_branch(inode, indirect_blks, &count, goal, |
| 683 | offsets + (partial - chain), partial); |

ext2_alloc_branch 完成如下的功能：

■ 设备上具体记录块的分配（目标记录块和可能用于间接映射的中间记录块）。

■ 映射的建立

| 684 | |
| 685 | if (err) { |
| 686 | mutex_unlock(&ei->truncate_mutex); |
| 687 | goto cleanup; |
| 688 | } |
| 689 | |
| 690 | if (ext2_use_xip(inode->i_sb)) { |
| 691 | /* |
| 692 | * we need to clear the block |
| 693 | */ |
| 694 | err = ext2_clear_xip_target (inode, |
| 695 | le32_to_cpu(chain[depth-1].key)); |
| 696 | if (err) { |
| 697 | mutex_unlock(&ei->truncate_mutex); |
| 698 | goto cleanup; |
| 699 | } |

| | |
|---|---|
| 700 | } |
| 701 | |
| 702 | ext2_splice_branch(inode, iblock, partial, indirect_blks, count); |

将断裂的键给安装上。

| | |
|---|---|
| 703 | mutex_unlock(&ei->truncate_mutex); |
| 704 | set_buffer_new(bh_result); |
| 705 got_it: | |
| 706 | map_bh(bh_result, inode->i_sb, le32_to_cpu(chain[depth-1].key)); |
| 707 | if (count > blocks_to_boundary) |
| 708 | set_buffer_boundary(bh_result); |
| 709 | err = count; |
| 710 | /* Clean up and exit */ |
| 711 | partial = chain + depth - 1;      /* the whole chain */ |
| 712 cleanup: | |
| 713 | while (partial > chain) { |
| 714 | brelse(partial->bh); |
| 715 | partial--; |
| 716 | } |
| 717 | return err; |
| 718 } | |

### 4.3.2.1  ext2_block_to_path

| | |
|---|---|
| 132 static int ext2_block_to_path(struct inode *inode, | |
| 133 | long i_block, int offsets[4], int *boundary) |
| 134 { | |
| 135 | int ptrs = EXT2_ADDR_PER_BLOCK(inode->i_sb); |
| 136 | int ptrs_bits = EXT2_ADDR_PER_BLOCK_BITS(inode->i_sb); |
| 137 | const long direct_blocks = EXT2_NDIR_BLOCKS, |
| 138 | indirect_blocks = ptrs, |
| 139 | double_blocks = (1 << (ptrs_bits * 2)); |

```c
140        int n = 0;
141        int final = 0;
142
143        if (i_block < 0) {
144                ext2_msg(inode->i_sb, KERN_WARNING,
145                        "warning: %s: block < 0", __func__);
146        } else if (i_block < direct_blocks) {
147                offsets[n++] = i_block;
148                final = direct_blocks;
149        } else if ( (i_block -= direct_blocks) < indirect_blocks) {
150                offsets[n++] = EXT2_IND_BLOCK;
151                offsets[n++] = i_block;
152                final = ptrs;
153        } else if ((i_block -= indirect_blocks) < double_blocks) {
154                offsets[n++] = EXT2_DIND_BLOCK;
155                offsets[n++] = i_block >> ptrs_bits;
156                offsets[n++] = i_block & (ptrs - 1);
157                final = ptrs;
158        } else if (((i_block -= double_blocks) >> (ptrs_bits * 2)) < ptrs) {
159                offsets[n++] = EXT2_TIND_BLOCK;
160                offsets[n++] = i_block >> (ptrs_bits * 2);
161                offsets[n++] = (i_block >> ptrs_bits) & (ptrs - 1);
162                offsets[n++] = i_block & (ptrs - 1);
163                final = ptrs;
164        } else {
165                ext2_msg(inode->i_sb, KERN_WARNING,
166                        "warning: %s: block is too big", __func__);
167        }
168        if (boundary)
169                *boundary = final - 1 - (i_block & (ptrs - 1));
170
```

```
171       return n;

172 }
```

该函数逻辑比较简单，根据文件中的块号计算在磁盘上映射的深度以及每层中对应数组的下标。

## 4.3.2.2 ext2_get_branch

在介绍该函数之前，先讲解一下对应的数据结构。

```
83 typedef struct {

84       __le32   *p;

85       __le32   key;

86       struct buffer_head *bh;

87 } Indirect;
```

Indirect 结构用来记录映射的中间层。该结构根据 offset 来填充。

p 指向本层记录块映射表中相应的表项

key 是该表项的值

bh 指向缓冲区的指针

辅助函数 add_chain 建立链接，verify_chain 检测链接是否正确。

```
89 static inline void add_chain(Indirect *p, struct buffer_head *bh, __le32 *v)

90 {

91       p->key = *(p->p = v);

92       p->bh = bh;

93 }

94

95 static inline int verify_chain(Indirect *from, Indirect *to)

96 {

97       while (from <= to && from->key == *from->p)

98           from++;

99       return (from > to);

100 }
```

```
203 static Indirect *ext2_get_branch(struct inode *inode,
```

| 204 | int depth, |
| 205 | int *offsets, |
| 206 | Indirect chain[4], |
| 207 | int *err) |
| 208 { |
| 209 | struct super_block *sb = inode->i_sb; |
| 210 | Indirect *p = chain; |
| 211 | struct buffer_head *bh; |
| 212 | |
| 213 | *err = 0; |
| 214 | /* i_data is not going away, no lock needed */ |
| 215 | add_chain (chain, NULL, EXT2_I(inode)->i_data + *offsets); |

第一层映射的填充。

| 216 | if (!p->key) |
| 217 | goto no_block; |

如果下标对应的数组元素值为 0，表示链接断开。

| 218 | while (--depth) { |
| 219 | bh = sb_bread(sb, le32_to_cpu(p->key)); |

读取中间目录的记录块

| 220 | if (!bh) |
| 221 | goto failure; |
| 222 | read_lock(&EXT2_I(inode)->i_meta_lock); |
| 223 | if (!verify_chain(chain, p)) |
| 224 | goto changed; |

如果间接映射的记录块有变，跳转到 changed。

| 225 | add_chain(++p, bh, (__le32*)bh->b_data + *++offsets); |

建立映射关系

| 226 | read_unlock(&EXT2_I(inode)->i_meta_lock); |
| 227 | if (!p->key) |
| 228 | goto no_block; |

如果对以的记录为 0，跳转到 no_block。

```
229      }
230      return NULL;
231
232 changed:
233      read_unlock(&EXT2_I(inode)->i_meta_lock);
234      brelse(bh);
235      *err = -EAGAIN;
236      goto no_block;
237 failure:
238      *err = -EIO;
239 no_block:
240      return p;
241 }
```

### 4.3.2.3  ext2_find_goal

```
299 static inline ext2_fsblk_t ext2_find_goal(struct inode *inode, long block,
300                      Indirect *partial)
301 {
302      struct ext2_block_alloc_info *block_i;
303
304      block_i = EXT2_I(inode)->i_block_alloc_info;
305
306      /*
307       * try the heuristic for sequential allocation,
308       * failing that at least try to get decent locality.
309       */
310      if (block_i && (block == block_i->last_alloc_logical_block + 1)
311          && (block_i->last_alloc_physical_block != 0)) {
312          return block_i->last_alloc_physical_block + 1;
313      }
```

参数 block 为文件内的逻辑块号，goal 用来返回所建议的设备上的块号。从本文件的角度当然希望所有的记录块在设备上都是连续的。为此目的，ext2_block_alloc_info 中的 last_alloc_logical_block 和 last_alloc_physical_block 用于次目的。last_alloc_logical_block 用于记录最后一次分配的逻辑块号，last_alloc_physical_block 用于记录设备上最后一次分配的块号。在正常的情况下对文件的扩充是顺序的，所以每次的文件内块号与前一次的连续，而理想的设备上块号也同样连续，二种平行的推进。当然，这只是从特定的文件角度提供的建议块号，能否实现还要看条件是否允许，不过内核会尽量满足要求，不能满足也会尽可能靠近建议值的块号分配。

但是，文件内逻辑块号也有可能不连续，也就是说文件的扩充是跳跃的，也的逻辑块号与文件原有的最后一个逻辑块号之间有"空洞"，这时调用 ext2_find_near 函数。

| |
|---|
| 314 |
| 315　　　return ext2_find_near(inode, partial); |

调用 ext2_find_near 来完成分配块的查找。

| |
|---|
| 316 } |


## ext2_find_near

| |
|---|
| 263 static ext2_fsblk_t ext2_find_near(struct inode *inode, Indirect *ind) |
| 264 { |
| 265　　　struct ext2_inode_info *ei = EXT2_I(inode); |
| 266　　　__le32 *start = ind->bh ? (__le32 *) ind->bh->b_data : ei->i_data; |
| 267　　　__le32 *p; |
| 268　　　ext2_fsblk_t bg_start; |
| 269　　　ext2_fsblk_t colour; |
| 270 |
| 271　　　/* Try to find previous block */ |
| 272　　　for (p = ind->p - 1; p >= start; p--) |
| 273　　　　　if (*p) |
| 274　　　　　　　return le32_to_cpu(*p); |

首先将起点 start 设置成指向当前映射表的（映射过程中首次发现断裂的那个映射表）的起点，然后在当前映射表内往回搜索。如果要分配的是空洞后面的第一个记录块，那就需要往回找到空洞之前的表项所对应的物理块号，并以此为建议块号。当然，这个物理块号已经被使用，这个要求无法满足（记住，该函数给出的只是建议块号，便于内核在该块号的附件范围查找）。内核在分配物理记录块时会在位图中从这里开始往前搜索，就近分配物理记录块。

```
275
276    /* No such thing, so let's try location of indirect block */
277    if (ind->bh)
278            return ind->bh->b_blocknr;
```

当空洞在间接映射表的开头处时，往回搜索在本映射中找不到空洞之前的表项，就以间接映射表本身所在的记录块号作为建议块号。

```
279
280    /*
281     * It is going to be refered from inode itself? OK, just put it into
282     * the same cylinder group then.
283     */
284    bg_start = ext2_group_first_block_no(inode->i_sb, ei->i_block_group);
285    colour = (current->pid % 16) *
286            (EXT2_BLOCKS_PER_GROUP(inode->i_sb) / 16);
287    return bg_start + colour;
```

空洞在文件的开头，以索引节点所在块组的第一个数据记录块作为建议号。

```
288 }
```

## 4.3.2.4  ext2_blks_to_allocate

计算数据项需要分配的块数。

```
330 static int
331 ext2_blks_to_allocate(Indirect * branch, int k, unsigned long blks,
332            int blocks_to_boundary)
```

```
333 {
334     unsigned long count = 0;
335
336     /*
337      * Simple case, [t,d]Indirect block(s) has not allocated yet
338      * then it's clear blocks on that path have not allocated
339      */
340     if (k > 0) {
```

k>0 表示需要分配中间目录块。

```
341         /* right now don't hanel cross boundary allocation */
342         if (blks < blocks_to_boundary + 1)
343             count += blks;
344         else
345             count += blocks_to_boundary + 1;
346     return count;
```

返回直接分配的块数。

```
347     }
348
349     count++;
350     while (count < blks && count <= blocks_to_boundary
351         && le32_to_cpu(*(branch[0].p + count)) == 0) {
352         count++;
```

直接分配，返回连续分配的块数。

```
353     }
354     return count;
355 }
```

## 4.3.2.5  ext2_alloc_branch

```
443 static int ext2_alloc_branch(struct inode *inode,
444                 int indirect_blks, int *blks, ext2_fsblk_t goal,
445                 int *offsets, Indirect *branch)
```

参数：

- indirect_blks：表示还有几次映射需要建立，实际上就是需要分配几个记录块。

- branch：指向 chain[]中从映射断裂后开始的那一部分。

- offsets：指向 offsets 中的相应部分

```
446 {
447     int blocksize = inode->i_sb->s_blocksize;
448     int i, n = 0;
449     int err = 0;
450     struct buffer_head *bh;
451     int num;
452     ext2_fsblk_t new_blocks[4];
```

保存分配的块号

```
453     ext2_fsblk_t current_block;
454
455     num = ext2_alloc_blocks(inode, goal, indirect_blks,
456                 *blks, new_blocks, &err);
457     if (err)
458         return err;
459
460     branch[0].key = cpu_to_le32(new_blocks[0]);
461     /*
462      * metadata blocks and data blocks are allocated.
463      */
464     for (n = 1; n <= indirect_blks;   n++) {
465         /*
466          * Get buffer_head for parent block, zero it out
467          * and set the pointer to new one, then send
468          * parent to disk.
469          */
```

| 470 | bh = sb_getblk(inode->i_sb, new_blocks[n-1]); |

在内存中分配缓存。

| 471 | branch[n].bh = bh; |
| 472 | lock_buffer(bh); |
| 473 | memset(bh->b_data, 0, blocksize); |

设置为 0

| 474 | branch[n].p = (__le32 *) bh->b_data + offsets[n]; |
| 475 | branch[n].key = cpu_to_le32(new_blocks[n]); |
| 476 | *branch[n].p = branch[n].key; |

设置 Indirect 中的各个元素。

| 477 | if ( n == indirect_blks) { |

表示直接数据块

| 478 | current_block = new_blocks[n]; |
| 479 | /* |
| 480 | * End of chain, update the last new metablock of |
| 481 | * the chain to point to the new allocated |
| 482 | * data blocks numbers |
| 483 | */ |
| 484 | for (i=1; i < num; i++) |
| 485 | *(branch[n].p + i) = cpu_to_le32(++current_block); |

设置连续的数据块

| 486 | } |
| 487 | set_buffer_uptodate(bh); |
| 488 | unlock_buffer(bh); |
| 489 | mark_buffer_dirty_inode(bh, inode); |
| 490 | /* We used to sync bh here if IS_SYNC(inode). |
| 491 | * But we now rely upon generic_write_sync() |
| 492 | * and b_inode_buffers.   But not for directories. |
| 493 | */ |
| 494 | if (S_ISDIR(inode->i_mode) && IS_DIRSYNC(inode)) |
| 495 | sync_dirty_buffer(bh); |

设置相应项的标志，将内存中的修改写回硬盘。

```
496        }
497        *blks = num;
498        return err;
499 }
```

### 4.3.2.6  ext2_alloc_blocks

```
357 /**
358  *   ext2_alloc_blocks: multiple allocate blocks needed for a branch
359  *   @indirect_blks: the number of blocks need to allocate for indirect
360  *           blocks
361  *
362  *   @new_blocks: on return it will store the new block numbers for
363  *   the indirect blocks(if needed) and the first direct block,
364  *   @blks:   on return it will store the total number of allocated
365  *       direct blocks
366  */
367 static int ext2_alloc_blocks(struct inode *inode,
368                 ext2_fsblk_t goal, int indirect_blks, int blks,
369                 ext2_fsblk_t new_blocks[4], int *err)
370 {
```

参数:

goal：建议的块号

indirect_blks：中间目录需要分配的块数

blks：数据记录块的块数

new_blocks：指向 offset 数组中链接断开处的元素。

```
371        int target, i;
372        unsigned long count = 0;
373        int index = 0;
374        ext2_fsblk_t current_block = 0;
375        int ret = 0;
```

| 376 | |
|---|---|
| 377 | /* |
| 378 | * Here we try to allocate the requested multiple blocks at once, |
| 379 | * on a best-effort basis. |
| 380 | * To build a branch, we should allocate blocks for |
| 381 | * the indirect blocks(if not allocated yet), and at least |
| 382 | * the first direct block of this branch.   That's the |
| 383 | * minimum number of blocks need to allocate(required) |
| 384 | */ |
| 385 | target = blks + indirect_blks; |

target 总共需要分配的块数。

| 386 | |
|---|---|
| 387 | while (1) { |

调用 ext2_new_blocks 循环分配需要的块数，知道满足要求为止。

| 388 | count = target; |
|---|---|

count 表示需要分配的块数

| 389 | /* allocating blocks for indirect blocks and direct blocks */ |
|---|---|
| 390 | current_block = ext2_new_blocks(inode,goal,&count,err); |

分配块数。返回值 current_block 表示分配块的起始块号。参数 count 是一个 in/out 参数，[in]count 表示需要分配的块数，[out]count 表示实际分配到的块数。

| 391 | if (*err) |
|---|---|
| 392 | goto failed_out; |
| 393 | |
| 394 | target -= count; |

target 表示剩余的需要分配的块数

| 395 | /* allocate blocks for indirect blocks */ |
|---|---|
| 396 | while (index < indirect_blks && count) { |
| 397 | new_blocks[index++] = current_block++; |
| 398 | count--; |
| 399 | } |

设置间接分配块的块号。

```
400
401        if (count > 0)
402            break;
403    }
```

count>0 表示需要分配的块数以完成，跳出循环。

```
404
405    /* save the new block number for the first direct block */
406    new_blocks[index] = current_block;
```

数据块的起始块号

```
407
408    /* total number of blocks allocated for direct blocks */
409    ret = count;
```

数据块分配的总块数。

```
410    *err = 0;
411    return ret;
412 failed_out:
413    for (i = 0; i <index; i++)
414        ext2_free_blocks(inode, new_blocks[i], 1);
415    return ret;
416 }
```

ext2_alloc_blocks 函数完成实际的分配任务，分配的最大的块的数目是 indirect_blks+blks，最小的块的数目是 indirect_blcks + 1，也就是该函数努力分配块，但是当分配不到更多的块时，至少要保证分配中间目录的块和 1 个数据块。

### 4.3.2.7 ext2_new_blocks

fs/ext2/balloc.c

分配时首先满足"客户"需要，如果所建议的记录块还空闲着就把它分配出去。否则，若果所建议的块号已经分配，就试图在它附件 32 个记录块中的范围内分配。还不行就向前在本块组的位图中搜索，先找位图中整个字节都是 8 个记录块的空闲区间，若达不到目的再降格以求。最后，如果实在找不到，则在整个设备的范围内寻找和分配。

```
1204 /*
1205  * ext2_new_blocks() -- core block(s) allocation function
1206  * @inode:        file inode
1207  * @goal:         given target block(filesystem wide)
1208  * @count:        target number of blocks to allocate
1209  * @errp:         error code
1210  *
1211  * ext2_new_blocks uses a goal block to assist allocation.   If the goal is
1212  * free, or there is a free block within 32 blocks of the goal, that block
1213  * is allocated.   Otherwise a forward search is made for a free block; within
1214  * each block group the search first looks for an entire free byte in the block
1215  * bitmap, and then for any free bit if that fails.
1216  * This function also updates quota and i_blocks field.
1217  */
1218 ext2_fsblk_t ext2_new_blocks(struct inode *inode, ext2_fsblk_t goal,
1219                 unsigned long *count, int *errp)
1220 {
1221     struct buffer_head *bitmap_bh = NULL;
1222     struct buffer_head *gdp_bh;
1223     int group_no;
1224     int goal_group;
```

```
1225        ext2_grpblk_t grp_target_blk;      /* blockgroup relative goal block */
1226        ext2_grpblk_t grp_alloc_blk;        /* blockgroup-relative allocated
block*/
1227        ext2_fsblk_t ret_block;         /* filesyetem-wide allocated block */
1228        int bgi;                      /* blockgroup iteration index */
1229        int performed_allocation = 0;
1230        ext2_grpblk_t free_blocks;    /* number of free blocks in a group */
1231        struct super_block *sb;
1232        struct ext2_group_desc *gdp;
1233        struct ext2_super_block *es;
1234        struct ext2_sb_info *sbi;
1235        struct ext2_reserve_window_node *my_rsv = NULL;
1236        struct ext2_block_alloc_info *block_i;
1237        unsigned short windowsz = 0;
1238        unsigned long ngroups;
1239        unsigned long num = *count;
1240        int ret;
1241
1242        *errp = -ENOSPC;
1243        sb = inode->i_sb;
1244        if (!sb) {
1245            printk("ext2_new_blocks: nonexistent device");
1246            return 0;
1247        }
1248
1249        /*
1250         * Check quota for allocation of this block.
1251         */
1252        ret = dquot_alloc_block(inode, num);
1253        if (ret) {
1254            *errp = ret;
1255            return 0;
```

```
1256        }
1257
1258        sbi = EXT2_SB(sb);
1259        es = EXT2_SB(sb)->s_es;
1260        ext2_debug("goal=%lu.\n", goal);
1261        /*
1262         * Allocate a block from reservation only when
1263         * filesystem is mounted with reservation(default,-o reservation), and
1264         * it's a regular file, and
1265         * the desired window size is greater than 0 (One could use ioctl
1266         * command EXT2_IOC_SETRSVSZ to set the window size to 0 to turn off
1267         * reservation on that particular file)
1268         */
1269        block_i = EXT2_I(inode)->i_block_alloc_info;
1270        if (block_i) {
1271            windowsz = block_i->rsv_window_node.rsv_goal_size;
1272            if (windowsz > 0)
1273                my_rsv = &block_i->rsv_window_node;
1274        }
1275
1276        if (!ext2_has_free_blocks(sbi)) {
1277            *errp = -ENOSPC;
1278            goto out;
1279        }
1280
1281        /*
1282         * First, test whether the goal block is free.
1283         */
1284        if (goal < le32_to_cpu(es->s_first_data_block) ||
1285            goal >= le32_to_cpu(es->s_blocks_count))
```

| 1286 | goal = le32_to_cpu(es->s_first_data_block); |
|---|---|

按需调整 goal 号

| 1287 | group_no = (goal - le32_to_cpu(es->s_first_data_block)) / |
|---|---|
| 1288 | EXT2_BLOCKS_PER_GROUP(sb); |
| 1289 | goal_group = group_no; |

计算 goal 所在的块组号。

| 1290 retry_alloc: | |
|---|---|
| 1291 | gdp = ext2_get_group_desc(sb, group_no, &gdp_bh); |

得到 goal 所在的块组的组描述符

| 1292 | if (!gdp) |
|---|---|
| 1293 | goto io_error; |
| 1294 | |
| 1295 | free_blocks = le16_to_cpu(gdp->bg_free_blocks_count); |
| 1296 | /* |
| 1297 | * if there is not enough free blocks to make a new resevation |
| 1298 | * turn off reservation for this allocation |
| 1299 | */ |
| 1300 | if (my_rsv && (free_blocks < windowsz |
| 1301 | && (free_blocks > 0) |
| 1302 | && (rsv_is_empty(&my_rsv->rsv_window))) |
| 1303 | my_rsv = NULL; |
| 1304 | |
| 1305 | if (free_blocks > 0) { |
| 1306 | grp_target_blk = ((goal - le32_to_cpu(es->s_first_data_block)) % |
| 1307 | EXT2_BLOCKS_PER_GROUP(sb)); |
| 1308 | bitmap_bh = read_block_bitmap(sb, group_no); |
| 1309 | if (!bitmap_bh) |
| 1310 | goto io_error; |
| 1311 | grp_alloc_blk = ext2_try_to_allocate_with_rsv(sb, group_no, |
| 1312 | bitmap_bh, grp_target_blk, |
| 1313 | my_rsv, &num); |

该函数进行块分配，返回值 grp_alloc_blk 表示分配的块号在本组内的序号。num 是一个 in/out 参数，[in]count 表示需要分配的块数目，[out]count 表示实际分配的块数。grp_target_blk 表示建议的分配块号。

```
1314            if (grp_alloc_blk >= 0)
1315                goto allocated;
1316        }
```

如果本组内有空闲块，则在本组内分配。

```
1317
1318    ngroups = EXT2_SB(sb)->s_groups_count;
1319    smp_rmb();
1320
1321    /*
1322     * Now search the rest of the groups.   We assume that
1323     * group_no and gdp correctly point to the last group visited.
1324     */
1325    for (bgi = 0; bgi < ngroups; bgi++) {
1326        group_no++;
1327        if (group_no >= ngroups)
1328            group_no = 0;
1329        gdp = ext2_get_group_desc(sb, group_no, &gdp_bh);
1330        if (!gdp)
1331            goto io_error;
1332
1333        free_blocks = le16_to_cpu(gdp->bg_free_blocks_count);
1334        /*
1335         * skip this group if the number of
1336         * free blocks is less than half of the reservation
1337         * window size.
1338         */
1339        if (my_rsv && (free_blocks <= (windowsz/2)))
1340            continue;
```

```
1341
1342        brelse(bitmap_bh);
1343        bitmap_bh = read_block_bitmap(sb, group_no);
1344        if (!bitmap_bh)
1345            goto io_error;
1346        /*
1347         * try to allocate block(s) from this group, without a goal(-1).
1348         */
1349        grp_alloc_blk = ext2_try_to_allocate_with_rsv(sb, group_no,
1350                        bitmap_bh, -1, my_rsv, &num);
```

该函数进行块分配，返回值 grp_alloc_blk 表示分配的块号在本组内的序号。num 是一个 in/out 参数，[in]count 表示需要分配的块数目，[out]count 表示实际分配的块数。grp_target_blk 表示建议的分配块号。

```
1351        if (grp_alloc_blk >= 0)
1352            goto allocated;
1353    }
```

本组内分配失败，则遍历所有块组分配。

```
1354    /*
1355     * We may end up a bogus ealier ENOSPC error due to
1356     * filesystem is "full" of reservations, but
1357     * there maybe indeed free blocks avaliable on disk
1358     * In this case, we just forget about the reservations
1359     * just do block allocation as without reservations.
1360     */
1361    if (my_rsv) {
1362        my_rsv = NULL;
1363        windowsz = 0;
1364        group_no = goal_group;
1365        goto retry_alloc;
1365        goto retry_alloc;
1366    }
```

| 1367 | /* No space left on the device */ |
| 1368 | *errp = -ENOSPC; |
| 1369 | goto out; |
| 1370 | |
| 1371 allocated: | |
| 1372 | |
| 1373 | ext2_debug("using block group %d(%d)\n", |
| 1374 | group_no, gdp->bg_free_blocks_count); |
| 1375 | |
| 1376 | ret_block = grp_alloc_blk + ext2_group_first_block_no(sb, group_no); |

返回分配的块的起始块号

| 1377 | |
| 1378 | if (in_range(le32_to_cpu(gdp->bg_block_bitmap), ret_block, num) || |
| 1379 | in_range(le32_to_cpu(gdp->bg_inode_bitmap), ret_block, num) || |
| 1380 | in_range(ret_block, le32_to_cpu(gdp->bg_inode_table), |
| 1381 | EXT2_SB(sb)->s_itb_per_group) || |
| 1382 | in_range(ret_block + num - 1, le32_to_cpu(gdp->bg_inode_table), |
| 1383 | EXT2_SB(sb)->s_itb_per_group)) { |
| 1384 | ext2_error(sb, "ext2_new_blocks", |
| 1385 | "Allocating block in system zone - " |
| 1386 | "blocks from "E2FSBLK", length %lu", |
| 1387 | ret_block, num); |
| 1388 | /* |
| 1389 | * ext2_try_to_allocate marked the blocks we allocated as in |
| 1390 | * use.   So we may want to selectively mark some of the blocks |
| 1391 | * as free |
| 1392 | */ |
| 1393 | goto retry_alloc; |
| 1394 | } |

块的合理性检查

```
1395
1396        performed_allocation = 1;
1397
1398        if (ret_block + num - 1 >= le32_to_cpu(es->s_blocks_count)) {
1399            ext2_error(sb, "ext2_new_blocks",
1400                    "block("E2FSBLK") >= blocks count(%d) - "
1401                    "block_group = %d, es == %p ", ret_block,
1402                le32_to_cpu(es->s_blocks_count), group_no, es);
1403            goto out;
1404        }
1405
1406        group_adjust_blocks(sb, group_no, gdp, gdp_bh, -num);
```

调整相关数据结构中空闲块的记录字段

```
1407        percpu_counter_sub(&sbi->s_freeblocks_counter, num);
1408
1409        mark_buffer_dirty(bitmap_bh);
```

设置 bitmap 对应的缓存区为脏，将数据回写到硬盘。

```
1410        if (sb->s_flags & MS_SYNCHRONOUS)
1411            sync_dirty_buffer(bitmap_bh);
1412
1413        *errp = 0;
1414        brelse(bitmap_bh);
1415        dquot_free_block(inode, *count-num);
1416        *count = num;
1417        return ret_block;
```

返回实际分配的块数目，连续块的起始块号。

```
1418
1419 io_error:
1420        *errp = -EIO;
1421 out:
1422        /*
```

```
1423        * Undo the block allocation
1424        */
1425       if (!performed_allocation)
1426            dquot_free_block(inode, *count);
1427       brelse(bitmap_bh);
1428       return 0;
1429 }
```

### 4.3.2.8 ext2_try_to_allocate_with_rsv

```
1074 /**
1075    * ext2_try_to_allocate_with_rsv()
1076    * @sb:            superblock
1077    * @group:         given allocation block group
1078    * @bitmap_bh:      bufferhead holds the block bitmap
1079    * @grp_goal:       given target block within the group
1080    * @count:         target number of blocks to allocate
1081    * @my_rsv:       reservation window
1082    *
1083    * This is the main function used to allocate a new block and its
reservation
1084    * window.
1085    *
1086    * Each time when a new block allocation is need, first try to allocate from
1087    * its own reservation.   If it does not have a reservation window, instead
of
1088    * looking for a free bit on bitmap first, then look up the reservation list to
1089    * see if it is inside somebody else's reservation window, we try to allocate
a
1090    * reservation window for it starting from the goal first. Then do the block
1091    * allocation within the reservation window.
```

```
1092   *
1093   * This will avoid keeping on searching the reservation list again and
1094   * again when somebody is looking for a free block (without
1095   * reservation), and there are lots of free blocks, but they are all
1096   * being reserved.
1097   *
1098   * We use a red-black tree for the per-filesystem reservation list.
1099   */
1100 static ext2_grpblk_t
1101 ext2_try_to_allocate_with_rsv(struct super_block *sb, unsigned int group,
1102                 struct buffer_head *bitmap_bh, ext2_grpblk_t grp_goal,
1103                 struct ext2_reserve_window_node * my_rsv,
1104                 unsigned long *count)
1105 {
1106       ext2_fsblk_t group_first_block, group_last_block;
1107       ext2_grpblk_t ret = 0;
1108       unsigned long num = *count;
1109
/*
1111           * we don't deal with reservation when
1112           * filesystem is mounted without reservation
1113           * or the file is not a regular file
1114           * or last attempt to allocate a block with reservation turned on failed
1115           */
1116       if (my_rsv == NULL) {
1117             return ext2_try_to_allocate(sb, group, bitmap_bh,
1118                             grp_goal, count, NULL);
1119       }
1120       /*
1121           * grp_goal is a group relative block number (if there is a goal)
1122           * 0 <= grp_goal < EXT2_BLOCKS_PER_GROUP(sb)
```

```
1123          * first block is a filesystem wide block number
1124          * first block is the block number of the first block in this group
1125          */
1126         group_first_block = ext2_group_first_block_no(sb, group);
1127         group_last_block = group_first_block +
(EXT2_BLOCKS_PER_GROUP(sb) - 1);
1128
1129         /*
1130          * Basically we will allocate a new block from inode's reservation
1131          * window.
1132          *
1133          * We need to allocate a new reservation window, if:
1134          * a) inode does not have a reservation window; or
1135          * b) last attempt to allocate a block from existing reservation
1136          *     failed; or
1137          * c) we come here with a goal and with a reservation window
1138          *
1139          * We do not need to allocate a new reservation window if we come
here
1140          * at the beginning with a goal and the goal is inside the window, or
1141          * we don't have a goal but already have a reservation window.
1142          * then we could go to allocate from the reservation window directly.
1143          */
1144         while (1) {
1145             if (rsv_is_empty(&my_rsv->rsv_window) || (ret < 0) ||
1146                 !goal_in_my_reservation(&my_rsv->rsv_window,
1147                             grp_goal, group, sb)) {
1148                 if (my_rsv->rsv_goal_size < *count)
1149                     my_rsv->rsv_goal_size = *count;
1150                 ret = alloc_new_reservation(my_rsv, grp_goal, sb,
1151                             group, bitmap_bh);
1152                 if (ret < 0)
```

```
1153                    break;              /* failed */
1154
1155               if (!goal_in_my_reservation(&my_rsv->rsv_window,
1156                            grp_goal, group, sb))
1157                    grp_goal = -1;
1158          } else if (grp_goal >= 0) {
1159              int curr = my_rsv->rsv_end -
1160                      (grp_goal + group_first_block) + 1;
1161
1162              if (curr < *count)
1163                  try_to_extend_reservation(my_rsv, sb,
1164                              *count - curr);
1165          }
1166
1167          if ((my_rsv->rsv_start > group_last_block) ||
1168                  (my_rsv->rsv_end < group_first_block)) {
1169              rsv_window_dump(&EXT2_SB(sb)->s_rsv_window_root, 1);
1170              BUG();
1171          }
1172          ret = ext2_try_to_allocate(sb, group, bitmap_bh, grp_goal,
1173                          &num, &my_rsv->rsv_window);
```

调用 ext2_try_to_allocate 分配块。

```
1174          if (ret >= 0) {
1175              my_rsv->rsv_alloc_hit += num;
1176              *count = num;
1177              break;              /* succeed */
1178          }
1179          num = *count;
1180      }
1181      return ret;
1182 }
```

## 4.3.2.9 ext2_try_to_allocate

```
672 static int
673 ext2_try_to_allocate(struct super_block *sb, int group,
674                 struct buffer_head *bitmap_bh, ext2_grpblk_t grp_goal,
675                 unsigned long *count,
676                 struct ext2_reserve_window *my_rsv)
677 {
678     ext2_fsblk_t group_first_block;
679         ext2_grpblk_t start, end;
680     unsigned long num = 0;
681
682     /* we do allocation within the reservation window if we have a window */
683     if (my_rsv) {
684         group_first_block = ext2_group_first_block_no(sb, group);
685         if (my_rsv->_rsv_start >= group_first_block)
686             start = my_rsv->_rsv_start - group_first_block;
687         else
688             /* reservation window cross group boundary */
689             start = 0;
690         end = my_rsv->_rsv_end - group_first_block + 1;
691         if (end > EXT2_BLOCKS_PER_GROUP(sb))
692             /* reservation window crosses group boundary */
693             end = EXT2_BLOCKS_PER_GROUP(sb);
694         if ((start <= grp_goal) && (grp_goal < end))
695             start = grp_goal;
696         else
697             grp_goal = -1;
698     } else {
```

```
699            if (grp_goal > 0)
700                 start = grp_goal;
701            else
702                 start = 0;
703            end = EXT2_BLOCKS_PER_GROUP(sb);
704       }
705
706      BUG_ON(start > EXT2_BLOCKS_PER_GROUP(sb));
707
708 repeat:
709       if (grp_goal < 0) {
710            grp_goal = find_next_usable_block(start, bitmap_bh, end);
711            if (grp_goal < 0)
712                 goto fail_access;
713            if (!my_rsv) {
714                 int i;
715
716                 for (i = 0; i < 7 && grp_goal > start &&
717                           !ext2_test_bit(grp_goal - 1,
718                                     bitmap_bh->b_data);
719                                i++, grp_goal--)
720                      ;
721            }
722       }
723      start = grp_goal;
724
725      if (ext2_set_bit_atomic(sb_bgl_lock(EXT2_SB(sb), group), grp_goal,
726                                     bitmap_bh->b_data)) {
```

设置块位图中对应的 bit 位，表示分配块号。

```
727           /*
728            * The block was allocated by another thread, or it was
```

```
729          * allocated and then freed by another thread
730          */
731         start++;
732         grp_goal++;
733         if (start >= end)
734             goto fail_access;
735         goto repeat;
736     }
737     num++;
738     grp_goal++;
739     while (num < *count && grp_goal < end
740         && !ext2_set_bit_atomic(sb_bgl_lock(EXT2_SB(sb), group),
741                     grp_goal, bitmap_bh->b_data)) {
742         num++;
743         grp_goal++;
744     }
```

设置块位图中的 bit 位表示分配块

```
745     *count = num;
746     return grp_goal - num;
747 fail_access:
748     *count = num;
749     return -1;
750 }
```

# 5. ext2 文件系统的具体实现

本小结中以 ext2 文件系统为例，具体分析文件系统的实现，主要内容包含如下几部分：

- 从文件路径到索引节点(inode)
- 文件系统的安装与拆卸
- 文件的打开与关闭

- 文件的读与写
- 文件的映射

## 5.1 文件路径到索引节点

文件打开、读写都涉及到从文件路径到索引节点的转换，所以这个部分是其他内容的基础。

ext2 文件系统在硬盘上的布局如下：

| 启动块 | 块组 0 | 块组 1 | …… | 块组 n |
|---|---|---|---|---|

每一个块组的划分如下：

| 超级块 | 组描述符 | 数据位图 | inode 位图 | inode 表 | 数据块 |
|---|---|---|---|---|---|

在 linux 下面，目录也是作为一种特定的文件来实现的。

### 5.1.1 相关数据结构

#### 5.1.1.1 ext2 文件节点ext2_inode

vfs 的 inode 结构是内存中用来表示索引节点的数据结构，它要兼容各种文件系统，与磁盘上 indoe 的索引节点结构不同。对于 ext2 文件系统而言，磁盘上 indoe 结构就是 ext2_inode。

include/linux/ext2_fs.h

```
239 /*
240    * Structure of an inode on the disk
241    */
242 struct ext2_inode {
243        __le16  i_mode;        /* File mode */
244        __le16  i_uid;         /* Low 16 bits of Owner Uid */
245        __le32  i_size;        /* Size in bytes */
246        __le32  i_atime;       /* Access time */
247        __le32  i_ctime;       /* Creation time */
248        __le32  i_mtime;       /* Modification time */
249        __le32  i_dtime;       /* Deletion Time */
250        __le16  i_gid;         /* Low 16 bits of Group Id */
251        __le16  i_links_count; /* Links count */
```

```c
252        __le32  i_blocks;     /* Blocks count */
253        __le32  i_flags;      /* File flags */
254        union {
255            struct {
256                    __le32  l_i_reserved1;
257            } linux1;
258            struct {
259                    __le32  h_i_translator;
260            } hurd1;
261            struct {
262                    __le32  m_i_reserved1;
263            } masix1;
264        } osd1;               /* OS dependent 1 */
265        __le32  i_block[EXT2_N_BLOCKS];/* Pointers to blocks */
266        __le32  i_generation;    /* File version (for NFS) */
267        __le32  i_file_acl; /* File ACL */
268        __le32  i_dir_acl;   /* Directory ACL */
269        __le32  i_faddr;      /* Fragment address */
270        union {
271            struct {
272                    __u8     l_i_frag;    /* Fragment number */
273                    __u8     l_i_fsize;   /* Fragment size */
274                    __u16    i_pad1;
275                    __le16  l_i_uid_high;    /* these 2 fields     */
276                    __le16  l_i_gid_high;    /* were reserved2[0] */
277                    __u32    l_i_reserved2;
278            } linux2;
279            struct {
280                    __u8     h_i_frag;    /* Fragment number */
281                    __u8     h_i_fsize;   /* Fragment size */
282                    __le16   h_i_mode_high;
```

```
283              __le16   h_i_uid_high;
284              __le16   h_i_gid_high;
285              __le32   h_i_author;
286          } hurd2;
287      struct {
288              __u8     m_i_frag;     /* Fragment number */
289              __u8     m_i_fsize;   /* Fragment size */
290              __u16    m_pad1;
291              __u32    m_i_reserved2[2];
292          } masix2;
293      } osd2;              /* OS dependent 2 */
294 };
```

```
159 /*
160   * Constants relative to the data blocks
161   */
162 #define EXT2_NDIR_BLOCKS          12
163 #define EXT2_IND_BLOCK            EXT2_NDIR_BLOCKS
164 #define EXT2_DIND_BLOCK           (EXT2_IND_BLOCK + 1)
165 #define EXT2_TIND_BLOCK           (EXT2_DIND_BLOCK + 1)
166 #define EXT2_N_BLOCKS             (EXT2_TIND_BLOCK + 1)
```

重点看 i_block 结构，i_block 是从文件逻辑块号到磁盘块号转换的关键点。

i_block 是一个数组，通过文件的逻辑块号作为数组下表，索引到该下表的元素，也就是逻辑块号对应的磁盘块号。（为了尽量简单，这里没有考虑间接索引的情况）。

```
549 /*
550   * The new version of the directory entry.   Since EXT2 structures are
551   * stored in intel byte order, and the name_len field could never be
552   * bigger than 255 chars, it's safe to reclaim the extra byte for the
553   * file_type field.
554   */
```

```
555 struct ext2_dir_entry_2 {
556     __le32   inode;              /* Inode number */
557     __le16   rec_len;           /* Directory entry length */
558     __u8     name_len;           /* Name length */
559     __u8     file_type;
560     char     name[EXT2_NAME_LEN];     /* File name */
561 };
```

- ■　　inode:索引节点号
- ■　　rec_len：文件实际长度
- ■　　name_len：文件名长度
- ■　　file_type：文件类型
- ■　　name：保存文件名的空间

## 5.1.2 文件路径到索引节点的实现概述

以/home/zenhumany/test.txt 为例，看看是如何通过路径名找到 test.txt 对应的索引节点的。

首先，系统先根据"/"的索引节点（根 inode，至于这一个索引节点哪来的，下面介绍）可以找到"/"在硬盘上的内容，根据索引节点中的 i_block 数组，可以找到文件在硬盘上的所有块号，也就得到了该文件的所有内容。

找到硬盘上的内容后，将其读入内存。前面说过，目录也是文件，是一种特需的文件，其中保存的是目录项。对 ext2 文件系统而言，里面保存的是类型为 ext2_dir_entry_2 的实例。通过根索引节点可以的到根文件（目录文件）的内容，在其中搜索名为 home 的目录项，可以得到 home 对应的 ext2_dir_entry_2 目录项，根据其 inode 元素，可以得到 home 对应的文件内容。

在 home 的文件内容中，查找 zenhumany 的目录项，可以得到 zenhumany 对应的 ext2_dir_entry_2 实例，根据 indoe 元素，可以得到 zenhumany 的文件内容，在其中搜索 test.txt，可以得到 test.txt 的 ext2_dir_entry_2 实例，根据 inode,可以得到 test.txt 对应的索引号，也就可以找到 test.txt 的文件内容了。

前面留了一个疑问，"/"的索引节点是如何来的。要访问一个文件，先要访问一个目录，才能根据文件名从目录中找到给文件的目录项，进而找到其 i 节点。可是目录自身也是文件，它本身的目录项又在另一个目录项中，这样一来就成先有鸡还是先有蛋的问题了。

linux 解决办法是："/"目录，或者说根设备上的根目录的目录项不在其他目录中，可以通过在一个固定的位置上或者通过一个固定的算法找到，从这个目录出发，可以找到系统中任何一个文件。linux 中每一个"文件系统"，即每一格式化成某种文件系统的存储设备上都有一个根目录，同时都有一个"超级块"，根目录的位置以及文件系统中的其他一些参数就记录在超级块中。超级块在设备上的逻辑位置都是固定的。

## 5.1.3 具体实现

fs/namei.c

```
1054 /* Returns 0 and nd will be valid on success; Retuns error, otherwise. */
1055 static int do_path_lookup(int dfd, const char *name,
1056                          unsigned int flags, struct nameidata *nd)
1057 {
1058     int retval = path_init(dfd, name, flags, nd);
1059     if (!retval)
1060         retval = path_walk(name, nd);
1061     if (unlikely(!retval && !audit_dummy_context() && nd->path.dentry &&
1062                 nd->path.dentry->d_inode))
1063         audit_inode(name, nd->path.dentry);
1064     if (nd->root.mnt) {
1065         path_put(&nd->root);
1066         nd->root.mnt = NULL;
1067     }
1068     return retval;
1069 }
1070
1071 int path_lookup(const char *name, unsigned int flags,
```

```
1072                 struct nameidata *nd)
1073 {
1074      return do_path_lookup(AT_FDCWD, name, flags, nd);
1075 }
```

由路径到索引节点实现的入口函数是 path_lookup，其调用 do_path_lookup 来实现具体功能。

do_path_lookup 调用连个函数来实现具体功能

- path_init：确定搜索路径的起始 dentry 结构。
- path_walk：根据 path_init 中提供的起始 dentry 结构，查找目录项的对应结构。

### 5.1.3.1 path_init

参数 nd 的类型为 nameidata，它是一个临时变量，用来在查找 inode 的过程中。

```
18 struct nameidata {
19      struct path path;
20      struct qstr last;
21      struct path root;
22      unsigned int      flags;
23      int        last_type;
24      unsigned      depth;
25      char *saved_names[MAX_NESTED_LINKS + 1];
26
27      /* Intent data */
28      union {
29          struct open_intent open;
30      } intent;
31 };
```

path：当前路径的 path 实例

root：进程更目录的 path 实例。

```
1002 static int path_init(int dfd, const char *name, unsigned int flags, struct
nameidata *nd)
1003 {
1004     int retval = 0;
1005     int fput_needed;
1006     struct file *file;
1007
1008     nd->last_type = LAST_ROOT; /* if there are only slashes... */
1009     nd->flags = flags;
1010     nd->depth = 0;
1011     nd->root.mnt = NULL;
1012
1013     if (*name=='/') {
1014         set_root(nd);
1015         nd->path = nd->root;
1016         path_get(&nd->root);
1017     } else if (dfd == AT_FDCWD) {
1018         struct fs_struct *fs = current->fs;
1019         read_lock(&fs->lock);
1020         nd->path = fs->pwd;
1021         path_get(&fs->pwd);
1022         read_unlock(&fs->lock);
1023     } else {
1024         struct dentry *dentry;
1025
1026         file = fget_light(dfd, &fput_needed);
1027         retval = -EBADF;
1028         if (!file)
1029             goto out_fail;
1030
1031         dentry = file->f_path.dentry;
1032
```

```
1033            retval = -ENOTDIR;
1034            if (!S_ISDIR(dentry->d_inode->i_mode))
1035                    goto fput_fail;
1036
1037            retval = file_permission(file, MAY_EXEC);
1038            if (retval)
1039                    goto fput_fail;
1040
1041            nd->path = file->f_path;
1042            path_get(&file->f_path);
1043
1044            fput_light(file, fput_needed);
1045        }
1046        return 0;
1047
1048 fput_fail:
1049        fput_light(file, fput_needed);
1050 out_fail:
1051        return retval;
1052 }
```

当路径名是以"/"开始的，则认为是从当前进程的根目录开始操作。将 nd->path 设置为 fs->root。否则，将 nd->path 设置为 fs->pwd。

### 5.1.3.2 path_walk

```
977 static int path_walk(const char *name, struct nameidata *nd)
978 {
979      struct path save = nd->path;
980      int result;
981
```

```
982        current->total_link_count = 0;

983

984        /* make sure the stuff we saved doesn't go away */

985        path_get(&save);

986

987        result = link_path_walk(name, nd);

988        if (result == -ESTALE) {

989            /* nd->path had been dropped */

990            current->total_link_count = 0;

991            nd->path = save;

992            path_get(&nd->path);

993            nd->flags |= LOOKUP_REVAL;

994            result = link_path_walk(name, nd);

995        }

996

997        path_put(&save);

998

999        return result;

1000 }
```

987 行：调用 link_path_walk 实现具体功能。

## link_path_walk：

该函数的逻辑比较简单，就是在一个 for 循环中处理每一个目录项，而每一次循环中针对目录项的不同做不同的处理：

■    目录项为".."，由 follow_dotot 处理。

■    目录项为正常目录，由 do_lookup 处理。

```
806 /*

807  * Name resolution.

808  * This is the basic name resolution function, turning a pathname into

809  * the final dentry. We expect 'base' to be positive and a directory.

810  *

811  * Returns 0 and nd will have valid dentry and mnt on success.
```

```
812    * Returns error and drops reference to input namei data on failure.
813   */
814 static int link_path_walk(const char *name, struct nameidata *nd)
815 {
816       struct path next;
817       struct inode *inode;
818       int err;
819       unsigned int lookup_flags = nd->flags;
820
821       while (*name=='/')
822           name++;
823       if (!*name)
824           goto return_reval;
825
826       inode = nd->path.dentry->d_inode;
827       if (nd->depth)
828           lookup_flags = LOOKUP_FOLLOW | (nd->flags & LOOKUP_CONTINUE);
```

如果是 nd->depth > 0，则表示 link_path_walk 为嵌套调用。

```
829
830       /* At this point we know we have a real path component. */
```

如果路径名是

```
831       for(;;) {
832           unsigned long hash;
833           struct qstr this;
834           unsigned int c;
835
836           nd->flags |= LOOKUP_CONTINUE;
837           err = exec_permission(inode);
838           if (err)
839               break;
```

| 840 | |
|---|---|
| 841 | this.name = name; |

this.name 中保存本次循环中 name 的起始地址。

| 842 | c = *(const unsigned char *)name; |
|---|---|
| 843 | |
| 844 | hash = init_name_hash(); |
| 845 | do { |
| 846 | name++; |
| 847 | hash = partial_name_hash(c, hash); |
| 848 | c = *(const unsigned char *)name; |
| 849 | } while (c && (c != '/')); |

向前移动指针 name，知道 c 为空或者遇到"/"。

| 850 | this.len = name - (const char *) this.name; |
|---|---|
| 851 | this.hash = end_name_hash(hash); |

保存本次查找目录的目录长度。

| 852 | |
|---|---|
| 853 | /* remove trailing slashes? */ |
| 854 | if (!c) |
| 855 | goto last_component; |
| 856 | while (*++name == '/'); |
| 857 | if (!*name) |
| 858 | goto last_with_slashes; |
| 859 | |
| 860 | /* |
| 861 | * "." and ".." are special - ".." especially so because it has |
| 862 | * to be able to know about the current root directory and |
| 863 | * parent relationships. |
| 864 | */ |
| 865 | if (this.name[0] == '.') switch (this.len) { |
| 866 | default: |

```
867                 break;
868         case 2:
869             if (this.name[1] != '.')
870                 break;
871             follow_dotdot(nd);
872             inode = nd->path.dentry->d_inode;
873             /* fallthrough */
874         case 1:
875             continue;
876         }
```

查看本次目录需要查找的是否为"."或者"..",如果是".",表示当前目录，则本次循环接触，跳转到下一次循环。如果是"..",表示父目录，调用函数 follow_dotdot，查找父目录索引节点。

如果以"."开头，但有不是上面两种情况，则路径有问题，跳出循环。

```
877         /* This does the actual lookups.. */
878         err = do_lookup(nd, &this, &next);
```

如果是正常的目录，调用 do_lookup 实现查找功能。

```
879         if (err)
880             break;
881
882         err = -ENOENT;
883         inode = next.dentry->d_inode;
884         if (!inode)
885             goto out_dput;
```

查找到了目录项，对应的目录项没有索引节点，则目录项有问题，跳转到 out_dput。

```
886
887         if (inode->i_op->follow_link) {
888             err = do_follow_link(&next, nd);
889             if (err)
890                 goto return_err;
```

| 891 | err = -ENOENT; |
| 892 | inode = nd->path.dentry->d_inode; |
| 893 | if (!inode) |
| 894 | break; |

如果目录项是一个链接，则调用 do_follow_link 函数。

| 895 | } else |
| 896 | path_to_nameidata(&next, nd); |
| 897 | err = -ENOTDIR; |
| 898 | if (!inode->i_op->lookup) |
| 899 | break; |
| 900 | continue; |

进入下一轮循环。

| 901 | /* here ends the main loop */ |
| 902 | |
| 903 | last_with_slashes: |
| 904 | lookup_flags |= LOOKUP_FOLLOW | LOOKUP_DIRECTORY; |

如果是以"/"结尾，则设置查找参数为 LOOKUP_FOLLOW | LOOKUP_DIRECTORY

| 905 | last_component: |
| 906 | /* Clear LOOKUP_CONTINUE iff it was previously unset */ |
| 907 | nd->flags &= lookup_flags | ~LOOKUP_CONTINUE; |
| 908 | if (lookup_flags & LOOKUP_PARENT) |
| 909 | goto lookup_parent; |

908-909 行：如果要查找的是父节点，则最后一个目录项不用查找了，直接跳转到 lookup_parent

| 910 | if (this.name[0] == '.') switch (this.len) { |
| 911 | default: |
| 912 | break; |
| 913 | case 2: |
| 914 | if (this.name[1] != '.') |
| 915 | break; |

| 916 | follow_dotdot(nd); |
| 917 | inode = nd->path.dentry->d_inode; |
| 918 | /* fallthrough */ |
| 919 | case 1: |
| 920 | goto return_reval; |

919-920 行：和中间目录项处理不同，中间目录项时跳转到下一次循环，这里跳出循环。

| 921 | } |
| 922 | err = do_lookup(nd, &this, &next); |
| 923 | if (err) |
| 924 | break; |
| 925 | inode = next.dentry->d_inode; |
| 926 | if (follow_on_final(inode, lookup_flags)) { |
| 927 | err = do_follow_link(&next, nd); |
| 928 | if (err) |
| 929 | goto return_err; |
| 930 | inode = nd->path.dentry->d_inode; |
| 931 | } else |
| 932 | path_to_nameidata(&next, nd); |
| 933 | err = -ENOENT; |
| 934 | if (!inode) |
| 935 | break; |
| 936 | if (lookup_flags & LOOKUP_DIRECTORY) { |
| 937 | err = -ENOTDIR; |
| 938 | if (!inode->i_op->lookup) |
| 939 | break; |
| 940 | } |
| 941 | goto return_base; |

906-941 行：处理最后一个目录项。

| 942 lookup_parent: |

| 943 | nd->last = this; |

```
944            nd->last_type = LAST_NORM;
945            if (this.name[0] != '.')
946                goto return_base;
947            if (this.len == 1)
948                nd->last_type = LAST_DOT;
949            else if (this.len == 2 && this.name[1] == '.')
950                nd->last_type = LAST_DOTDOT;
951            else
952                goto return_base;
```

942 行：如果要查找的是父节点，则

```
953 return_reval:
954            /*
955             * We bypassed the ordinary revalidation routines.
956             * We may need to check the cached dentry for staleness.
957             */
958            if (nd->path.dentry && nd->path.dentry->d_sb &&
959                (nd->path.dentry->d_sb->s_type->fs_flags &
FS_REVAL_DOT)) {
960                err = -ESTALE;
961                /* Note: we do not d_invalidate() */
962                if (!nd->path.dentry->d_op->d_revalidate(
963                        nd->path.dentry, nd))
964                    break;
965            }
966 return_base:
967            return 0;
```

正确返回

```
968 out_dput:
969            path_put_conditional(&next, nd);
970            break;
971        }
```

| | |
|---|---|
| 972 | path_put(&nd->path); |
| 973 return_err: | |
| 974 | return err; |

错误返回。

| | |
|---|---|
| 975 } | |

### 5.1.3.2.1 回退到父目录follow_dotdot

| | |
|---|---|
| 670 static __always_inline void follow_dotdot(struct nameidata *nd) | |
| 671 { | |
| 672 | set_root(nd); |

设置 nd->root 为当前进程的 root。

| | |
|---|---|
| 673 | |
| 674 | while(1) { |
| 675 | struct dentry *old = nd->path.dentry; |
| 676 | |
| 677 | if (nd->path.dentry == nd->root.dentry && |
| 678 | nd->path.mnt == nd->root.mnt) { |
| 679 | break; |
| 680 | } |
| 681 | if (nd->path.dentry != nd->path.mnt->mnt_root) { |
| 682 | /* rare case of legitimate dget_parent()... */ |
| 683 | nd->path.dentry = dget_parent(nd->path.dentry); |
| 684 | dput(old); |
| 685 | break; |
| 686 | } |
| 687 | if (!follow_up(&nd->path)) |
| 688 | break; |

当 nd->path.dentry 所在的设备为根设备时，follow_up 返回 0，则执行到 688 行，跳出循环。

当 nd->path.dentry 所在的设备不是跟设备时，nd->path.dentry 为父设备中的安装点，nd->path.mnt 为父设备的 vfsmount。此时 follow_up 返回 1,688 行不执行，跳转到 674 行继续循环。

为啥 follow_up 会做如此的设计呢。因为 nd->path.dentry 已经为其所在设备的根设备了，并且 nd->path.mnt 所代表的设备也就是根设备了，所以没法再往上会退了，故返回 0，跳出循环。follow_up 返回 1，也就是 nd->path.dentry 所在的设备不是根设备，由 vfsmount 中 mnt_mountpoint（父设备安装点）与 mnt_mount(子设备根节点)表示同一层级，而 follow_dotdot 需要回退一级，所以 follow_up 返回 1，进而 where 循环继续运行，回退到 mnt_mountpoint 中的父节点。

| 689 | } |
|---|---|
| 690 | follow_mount(&nd->path); |

follow_mount 处理 nd->path.dentry 为安装点的情况。

| 691 | } |

该函数三种情况来处理：

- 第一种情况，已经到达节点 nd->path.dentry 就是本进程的根节点，且到达节点的安装点 nd->path.mnt 与本进程根节点的安装点相同，此时就不能在往上跑了，保持 nd->path.dentry 不变。（677-680）
- 第二种情况，已经到达的节点 nd->path.dentry 与其父节点在同一设备上。且 nd->path.dentry 不是根节点，这种情况处理比较简单，去 path.dentry 的父目录项即可。681 行为何能保证这一点能（因为 nd->path 所在的安装点就是 nd->path.mnt，所以 nd->path 与 nd->paht.mnt->mnt_root 一定在同一设备上）。（681-686）
- 第三中情况，已经到达的节点就是该节点所在设备的根节点，在往上跑一层就要跑到另一个设备上去了。这种情况交由 follow_up 处理。（687 行）

当将一个存储设备"安装"到另一个设备上的某个节点时，内核会分配和设置一个 vfsmount 设备，通过这个结构上的 mnt_mountpoint 和 mnt_root 将这两个设备节点关联起来。vfsmount 结构中 mnt_parent 指向"父设备"。mnt_mountpoint 指向安装点，mnt_root 指向安装设备（子设备）的根目录。

| 50 struct vfsmount { |
|---|

```
51      struct list_head mnt_hash;

52      struct vfsmount *mnt_parent;      /* fs we are mounted on */

53      struct dentry *mnt_mountpoint;   /* dentry of mountpoint */

54      struct dentry *mnt_root;      /* root of the mounted tree */

55      struct super_block *mnt_sb; /* pointer to superblock */
```

### follow_up 函数

```
599 int follow_up(struct path *path)
```

```
600 {

601      struct vfsmount *parent;

602      struct dentry *mountpoint;

603      spin_lock(&vfsmount_lock);

604      parent = path->mnt->mnt_parent;
```

parent 父设备安装点

```
605      if (parent == path->mnt) {

606          spin_unlock(&vfsmount_lock);

607          return 0;

608      }
```

如果设备 mnt 与 mnt->mnt_parent 相同，则该设备为根设备，所以返回 0。

```
609      mntget(parent);

610      mountpoint = dget(path->mnt->mnt_mountpoint);

611      spin_unlock(&vfsmount_lock);

612      dput(path->dentry);

613      path->dentry = mountpoint;
```

设置 path->dentry 为父设备中的安装点

```
614      mntput(path->mnt);

615      path->mnt = parent;
```

设置 path->mnt 为父设备的 vfsmount。

```
616      return 1;
```

返回 1。

617 }

### follow_mount 函数

639 static void follow_mount(struct path *path)

640 {

641    while (d_mountpoint(path->dentry)) {

642        struct vfsmount *mounted = lookup_mnt(path);

643        if (!mounted)

644            break;

645        dput(path->dentry);

646        mntput(path->mnt);

647        path->mnt = mounted;

648        path->dentry = dget(mounted->mnt_root);

649    }

650 }

d_mountpoint 函数检测 path->dentry 是否为安装点，如果是，调用 lookup_mnt 搜索目录项高速缓存中已安装文件系统的根目录，并把 nd->dentry 和 nd->mnt 更新为相应已安装文件系统的安装点和安装系统对象地址。然后重复整个操作（几个文件系统可以安装在同一个安装点上）。本质上讲，由于进程可能从某个文件系统的目录开始路径名的查找，而该目录被另一个安装在其父目录上的文件系统所隐藏，那么当需要回到父目录时，则调用 follow_mount 函数。

### 5.1.3.2.2 目录查找函数 do_lookup

该函数是实现由中间目录名到 inode 节点转换的具体函数。

693 /*

694    *   It's more convoluted than I'd like it to be, but... it's still fairly

695    *   small and for now I'd prefer to have fast path as straight as possible.

696    *   It _is_ time-critical.

```
697    */
698 static int do_lookup(struct nameidata *nd, struct qstr *name,
699                struct path *path)
700 {
701     struct vfsmount *mnt = nd->path.mnt;
702     struct dentry *dentry, *parent;
703     struct inode *dir;
704     /*
705      * See if the low-level filesystem might want
706      * to use its own hash..
707      */
708     if (nd->path.dentry->d_op && nd->path.dentry->d_op->d_hash) {
709         int err = nd->path.dentry->d_op->d_hash(nd->path.dentry, name);
710         if (err < 0)
711             return err;
712     }
713
714     dentry = __d_lookup(nd->path.dentry, name);
```

在 nd->path.dentry 对应的缓存中查找名称为 name 的子 dentry 项。

```
715     if (!dentry)
716         goto need_lookup;
```

如果在缓存中没有找到，则调整到 need_lookup 函数，此时需要从磁盘中读入相关信息，建立 dentry 缓存。

```
717     if (dentry->d_op && dentry->d_op->d_revalidate)
718         goto need_revalidate;
```

如果 dentry 的 d_op-> d_revalidate 非空，则调整到验证代码。

```
719 done:
720     path->mnt = mnt;
721     path->dentry = dentry;
```

查找到 dentry 后，设置 path 相关变量。

| 722 | __follow_mount(path); |
|---|---|

　　如果 path->dentry 为安装点，则调用__follow_mount 函数前进到子设备的根目录。

| 723 | return 0; |
|---|---|
| 724 | |
| 725 | need_lookup: |
| 726 | parent = nd->path.dentry; |
| 727 | dir = parent->d_inode; |
| 728 | |
| 729 | mutex_lock(&dir->i_mutex); |
| 730 | /* |
| 731 | * First re-do the cached lookup just in case it was created |
| 732 | * while we waited for the directory semaphore.. |
| 733 | * |
| 734 | * FIXME! This could use version numbering or similar to |
| 735 | * avoid unnecessary cache lookups. |
| 736 | * |
| 737 | * The "dcache_lock" is purely to protect the RCU list walker |
| 738 | * from concurrent renames at this point (we mustn't get false |
| 739 | * negatives from the RCU list walk here, unlike the optimistic |
| 740 | * fast walk). |
| 741 | * |
| 742 | * so doing d_lookup() (with seqlock), instead of lockfree __d_lookup |
| 743 | */ |
| 744 | dentry = d_lookup(parent, name); |

　　再次在缓存中查找。

| 745 | if (!dentry) { |
|---|---|
| 746 | struct dentry *new; |
| 747 | |
| 748 | /* Don't create child dentry for a dead directory. */ |
| 749 | dentry = ERR_PTR(-ENOENT); |

| | |
|---|---|
| 750 | if (IS_DEADDIR(dir)) |
| 751 | goto out_unlock; |
| 752 | |
| 753 | new = d_alloc(parent, name); |

在内存中分配一个新的 dentry 结构。

| | |
|---|---|
| 754 | dentry = ERR_PTR(-ENOMEM); |
| 755 | if (new) { |
| 756 | dentry = dir->i_op->lookup(dir, new, nd); |

根据具体的文件系统，选择具体的查找函数，对于 ext2 而言，就是 ext2_lookup 函数。

| | |
|---|---|
| 757 | if (dentry) |
| 758 | dput(new); |
| 759 | else |
| 760 | dentry = new; |
| 761 | } |
| 762 out_unlock: | |
| 763 | mutex_unlock(&dir->i_mutex); |
| 764 | if (IS_ERR(dentry)) |
| 765 | goto fail; |
| 766 | goto done; |
| 767 | } |
| 768 | |
| 769 | /* |
| 770 | * Uhhuh! Nasty case: the cache was re-populated while |
| 771 | * we waited on the semaphore. Need to revalidate. |
| 772 | */ |
| 773 | mutex_unlock(&dir->i_mutex); |
| 774 | if (dentry->d_op && dentry->d_op->d_revalidate) { |
| 775 | dentry = do_revalidate(dentry, nd); |
| 776 | if (!dentry) |
| 777 | dentry = ERR_PTR(-ENOENT); |

| | |
|---|---|
| 778 | } |

如果是在缓存中找到的 dentry，需要验证。

| | |
|---|---|
| 779 | if (IS_ERR(dentry)) |
| 780 | goto fail; |
| 781 | goto done; |
| 782 | |
| 783 | need_revalidate: |
| 784 | dentry = do_revalidate(dentry, nd); |

验证

| | |
|---|---|
| 785 | if (!dentry) |
| 786 | goto need_lookup; |

验证失败，跳转到真正查找的地方

| | |
|---|---|
| 787 | if (IS_ERR(dentry)) |
| 788 | goto fail; |
| 789 | goto done; |
| 790 | |
| 791 | fail: |
| 792 | return PTR_ERR(dentry); |
| 793 | } |

do_lookup 函数现在缓存中查找（d_lookup）name 对应的 dentry 是否存在，如果存在，在返回该 dentry。 否则，就需要先分配一个 dentry 结构，然后调用具体的文件系统查找函数，从磁盘中读取相关内容建立内存中的 dentry 数据结构。

## 缓存中搜索 dentry 结构

| | |
|---|---|
| 1359 | struct dentry * d_lookup(struct dentry * parent, struct qstr * name) |
| 1360 | { |
| 1361 | struct dentry * dentry = NULL; |
| 1362 | unsigned long seq; |
| 1363 | |

```
1364          do {
1365                    seq = read_seqbegin(&rename_lock);
1366                    dentry = __d_lookup(parent, name);
1367                    if (dentry)
1368              break;
1369      } while (read_seqretry(&rename_lock, seq));
1370      return dentry;
1371 }
1372 EXPORT_SYMBOL(d_lookup);
1373
1374 struct dentry * __d_lookup(struct dentry * parent, struct qstr * name)
1375 {
```

参数 parent 表示父目录，name 表示需要在父目录中查找的名称。

```
1376      unsigned int len = name->len;
1377      unsigned int hash = name->hash;
1378      const unsigned char *str = name->name;
1379      struct hlist_head *head = d_hash(parent,hash);
1380      struct dentry *found = NULL;
1381      struct hlist_node *node;
1382      struct dentry *dentry;
1383
1384      rcu_read_lock();
1385
1386      hlist_for_each_entry_rcu(dentry, node, head, d_hash) {
1387          struct qstr *qstr;
1388
1389          if (dentry->d_name.hash != hash)
1390              continue;
1391          if (dentry->d_parent != parent)
1392              continue;
```

判断当前的 dentry 是否满足要求，满足要求的条件是：

- hash 值与传递进来的 hash 值相同
- dentry 是 parent 的子目录

```
1393
1394        spin_lock(&dentry->d_lock);
1395
1396        /*
1397         * Recheck the dentry after taking the lock - d_move may have
1398         * changed things.   Don't bother checking the hash because we're
1399         * about to compare the whole name anyway.
1400         */
1401        if (dentry->d_parent != parent)
1402            goto next;
1403
1404        /* non-existing due to RCU? */
1405        if (d_unhashed(dentry))
1406            goto next;
1407
1408        /*
1409         * It is safe to compare names since d_move() cannot
1410         * change the qstr (protected by d_lock).
1411         */
1412        qstr = &dentry->d_name;
1413        if (parent->d_op && parent->d_op->d_compare) {
1414            if (parent->d_op->d_compare(parent, qstr, name))
1415                goto next;
1416        } else {
1417            if (qstr->len != len)
1418                goto next;
1419            if (memcmp(qstr->name, str, len))
1420                goto next;
```

| |
|---|
| 1421　　　　} |

名称，长度的进一步检验

| |
|---|
| 1422 |
| 1423　　　　atomic_inc(&dentry->d_count); |
| 1424　　　　found = dentry; |
| 1425　　　　spin_unlock(&dentry->d_lock); |
| 1426　　　　break; |
| 1427 next: |
| 1428　　　　spin_unlock(&dentry->d_lock); |
| 1429　　} |
| 1430　　rcu_read_unlock(); |
| 1431 |
| return found; |
| 1433 } |

# 分配一个新的 dentry

| |
|---|
| 915 /** |
| 916　　* d_alloc　-　　allocate a dcache entry |
| 917　　* @parent: parent of entry to allocate |
| 918　　* @name: qstr of the name |
| 919　　* |
| 920　　* Allocates a dentry. It returns %NULL if there is insufficient memory |
| 921　　* available. On a success the dentry is returned. The name passed in is |
| 922　　* copied and the copy passed in may be reused after this call. |
| 923　　*/ |
| 924 |
| 925 struct dentry *d_alloc(struct dentry * parent, const struct qstr *name) |
| 926 { |

参数 parent 为新分配目录项的父目录。

name 为新分配目录项的名称。

| |
|---|
| 927　　struct dentry *dentry; |

```
928        char *dname;
929
930        dentry = kmem_cache_alloc(dentry_cache, GFP_KERNEL);
```

从 dentry_cache 缓存中分配 dentry

```
931        if (!dentry)
932            return NULL;
933
934        if (name->len > DNAME_INLINE_LEN-1) {
935            dname = kmalloc(name->len + 1, GFP_KERNEL);
936            if (!dname) {
937                kmem_cache_free(dentry_cache, dentry);
938                return NULL;
939            }
940        } else   {
941            dname = dentry->d_iname;
942        }
943        dentry->d_name.name = dname;
944
945        dentry->d_name.len = name->len;
946        dentry->d_name.hash = name->hash;
947        memcpy(dname, name->name, name->len);
948        dname[name->len] = 0;
949
950        atomic_set(&dentry->d_count, 1);
951        dentry->d_flags = DCACHE_UNHASHED;
952        spin_lock_init(&dentry->d_lock);
953        dentry->d_inode = NULL;
954        dentry->d_parent = NULL;
955        dentry->d_sb = NULL;
956        dentry->d_op = NULL;
957        dentry->d_fsdata = NULL;
```

| | |
|---|---|
| 958 | dentry->d_mounted = 0; |
| 959 | INIT_HLIST_NODE(&dentry->d_hash); |
| 960 | INIT_LIST_HEAD(&dentry->d_lru); |
| 961 | INIT_LIST_HEAD(&dentry->d_subdirs); |
| 962 | INIT_LIST_HEAD(&dentry->d_alias); |

初始化 dentry 结构相关变量

| | |
|---|---|
| 963 | |
| 964 | if (parent) { |
| 965 | dentry->d_parent = dget(parent); |

设置 dentry 父节点

| | |
|---|---|
| 967 | } else { |
| 968 | INIT_LIST_HEAD(&dentry->d_u.d_child); |
| 969 | } |
| 970 | |
| 971 | spin_lock(&dcache_lock); |
| 972 | if (parent) |
| 973 | list_add(&dentry->d_u.d_child, &parent->d_subdirs); |
| 974 | dentry_stat.nr_dentry++; |
| 975 | spin_unlock(&dcache_lock); |
| 976 | |
| 977 | return dentry; |
| 978 | } |
| 979 | EXPORT_SYMBOL(d_alloc); |

## 5.1.3.3  ext2 目录查找函数ext2_lookup

在 do_lookup 的 756 行中有如下的调用，

| | |
|---|---|
| 756 | dentry = dir->i_op->lookup(dir, new, nd); |

最终会转换为 ext2_lookup 函数。

fs/ext2/namei.c

| | |
|---|---|
| 54 | /* |

```
55  * Methods themselves.
56  */
57
58 static struct dentry *ext2_lookup(struct inode * dir, struct dentry *dentry,
struct nameidata *nd)
59 {
60      struct inode * inode;
61      ino_t ino;
62
63      if (dentry->d_name.len > EXT2_NAME_LEN)
64          return ERR_PTR(-ENAMETOOLONG);
65
66      ino = ext2_inode_by_name(dir, &dentry->d_name);
```

有名称查找对应的 inode 索引。

```
67      inode = NULL;
68      if (ino) {
69          inode = ext2_iget(dir->i_sb, ino);
```

通过 indoe 索引，从磁盘读入相关信息，在内存中建立 inode 结构。

```
70          if (unlikely(IS_ERR(inode))) {
71              if (PTR_ERR(inode) == -ESTALE) {
72                  ext2_error(dir->i_sb, __func__,
73                          "deleted inode referenced: %lu",
74                          (unsigned long) ino);
75                  return ERR_PTR(-EIO);
76              } else {
77                  return ERR_CAST(inode);
78              }
79          }
80      }
81      return d_splice_alias(inode, dentry);
```

安装 inode 到 dentry 结构中。

```
82 }
```

### 5.1.3.3.1  name到索引节点号ext2_inode_by_name

```
437 ino_t ext2_inode_by_name(struct inode *dir, struct qstr *child)

438 {

439     ino_t res = 0;

440     struct ext2_dir_entry_2 *de;

441     struct page *page;

442

443     de = ext2_find_entry (dir, child, &page);

444     if (de) {

445         res = le32_to_cpu(de->inode);
```

得到 inode 节点的索引号

```
446         ext2_put_page(page);

447     }

448     return res;

449 }
```

## 查找磁盘上目录项的相关内容

```
352 /*

353   *   ext2_find_entry()

354   *

355   * finds an entry in the specified directory with the wanted name. It

356   * returns the page in which the entry was found (as a parameter -
res_page),

357   * and the entry itself. Page is returned mapped and unlocked.

358   * Entry is guaranteed to be valid.

359   */

360 struct ext2_dir_entry_2 *ext2_find_entry (struct inode * dir,

361                 struct qstr *child, struct page ** res_page)
```

```
362 {

363     const char *name = child->name;

364     int namelen = child->len;

365     unsigned reclen = EXT2_DIR_REC_LEN(namelen);

366     unsigned long start, n;

367     unsigned long npages = dir_pages(dir);

368     struct page *page = NULL;

369     struct ext2_inode_info *ei = EXT2_I(dir);

370     ext2_dirent * de;

371     int dir_has_error = 0;

372

373     if (npages == 0)

374         goto out;

375

376     /* OFFSET_CACHE */

377     *res_page = NULL;

378

379     start = ei->i_dir_start_lookup;

380     if (start >= npages)

381         start = 0;

382     n = start;

383     do {

384         char *kaddr;

385         page = ext2_get_page(dir, n, dir_has_error);
```

从父目录文件中读取一页内容，其中可能涉及到磁盘读取操作。

```
386         if (!IS_ERR(page)) {

387             kaddr = page_address(page);

388             de = (ext2_dirent *) kaddr;
```

Kaddr 可以看做是一个 ext2_dir_entry_2 数组

```
389             kaddr += ext2_last_byte(dir, n) - reclen;

390             while ((char *) de <= kaddr) {
```

```
391                    if (de->rec_len == 0) {
392                        ext2_error(dir->i_sb, __func__,
393                            "zero-length directory entry");
394                        ext2_put_page(page);
395                        goto out;
396                    }
397                    if (ext2_match (namelen, name, de))
```

检测是否为需要的元素

```
398                        goto found;
399                    de = ext2_next_entry(de);
400                }
401                ext2_put_page(page);
402            } else
403                dir_has_error = 1;
404
405            if (++n >= npages)
406                n = 0;
407            /* next page is past the blocks we've got */
408            if (unlikely(n > (dir->i_blocks >> (PAGE_CACHE_SHIFT - 9)))) {
409                ext2_error(dir->i_sb, __func__,
410                    "dir %lu size %lld exceeds block count %llu",
411                    dir->i_ino, dir->i_size,
412                    (unsigned long long)dir->i_blocks);
413                goto out;
414            }
415    } while (n != start);
416 out:
417    return NULL;
418
419 found:
420    *res_page = page;
```

```
421     ei->i_dir_start_lookup = n;

422     return de;

423 }
```

### 5.1.3.3.2 索引号到索引节点 ext2_iget

```
1219 struct inode *ext2_iget (struct super_block *sb, unsigned long ino)

1220 {

1221     struct ext2_inode_info *ei;

1222     struct buffer_head * bh;

1223     struct ext2_inode *raw_inode;

1224     struct inode *inode;

1225     long ret = -EIO;

1226     int n;

1227

1228     inode = iget_locked(sb, ino);

1229     if (!inode)

1230         return ERR_PTR(-ENOMEM);

1231     if (!(inode->i_state & I_NEW))

1232         return inode;

1233

1234     ei = EXT2_I(inode);
```

从 inode 中提取出 ext2_inode 信息。

```
1235     ei->i_block_alloc_info = NULL;

1236

1237     raw_inode = ext2_get_inode(inode->i_sb, ino, &bh);
```

该函数是重点，获得原始的 inode 信息。

```
1238     if (IS_ERR(raw_inode)) {

1239         ret = PTR_ERR(raw_inode);

1240         goto bad_inode;

1241     }
```

```
1242
1243        inode->i_mode = le16_to_cpu(raw_inode->i_mode);
1244        inode->i_uid = (uid_t)le16_to_cpu(raw_inode->i_uid_low);
1245        inode->i_gid = (gid_t)le16_to_cpu(raw_inode->i_gid_low);
1246        if (!(test_opt (inode->i_sb, NO_UID32))) {
1247            inode->i_uid |= le16_to_cpu(raw_inode->i_uid_high) << 16;
1248            inode->i_gid |= le16_to_cpu(raw_inode->i_gid_high) << 16;
1249        }
1250        inode->i_nlink = le16_to_cpu(raw_inode->i_links_count);
1251        inode->i_size = le32_to_cpu(raw_inode->i_size);
1252        inode->i_atime.tv_sec = (signed)le32_to_cpu(raw_inode->i_atime);
1253        inode->i_ctime.tv_sec = (signed)le32_to_cpu(raw_inode->i_ctime);
1254        inode->i_mtime.tv_sec = (signed)le32_to_cpu(raw_inode->i_mtime);
1255        inode->i_atime.tv_nsec = inode->i_mtime.tv_nsec =
inode->i_ctime.tv_nsec = 0;
1256        ei->i_dtime = le32_to_cpu(raw_inode->i_dtime);
1257        /* We now have enough fields to check if the inode was active or not.
1258         * This is needed because nfsd might try to access dead inodes
1259         * the test is that same one that e2fsck uses
1260         * NeilBrown 1999oct15
1261         */
1262        if (inode->i_nlink == 0 && (inode->i_mode == 0 || ei->i_dtime)) {
1263            /* this inode is deleted */
1264            brelse (bh);
1265            ret = -ESTALE;
1266            goto bad_inode;
1267        }
1268        inode->i_blocks = le32_to_cpu(raw_inode->i_blocks);
1269        ei->i_flags = le32_to_cpu(raw_inode->i_flags);
1270        ei->i_faddr = le32_to_cpu(raw_inode->i_faddr);
1271        ei->i_frag_no = raw_inode->i_frag;
1272        ei->i_frag_size = raw_inode->i_fsize;
```

```
1273        ei->i_file_acl = le32_to_cpu(raw_inode->i_file_acl);
1274        ei->i_dir_acl = 0;
1275        if (S_ISREG(inode->i_mode))
1276            inode->i_size |= ((__u64)le32_to_cpu(raw_inode->i_size_high))
<< 32;
1277        else
1278            ei->i_dir_acl = le32_to_cpu(raw_inode->i_dir_acl);
1279        ei->i_dtime = 0;
1280        inode->i_generation = le32_to_cpu(raw_inode->i_generation);
1281        ei->i_state = 0;
1282        ei->i_block_group = (ino - 1) /
EXT2_INODES_PER_GROUP(inode->i_sb);
1283        ei->i_dir_start_lookup = 0;
1284
1285        /*
1286         * NOTE! The in-memory inode i_data array is in little-endian order
1287         * even on big-endian machines: we do NOT byteswap the block
numbers!
1288         */
1289        for (n = 0; n < EXT2_N_BLOCKS; n++)
1290            ei->i_data[n] = raw_inode->i_block[n];
```

将磁盘中 i_block 数组拷贝到 inode 中。

```
1291
1292        if (S_ISREG(inode->i_mode)) {
1293            inode->i_op = &ext2_file_inode_operations;
1294            if (ext2_use_xip(inode->i_sb)) {
1295                inode->i_mapping->a_ops = &ext2_aops_xip;
1296                inode->i_fop = &ext2_xip_file_operations;
1297            } else if (test_opt(inode->i_sb, NOBH)) {
1298                inode->i_mapping->a_ops = &ext2_nobh_aops;
1299                inode->i_fop = &ext2_file_operations;
1300            } else {
```

```
1301              inode->i_mapping->a_ops = &ext2_aops;
1302              inode->i_fop = &ext2_file_operations;
1303          }
1304      } else if (S_ISDIR(inode->i_mode)) {
1305          inode->i_op = &ext2_dir_inode_operations;
1306          inode->i_fop = &ext2_dir_operations;
1307          if (test_opt(inode->i_sb, NOBH))
1308              inode->i_mapping->a_ops = &ext2_nobh_aops;
1309          else
1310              inode->i_mapping->a_ops = &ext2_aops;
1311      } else if (S_ISLNK(inode->i_mode)) {
1312          if (ext2_inode_is_fast_symlink(inode)) {
1313              inode->i_op = &ext2_fast_symlink_inode_operations;
1314              nd_terminate_link(ei->i_data, inode->i_size,
1315                  sizeof(ei->i_data) - 1);
1316          } else {
1317              inode->i_op = &ext2_symlink_inode_operations;
1318              if (test_opt(inode->i_sb, NOBH))
1319                  inode->i_mapping->a_ops = &ext2_nobh_aops;
1320              else
1321                  inode->i_mapping->a_ops = &ext2_aops;
1322          }
1323      } else {
1324          inode->i_op = &ext2_special_inode_operations;
1325          if (raw_inode->i_block[0])
1326              init_special_inode(inode, inode->i_mode,
1327                  old_decode_dev(le32_to_cpu(raw_inode->i_block[0])));
1328          else
1329              init_special_inode(inode, inode->i_mode,
1330                  new_decode_dev(le32_to_cpu(raw_inode->i_block[1])));
1331      }
```

| | |
|---|---|
| 1332 | brelse (bh); |
| 1333 | ext2_set_inode_flags(inode); |
| 1334 | unlock_new_inode(inode); |
| 1335 | return inode; |
| 1336 | |
| 1337 | bad_inode: |
| 1338 | iget_failed(inode); |
| 1339 | return ERR_PTR(ret); |
| 1340 | } |

本函数的重点是 ext2_get_inode 函数，它根据 inode 索引节点号的到磁盘上的原始 indoe 信息，然后用此初始化内存中的 inode 结构。

| | |
|---|---|
| 1140 | static struct ext2_inode *ext2_get_inode(struct super_block *sb, ino_t ino, |
| 1141 | struct buffer_head **p) |
| 1142 | { |
| 1143 | struct buffer_head * bh; |
| 1144 | unsigned long block_group; |
| 1145 | unsigned long block; |
| 1146 | unsigned long offset; |
| 1147 | struct ext2_group_desc * gdp; |
| 1148 | |
| 1149 | *p = NULL; |
| 1150 | if ((ino != EXT2_ROOT_INO && ino < EXT2_FIRST_INO(sb)) \|\| |
| 1151 | ino > le32_to_cpu(EXT2_SB(sb)->s_es->s_inodes_count)) |
| 1152 | goto Einval; |
| 1153 | |
| 1154 | block_group = (ino - 1) / EXT2_INODES_PER_GROUP(sb); |

计算块组号

| | |
|---|---|
| 1155 | gdp = ext2_get_group_desc(sb, block_group, NULL); |

得到块组

| | |
|---|---|
| 1156 | if (!gdp) |

```
1157          goto Egdp;
1158      /*
1159       * Figure out the offset within the block group inode table
1160       */
1161      offset = ((ino - 1) % EXT2_INODES_PER_GROUP(sb)) *
EXT2_INODE_SIZE(sb);
```

节点在块内偏移

```
1162      block = le32_to_cpu(gdp->bg_inode_table) +
1163          (offset >> EXT2_BLOCK_SIZE_BITS(sb));
```

节点在磁盘上的块号

```
1164      if (!(bh = sb_bread(sb, block)))
1165          goto Eio;
```

调用底层驱动函数，读取指定块。

```
1166
1167      *p = bh;
1168      offset &= (EXT2_BLOCK_SIZE(sb) - 1);
1169      return (struct ext2_inode *) (bh->b_data + offset);
```

返回

```
1170
1171 Einval:
1172      ext2_error(sb, "ext2_get_inode", "bad inode number: %lu",
1173              (unsigned long) ino);
1174      return ERR_PTR(-EINVAL);
1175 Eio:
1176      ext2_error(sb, "ext2_get_inode",
1177              "unable to read inode block - inode=%lu, block=%lu",
1178              (unsigned long) ino, block);
1179 Egdp:
1180      return ERR_PTR(-EIO);
1181 }
```

### 5.1.3.4 符号链接处理do_follow_link

在由名称到节点的转换过程中，还需检测最后得到的目录项是否为链接，如果是链接的话，调用 do_follow_walk 函数处理。

考虑如下的情况，如果/home/henry 是指向/home/zenhumany 的一个符号链接，那么，/home/henry/test.txt 必须由内核解析为/home/zenhumany/test.txt。在此，内核指向两个操作。

第一个查找操作：解析到/home/henry 时，发现 henry 是一个链接，内核必须将其链接内容/home/zenhumany 提取出来。进入第二个查找操作。

第二个查找操作，解析/home/zenhumany 目录，当得到 zenhumany 的 dentry 的目录项后，转入到第一个查找操作。

第一个查找操作从 zenhumany 的 dentry 中继续查找 test.txt 的目录项。

然而，难以驾驭的递归本质上是危险的。例如，假定一个符号链接指向自己。当然，解析含有这样符号链接的路径名可能导致无休止的递归调用流，这又依次引发内核栈的溢出。当前进程的描述符中的 link_count 字段用来避免这种问题：每次递归执行前增加这个字段的值，执行之后减少其值。如果该字段的值达到 6，整个循环操作就以错误码结束。因此，符号链接嵌套的层数不超过 5。

另外，当前进程的描述符中的 total_link_count 字段记录在原查找操作中有多少符号链接（甚至非嵌套的）被跟踪。如果这个计数器的值到 40，则查找操作中止。没有这个计数器，怀有恶意的用户就可能创建一个病态的路径名，让其中包含很多连续的符号链接，使内核在无休止的查找操作中冻结。

```
570 static inline int do_follow_link(struct path *path, struct nameidata *nd)
571 {
572     void *cookie;
573     int err = -ELOOP;
574     if (current->link_count >= MAX_NESTED_LINKS)
575         goto loop;
576     if (current->total_link_count >= 40)
577         goto loop;
578     BUG_ON(nd->depth >= MAX_NESTED_LINKS);
```

| | |
|---|---|
| 579 | cond_resched(); |
| 580 | err = security_inode_follow_link(path->dentry, nd); |
| 581 | if (err) |
| 582 | goto loop; |
| 583 | current->link_count++; |

增加符号链接计数，防止嵌套的层数过的导致内核溢出。

| | |
|---|---|
| 584 | current->total_link_count++; |

原查找路径中总的符号链接数。

| | |
|---|---|
| 585 | nd->depth++; |
| 586 | err = __do_follow_link(path, nd, &cookie); |

调用__do_follow_link 处理。

| | |
|---|---|
| 587 | if (!IS_ERR(cookie) && path->dentry->d_inode->i_op->put_link) |
| 588 | path->dentry->d_inode->i_op->put_link(path->dentry, nd, cookie); |
| 589 | path_put(path); |
| 590 | current->link_count--; |
| 591 | nd->depth--; |

减少链接引用计数。

| | |
|---|---|
| 592 | return err; |
| 593 loop: | |
| 594 | path_put_conditional(path, nd); |
| 595 | path_put(&nd->path); |
| 596 | return err; |
| 597 } | |

| | |
|---|---|
| 532 static __always_inline int | |
| 533 __do_follow_link(struct path *path, struct nameidata *nd, void **p) | |
| 534 { | |
| 535 | int error; |
| 536 | struct dentry *dentry = path->dentry; |
| 537 | |

```
538        touch_atime(path->mnt, dentry);

539        nd_set_link(nd, NULL);

540

541        if (path->mnt != nd->path.mnt) {

542            path_to_nameidata(path, nd);

543            dget(dentry);

544        }

545        mntget(path->mnt);

546        nd->last_type = LAST_BIND;

547        *p = dentry->d_inode->i_op->follow_link(dentry, nd);

548        error = PTR_ERR(*p);

549        if (!IS_ERR(*p)) {

550            char *s = nd_get_link(nd);
```

得到符号链接所表示的目录。

```
551            error = 0;

552            if (s)

553                error = __vfs_follow_link(nd, s);

554            else if (nd->last_type == LAST_BIND) {

555                error = force_reval_path(&nd->path, nd);

556                if (error)

557                    path_put(&nd->path);

558            }

559        }

560        return error;

561 }
```

```
498 static __always_inline int __vfs_follow_link(struct nameidata *nd, const char
*link)

499 {

500     if (IS_ERR(link))

501         goto fail;
```

```
502
503     if (*link == '/') {
504          set_root(nd);
505          path_put(&nd->path);
506          nd->path = nd->root;
507          path_get(&nd->root);
508     }
509
510     return link_path_walk(link, nd);
```

从符号链接的目录开始调用 link_path_walk。大家记得 do_follow_link 是从 link_path_walk 中调用的，现在回调 link_path_walk 函数，构成递归调用。

```
511 fail:
512     path_put(&nd->path);
513     return PTR_ERR(link);
 514 }
```

## 5.2 文件系统的安装与拆卸

### 5.2.1 概述

在一个块设备上按一定的格式建立起文件系统的时候，或者系统引导之初，设备上的文件和节点都还是不可访问的。也就是说，还不能按一定的路径名访问其中的特定的节点或文件（虽然作为"设备"是可访问的）。只有把它"安装"到计算机系统的文件系统中的某个节点上，才能使设备上的文件和节点成为可访问的。经过安装以后，设备上的"文件系统"就成为整个文件系统的一部分，或者说一个子系统。一般而言，文件系统就像一颗倒立的树，不过由于可能存在着的节点间的"连接"和"符号连接"而不一定是严格意义上的"树"。最初时，整个系统中只有一个节点，那就是整个文件系统的"根"节点，这个节点存在于内存中，而不在任何具体的设备上。系统在初始化时将一个"根设备"安装到节点"/"上，这个设备上的文件系统就成了整个系统中原始的、基本的文件系统（所以才称为根设备）。此后，就可以由超级用户进程通过系统调用 mount( )把其他的子系统安装到已经存在于文件系统中的空

闲节点上，使整个文件系统得以扩展，当不再需要使用某个子系统时，或者在关闭系统之前，则通过系统调用 umount 把已经安装的设备逐个"拆卸"下来。

系统调用 mount 将一个可访问的块设备（代表设备本身）安装到一个可访问的节点（设备安装点）上。所谓"可访问"是指该节点或者文件已经存在于已安装的文件系统中，可以通过路径名寻访。Uinx（以及 linux）将设备看作一种特殊的文件，并在文件系统中有代表着具体设备的节点，称为"设备"文件，通常都在目录"/dev"中。例如 IDE 硬盘上的第一个分区就是/dev/hda1。每个设备文件实际上只是一个索引节点，节点中提供了设备的"设备号"，由"主设备号"和"次设备号"两部分构成。其中主设备号指明了设备的种类，或者更确切地说是指明了应该使用哪一组驱动程序。同一个物理的设备，如果有两组不同的驱动程序，在逻辑上就被视作两种不同的设备而在文件系统中有两个不同的"设备文件"。次设备号则指明了该设备时同种设备中的第几个。所以，只要找到代表着某个设备的索引节点，就知道该怎样读/写这个设备了。既然是一个"可访问"的块设备，那为什么还要安装呢？答案是在安装之前可访问的只是这个设备，通常是作为一个线性的无结构的字节流来访问的，称为"原始设备"（raw device）。而在设备上的文件系统则是不可访问的。经过安装之后，设备上的文件系统就称为可访问的了。

这里有一个文件就是：系统调用 mount 要求被安装的块设备在安装之前就是可访问的，那根设备怎么办？在安装根设备之前，系统中只有一个"/"节点，根本就不存在可访问的块设备啊。其实，明白了"可访问的块设备"的含义就好解决这个问题了。"可访问的块设备"只是提供了设备的主设备号和次设备号，如果可以通过其他途径得到设备的主设备号和次设备号（比喻通过内核启动参数），那么同样可以安装设备上的文件系统。根设备的安装就是通过这种方式来的。事实上，根据情况的不同，内核中有三个函数是用于设备安装的，那就是 sys_mount，mount_root 以及 kem_mount。

## 5.2.2 sys_mount 文件系统的安装

```
2130 SYSCALL_DEFINE5(mount, char __user *, dev_name, char __user *, dir_name,
2131          char __user *, type, unsigned long, flags, void __user *, data)
```

参数：

■ dev_name：设备名称（如/dev/hda1）

■ dir_name：安装点路径

■ type：安装类型

```
2132 {
2133     int ret;
2134     char *kernel_type;
2135     char *kernel_dir;
2136     char *kernel_dev;
2137     unsigned long data_page;
2138
2139     ret = copy_mount_string(type, &kernel_type);
2140     if (ret < 0)
2141         goto out_type;
2142
2143     kernel_dir = getname(dir_name);
2144     if (IS_ERR(kernel_dir)) {
2145         ret = PTR_ERR(kernel_dir);
2146         goto out_dir;
2147     }
2148
2149     ret = copy_mount_string(dev_name, &kernel_dev);
2150     if (ret < 0)
2151         goto out_dev;
2152
2153     ret = copy_mount_options(data, &data_page);
2154     if (ret < 0)
2155         goto out_data;
```

将对应的参数从用户空间拷贝到内核空间。

```
2156
2157     ret = do_mount(kernel_dev, kernel_dir, kernel_type, flags,
2158         (void *) data_page);
```

do_mount 实现具体功能。

```
2159
```

```
2160      free_page(data_page);
2161 out_data:
2162      kfree(kernel_dev);
2163 out_dev:
2164      putname(kernel_dir);
2165 out_dir:
2166      kfree(kernel_type);
2167 out_type:
2168      return ret;
2169 }
```

### 5.2.2.1  do_mount

```
1950 long do_mount(char *dev_name, char *dir_name, char *type_page,
1951              unsigned long flags, void *data_page)
1952 {
1953      struct path path;
1954      int retval = 0;
1955      int mnt_flags = 0;
1956
1957      /* Discard magic */
1958      if ((flags & MS_MGC_MSK) == MS_MGC_VAL)
1959          flags &= ~MS_MGC_MSK;
1960
1961      /* Basic sanity checks */
1962
1963      if (!dir_name || !*dir_name || !memchr(dir_name, 0, PAGE_SIZE))
1964          return -EINVAL;
1965
1966      if (data_page)
1967          ((char *)data_page)[PAGE_SIZE - 1] = 0;
```

| | |
|---|---|
| 1968 | |
| 1969 | /* ... and get the mountpoint */ |
| 1970 | retval = kern_path(dir_name, LOOKUP_FOLLOW, &path); |

得到安装点

| | |
|---|---|
| 1971 | if (retval) |
| 1972 | return retval; |
| 1973 | |
| 1974 | retval = security_sb_mount(dev_name, &path, |
| 1975 | type_page, flags, data_page); |

安全性检查

| | |
|---|---|
| 1976 | if (retval) |
| 1977 | goto dput_out; |
| 1978 | |
| 1979 | /* Default to relatime unless overriden */ |
| 1980 | if (!(flags & MS_NOATIME)) |
| 1981 | mnt_flags |= MNT_RELATIME; |
| 1982 | |
| 1983 | /* Separate the per-mountpoint flags */ |
| 1984 | if (flags & MS_NOSUID) |
| 1985 | mnt_flags |= MNT_NOSUID; |
| 1986 | if (flags & MS_NODEV) |
| 1987 | mnt_flags |= MNT_NODEV; |
| 1988 | if (flags & MS_NOEXEC) |
| 1989 | mnt_flags |= MNT_NOEXEC; |
| 1990 | if (flags & MS_NOATIME) |
| 1991 | mnt_flags |= MNT_NOATIME; |
| 1992 | if (flags & MS_NODIRATIME) |
| 1993 | mnt_flags |= MNT_NODIRATIME; |
| 1994 | if (flags & MS_STRICTATIME) |
| 1995 | mnt_flags &= ~(MNT_RELATIME | MNT_NOATIME); |
| 1996 | if (flags & MS_RDONLY) |

```
1997            mnt_flags |= MNT_READONLY;

1998

1999      flags &= ~(MS_NOSUID | MS_NOEXEC | MS_NODEV | MS_ACTIVE |

2000               MS_NOATIME | MS_NODIRATIME | MS_RELATIME| MS_KERNMOUNT |

2001               MS_STRICTATIME);

2002

2003      if (flags & MS_REMOUNT)

2004          retval = do_remount(&path, flags & ~MS_REMOUNT, mnt_flags,

2005                        data_page);

2006      else if (flags & MS_BIND)

2007          retval = do_loopback(&path, dev_name, flags & MS_REC);

2008      else if (flags & (MS_SHARED | MS_PRIVATE | MS_SLAVE | MS_UNBINDABLE))

2009          retval = do_change_type(&path, flags);

2010      else if (flags & MS_MOVE)

2011          retval = do_move_mount(&path, dev_name);

2012      else

2013          retval = do_new_mount(&path, type_page, flags, mnt_flags,

2014                        dev_name, data_page);
```

根据不同的安装方式，选择不同的调用函数。

MS_REMOUNT 表示所要求的只是改变一个原已安装的设备的安装方式。例如，原来是按"只读"方式安装的，而现在要改为"可写"方式。这种操作成为"重安装"。

MS_MOVE：改变安装路径。

在这里，只关系 do_new_mount 函数。

```
2015 dput_out:

2016      path_put(&path);

2017      return retval;

2018 }
```

## 5.2.2.2  do_new_mount

```
1681 static int do_new_mount(struct path *path, char *type, int flags,
1682                 int mnt_flags, char *name, void *data)
1683 {
1684     struct vfsmount *mnt;
1685
1686     if (!type)
1687         return -EINVAL;
1688
1689     /* we need capabilities... */
1690     if (!capable(CAP_SYS_ADMIN))
1691         return -EPERM;
1692
1693     lock_kernel();
1694     mnt = do_kern_mount(type, flags, name, data);
```

调用 do_kern_mount 完成具体功能。

```
1695     unlock_kernel();
1696     if (IS_ERR(mnt))
1697         return PTR_ERR(mnt);
1698
1699     return do_add_mount(mnt, path, mnt_flags, NULL);
```

将一个 mount 挂入安装树。

```
1700 }
```

## 5.2.2.2.1  do_kern_mount

```
1016 struct vfsmount *
1017 do_kern_mount(const char *fstype, int flags, const char *name, void *data)
1018 {
1019     struct file_system_type *type = get_fs_type(fstype);
```

根据文件系统的名称获得 file_system_type 结构

| | |
|---|---|
| 1020 | struct vfsmount *mnt; |
| 1021 | if (!type) |
| 1022 | return ERR_PTR(-ENODEV); |
| 1023 | mnt = vfs_kern_mount(type, flags, name, data); |

调用 vfs_kern_mount。

| | |
|---|---|
| 1024 | if (!IS_ERR(mnt) && (type->fs_flags & FS_HAS_SUBTYPE) && |
| 1025 | !mnt->mnt_sb->s_subtype) |
| 1026 | mnt = fs_set_subtype(mnt, fstype); |
| 1027 | put_filesystem(type); |
| 1028 | return mnt; |
| 1029 } | |
| 1030 EXPORT_SYMBOL_GPL(do_kern_mount); | |

## vfs_kern_mount

927 struct vfsmount *

928 vfs_kern_mount(struct file_system_type *type, int flags, const char *name, void *data)

929 {

| | |
|---|---|
| 930 | struct vfsmount *mnt; |
| 931 | char *secdata = NULL; |
| 932 | int error; |
| 933 | |
| 934 | if (!type) |
| 935 | return ERR_PTR(-ENODEV); |
| 936 | |
| 937 | error = -ENOMEM; |
| 938 | mnt = alloc_vfsmnt(name); |

分配一个 vfsmount 结构。

| | |
|---|---|
| 939 | if (!mnt) |
| 940 | goto out; |

| 941 | |
|---|---|
| 942 | if (flags & MS_KERNMOUNT) |
| 943 | mnt->mnt_flags = MNT_INTERNAL; |
| 944 | |
| 945 | if (data && !(type->fs_flags & FS_BINARY_MOUNTDATA)) { |
| 946 | secdata = alloc_secdata(); |
| 947 | if (!secdata) |
| 948 | goto out_mnt; |
| 949 | |
| 950 | error = security_sb_copy_data(data, secdata); |
| 951 | if (error) |
| 952 | goto out_free_secdata; |
| 953 | } |
| 954 | |
| 955 | error = type->get_sb(type, flags, name, data, mnt); |

调用 ext2_get_sb 函数读入设备的超级块内容。

| 956 | if (error < 0) |
|---|---|
| 957 | goto out_free_secdata; |
| 958 | BUG_ON(!mnt->mnt_sb); |
| 959 | WARN_ON(!mnt->mnt_sb->s_bdi); |
| 960 | |
| 961 | error = security_sb_kern_mount(mnt->mnt_sb, flags, secdata); |
| 962 | if (error) |
| 963 | goto out_sb; |
| 964 | |
| 965 | /* |
| 966 | * filesystems should never set s_maxbytes larger than MAX_LFS_FILESIZE |
| 967 | * but s_maxbytes was an unsigned long long for many releases. Throw |
| 968 | * this warning for a little while to try and catch filesystems that |
| 969 | * violate this rule. This warning should be either removed or |

| |
|---|
| 970　　　 * converted to a BUG() in 2.6.34. |
| 971　　　 */ |
| 972　　　 WARN((mnt->mnt_sb->s_maxbytes < 0), "%s set sb->s_maxbytes to " |
| 973　　　　　 "negative value (%lld)\n", type->name, mnt->mnt_sb->s_maxbytes); |
| 974 |
| 975　　　 mnt->mnt_mountpoint = mnt->mnt_root; |
| 976　　　 mnt->mnt_parent = mnt; |

设置 mnt 的相关字段。

| |
|---|
| 977　　　 up_write(&mnt->mnt_sb->s_umount); |
| 978　　　 free_secdata(secdata); |
| 979　　　 return mnt; |
| 980 out_sb: |
| 981　　　 dput(mnt->mnt_root); |
| 982　　　 deactivate_locked_super(mnt->mnt_sb); |
| 983 out_free_secdata: |
| 984　　　 free_secdata(secdata); |
| 985 out_mnt: |
| 986　　　 free_vfsmnt(mnt); |
| 987 out: |
| 988　　　 return ERR_PTR(error); |
| 989 } |

对于 ext2 文件系统，函数会调用 ext2_get_sb。

| |
|---|
| 1368 static int ext2_get_sb(struct file_system_type *fs_type, |
| 1369　　　 int flags, const char *dev_name, void *data, struct vfsmount *mnt) |
| 1370 { |
| 1371　　　 return get_sb_bdev(fs_type, flags, dev_name, data, ext2_fill_super, mnt); |
| 1372 } |

## get_sb_bdev

```
784 int get_sb_bdev(struct file_system_type *fs_type,
785      int flags, const char *dev_name, void *data,
786      int (*fill_super)(struct super_block *, void *, int),
787      struct vfsmount *mnt)
788 {
789      struct block_device *bdev;
790      struct super_block *s;
791      fmode_t mode = FMODE_READ;
792      int error = 0;
793
794      if (!(flags & MS_RDONLY))
795          mode |= FMODE_WRITE;
796
797      bdev = open_bdev_exclusive(dev_name, mode, fs_type);
```

该函数根据设备名找到设备对应的 block_device 结构。

```
798      if (IS_ERR(bdev))
799          return PTR_ERR(bdev);
800
801      /*
802       * once the super is inserted into the list by sget, s_umount
803       * will protect the lockfs code from trying to start a snapshot
804       * while we are mounting
805       */
806      mutex_lock(&bdev->bd_fsfreeze_mutex);
807      if (bdev->bd_fsfreeze_count > 0) {
808          mutex_unlock(&bdev->bd_fsfreeze_mutex);
809          error = -EBUSY;
810          goto error_bdev;
811      }
```

```
812        s = sget(fs_type, test_bdev_super, set_bdev_super, bdev);
```

调用 sget()搜索文件系统的超级块对象链表。 如果找到一个与块设备相关的超级块，则返回它的地址。否则，分配并初始化一个新的超级块对象，把它插入到文件系统链表和超级块全局链表中，并返回其地址。

```
813        mutex_unlock(&bdev->bd_fsfreeze_mutex);

814        if (IS_ERR(s))

815            goto error_s;

816

817        if (s->s_root) {
```

如果不是新超级块。

```
818            if ((flags ^ s->s_flags) & MS_RDONLY) {

819                deactivate_locked_super(s);

820                error = -EBUSY;

821                goto error_bdev;

822            }

823

824            close_bdev_exclusive(bdev, mode);

825        } else {

826            char b[BDEVNAME_SIZE];

827

828            s->s_flags = flags;

829            s->s_mode = mode;

830            strlcpy(s->s_id, bdevname(bdev, b), sizeof(s->s_id));

831            sb_set_blocksize(s, block_size(bdev));

832            error = fill_super(s, data, flags & MS_SILENT ? 1 : 0);
```

如果是新超级块，则调用 ext2_fill_super 函数来从硬盘读入超级块。ext2_fill_super 在 4.1 中已经介绍。

```
833            if (error) {

834                deactivate_locked_super(s);

835                goto error;

836            }
```

```
837
838          s->s_flags |= MS_ACTIVE;
839          bdev->bd_super = s;
840     }
841
842     simple_set_mnt(mnt, s);
843     return 0;
844
845 error_s:
846     error = PTR_ERR(s);
847 error_bdev:
848     close_bdev_exclusive(bdev, mode);
849 error:
850     return error;
851 }
```

### 5.2.2.2.2  do_add_mount

　　在 do_new_mount 中，调用 do_kern_mount 后，待安装设备的 super_block 已经解决了，这边已经没有问题了。现在会来看看安装点这边。在读入超级块之前，已经通过 path_init 和 path_walk 找到了安装点的 dentry 结构、inode 结构以及 vfsmount 结构。现在，需要考虑一个问题：从设备上读入超级块是一个漫长的过程，当前进程在等待从设备上读入的过程中几乎肯定要睡眠，这样就可能会发生另一个设备捷足先登抢先将另一个设备安装到了同一个安装点上。如何探测这种情况，在 do_add_mount 函数有处理。

```
1706 int do_add_mount(struct vfsmount *newmnt, struct path *path,
1707          int mnt_flags, struct list_head *fslist)
1708 {
1709     int err;
1710
1711     mnt_flags &= ~(MNT_SHARED | MNT_WRITE_HOLD |
MNT_INTERNAL);
```

```
1712
1713        down_write(&namespace_sem);
1714        /* Something was mounted here while we slept */
1715        while (d_mountpoint(path->dentry) &&
1716                follow_down(path))
1717                ;
```

处理安装点是否已经被安装。如果安装点上已经有设备，从 follow_down
上可以看出，调用 follow_down 函数前进到已安装设备上的根目录，并且要
通过 while 循环进一步检测新的安装点（是否已有设备安装在其上），直到
尽头，即前进到不再有设备安装的某个设备上的根节点为准。已安装设备的
根目录下一般有内容，是否可以把一个设备安装到一个非空的目录节点上
呢？这个是可以的。

```
1718        err = -EINVAL;
1719        if (!(mnt_flags & MNT_SHRINKABLE) && !check_mnt(path->mnt))
1720            goto unlock;
1721
1722        /* Refuse the same filesystem on the same mount point */
1723        err = -EBUSY;
1724        if (path->mnt->mnt_sb == newmnt->mnt_sb &&
1725            path->mnt->mnt_root == path->dentry)
1726            goto unlock;
1727
1728        err = -EINVAL;
1729        if (S_ISLNK(newmnt->mnt_root->d_inode->i_mode))
1730            goto unlock;
```

如果要安装的文件系统已经被安装或者安装点是一个符号链接，则返回
错误。

```
1731
1732        newmnt->mnt_flags = mnt_flags;
1733        if ((err = graft_tree(newmnt, path)))
```

将新安装的文件系统对象插入到 namespace 链表、散列表中。

```
1734            goto unlock;
```

```
1735
1736     if (fslist) /* add to the specified expiration list */
1737          list_add_tail(&newmnt->mnt_expire, fslist);
1738
1739     up_write(&namespace_sem);
1740     return 0;
1741
1742 unlock:
1743     up_write(&namespace_sem);
1744     mntput(newmnt);
1745     return err;
1746 }
```

d_mountpoint 函数检测目录 dentry 是否已经有设备安装在其上。

include/linux/dcache.h

```
387 static inline int d_mountpoint(struct dentry *dentry)
388 {
389      return dentry->d_mounted;
390 }
```

### 5.2.2.2.2.1    follow_down函数

fs/namei.c

```
652 /* no need for dcache_lock, as serialization is taken care in
653  * namespace.c
654  */
655 int follow_down(struct path *path)
656 {
657     struct vfsmount *mounted;
658
659     mounted = lookup_mnt(path);
660     if (mounted) {
```

| | |
|---|---|
| 661 | dput(path->dentry); |
| 662 | mntput(path->mnt); |
| 663 | path->mnt = mounted; |
| 664 | path->dentry = dget(mounted->mnt_root); |

前进到安装点的根目录。

| | |
|---|---|
| 665 | return 1; |
| 666 | } |
| 667 | return 0; |
| 668 | } |

## lookup_mnt

| | |
|---|---|
| 439 | struct vfsmount *lookup_mnt(struct path *path) |
| 440 | { |
| 441 | struct vfsmount *child_mnt; |
| 442 | spin_lock(&vfsmount_lock); |
| 443 | if ((child_mnt = __lookup_mnt(path->mnt, path->dentry, 1))) |

根据 mnt，dentry 查找安装点 vfsmount 的子 vfsmount。

| | |
|---|---|
| 444 | mntget(child_mnt); |
| 445 | spin_unlock(&vfsmount_lock); |
| 446 | return child_mnt; |
| 447 | } |

## __lookup_mnt

| | |
|---|---|
| 410 | /* |
| 411 | * find the first or last mount at @dentry on vfsmount @mnt depending on |
| 412 | * @dir. If @dir is set return the first mount else return the last mount. |
| 413 | */ |
| 414 | struct vfsmount *__lookup_mnt(struct vfsmount *mnt, struct dentry *dentry, |
| 415 | int dir) |
| 416 | { |
| 417 | struct list_head *head = mount_hashtable + hash(mnt, dentry); |

```
418        struct list_head *tmp = head;

419        struct vfsmount *p, *found = NULL;

420

421        for (;;) {

422            tmp = dir ? tmp->next : tmp->prev;

423            p = NULL;

424            if (tmp == head)

425                break;

426            p = list_entry(tmp, struct vfsmount, mnt_hash);

427            if (p->mnt_parent == mnt && p->mnt_mountpoint == dentry) {

428                found = p;

429                break;

430            }
```

子安装点必须满足两个条件：

- 子 vfsmount 的 mnt_parent 必须指向父 vfsmount。
- 子 vfsmount 的 mnt_mountpoint 必须指向安装点。

```
431    }

432    return found;

433 }
```

### 427-430 行代码解析

把一个设备安装到一个目录节点时要用一个 vfsmount 数据结构作为"连接件"，该数据结构定义如下。

```
50 struct vfsmount {

51        struct list_head mnt_hash;

52        struct vfsmount *mnt_parent;      /* fs we are mounted on */

53        struct dentry *mnt_mountpoint;   /* dentry of mountpoint */

54        struct dentry *mnt_root;      /* root of the mounted tree */

55        struct super_block *mnt_sb; /* pointer to superblock */

56        struct list_head mnt_mounts;      /* list of children, anchored here */

57        struct list_head mnt_child; /* and going through their mnt_child */
```

```
58      int mnt_flags;

59      /* 4 bytes hole on 64bits arches */

60      const char *mnt_devname;      /* Name of device e.g. /dev/dsk/hda1 */

61      struct list_head mnt_list;

62      struct list_head mnt_expire;      /* link in fs-specific expiry list */

63      struct list_head mnt_share; /* circular list of shared mounts */

64      struct list_head mnt_slave_list;/* list of slave mounts */

65      struct list_head mnt_slave; /* slave list entry */

66      struct vfsmount *mnt_master;      /* slave is on master->mnt_slave_list
*/

67      struct mnt_namespace *mnt_ns;     /* containing namespace */

68      int mnt_id;            /* mount identifier */

69      int mnt_group_id;         /* peer group identifier */

70      /*

71       * We put mnt_count & mnt_expiry_mark at the end of struct vfsmount

72       * to let these frequently modified fields in a separate cache line

73       * (so that reads of mnt_flags wont ping-pong on SMP machines)

74       */

75      atomic_t mnt_count;

76      int mnt_expiry_mark;          /* true if marked for expiry */

77      int mnt_pinned;

78      int mnt_ghosts;

79 #ifdef CONFIG_SMP

80      int __percpu *mnt_writers;

81 #else

82      int mnt_writers;

83 #endif

84 };
```

- ■ mnt_mountpoint，mount_root：mnt_mountpoing 指向安装点的 dentry
  数据结构,mount_root 指向所安装设备上的根目录的 dentry 数据结构。
  而则在安装点和安装设备上建立连接。

- dentry 结构中没有直接指向 vfsmount 数据结构的指针，有一个队列头 d_subdirs，这是因为安装点和设备之间是一对多的关系，在同一个安装点上可以安装多个设备。

- mnt_sb 指向所安装设备的超级块 super_block 数据结构。反之，在所安装设备的 super_block 数据结构中并没有直接指向 vfsmount 的数据结构，而是有个队列头 s_instalnces，因为设备与安装点也是一对多的关系，同一个设备可以安装到多个安装点上。

- 指针 mnt_parent 指向安装点所在设备当初安装时的 vfsmount 数据结构，就是上一层的 vfsmount。

所安装设备的 super_block 数据结构与作为"连接件"的 vfsmount 之间可以是一对多的关系，这个很容易理解。可是，安装点 dentry 与 vfsmount 也是一对多的关系，这个就不容易理解了。

通过一个例子来看这个问题：

```
        /d11                    /d12
     (/dev/hda1)             (/dev/hda1)
          ↑                      ↑
           \                    /
            \                  /
              /d2
           (/dev/hda2)
            ↑      ↑
           /        \
          /          \
    /d11/d2            /d12/d2
  (/dev/hda3)        (/dev/hda4)
```

如上图所示：

假定有/dev/hda1，/dev/hda2，/dev/hda3，/dev/hda4四个设备，/dev/hda1 为根设备（这四个设备文件节点都在/dev/hda1 上的/dev 目录下），并且，在/dev/hda1 的根目录下有两个空闲的目录节点/d11 和/d12，而在/dev/hda2 的根目录下则有个空闲的目录节点 d2。现在将/dev/hda2 分别安装到/d11 和/d12 上去，这当然是可以的。这样，/dll/d2 和/d12/d2 两个路径通往同一个物理的目录节点。然后，把/dev/hda3 安装到/d11/d2 上，这样/d11/d2 就代表了着/dev/hda3 了。可是/d12/d2 呢？显然应该是空的，因为/d12 代表着一棵独立的子树。再把、/dev/hda4 安装到/d12/d2 上，这也是允许的。

现在，/dev/hda2 上的目录节点 d2 就安装两个设备了，从而有两个 vfsmount 数据结构在其 dentry 结构的 d_ subdirs 队列中。

现在有个问题，当沿路径名搜索时，发现 d2 是个安装点而要前进到所安装的设备上时，怎么知道是前进到/dev/hda3 还是/dev/hda4 呢？显然需要根据搜索路径时的上下文来决定。具体地，要看是顺着/dev/hda2 的那一次安装（每一次安装都有一个 vfsmount）（/d11/d2 或/d12/d2）搜索下来，而上一层的 vfsmount 数据结构就是这个上下文。

### 5.2.2.2.2.2  graft_tree

```
1423 static int graft_tree(struct vfsmount *mnt, struct path *path)
1424 {
1425     int err;
1426     if (mnt->mnt_sb->s_flags & MS_NOUSER)
1427         return -EINVAL;
1428
1429     if (S_ISDIR(path->dentry->d_inode->i_mode) !=
1430             S_ISDIR(mnt->mnt_root->d_inode->i_mode))
1431         return -ENOTDIR;
1432
1433     err = -ENOENT;
1434     mutex_lock(&path->dentry->d_inode->i_mutex);
1435     if (cant_mount(path->dentry))
1436         goto out_unlock;
1437
1438     err = security_sb_check_sb(mnt, path);
1439     if (err)
1440         goto out_unlock;
1441
1442     err = -ENOENT;
1443     if (!d_unlinked(path->dentry))
1444         err = attach_recursive_mnt(mnt, path, NULL);
1445 out_unlock:
1446     mutex_unlock(&path->dentry->d_inode->i_mutex);
```

```
1447        if (!err)
1448            security_sb_post_addmount(mnt, path);
1449        return err;
1450 }
```

```
1376 static int attach_recursive_mnt(struct vfsmount *source_mnt,
1377                struct path *path, struct path *parent_path)
1378 {
1379        LIST_HEAD(tree_list);
1380        struct vfsmount *dest_mnt = path->mnt;
1381        struct dentry *dest_dentry = path->dentry;
1382        struct vfsmount *child, *p;
1383        int err;
1384
1385        if (IS_MNT_SHARED(dest_mnt)) {
1386            err = invent_group_ids(source_mnt, true);
1387            if (err)
1388                goto out;
1389        }
1390        err = propagate_mnt(dest_mnt, dest_dentry, source_mnt, &tree_list);
1391        if (err)
1392            goto out_cleanup_ids;
1393
1394        spin_lock(&vfsmount_lock);
1395
1396        if (IS_MNT_SHARED(dest_mnt)) {
1397            for (p = source_mnt; p; p = next_mnt(p, source_mnt))
1398                set_mnt_shared(p);
1399        }
1400        if (parent_path) {
1401            detach_mnt(source_mnt, parent_path);
```

```
1402            attach_mnt(source_mnt, path);
1403            touch_mnt_namespace(parent_path->mnt->mnt_ns);
1404        } else {
1405            mnt_set_mountpoint(dest_mnt, dest_dentry, source_mnt);
1406            commit_tree(source_mnt);
1407        }
1408
1409        list_for_each_entry_safe(child, p, &tree_list, mnt_hash) {
1410            list_del_init(&child->mnt_hash);
1411            commit_tree(child);
1412        }
1413        spin_unlock(&vfsmount_lock);
1414        return 0;
1415
1416    out_cleanup_ids:
1417        if (IS_MNT_SHARED(dest_mnt))
1418            cleanup_group_ids(source_mnt, NULL);
1419    out:
1420        return err;
1421 }
```

安装块设备到安装点。

```
481 void mnt_set_mountpoint(struct vfsmount *mnt, struct dentry *dentry,
482                struct vfsmount *child_mnt)
483 {
484     child_mnt->mnt_parent = mntget(mnt);
485     child_mnt->mnt_mountpoint = dget(dentry);
486     dentry->d_mounted++;
487 }
```

## 5.2.3  根文件系统的安装

在 5.2 的开头讲了根文件系统的安装比较特殊。下面来看看具体的实现。根文件系统的安装时系统初始化中重要的一部分，这个过程比较复杂。

前面讲过，根文件系统的设备号是通过内核参数传递过来的。

当系统启动时，内核就要在变量 ROOT_DEV 中寻找包含根文件系统的磁盘主设备号：

//init/Do_mounts.c

dev_t ROOT_DEV;

当编译内核时，或者向最初的启动装入程序传递一个合适的"root"选项时，根文件系统可以被指定为/dev 目录下的一个设备文件。类似地，根文件系统的安装标志存放在

root_mountflags 变量中：

//init/Do_mounts.c

int root_mountflags = MS_RDONLY | MS_SILENT;

用户可以指定这些标志，或者通过对已编译的内核映像使用 rdev 外部程序，或者向最初的启动装入程序传递一个合适的 rootflags 选项来达到。

安装根文件系统分两个阶段：

（1）内核安装特殊 rootfs 文件系统，该文件系统仅提供一个作为初始安装点的空目录。

（2）内核在空目录上安装实际根文件系统。

为什么内核不怕麻烦，要在安装实际根文件系统之前安装 rootfs 文件系统呢？这是因为，rootfs 文件系统允许内核容易地改变实际根文件系统。实际上，在大多数情况下，系统初始化是内核会逐个地安装和卸载几个根文件系统。例如，一个发布版的初始启动光盘可能把具有一组最小驱动程序的内核装人 RAM 中，内核把存放在 ramdisk 中的一个最小的文件系统作为根安装。接下来，在这个初始根文件系统中的程序探测系统的硬件（例如，它们判断硬盘是否是 EIDE、SCSI 等等），装入所有必需的内核模块，并从物理块设备重新安装根文件系统。

## 5.2.3.1  安装rootfs文件系统

在内核初始化过程中，会调用 mnt_init 函数。

```
2321 void __init mnt_init(void)
2322 {
```

```
2323        unsigned u;

2324        int err;

2325

2326        init_rwsem(&namespace_sem);

2327

2328        mnt_cache = kmem_cache_create("mnt_cache", sizeof(struct
vfsmount),

2329                0, SLAB_HWCACHE_ALIGN | SLAB_PANIC, NULL);

2330

2331        mount_hashtable = (struct list_head
*)__get_free_page(GFP_ATOMIC);

2332

2333        if (!mount_hashtable)

2334            panic("Failed to allocate mount hash table\n");

2335

2336        printk("Mount-cache hash table entries: %lu\n", HASH_SIZE);

2337

2338        for (u = 0; u < HASH_SIZE; u++)

2339            INIT_LIST_HEAD(&mount_hashtable[u]);

2340

2341        err = sysfs_init();

2342        if (err)

2343            printk(KERN_WARNING "%s: sysfs_init error: %d\n",

2344                __func__, err);

2345        fs_kobj = kobject_create_and_add("fs", NULL);

2346        if (!fs_kobj)

2347            printk(KERN_WARNING "%s: kobj create error\n", __func__);

2348        init_rootfs();

2349        init_mount_tree();
```

init_rootfs 和 init_mount_tree 函数实现 rootfs 文件系统的注册等初始化
步骤。

```
2350 }
```

fs/ramfs/inode.c

289 static struct file_system_type rootfs_fs_type = {

290     .name         = "rootfs",

291     .get_sb      = rootfs_get_sb,

292     .kill_sb    = kill_litter_super,

293 };

**init_rootfs**

308 int __init init_rootfs(void)

309 {

310     int err;

311

312     err = bdi_init(&ramfs_backing_dev_info);

313     if (err)

314         return err;

315

316     err = register_filesystem(&rootfs_fs_type);

注册 rootfs 文件系统。

317     if (err)

318         bdi_destroy(&ramfs_backing_dev_info);

319

320     return err;

321 }


**init_mount_tree**

2298 static void __init init_mount_tree(void)

2299 {

2300     struct vfsmount *mnt;

2301     struct mnt_namespace *ns;

2302     struct path root;

2303

| | |
|---|---|
| 2304 | mnt = do_kern_mount("rootfs", 0, "rootfs", NULL); |

安装 rootfs 文件系统。

| | |
|---|---|
| 2305 | if (IS_ERR(mnt)) |
| 2306 | panic("Can't create rootfs"); |
| 2307 | ns = create_mnt_ns(mnt); |
| 2308 | if (IS_ERR(ns)) |
| 2309 | panic("Can't allocate initial namespace"); |
| 2310 | |

| | |
|---|---|
| 2311 | init_task.nsproxy->mnt_ns = ns; |
| 2312 | get_mnt_ns(ns); |
| 2313 | |
| 2314 | root.mnt = ns->root; |
| 2315 | root.dentry = ns->root->mnt_root; |
| 2316 | |
| 2317 | set_fs_pwd(current->fs, &root); |
| 2318 | set_fs_root(current->fs, &root); |
| 2319 } | |

### 5.2.3.2 实际根文件系统安装

在内核初始化的最后阶段，会调用 prepare_namespace 函数进行实际根文件系统的安装。

init/do_mounts.c

| | |
|---|---|
| 366 void __init prepare_namespace(void) | |
| 367 { | |
| 368 | int is_floppy; |
| 369 | |
| 370 | if (root_delay) { |
| 371 | printk(KERN_INFO "Waiting %dsec before mounting root |
| device...\n", | |
| 372 | root_delay); |

```
373            ssleep(root_delay);
374      }
375
376      /*
377       * wait for the known devices to complete their probing
378       *
379       * Note: this is a potential source of long boot delays.
380       * For example, it is not atypical to wait 5 seconds here
381       * for the touchpad of a laptop to initialize.
382       */
383      wait_for_device_probe();
384
385      md_run_setup();
386
387      if (saved_root_name[0]) {
388            root_device_name = saved_root_name;
```

把 root_device_name 变量置为从启动参数"root"中获取的设备文件名。同样，把 ROOT_DEV 变量置为同一设备文件的主设备号和次设备号。

```
389            if (!strncmp(root_device_name, "mtd", 3) ||
390                !strncmp(root_device_name, "ubi", 3)) {
391                mount_block_root(root_device_name, root_mountflags);
```

调用 mount_block_root()函数，将最常用的块设备作为 rootfs 文件系统的子文件系统。

```
392                goto out;
393            }
394            ROOT_DEV = name_to_dev_t(root_device_name);
395            if (strncmp(root_device_name, "/dev/", 5) == 0)
396                root_device_name += 5;
397      }
398
399      if (initrd_load())
```

```
400              goto out;
```

如果加载用户指定的设备失败，则根据内核配置参数（是否设置了加载 initrd 的参数），决定是否加载 initrd。

```
401
402     /* wait for any asynchronous scanning to complete */
403     if ((ROOT_DEV == 0) && root_wait) {
404          printk(KERN_INFO "Waiting for root device %s...\n",
405               saved_root_name);
406          while (driver_probe_done() != 0 ||
407               (ROOT_DEV = name_to_dev_t(saved_root_name)) == 0)
408               msleep(100);
409          async_synchronize_full();
410     }
411
412     is_floppy = MAJOR(ROOT_DEV) == FLOPPY_MAJOR;
413
414     if (is_floppy && rd_doload && rd_load_disk(0))
415          ROOT_DEV = Root_RAM0;
416
417     mount_root();//大部分时间会调用该函数。
```

如果上面都失败，调用 mount_root 进一步努力。mount_root 可以看做是 mount_block_root 封装，只是将文件的设备号变为文件路径。

```
418 out:
419     devtmpfs_mount("dev");
420     sys_mount(".", "/", NULL, MS_MOVE, NULL);
```

将当前的目录移动到根目录下。

```
421     sys_chroot(".");
422 }
```

### 5.2.3.2.1 mount_block_root

```
233 void __init mount_block_root(char *name, int flags)
234 {
235     char *fs_names = __getname_gfp(GFP_KERNEL
236         | __GFP_NOTRACK_FALSE_POSITIVE);
237     char *p;
238 #ifdef CONFIG_BLOCK
239     char b[BDEVNAME_SIZE];
240 #else
241     const char *b = name;
242 #endif
243
244     get_fs_names(fs_names);
245 retry:
246     for (p = fs_names; *p; p += strlen(p)+1) {
247         int err = do_mount_root(name, p, flags, root_mount_data);
```

完成实际的操作。

```
248         switch (err) {
249             case 0:
250                 goto out;
251             case -EACCES:
252                 flags |= MS_RDONLY;
253                 goto retry;
254             case -EINVAL:
255                 continue;
256         }
257             /*
258          * Allow the user to distinguish between failed sys_open
259          * and bad superblock on root device.
260          * and give them a list of the available devices
261          */
```

```
262 #ifdef CONFIG_BLOCK
263         __bdevname(ROOT_DEV, b);
264 #endif
265         printk("VFS: Cannot open root device \"%s\" or %s\n",
266                 root_device_name, b);
267         printk("Please append a correct \"root=\" boot option; here are the
available partitions:\n");
268
269         printk_all_partitions();
270 #ifdef CONFIG_DEBUG_BLOCK_EXT_DEVT
271         printk("DEBUG_BLOCK_EXT_DEVT is enabled, you need to
specify "
272                 "explicit textual name for \"root=\" boot option.\n");
273 #endif
274         panic("VFS: Unable to mount root fs on %s", b);
275     }
276
277     printk("List of all partitions:\n");
278     printk_all_partitions();
279     printk("No filesystem could mount root, tried: ");
280     for (p = fs_names; *p; p += strlen(p)+1)
281         printk(" %s", p);
282     printk("\n");
283 #ifdef CONFIG_BLOCK
284     __bdevname(ROOT_DEV, b);
285 #endif
286     panic("VFS: Unable to mount root fs on %s", b);
287 out:
288     putname(fs_names);
289 }
```

**do_mount_root**

```
218 static int __init do_mount_root(char *name, char *fs, int flags, void *data)
219 {
220      int err = sys_mount(name, "/root", fs, flags, data);
```

将设备安装到 rootfs 的/root 目录下。

```
221      if (err)
222          return err;
223
224      sys_chdir("/root");
```

切换当前目录为/root。

```
225      ROOT_DEV = current->fs->pwd.mnt->mnt_sb->s_dev;
226      printk("VFS: Mounted root (%s filesystem)%s on device %u:%u.\n",
227              current->fs->pwd.mnt->mnt_sb->s_type->name,
228              current->fs->pwd.mnt->mnt_sb->s_flags & MS_RDONLY ?
229              " readonly" : "", MAJOR(ROOT_DEV), MINOR(ROOT_DEV));
230      return 0;
231 }
```

### 5.2.3.2.2 mount_root

```
334 void __init mount_root(void)
335 {
336 #ifdef CONFIG_ROOT_NFS
337      if (MAJOR(ROOT_DEV) == UNNAMED_MAJOR) {
338          if (mount_nfs_root())
339              return;
340
341          printk(KERN_ERR "VFS: Unable to mount root fs via NFS, trying
floppy.\n");
342          ROOT_DEV = Root_FD0;
343      }
344 #endif
```

```
345 #ifdef CONFIG_BLK_DEV_FD
346     if (MAJOR(ROOT_DEV) == FLOPPY_MAJOR) {
347         /* rd_doload is 2 for a dual initrd/ramload setup */
348         if (rd_doload==2) {
349             if (rd_load_disk(1)) {
350                 ROOT_DEV = Root_RAM1;
351                 root_device_name = NULL;
352             }
353         } else
354             change_floppy("root floppy");
355     }
356 #endif
357 #ifdef CONFIG_BLOCK
358     create_dev("/dev/root", ROOT_DEV);
```

在 rootfs 文件系统上创建文件/dev/root，文件表示一个设备文件，inode 中存放的是文件的设备号。

```
359     mount_block_root("/dev/root", root_mountflags);
```

调用 mount_block_root 函数安装设备文件到 rootfs 的/root 目录下。

```
360 #endif
361 }
```

## 5.2.4 sys_umount文件系统的拆卸

```
1135 SYSCALL_DEFINE2(umount, char __user *, name, int, flags)
1136 {
1137     struct path path;
1138     int retval;
1139     int lookup_flags = 0;
1140
1141     if (flags & ~(MNT_FORCE | MNT_DETACH | MNT_EXPIRE |
UMOUNT_NOFOLLOW))
```

```
1142        return -EINVAL;

1143

1144    if (!(flags & UMOUNT_NOFOLLOW))

1145        lookup_flags |= LOOKUP_FOLLOW;

1146

1147    retval = user_path_at(AT_FDCWD, name, lookup_flags, &path);

1148    if (retval)

1149        goto out;

1150    retval = -EINVAL;

1151    if (path.dentry != path.mnt->mnt_root)

1152        goto dput_and_out;

1153    if (!check_mnt(path.mnt))

1154        goto dput_and_out;

1155

1156    retval = -EPERM;

1157    if (!capable(CAP_SYS_ADMIN))

1158        goto dput_and_out;

1159

1160    retval = do_umount(path.mnt, flags);

1161 dput_and_out:

1162    /* we mustn't call path_put() as that would clear mnt_expiry_mark */

1163    dput(path.dentry);

1164    mntput_no_expire(path.mnt);

1165 out:

1166    return retval;

1167 }
```

### 5.2.4.1  do_umount

```
1043 static int do_umount(struct vfsmount *mnt, int flags)
```

```
1044 {
1045     struct super_block *sb = mnt->mnt_sb;
1046     int retval;
1047     LIST_HEAD(umount_list);
1048
1049     retval = security_sb_umount(mnt, flags);
1050     if (retval)
1051         return retval;
1052
1053     /*
1054      * Allow userspace to request a mountpoint be expired rather than
1055      * unmounting unconditionally. Unmount only happens if:
1056      *  (1) the mark is already set (the mark is cleared by mntput())
1057      *  (2) the usage count == 1 [parent vfsmount] + 1 [sys_umount]
1058      */
1059     if (flags & MNT_EXPIRE) {
1060         if (mnt == current->fs->root.mnt ||
1061             flags & (MNT_FORCE | MNT_DETACH))
1062             return -EINVAL;
1063
1064         if (atomic_read(&mnt->mnt_count) != 2)
1065             return -EBUSY;
1066
1067         if (!xchg(&mnt->mnt_expiry_mark, 1))
1068             return -EAGAIN;
1069     }
1070
1071     /*
1072      * If we may have to abort operations to get out of this
1073      * mount, and they will themselves hold resources we must
1074      * allow the fs to do things. In the Unix tradition of
```

```
1075        * 'Gee thats tricky lets do it in userspace' the umount_begin
1076        * might fail to complete on the first run through as other tasks
1077        * must return, and the like. Thats for the mount program to worry
1078        * about for the moment.
1079        */
1080
1081       if (flags & MNT_FORCE && sb->s_op->umount_begin) {
1082           sb->s_op->umount_begin(sb);
1083       }
1084
1085       /*
1086        * No sense to grab the lock for this test, but test itself looks
1087        * somewhat bogus. Suggestions for better replacement?
1088        * Ho-hum... In principle, we might treat that as umount + switch
1089        * to rootfs. GC would eventually take care of the old vfsmount.
1090        * Actually it makes sense, especially if rootfs would contain a
1091        * /reboot - static binary that would close all descriptors and
1092        * call reboot(9). Then init(8) could umount root and exec /reboot.
1093        */
1094       if (mnt == current->fs->root.mnt && !(flags & MNT_DETACH)) {
1095           /*
1096            * Special case for "unmounting" root ...
1097            * we just try to remount it readonly.
1098            */
1099           down_write(&sb->s_umount);
1100           if (!(sb->s_flags & MS_RDONLY))
1101               retval = do_remount_sb(sb, MS_RDONLY, NULL, 0);
1102           up_write(&sb->s_umount);
1103           return retval;
1104       }
```

如果要卸载的文件系统是根文件系统，且用户并不要求真正地把它卸载下来则调用 do_remount_sb()重新安装根文件系统为只读并终止

```
1105
1106        down_write(&namespace_sem);
1107        spin_lock(&vfsmount_lock);
1108        event++;
1109
1110        if (!(flags & MNT_DETACH))
1111            shrink_submounts(mnt, &umount_list);
1112
1113        retval = -EBUSY;
1114        if (flags & MNT_DETACH || !propagate_mount_busy(mnt, 2)) {
1115            if (!list_empty(&mnt->mnt_list))
1116                umount_tree(mnt, 1, &umount_list);
```

如果已安装文件系统不包含任何子安装文件系统的安装点，或者用户要求强制卸载文件系统，则调用 umount_tree()卸载文件系统（及其所有子文件系统）

```
1117            retval = 0;
1118        }
1119        spin_unlock(&vfsmount_lock);
1120        if (retval)
1121            security_sb_umount_busy(mnt);
1122        up_write(&namespace_sem);
1123        release_mounts(&umount_list);
1124        return retval;
1125 }
```

## 5.3 文件的打开与关闭

### 5.3.1 文件的打开do_sys_open

```
1066 SYSCALL_DEFINE3(open, const char __user *, filename, int, flags, int, mode)
1067 {
1068     long ret;
1069
1070     if (force_o_largefile())
1071         flags |= O_LARGEFILE;
1072
1073     ret = do_sys_open(AT_FDCWD, filename, flags, mode);
1074     /* avoid REGPARM breakage on x86: */
1075     asmlinkage_protect(3, ret, filename, flags, mode);
1076     return ret;
1077 }
```

```
1044 long do_sys_open(int dfd, const char __user *filename, int flags, int mode)
1045 {
1046     char *tmp = getname(filename);
1047     int fd = PTR_ERR(tmp);
1048
1049     if (!IS_ERR(tmp)) {
1050         fd = get_unused_fd_flags(flags);
```

取得一个为使用的文件号

```
1051         if (fd >= 0) {
1052             struct file *f = do_filp_open(dfd, tmp, flags, mode, 0);
```

具体功能函数

```
1053             if (IS_ERR(f)) {
```

```
1054                    put_unused_fd(fd);
1055                    fd = PTR_ERR(f);
1056            } else {
1057                    fsnotify_open(f->f_path.dentry);
1058                    fd_install(fd, f);
```

安装 file 结构到 fd 对应的位置。

```
1059            }
1060        }
1061        putname(tmp);
1062    }
1063    return fd;
1064 }
```

## do_filp_open

```
1755 /*
1756  * Note that the low bits of the passed in "open_flag"
1757  * are not the same as in the local variable "flag". See
1758  * open_to_namei_flags() for more details.
1759  */
1760 struct file *do_filp_open(int dfd, const char *pathname,
1761            int open_flag, int mode, int acc_mode)
1762 {
1763    struct file *filp;
1764    struct nameidata nd;
1765    int error;
1766    struct path path;
1767    int count = 0;
1768    int flag = open_to_namei_flags(open_flag);
1769    int force_reval = 0;
1770
1771    if (!(open_flag & O_CREAT))
```

```
1772            mode = 0;

1773

1774        /*

1775         * O_SYNC is implemented as __O_SYNC|O_DSYNC.   As many
places only

1776         * check for O_DSYNC if the need any syncing at all we enforce it's

1777         * always set instead of having to deal with possibly weird behaviour

1778         * for malicious applications setting only __O_SYNC.

1779         */

1780        if (open_flag & __O_SYNC)

1781            open_flag |= O_DSYNC;

1782

1783        if (!acc_mode)

1784            acc_mode = MAY_OPEN | ACC_MODE(open_flag);

1785

1786        /* O_TRUNC implies we need access checks for write permissions */

1787        if (open_flag & O_TRUNC)

1788            acc_mode |= MAY_WRITE;

1789

1790        /* Allow the LSM permission hook to distinguish append

1791           access from general write access. */

1792        if (open_flag & O_APPEND)

1793            acc_mode |= MAY_APPEND;

1794

1795        /* find the parent */

1796 reval:

1797        error = path_init(dfd, pathname, LOOKUP_PARENT, &nd);

1798        if (error)

1799            return ERR_PTR(error);

1800        if (force_reval)

1801            nd.flags |= LOOKUP_REVAL;

1802
```

```
1803        current->total_link_count = 0;

1804        error = link_path_walk(pathname, &nd);
```

调用 path_init，link_path_walk 到父目录。

```
1805        if (error) {

1806             filp = ERR_PTR(error);

1807             goto out;

1808        }

1809        if (unlikely(!audit_dummy_context()) && (open_flag & O_CREAT))

1810             audit_inode(pathname, nd.path.dentry);

1811

1812        /*

1813         * We have the parent and last component.

1814         */

1815

1816        error = -ENFILE;

1817        filp = get_empty_filp();

1818        if (filp == NULL)

1819             goto exit_parent;

1820        nd.intent.open.file = filp;

1821        filp->f_flags = open_flag;

1822        nd.intent.open.flags = flag;

1823        nd.intent.open.create_mode = mode;

1824        nd.flags &= ~LOOKUP_PARENT;

1825        nd.flags |= LOOKUP_OPEN;

1826        if (open_flag & O_CREAT) {

1827             nd.flags |= LOOKUP_CREATE;

1828             if (open_flag & O_EXCL)

1829                  nd.flags |= LOOKUP_EXCL;

1830        }

1831        if (open_flag & O_DIRECTORY)

1832             nd.flags |= LOOKUP_DIRECTORY;
```

| 1833 | if (!(open_flag & O_NOFOLLOW)) |
| 1834 | nd.flags |= LOOKUP_FOLLOW; |

设置相应的处理参数。

| 1835 | filp = do_last(&nd, &path, open_flag, acc_mode, mode, pathname); |

处理最后一个节点。

| 1836 | while (unlikely(!filp)) { /* trailing symlink */ |

处理符号链接。

| 1837 | struct path holder; |
| 1838 | struct inode *inode = path.dentry->d_inode; |
| 1839 | void *cookie; |
| 1840 | error = -ELOOP; |
| 1841 | /* S_ISDIR part is a temporary automount kludge */ |
| 1842 | if (!(nd.flags & LOOKUP_FOLLOW) && !S_ISDIR(inode->i_mode)) |
| 1843 | goto exit_dput; |
| 1844 | if (count++ == 32) |
| 1845 | goto exit_dput; |
| 1846 | /* |
| 1847 | * This is subtle. Instead of calling do_follow_link() we do |
| 1848 | * the thing by hands. The reason is that this way we have zero |
| 1849 | * link_count and path_walk() (called from ->follow_link) |
| 1850 | * honoring LOOKUP_PARENT.   After that we have the parent and |
| 1851 | * last component, i.e. we are in the same situation as after |
| 1852 | * the first path_walk().   Well, almost - if the last component |
| 1853 | * is normal we get its copy stored in nd->last.name and we will |
| 1854 | * have to putname() it when we are done. Procfs-like symlinks |
| 1855 | * just set LAST_BIND. |
| 1856 | */ |
| 1857 | nd.flags |= LOOKUP_PARENT; |
| 1858 | error = security_inode_follow_link(path.dentry, &nd); |

```
1859        if (error)
1860            goto exit_dput;
1861        error = __do_follow_link(&path, &nd, &cookie);
1862        if (unlikely(error)) {
1863            /* nd.path had been dropped */
1864            if (!IS_ERR(cookie) && inode->i_op->put_link)
1865                inode->i_op->put_link(path.dentry, &nd, cookie);
1866            path_put(&path);
1867            release_open_intent(&nd);
1868            filp = ERR_PTR(error);
1869            goto out;
1870        }
1871        holder = path;
1872        nd.flags &= ~LOOKUP_PARENT;
1873        filp = do_last(&nd, &path, open_flag, acc_mode, mode,
pathname);
```

调用 do_last 来处理函数。

```
1874        if (inode->i_op->put_link)
1875            inode->i_op->put_link(holder.dentry, &nd, cookie);
1876        path_put(&holder);
1877    }
1878 out:
1879    if (nd.root.mnt)
1880        path_put(&nd.root);
1881    if (filp == ERR_PTR(-ESTALE) && !force_reval) {
1882        force_reval = 1;
1883        goto reval;
1884    }
1885    return filp;
1886
1887 exit_dput:
```

```
1888        path_put_conditional(&path, &nd);
1889        if (!IS_ERR(nd.intent.open.file))
1890            release_open_intent(&nd);
1891 exit_parent:
1892        path_put(&nd.path);
1893        filp = ERR_PTR(error);
1894        goto out;
1895 }
```

## do_last 函数

```
1611 static struct file *do_last(struct nameidata *nd, struct path *path,
1612                    int open_flag, int acc_mode,
1613                    int mode, const char *pathname)
1614 {
1615        struct dentry *dir = nd->path.dentry;
1616        struct file *filp;
1617        int error = -EISDIR;
1618
1619        switch (nd->last_type) {
1620        case LAST_DOTDOT:
1621            follow_dotdot(nd);
1622            dir = nd->path.dentry;
1623            if (nd->path.mnt->mnt_sb->s_type->fs_flags & FS_REVAL_DOT) {
1624                if (!dir->d_op->d_revalidate(dir, nd)) {
1625                    error = -ESTALE;
1626                    goto exit;
1627                }
1628            }
1629            /* fallthrough */
1630        case LAST_DOT:
```

```
1631        case LAST_ROOT:
1632            if (open_flag & O_CREAT)
1633                goto exit;
1634            /* fallthrough */
1635        case LAST_BIND:
1636            audit_inode(pathname, dir);
1637            goto ok;
1638        }
```

特需节点处理

```
1639
1640        /* trailing slashes? */
1641        if (nd->last.name[nd->last.len]) {
1642            if (open_flag & O_CREAT)
1643                goto exit;
1644            nd->flags |= LOOKUP_DIRECTORY | LOOKUP_FOLLOW;
1645        }
1646
1647        /* just plain open? */
1648        if (!(open_flag & O_CREAT)) {
1649            error = do_lookup(nd, &nd->last, path);
1650            if (error)
1651                goto exit;
1652            error = -ENOENT;
1653            if (!path->dentry->d_inode)
1654                goto exit_dput;
1655            if (path->dentry->d_inode->i_op->follow_link)
1656                return NULL;
1657            error = -ENOTDIR;
1658            if (nd->flags & LOOKUP_DIRECTORY) {
1659                if (!path->dentry->d_inode->i_op->lookup)
1660                    goto exit_dput;
```

```
1661                }
1662                path_to_nameidata(path, nd);
1663                audit_inode(pathname, nd->path.dentry);
1664                goto ok;
1665        }
```

处理非创建的打开

```
1666
1667     /* OK, it's O_CREAT */
1668     mutex_lock(&dir->d_inode->i_mutex);
1669
1670     path->dentry = lookup_hash(nd);
1671     path->mnt = nd->path.mnt;
1672
1673     error = PTR_ERR(path->dentry);
1674     if (IS_ERR(path->dentry)) {
1675         mutex_unlock(&dir->d_inode->i_mutex);
1676         goto exit;
1677     }
1678
1679     if (IS_ERR(nd->intent.open.file)) {
1680         error = PTR_ERR(nd->intent.open.file);
1681         goto exit_mutex_unlock;
1682     }
1683
1684     /* Negative dentry, just create the file */
1685     if (!path->dentry->d_inode) {
1686         /*
1687          * This write is needed to ensure that a
1688          * ro->rw transition does not occur between
1689          * the time when the file is created and when
1690          * a permanent write count is taken through
```

```
1691                * the 'struct file' in nameidata_to_filp().
1692                */
1693               error = mnt_want_write(nd->path.mnt);
1694               if (error)
1695                    goto exit_mutex_unlock;
1696               error = __open_namei_create(nd, path, open_flag, mode);
1697               if (error) {
1698                    mnt_drop_write(nd->path.mnt);
1699                    goto exit;
1700               }
1701               filp = nameidata_to_filp(nd);
1702               mnt_drop_write(nd->path.mnt);
1703               if (!IS_ERR(filp)) {
1704                    error = ima_file_check(filp, acc_mode);
1705                    if (error) {
1706                         fput(filp);
1707                         filp = ERR_PTR(error);
1708                    }
1709               }
1710               return filp;
1711          }
```

创建文件处理

```
1712
1713     /*
1714      * It already exists.
1715      */
1716     mutex_unlock(&dir->d_inode->i_mutex);
1717     audit_inode(pathname, path->dentry);
1718
1719     error = -EEXIST;
1720     if (open_flag & O_EXCL)
```

```
1721            goto exit_dput;
1722
1723       if (__follow_mount(path)) {
1724            error = -ELOOP;
1725            if (open_flag & O_NOFOLLOW)
1726                goto exit_dput;
1727       }
1728
1729       error = -ENOENT;
1730       if (!path->dentry->d_inode)
1731            goto exit_dput;
1732
1733       if (path->dentry->d_inode->i_op->follow_link)
1734            return NULL;
1735
1736       path_to_nameidata(path, nd);
1737       error = -EISDIR;
1738       if (S_ISDIR(path->dentry->d_inode->i_mode))
1739            goto exit;
1740 ok:
1741       filp = finish_open(nd, open_flag, acc_mode);
```

inode 创建成功后，分配 file 结构。

```
1742       return filp;
1743
1744 exit_mutex_unlock:
1745       mutex_unlock(&dir->d_inode->i_mutex);
1746 exit_dput:
1747       path_put_conditional(path, nd);
1748 exit:
1749       if (!IS_ERR(nd->intent.open.file))
1750            release_open_intent(nd);
```

```
1751        path_put(&nd->path);

1752        return ERR_PTR(error);

1753 }
```

## 5.3.1.1 __open_namei_create

```
1497 /*

1498  * Be careful about ever adding any more callers of this

1499  * function.   Its flags must be in the namei format, not

1500  * what get passed to sys_open().

1501  */

1502 static int __open_namei_create(struct nameidata *nd, struct path *path,

1503                          int open_flag, int mode)

1504 {

1505        int error;

1506        struct dentry *dir = nd->path.dentry;

1507

1508        if (!IS_POSIXACL(dir->d_inode))

1509            mode &= ~current_umask();

1510        error = security_path_mknod(&nd->path, path->dentry, mode, 0);

1511        if (error)

1512            goto out_unlock;

1513        error = vfs_create(dir->d_inode, path->dentry, mode, nd);

1514 out_unlock:

1515        mutex_unlock(&dir->d_inode->i_mutex);

1516        dput(nd->path.dentry);

1517        nd->path.dentry = path->dentry;

1518        if (error)

1519            return error;

1520        /* Don't check for write permission, don't truncate */

1521        return may_open(&nd->path, 0, open_flag & ~O_TRUNC);
```

```
1522 }
```

## vfs_create

```
1403 int vfs_create(struct inode *dir, struct dentry *dentry, int mode,
1404         struct nameidata *nd)
1405 {
1406     int error = may_create(dir, dentry);
1407
1408     if (error)
1409         return error;
1410
1411     if (!dir->i_op->create)
1412         return -EACCES; /* shouldn't it be ENOSYS? */
1413     mode &= S_IALLUGO;
1414     mode |= S_IFREG;
1415     error = security_inode_create(dir, dentry, mode);
1416     if (error)
1417         return error;
1418     error = dir->i_op->create(dir, dentry, mode, nd);
1419     if (!error)
1420         fsnotify_create(dir, dentry);
```

文件系统创建通知。

```
1421     return error;
1422 }
```

## ext2_create

```
93 /*
94    * By the time this is called, we already have created
95    * the directory cache entry for the new file, but it
96    * is so far negative - it has no inode.
97    *
```

```
 98    * If the create succeeds, we fill in the inode information
 99    * with d_instantiate().
100    */
101 static int ext2_create (struct inode * dir, struct dentry * dentry, int mode,
struct nameidata *nd)
102 {
103      struct inode *inode;
104
105      dquot_initialize(dir);
106
107      inode = ext2_new_inode(dir, mode);
```

在磁盘上分配一个 inode 节点，该函数在前面已分析。

```
108      if (IS_ERR(inode))
109          return PTR_ERR(inode);
110
111      inode->i_op = &ext2_file_inode_operations;
112      if (ext2_use_xip(inode->i_sb)) {
113          inode->i_mapping->a_ops = &ext2_aops_xip;
114          inode->i_fop = &ext2_xip_file_operations;
115      } else if (test_opt(inode->i_sb, NOBH)) {
116          inode->i_mapping->a_ops = &ext2_nobh_aops;
117          inode->i_fop = &ext2_file_operations;
118      } else {
119          inode->i_mapping->a_ops = &ext2_aops;
120          inode->i_fop = &ext2_file_operations;
121      }
```

设置 inode 的相关字段

```
122      mark_inode_dirty(inode);
```

设置 inode 为脏，回写到磁盘。

```
123      return ext2_add_nondir(dentry, inode);
```

将目录项写入到父目录项中，同事将父目录的 inode 设置为脏。

124 }

## 5.3.2　文件关闭sys_close

1135 SYSCALL_DEFINE1(close, unsigned int, fd)

1136 {

1137　　struct file * filp;

1138　　struct files_struct *files = current->files;

1139　　struct fdtable *fdt;

1140　　int retval;

1141

1142　　spin_lock(&files->file_lock);

1143　　fdt = files_fdtable(files);

1144　　if (fd >= fdt->max_fds)

1145　　　goto out_unlock;

1146　　filp = fdt->fd[fd];

1147　　if (!filp)

1148　　　goto out_unlock;

1149　　rcu_assign_pointer(fdt->fd[fd], NULL);

1150　　FD_CLR(fd, fdt->close_on_exec);

1151　　__put_unused_fd(files, fd);

1152　　spin_unlock(&files->file_lock);

1153　　retval = filp_close(filp, files);

1154

1155　　/* can't restart close syscall because file table entry was cleared */

1156　　if (unlikely(retval == -ERESTARTSYS ||

1157　　　　　retval == -ERESTARTNOINTR ||

1158　　　　　retval == -ERESTARTNOHAND ||

1159　　　　　retval == -ERESTART_RESTARTBLOCK))

1160　　　retval = -EINTR;

1161

1162　　return retval;

```
1163
1164 out_unlock:
1165     spin_unlock(&files->file_lock);
1166     return -EBADF;
1167 }
1168 EXPORT_SYMBOL(sys_close);
```

```
1110 int filp_close(struct file *filp, fl_owner_t id)
1111 {
1112     int retval = 0;
1113
1114     if (!file_count(filp)) {
1115         printk(KERN_ERR "VFS: Close: file count is 0\n");
1116         return 0;
1117     }
1118
1119     if (filp->f_op && filp->f_op->flush)
1120         retval = filp->f_op->flush(filp, id);
1121
1122     dnotify_flush(filp, id);
1123     locks_remove_posix(filp, id);
1124     fput(filp);
1125     return retval;
1126 }
```

## 5.4 文件的读写

### 5.4.1 概述

在打开文件后，或者说建立起进程与文件之间的"连接"后，才能对文件进行读写。为了调高效率，复杂一些的操作系统对文件的读/写都是带缓存的，linux 也不例外。像 VFS 一样，Linux 文件系统的缓冲机制也是它的一

大特色。所谓缓冲，是指系统为最近刚读/写过的文件在内核中保留的一份副本，以便当再次需要已经缓存在副本中的内容时就不必再临时从设备上读入，而需要写的时候则可以先写入到副本中，待系统较为空闲时再重副本写入设备。在多进程的系统中，由于同一文件可能为多个进程所共享，缓冲的作用就更显著了。

然而，1>怎样实现缓冲，2>在哪一个层次上实现缓冲，却是一个值得仔细考虑的问题。在系统中处于最高层的是进程，这一层可以称为"应用层"，在用户空间运行的，在这里代表着目标文件的是"打开文件号"。在这一层中提供缓冲似乎最贴近文件内容的使用者，但是那样就需要用户的介入，从而不能做到对使用者的"透明"。并且缓冲的内容不能为其它进程所共享，显然是不妥的。在应用层以下是"文件层"，又可细分为 VFS 层和具体的文件系统层，在下面就是"设备层"了。这些层都在内核中，所以在这些层上实现缓冲都可以达到对用户透明的目标。设备层最贴近设备，即文件内容的"源头"的地方，在这里实现缓冲显然是可行的，事实上早期的 Unix 内核中的文件缓冲就是以数据块缓冲的形式在这一层上实现的。但是，设备层上的缓冲区离使用者的距离太远了一点，特别是当文件层有分为 VFS 层和具体文件系统子层是，每次读/写都要穿越这么多界面深入到设备层就难免会耗时。

在文件层中有三种主要的数据结构，就是 file、dentry、inode。

一个 file 代表着一个文件的上下文，不但不同的进程可以在同一个文件上建立不同的上下文，就是同一个进程也可以同事打开一个文件多次而建立起多个上下文。如果在 file 结构中设备一个缓冲区队列，那么缓存区中的内容虽然贴近这个特定上下文的使用者，却不便于在多个进程中共享，甚至不便于在同一个进程打开的不同上下文"共享"。这显然是不合适的，需要把这些缓存区公用的地方提取出来。
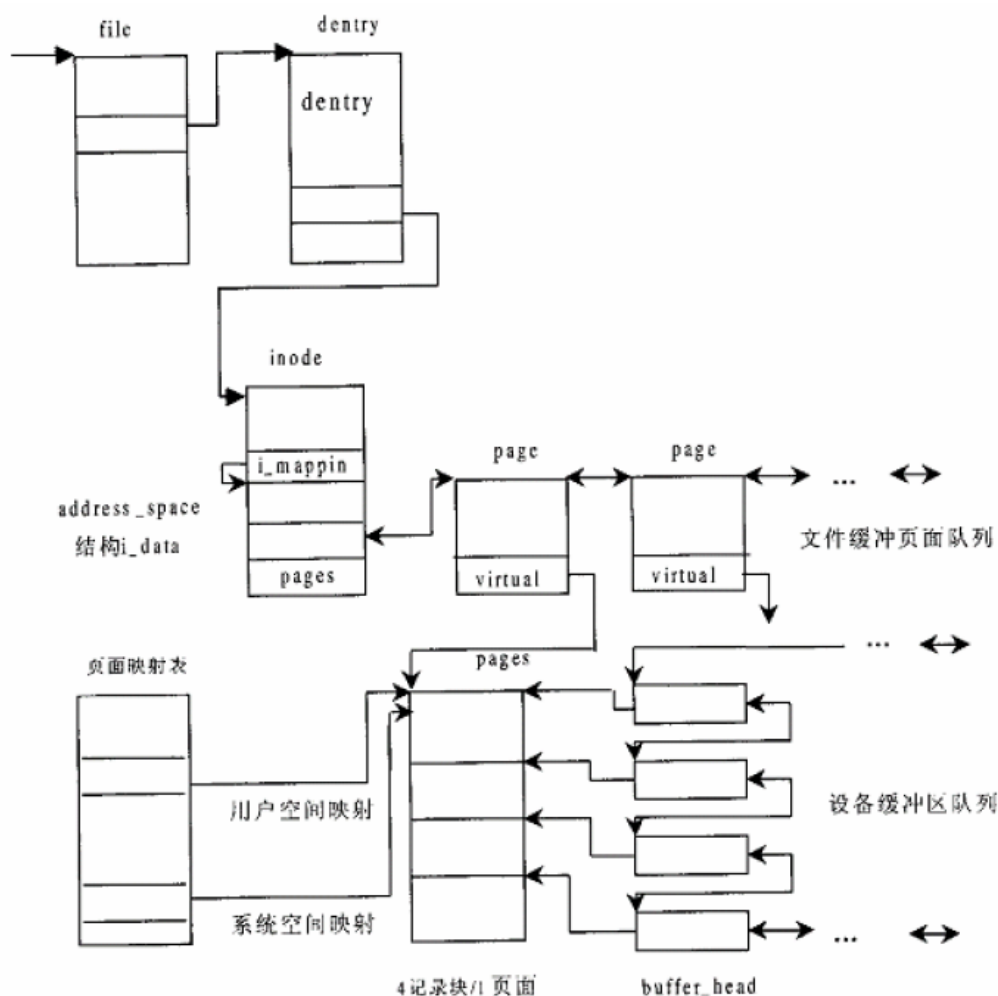
那么 dentry 结构如何？这个数据结构并不属于某一个上下文，也不属于某一个进程，可以为所有的进程和上下文共享。可是，dentry 结构与目标文件并不是一对一的关系，通过文件连接（硬连接），可以为已经存在的文件建立"别名"。一个 dentry 结构知识唯一的代表着这个文件系统中的一个节点，也就是一个路径名，但是多个节点可以同时代表同一个文件。

显然，在 inode 数据结构中设置一个缓冲区队列是最合适不过的了，首先，inode 结构与文件是一对一的关系，即使一个文件有多个路径名，最后也归纳到一个 inode 结构上。其次，一个文件中的内容是不能由其他文件共享的，在同一时间里，设备上的每一个记录块都只能至多属于一个文件（或者空闲），将载有同一文件内容的缓冲区都放在其所属文件的 inode 结构中

是很自然的是。因此，inode 结构中设置了一个指针 i_mapping，它指向一个 address_space 数据结构（通常这个数据结构就是 inode 结构中的 i_data0），缓存区队列就在这个数据结构中。

不过，挂在缓冲区队列中的并不是记录块而是内存页面。也就是说，文件的内容并不是以记录块为单位，而是已页面为单位进行缓冲的。如果记录块的大小为 1k 字节，那么一个页面就相当于 4 个记录块。为什么要这么做呢？这是为了将文件内容的缓冲与文件的内存映射结合在一起。在内存管理中讲过，一个进程可以通过 mmap 将一个文件映射到它的用户空间。建立了这样的映射后，就可以像访问内存一样访问这个文件。如果将文件的内容以页面为单位进行缓存，放在属于该文件的 inode 结构的缓冲队列中，那么只需要相应地设置进程的内存映射表，就可以自然地将这些缓冲页面映射到进程的用户空间中。这样，在按常规的文件访问一个文件时，可以通过 read 和 write 系统调用目标文件的 inode 结构访问这些缓冲页面；而通过内存映射机制访问这个文件时，就可以经由页面映射表直接读写这些缓冲着的页面。当目标页面不存在内存中时，常规文件通过系统调用 read,write 的底层将其从设备上读入，而内存映射机制则通过"缺页异常"的服务将目标页面从设备上读入。也就是所，同一个缓冲页面可以满足两方面的要求，文件系统的缓冲机制和文件的内存映射机制巧妙地结合在一起了。

可是，尽管以页面为单位的缓冲对于文件层来说确实是一个很好的选择，对于设备层则不那么合适了。对设备层而言，最自然的当然是以记录块为单位进行缓冲，因为设备的读/写都是以记录块为单位的。不过，从磁盘上读/写时主要的时间都花在准备工作（如磁头组的定位）上，一旦准备好了后读一个记录块与接连读几个记录块相差不大，而且每次只读写一个记录块反而是不经济的。所以每次读写若干连续的记录块、以页面为单位缓冲也并不成问题。另一方面，如果以页面为单位缓冲，而一个页面相当于若干个记录块，那么无论是对于缓冲页面还是对于记录块缓冲区，其控制和附加信息（如链接指针等）显然应该游离于页面之外，这些信息不应该映射到进程的用户空间。这个问题也不难解决，在"缓冲区头部"即 buffer_head 数据结构中有一个 b_data 指向缓冲区，而 buffer_head 结构本身不在缓冲区中。所以，在设备层中只要保持一些 buffer_head 结构，让它们的 b_data 指针分别指向缓冲区页面中的相应位置就可以了。这样，以页面为单位为文件内容建立缓冲是"一箭三雕"，如下示意图。

在这样一个结构框架中，一旦所欲访问的内容已经在缓冲页面队列中，读文件的效率就很高了。只要找到文件的 inode 结构（file->dentry->inode）就找到了缓冲区页面队列，从队列中找到相应的页面就可以读出了。

那么，写操作又如何呢？如前所述，一旦目标记录块已经存在与缓冲页面中，写操作只是把内容写到该缓冲页面中，所以从发动写操作的进程的角度看数度也是很快的。至于改变了内容的缓冲页面，则由系统负责在 CPU 较为空闲写入设备。

## 5.4.2 相关数据结构

### 5.4.2.1 inode

在 linux 系统中，磁盘上的一个文件由 indoe 索引节点唯一的标识。
include/linux/fs.h

```
724 struct inode {

……

754     struct address_space    *i_mapping;

755     struct address_space    i_data;

……

792 };
```

- i_mapping：指向缓存 address_space 对象的指针，该结构是文件系统缓存的关键。
- i_data：嵌入在 inode 中的 address_space 对象。

## 5.4.2.2  page

页帧代表系统内存的最小单位，对内存中的每个页都会创建 struct page 的一个实例。

内核设计需要尽量保证该结构尽量小。

由于 page 结构会由于很多地方，所以其中有很多字段在某种情况下有用，而在另外的情况下无用，为了保证 page 字段尽量小，内核使用联合的方式将同一内存出的内容做不同的解释。

```
27 /*

28   * Each physical page in the system has a struct page associated with

29   * it to keep track of whatever it is we are using the page for at the

30   * moment. Note that we have no way to track which tasks are using

31   * a page, though if it is a pagecache page, rmap structures can tell us

32   * who is mapping it.

33   */

34 struct page {

35      unsigned long flags;          /* Atomic flags, some possibly

36                        * updated asynchronously */

37      atomic_t _count;          /* Usage count, see below. */

38      union {

39          atomic_t _mapcount; /* Count of ptes mapped in mms,
```

```c
40                          * to show when page is mapped
41                          * & limit reverse map searches.
42                          */
43         struct {        /* SLUB */
44              u16 inuse;
45              u16 objects;
46         };
47     };
48     union {
49         struct {
50         unsigned long private;      /* Mapping-private opaque data:
51                                  * usually used for buffer_heads
52                                  * if PagePrivate set; used for
53                                  * swp_entry_t if PageSwapCache;
54                                  * indicates order in the buddy
55                                  * system if PG_buddy is set.
56                                  */
57         struct address_space *mapping;   /* If low bit clear, points to
58                                  * inode address_space, or NULL.
59                                  * If page mapped as anonymous
60                                  * memory, low bit is set, and
61                                  * it points to anon_vma object:
62                                  * see PAGE_MAPPING_ANON below.
63                                  */
64         };
65 #if USE_SPLIT_PTLOCKS
66         spinlock_t ptl;
67 #endif
68         struct kmem_cache *slab;      /* SLUB: Pointer to slab */
69         struct page *first_page;      /* Compound tail pages */
70     };
```

```c
71      union {
72          pgoff_t index;          /* Our offset within mapping. */
73          void *freelist;         /* SLUB: freelist req. slab lock */
74      };
75      struct list_head lru;          /* Pageout list, eg. active_list
76                              * protected by zone->lru_lock !
77                              */
78      /*
79       * On machines where all RAM is mapped into kernel address space,
80       * we can simply calculate the virtual address. On machines with
81       * highmem some memory is mapped into kernel virtual memory
82       * dynamically, so we need a place to store that address.
83       * Note that this field could be 16 bits on x86 ... ;)
84       *
85       * Architectures with slow multiplication can define
86       * WANT_PAGE_VIRTUAL in asm/page.h
87       */
88 #if defined(WANT_PAGE_VIRTUAL)
89      void *virtual;              /* Kernel virtual address (NULL if
90                              not kmapped, ie. highmem) */
91 #endif /* WANT_PAGE_VIRTUAL */
92 #ifdef CONFIG_WANT_PAGE_DEBUG_FLAGS
93      unsigned long debug_flags;   /* Use atomic bitops on this */
94 #endif
95
96 #ifdef CONFIG_KMEMCHECK
97      /*
98       * kmemcheck wants to track the status of each byte in a page; this
99       * is a pointer to such a status block. NULL if not tracked.
100      */
101     void *shadow;
```

```
102 #endif
103 };
```

- mapping 指定了页帧所在的地址空间 address_space。
- index 字段表示在所有者的地址空间中以页大小为单位的偏移，也就是所有者的磁盘映射中页中数据的位置。
- virtual 指向本页在内核中的虚拟地址

在页高速缓存中查找页使用 mapping、index 这两个字段。

## 5.4.2.3 address_space

页高速缓冲的核心数据结构是 address_space，它嵌入在页所有者的索引节点对象的数据结构。

```
623 struct address_space {
624     struct inode          *host;        /* owner: inode, block_device */
625     struct radix_tree_root   page_tree;   /* radix tree of all pages */
626     spinlock_t         tree_lock;   /* and lock protecting it */
627     unsigned int          i_mmap_writable;/* count VM_SHARED mappings */
628     struct prio_tree_root    i_mmap;        /* tree of private and shared mappings */
629     struct list_head      i_mmap_nonlinear;/*list VM_NONLINEAR mappings */
630     spinlock_t         i_mmap_lock;     /* protect tree, count, list */
631     unsigned int          truncate_count; /* Cover race condition with truncate */
632     unsigned long          nrpages;     /* number of total pages */
633     pgoff_t              writeback_index;/* writeback starts here */
634     const struct address_space_operations *a_ops;    /* methods */
635     unsigned long          flags;        /* error bits/gfp mask */
636     struct backing_dev_info *backing_dev_info; /* device readahead, etc */
637     spinlock_t         private_lock;    /* for use by the address_space */
```

```
638      struct list_head      private_list;    /* ditto */

639      struct address_space      *assoc_mapping; /* ditto */

640 } __attribute__((aligned(sizeof(long))));
```

- host：指向拥有该对象的索引节点的指针（如果存在）

- page_tree：拥有者页的基数的根。

- i_mmap_writable：地址空间中共享内存映射的个数

- i_mmap：radix 优先搜索树的根

## 5.4.2.4 buffer_head

```
52 /*

53   * Historically, a buffer_head was used to map a single block

54   * within a page, and of course as the unit of I/O through the

55   * filesystem and block layers.   Nowadays the basic I/O unit

56   * is the bio, and buffer_heads are used for extracting block

57   * mappings (via a get_block_t call), for tracking state within

58   * a page (via a page_mapping) and for wrapping bio submission

59   * for backward compatibility reasons (e.g. submit_bh).

60   */

61 struct buffer_head {

62      unsigned long b_state;        /* buffer state bitmap (see above) */

63      struct buffer_head *b_this_page;/* circular list of page's buffers */

64      struct page *b_page;            /* the page this bh is mapped to */

65

66      sector_t b_blocknr;        /* start block number */

67      size_t b_size;              /* size of mapping */

68      char *b_data;                /* pointer to data within the page */

69

70      struct block_device *b_bdev;

71      bh_end_io_t *b_end_io;        /* I/O completion */

72      void *b_private;            /* reserved for b_end_io */

73      struct list_head b_assoc_buffers; /* associated with another mapping */
```

```
74        struct address_space *b_assoc_map;   /* mapping this buffer is
75                                associated with */
76        atomic_t b_count;           /* users using this buffer_head */
77 };
```

- b_state：缓冲区状态标志
- b_this_page：指向缓冲区页的链表中的下一个元素的指针
- b_page：指向拥有该块的缓冲区页的描述符指针
- b_blocknr：与块设备相关的块号（起始逻辑块号）
- b_size：块大小
- b_data：块在缓冲区页内的位置
- b_bdev：指向块设备描述符的指针
- b_end_io：I/O 完成方法
- b_private：指向 I/O 完成方法的数据的指针
- b_assoc_buffers：为与某个索引节点相关的间接块的链表提供的指针
- b_count：块使用计数器。



各种数据结构关联图

## 5.4.3 文件的读sys_read

```
374 SYSCALL_DEFINE3(read, unsigned int, fd, char __user *, buf, size_t, count)
375 {
376     struct file *file;
377     ssize_t ret = -EBADF;
378     int fput_needed;
379
380     file = fget_light(fd, &fput_needed);
381     if (file) {
382         loff_t pos = file_pos_read(file);
383         ret = vfs_read(file, buf, count, &pos);
384         file_pos_write(file, pos);
385         fput_light(file, fput_needed);
386     }
387
388     return ret;
389 }
```

## 5.4.3.1 vfs_read

```
278 ssize_t vfs_read(struct file *file, char __user *buf, size_t count, loff_t *pos)
279 {
280     ssize_t ret;
281
282     if (!(file->f_mode & FMODE_READ))
283         return -EBADF;
284     if (!file->f_op || (!file->f_op->read && !file->f_op->aio_read))
285         return -EINVAL;
286     if (unlikely(!access_ok(VERIFY_WRITE, buf, count)))
287         return -EFAULT;
```

```
288
289    ret = rw_verify_area(READ, file, pos, count);
290    if (ret >= 0) {
291        count = ret;
292        if (file->f_op->read)
293            ret = file->f_op->read(file, buf, count, pos);
294        else
295            ret = do_sync_read(file, buf, count, pos);
296        if (ret > 0) {
297            fsnotify_access(file->f_path.dentry);
298            add_rchar(current, ret);
299        }
300        inc_syscr(current);
301    }
302
303    return ret;
304 }
```

## 5.4.3.2  do_sync_read

```
252 ssize_t do_sync_read(struct file *filp, char __user *buf, size_t len, loff_t
*ppos)
253 {
254    struct iovec iov = { .iov_base = buf, .iov_len = len };
255    struct kiocb kiocb;
256    ssize_t ret;
257
258    init_sync_kiocb(&kiocb, filp);
259    kiocb.ki_pos = *ppos;
260    kiocb.ki_left = len;
261    kiocb.ki_nbytes = len;
262
```

```
263     for (;;) {
264         ret = filp->f_op->aio_read(&kiocb, &iov, 1, kiocb.ki_pos);
265         if (ret != -EIOCBRETRY)
266             break;
267         wait_on_retry_sync_kiocb(&kiocb);
268     }
269
270     if (-EIOCBQUEUED == ret)
271         ret = wait_on_sync_kiocb(&kiocb);
272     *ppos = kiocb.ki_pos;
273     return ret;
274 }
```

### 5.4.3.3 generic_file_aio_read

该函数是所有文件系统实现同步和异步读取操作所使用的通用例程。

```
1260 ssize_t
1261 generic_file_aio_read(struct kiocb *iocb, const struct iovec *iov,
1262         unsigned long nr_segs, loff_t pos)
1263 {
1264     struct file *filp = iocb->ki_filp;
1265     ssize_t retval;
1266     unsigned long seg;
1267     size_t count;
1268     loff_t *ppos = &iocb->ki_pos;
1269
1270     count = 0;
1271     retval = generic_segment_checks(iov, &nr_segs, &count,
VERIFY_WRITE);
1272     if (retval)
1273         return retval;
1274
```

```
1275        /* coalesce the iovecs and go direct-to-BIO for O_DIRECT */
1276        if (filp->f_flags & O_DIRECT) {
1277            loff_t size;
1278            struct address_space *mapping;
1279            struct inode *inode;
1280
1281            mapping = filp->f_mapping;
1282            inode = mapping->host;
1283            if (!count)
1284                goto out; /* skip atime */
1285            size = i_size_read(inode);
1286            if (pos < size) {
1287                retval = filemap_write_and_wait_range(mapping, pos,
1288                        pos + iov_length(iov, nr_segs) - 1);
1289                if (!retval) {
1290                    retval = mapping->a_ops->direct_IO(READ, iocb,
1291                            iov, pos, nr_segs);
1292                }
1293                if (retval > 0)
1294                    *ppos = pos + retval;
1295                if (retval) {
1296                    file_accessed(filp);
1297                    goto out;
1298                }
1299            }
1300        }
1301
1302    for (seg = 0; seg < nr_segs; seg++) {
1303        read_descriptor_t desc;
1304
1305        desc.written = 0;
```

| | |
|---|---|
| 1306 | desc.arg.buf = iov[seg].iov_base; |
| 1307 | desc.count = iov[seg].iov_len; |
| 1308 | if (desc.count == 0) |
| 1309 | continue; |
| 1310 | desc.error = 0; |
| 1311 | do_generic_file_read(filp, ppos, &desc, file_read_actor); |
| 1312 | retval += desc.written; |
| 1313 | if (desc.error) { |
| 1314 | retval = retval ?: desc.error; |
| 1315 | break; |
| 1316 | } |
| 1317 | if (desc.count > 0) |
| 1318 | break; |
| 1319 | } |

for 循环中完成文件的读取操作。

```
1320 out:
1321     return retval;
1322 }
1323 EXPORT_SYMBOL(generic_file_aio_read);
```

### 5.4.3.4  do_generic_file_read

```
965 static void do_generic_file_read(struct file *filp, loff_t *ppos,
 966          read_descriptor_t *desc, read_actor_t actor)
 967 {
 968     struct address_space *mapping = filp->f_mapping;
```

地址空间

```
 969     struct inode *inode = mapping->host;
 970     struct file_ra_state *ra = &filp->f_ra;
 971     pgoff_t index;
 972     pgoff_t last_index;
 973     pgoff_t prev_index;
```

```
974        unsigned long offset;         /* offset into pagecache page */

975        unsigned int prev_offset;

976        int error;

977

978        index = *ppos >> PAGE_CACHE_SHIFT;
```

ppos 在缓存中对于的页面的索引号。

```
979        prev_index = ra->prev_pos >> PAGE_CACHE_SHIFT;

980        prev_offset = ra->prev_pos & (PAGE_CACHE_SIZE-1);

981        last_index = (*ppos + desc->count + PAGE_CACHE_SIZE-1) >>
PAGE_CACHE_SHIFT;

982        offset = *ppos & ~PAGE_CACHE_MASK;

983

984        for (;;) {

985            struct page *page;

986            pgoff_t end_index;

987            loff_t isize;

988            unsigned long nr, ret;

989

990            cond_resched();

991 find_page:

992            page = find_get_page(mapping, index);
```

在地址空间中查找对应的页。

```
993            if (!page) {
```

如果页不在缓存中，则调用 page_cache_sync_readahead 预读函数。

```
994                page_cache_sync_readahead(mapping,

995                        ra, filp,

996                        index, last_index - index);
```

从磁盘中读取相应内容。

```
997                page = find_get_page(mapping, index);

998                if (unlikely(page == NULL))

999                    goto no_cached_page;
```

| 1000 | } |
|------|---|

预读之后再次查询，如果还是没有找到，跳转到 no_cached_page 处。

| 1001 | if (PageReadahead(page)) { |
|------|---|
| 1002 | page_cache_async_readahead(mapping, |
| 1003 | ra, filp, page, |
| 1004 | index, last_index - index); |
| 1005 | } |

需要预读，则调用 page_cache_async_readhead 函数预读。

| 1006 | if (!PageUptodate(page)) { |
|------|---|

如果 page 已经存在于缓存中，但是 PG_uptodata 没有置为，说明数据没有更新，需要重磁盘读入。

| 1007 | if (inode->i_blkbits == PAGE_CACHE_SHIFT || |
|------|---|
| 1008 | !mapping->a_ops->is_partially_uptodate) |
| 1009 | goto page_not_up_to_date; |
| 1010 | if (!trylock_page(page)) |
| 1011 | goto page_not_up_to_date; |

跳转到 page_not_up_to_data。

| 1012 | if (!mapping->a_ops->is_partially_uptodate(page, |
|------|---|
| 1013 | desc, offset)) |
| 1014 | goto page_not_up_to_date_locked; |
| 1015 | unlock_page(page); |
| 1016 | } |
| 1017 page_ok: | |
| 1018 | /* |
| 1019 | * i_size must be checked after we know the page is Uptodate. |
| 1020 | * |
| 1021 | * Checking i_size after the check allows us to calculate |
| 1022 | * the correct value for "nr", which means the zero-filled |
| 1023 | * part of the page is not copied back to userspace (unless |
| 1024 | * another truncate extends the file - this is desired though). |
| 1025 | */ |

```
1026
1027            isize = i_size_read(inode);
1028            end_index = (isize - 1) >> PAGE_CACHE_SHIFT;
1029            if (unlikely(!isize || index > end_index)) {
1030                    page_cache_release(page);
1031                    goto out;
1032            }
1033
1034            /* nr is the maximum number of bytes to copy from this page */
1035            nr = PAGE_CACHE_SIZE;
1036            if (index == end_index) {
1037                    nr = ((isize - 1) & ~PAGE_CACHE_MASK) + 1;
1038                    if (nr <= offset) {
1039                            page_cache_release(page);
1040                            goto out;
1041                    }
1042            }
1043            nr = nr - offset;
1044
1045            /* If users can be writing to this page using arbitrary
1046             * virtual addresses, take care about potential aliasing
1047             * before reading the page on the kernel side.
1048             */
1049            if (mapping_writably_mapped(mapping))
1050                    flush_dcache_page(page);
1051
1052            /*
1053             * When a sequential read accesses a page several times,
1054             * only mark it as accessed the first time.
1055             */
1056            if (prev_index != index || offset != prev_offset)
```

```
1057            mark_page_accessed(page);
1058        prev_index = index;
1059
1060        /*
1061         * Ok, we have the page, and it's up-to-date, so
1062         * now we can copy it to user space...
1063         *
1064         * The actor routine returns how many bytes were actually used..
1065         * NOTE! This may not be the same as how much of a user buffer
1066         * we filled up (we may be padding etc), so we can only update
1067         * "pos" here (the actor routine has to update the user buffer
1068         * pointers and the remaining count).
1069         */
1070        ret = actor(desc, page, offset, nr);
1071        offset += ret;
1072        index += offset >> PAGE_CACHE_SHIFT;
1073        offset &= ~PAGE_CACHE_MASK;
1074        prev_offset = offset;
1075
1076        page_cache_release(page);
1077        if (ret == nr && desc->count)
1078            continue;
```

　　1018-1078：当要读取的内容已经存在与缓存中时，指向这段代码，代码比较简单，通过传过来的 actor 参数，将缓存中的内容读入用户传递过来的缓存区中。

```
1079        goto out;
1080
1081 page_not_up_to_date:
1082        /* Get exclusive access to the page ... */
1083        error = lock_page_killable(page);
1084        if (unlikely(error))
1085            goto readpage_error;
```

```
1086
1087 page_not_up_to_date_locked:
1088            /* Did it get truncated before we got the lock? */
1089            if (!page->mapping) {
1090                unlock_page(page);
1091                page_cache_release(page);
1092                continue;
1093            }
1094
1095            /* Did somebody else fill it already? */
1096            if (PageUptodate(page)) {
1097                unlock_page(page);
1098                goto page_ok;
1099            }
1100
1101 readpage:
1102            /* Start the actual read. The read will unlock the page. */
1103            error = mapping->a_ops->readpage(filp, page);
```

调用具体文件系统的函数从磁盘中读取内容到缓存。

```
1104
1105            if (unlikely(error)) {
1106                if (error == AOP_TRUNCATED_PAGE) {
1107                    page_cache_release(page);
1108                    goto find_page;
1109                }
1110                goto readpage_error;
1111            }
1112
if (!PageUptodate(page)) {
1114                error = lock_page_killable(page);
1115                if (unlikely(error))
```

```
1116                    goto readpage_error;
1117            if (!PageUptodate(page)) {
1118                if (page->mapping == NULL) {
1119                    /*
1120                     * invalidate_mapping_pages got it
1121                     */
1122                    unlock_page(page);
1123                    page_cache_release(page);
1124                    goto find_page;
1125                }
1126                unlock_page(page);
1127                shrink_readahead_size_eio(filp, ra);
1128                error = -EIO;
1129                goto readpage_error;
1130            }
1131            unlock_page(page);
1132        }
1133
1134        goto page_ok;
1135
1136 readpage_error:
1137        /* UHHUH! A synchronous read error occurred. Report it */
1138        desc->error = error;
1139        page_cache_release(page);
1140        goto out;
1141
1142 no_cached_page:
1143        /*
1144         * Ok, it wasn't cached, so we need to create a new
1145         * page..
1146         */
```

| | |
|---|---|
| 1147 | page = page_cache_alloc_cold(mapping); |

分配一个页面

| | |
|---|---|
| 1148 | if (!page) { |
| 1149 | desc->error = -ENOMEM; |
| 1150 | goto out; |
| 1151 | } |
| 1152 | error = add_to_page_cache_lru(page, mapping, |
| 1153 | index, GFP_KERNEL); |

将页面添加到对应的 address_space 中。

| | |
|---|---|
| 1154 | if (error) { |
| 1155 | page_cache_release(page); |
| 1156 | if (error == -EEXIST) |
| 1157 | goto find_page; |
| 1158 | desc->error = error; |
| 1159 | goto out; |
| 1160 | } |
| 1161 | goto readpage; |

跳转到 readpage 读取页面。

| | |
|---|---|
| 1162 | } |
| 1163 | |
| 1164 out: | |
| 1165 | ra->prev_pos = prev_index; |
| 1166 | ra->prev_pos <<= PAGE_CACHE_SHIFT; |
| 1167 | ra->prev_pos |= prev_offset; |
| 1168 | |
| 1169 | *ppos = ((loff_t)index << PAGE_CACHE_SHIFT) + offset; |
| 1170 | file_accessed(filp); |
| 1171 } | |

### 5.4.3.5  ext2_readpage

| |
|---|
| 745 static int ext2_readpage(struct file *file, struct page *page) |

```
746 {
747        return mpage_readpage(page, ext2_get_block);
```

调用 mpage_readpage 完成进一步操作，参数 ext2_get_block 函数获得硬盘上指定的块，该函数之前已经分析过。

```
748 }
```

### 5.4.3.6 mpage_readpage

```
407 /*
408   * This isn't called much at all
409   */
410 int mpage_readpage(struct page *page, get_block_t get_block)
411 {
412      struct bio *bio = NULL;
413      sector_t last_block_in_bio = 0;
414      struct buffer_head map_bh;
415      unsigned long first_logical_block = 0;
416
417      map_bh.b_state = 0;
418      map_bh.b_size = 0;
419      bio = do_mpage_readpage(bio, page, 1, &last_block_in_bio,
420              &map_bh, &first_logical_block, get_block);
421      if (bio)
422          mpage_bio_submit(READ, bio);
423      return 0;
424 }
425 EXPORT_SYMBOL(mpage_readpage);
```

初始化 buffer_head map_bh 相关字段，调用 do_mpage_readpage 函数创建一个 bio 请求，该请求指明了要读取数据库在磁盘的位置、数据库的数量以及拷贝该数据的目标位置——缓存区中的 page 信息，然后调用 mpage_bio_submit 函数处理请求。

### 5.4.3.7  do_mpage_readpage

struct bio {

sector_t bi_sector; //块 I/O 操作的第一个磁盘扇区

struct bio *bi_next; //链接到请求队列中的下一个 bio

struct block_device *bi_bdev;//指向块设备描述符的指针

unsigned long bi_flags; //bio 的状态标志

unsigned long bi_rw; //IO 操作标志，即这次 I.O 是读或写

unsigned short bi_vcnt; /* bio 的 bio_vec 数组中段的数目 */

unsigned short bi_idx; /* bio 的 bio_vec 数组中段的当前索引值 */

unsigned short bi_phys_segments; //合并之后 bio 中物理段的数目

unsigned short bi_hw_segments; //合并之后硬件段的数目

unsigned int bi_size; /* 需要传送的字节数 */

unsigned int bi_hw_front_size;// 硬件段合并算法使用

unsigned int bi_hw_back_size;// 硬件段合并算法使用

unsigned int bi_max_vecs; /* bio 的 bio vec 数组中允许的最大段数 */

struct bio_vec *bi_io_vec; /*指向 bio 的 bio_vec 数组中的段的指针 */

bio_end_io_t *bi_end_io; /* bio 的 I/O 操作结束时调用的方法 */

atomic_t bi_cnt; /* bio 的引用计数器 */

void *bi_private; //通用块层和块设备驱动程序的 I/O 完成方法使用的指针

bio_destructor_t *bi_destructor;//释放 bio 时调用的析构方法（通常是
bio_destructor()方法）

};

168 static struct bio *

169 do_mpage_readpage(struct bio *bio, struct page *page, unsigned nr_pages,

170          sector_t *last_block_in_bio, struct buffer_head *map_bh,

171          unsigned long *first_logical_block, get_block_t get_block)

172 {

173      struct inode *inode = page->mapping->host;

174      const unsigned blkbits = inode->i_blkbits;

175      const unsigned blocks_per_page = PAGE_CACHE_SIZE >> blkbits;

计算每页高速缓存能存放的块数目。

```
176     const unsigned blocksize = 1 << blkbits;
177     sector_t block_in_file;
178     sector_t last_block;
179     sector_t last_block_in_file;
180     sector_t blocks[MAX_BUF_PER_PAGE];
181     unsigned page_block;
182     unsigned first_hole = blocks_per_page;
183     struct block_device *bdev = NULL;
184     int length;
185     int fully_mapped = 1;
186     unsigned nblocks;
187     unsigned relative_block;
188
189     if (page_has_buffers(page))
190         goto confused;
191
192     block_in_file = (sector_t)page->index << (PAGE_CACHE_SHIFT -
blkbits);
193     last_block = block_in_file + nr_pages * blocks_per_page;
194     last_block_in_file = (i_size_read(inode) + blocksize - 1) >> blkbits;
195     if (last_block > last_block_in_file)
196         last_block = last_block_in_file;
```

缓存是以页为单位的，而记录块是以块为单位的，192 行记录 page 在文件内部的块号（也就是页面号乘以 4），last_block 记录所要读取内容的最后一个块号。如果 last_block 超过了文件系统的最大块号，则 last_block 赋值为文件系统最大块号。

```
197     page_block = 0;
199     /*
200      * Map blocks using the result from the previous get_blocks call first.
201      */
202     nblocks = map_bh->b_size >> blkbits;
```

| | |
|---|---|
| 203 | if (buffer_mapped(map_bh) && block_in_file > *first_logical_block && |
| 204 | block_in_file < (*first_logical_block + nblocks)) { |

如果页面已经被映射，且页中第一块的文件块号位于传递进来的 first_logical_block 参数和一个 buffer_heand 最后一个块之间，则对未映射部分进行处理。

| | |
|---|---|
| 205 | unsigned map_offset = block_in_file - *first_logical_block; |
| 206 | unsigned last = nblocks - map_offset; |
| 207 | |
| 208 | for (relative_block = 0; ; relative_block++) { |
| 209 | if (relative_block == last) { |
| 210 | clear_buffer_mapped(map_bh); |
| 211 | break; |
| 212 | } |
| 213 | if (page_block == blocks_per_page) |
| 214 | break; |
| 215 | blocks[page_block] = map_bh->b_blocknr + map_offset + |
| 216 | relative_block; |
| 217 | page_block++; |
| 218 | block_in_file++; |
| 219 | } |
| 220 | bdev = map_bh->b_bdev; |
| 221 | } |
| 222 | |
| 223 | /* |
| 224 | * Then do more get_blocks calls until we are done with this page. |
| 225 | */ |
| 226 | map_bh->b_page = page; |
| 227 | while (page_block < blocks_per_page) { |
| 228 | map_bh->b_state = 0; |
| 229 | map_bh->b_size = 0; |
| 230 | |
| 231 | if (block_in_file < last_block) { |

| 232 | map_bh->b_size = (last_block-block_in_file) << blkbits; |
|-----|--------------------------------------------------------|
| 233 | if (get_block(inode, block_in_file, map_bh, 0)) |

调用 get_block 将在文件中的逻辑块号转换为磁盘中的逻辑块号，保存在 buffer_head 中的 b_blocknr 字段中。

| 234 | goto confused; |
|-----|----------------|
| 235 | *first_logical_block = block_in_file; |
| 236 | } |
| 237 | |
| 238 | if (!buffer_mapped(map_bh)) { |
| 240 | if (first_hole == blocks_per_page) |
| 241 | first_hole = page_block; |
| 242 | page_block++; |
| 243 | block_in_file++; |

如果 buffer_mapped 没被映射，page_block++，block_in_file++，跳转到循环开头处继续。

| 244 | continue; |
|-----|----------|
| 245 | } |
| 246 | |
| 247 | /* some filesystems will copy data into the page during |
| 248 | * the get_block call, in which case we don't want to |
| 249 | * read it again.   map_buffer_to_page copies the data |
| 250 | * we just collected from get_block into the page's buffers |
| 251 | * so readpage doesn't have to repeat the get_block call |
| 252 | */ |
| 253 | if (buffer_uptodate(map_bh)) { |
| 254 | map_buffer_to_page(page, map_bh, page_block); |
| 255 | goto confused; |
| 256 | } |
| 257 | |
| 258 | if (first_hole != blocks_per_page) |
| 259 | goto confused;        /* hole -> non-hole */ |

```
260
261            /* Contiguous blocks? */
262            if (page_block && blocks[page_block-1] != map_bh->b_blocknr-1)
263                goto confused;
264            nblocks = map_bh->b_size >> blkbits;
265            for (relative_block = 0; ; relative_block++) {
266                if (relative_block == nblocks) {
267                    clear_buffer_mapped(map_bh);
268                    break;
269                } else if (page_block == blocks_per_page)
270                    break;
271                blocks[page_block] = map_bh->b_blocknr+relative_block;
272                page_block++;
273                block_in_file++;
274            }
275            bdev = map_bh->b_bdev;
276        }
277
    if (first_hole != blocks_per_page) {
278            zero_user_segment(page, first_hole << blkbits,
PAGE_CACHE_SIZE);
280            if (first_hole == 0) {
281                SetPageUptodate(page);
282                unlock_page(page);
283                goto out;
284            }
285        } else if (fully_mapped) {
286            SetPageMappedToDisk(page);
287        }
288
289        /*
290         * This page will go to BIO.   Do we need to send this BIO off first?
```

```
291         */
292         if (bio && (*last_block_in_bio != blocks[0] - 1))
293             bio = mpage_bio_submit(READ, bio);
294
295 alloc_new:
296         if (bio == NULL) {
297             bio = mpage_alloc(bdev, blocks[0] << (blkbits - 9),
298                     min_t(int, nr_pages, bio_get_nr_vecs(bdev)),
299                     GFP_KERNEL);
300             if (bio == NULL)
301                 goto confused;
302         }
303
304         length = first_hole << blkbits;
305         if (bio_add_page(bio, page, length, 0) < length) {
306             bio = mpage_bio_submit(READ, bio);
307             goto alloc_new;
308         }
```

分配一个 BIO，进行适当的初始化。

```
309
310         relative_block = block_in_file - *first_logical_block;
311         nblocks = map_bh->b_size >> blkbits;
312         if ((buffer_boundary(map_bh) && relative_block == nblocks) ||
313             (first_hole != blocks_per_page))
314             bio = mpage_bio_submit(READ, bio);
```

调用 mpage_bio_submit 处理读请求。完成读文件的实际操作。

```
315         else
316             *last_block_in_bio = blocks[blocks_per_page - 1];
317 out:
318         return bio;
319
```

```
320 confused:
321     if (bio)
322         bio = mpage_bio_submit(READ, bio);
323     if (!PageUptodate(page))
324             block_read_full_page(page, get_block);
325     else
326         unlock_page(page);
327     goto out;
328 }
```

## 5.4.4 sys_write 文件的写

fs/read_write.c

```
391 SYSCALL_DEFINE3(write, unsigned int, fd, const char __user *, buf,
392             size_t, count)
```

参数：

- fd：用户空间的文件句柄号
- buf：用户空间保存写入内容的缓冲区指针
- count：写入的内容的长度，以字节为单位

```
393 {
394     struct file *file;
395     ssize_t ret = -EBADF;
396     int fput_needed;
397
398     file = fget_light(fd, &fput_needed);
```

根据文件句柄号得到 file 结构。

```
399     if (file) {
400         loff_t pos = file_pos_read(file);
```

得到文件的上下文 pos，pos 即表示文件当前索引与文件起始处的偏移量。

```
401         ret = vfs_write(file, buf, count, &pos);
```

读文件。

| 402 | file_pos_write(file, pos); |
|---|---|

设置文件的上下文。

| 403 | fput_light(file, fput_needed); |
|---|---|
| 404 | } |
| 405 | |
| 406 | return ret; |
| 407 | } |

## 5.4.4.1  vfs_write

| 334 ssize_t vfs_write(struct file *file, const char __user *buf, size_t count, loff_t *pos) |
|---|
| 335 { |
| 336     ssize_t ret; |
| 337 |
| 338     if (!(file->f_mode & FMODE_WRITE)) |
| 339         return -EBADF; |
| 340     if (!file->f_op || (!file->f_op->write && !file->f_op->aio_write)) |
| 341         return -EINVAL; |
| 342     if (unlikely(!access_ok(VERIFY_READ, buf, count))) |
| 343         return -EFAULT; |

写操作必要条件检测

| 344 | |
|---|---|
| 345 |     ret = rw_verify_area(WRITE, file, pos, count); |

检测需要写入的内容是否被加锁

| 346 |     if (ret >= 0) { |
|---|---|
| 347 |         count = ret; |
| 348 |         if (file->f_op->write) |
| 349 |             ret = file->f_op->write(file, buf, count, pos); |

调用具体文件系统的写入函数进行写操作。对于 ext2 而言就是 do_sync_write 函数。

```
350            else
351                ret = do_sync_write(file, buf, count, pos);
352            if (ret > 0) {
353                fsnotify_modify(file->f_path.dentry);
354                add_wchar(current, ret);
355            }
356            inc_syscw(current);
357        }
358
359        return ret;
360 }
```

### 5.4.4.2  do_sync_write

```
308 ssize_t do_sync_write(struct file *filp, const char __user *buf, size_t len,
loff_t *ppos)
309 {
310        struct iovec iov = { .iov_base = (void __user *)buf, .iov_len = len };
311        struct kiocb kiocb;
312        ssize_t ret;
313
314        init_sync_kiocb(&kiocb, filp);
315        kiocb.ki_pos = *ppos;
316        kiocb.ki_left = len;
317        kiocb.ki_nbytes = len;
318
319        for (;;) {
320            ret = filp->f_op->aio_write(&kiocb, &iov, 1, kiocb.ki_pos);
```

调用该函数完成写入操作。

```
321          if (ret != -EIOCBRETRY)
322             break;
323          wait_on_retry_sync_kiocb(&kiocb);
324     }
325
326     if (-EIOCBQUEUED == ret)
327          ret = wait_on_sync_kiocb(&kiocb);
328     *ppos = kiocb.ki_pos;
329     return ret;
330 }
331
332 EXPORT_SYMBOL(do_sync_write);
```

### 5.4.4.3  generic_file_aio_write

mm/filemap.c

```
2424 ssize_t generic_file_aio_write(struct kiocb *iocb, const struct iovec *iov,
2425          unsigned long nr_segs, loff_t pos)
2426 {
2427     struct file *file = iocb->ki_filp;
2428     struct inode *inode = file->f_mapping->host;
2429     ssize_t ret;
2430
2431     BUG_ON(iocb->ki_pos != pos);
2432
2433     mutex_lock(&inode->i_mutex);
2434     ret = __generic_file_aio_write(iocb, iov, nr_segs, &iocb->ki_pos);
2435     mutex_unlock(&inode->i_mutex);
2436
2437     if (ret > 0 || ret == -EIOCBQUEUED) {
2438          ssize_t err;
```

```
2439
2440            err = generic_write_sync(file, pos, ret);
2441            if (err < 0 && ret > 0)
2442                    ret = err;
2443        }
2444        return ret;
2445 }
2446 EXPORT_SYMBOL(generic_file_aio_write);
```

## 5.4.4.4 __generic_file_aio_write

```
2316 ssize_t __generic_file_aio_write(struct kiocb *iocb, const struct iovec *iov,
2317                        unsigned long nr_segs, loff_t *ppos)
2318 {
2319      struct file *file = iocb->ki_filp;
2320      struct address_space * mapping = file->f_mapping;
2321      size_t ocount;         /* original count */
2322      size_t count;          /* after file limit checks */
2323      struct inode      *inode = mapping->host;
2324      loff_t        pos;
2325      ssize_t       written;
2326      ssize_t       err;
2327
2328      ocount = 0;
2329      err = generic_segment_checks(iov, &nr_segs, &ocount,
VERIFY_READ);
2330      if (err)
2331          return err;
2332
2333      count = ocount;
2334      pos = *ppos;
```

```
2335

2336        vfs_check_frozen(inode->i_sb, SB_FREEZE_WRITE);

2337

2338        /* We can write back this queue in page reclaim */

2339        current->backing_dev_info = mapping->backing_dev_info;

2340        written = 0;

2341

2342        err = generic_write_checks(file, &pos, &count,
S_ISBLK(inode->i_mode));

2343        if (err)

2344            goto out;

2345

2346        if (count == 0)

2347            goto out;

2348

2349        err = file_remove_suid(file);

2350        if (err)

2351            goto out;

2352

2353        file_update_time(file);

2354

2355        /* coalesce the iovecs and go direct-to-BIO for O_DIRECT */

2356        if (unlikely(file->f_flags & O_DIRECT)) {

2357            loff_t endbyte;

2358            ssize_t written_buffered;

2359

2360            written = generic_file_direct_write(iocb, iov, &nr_segs, pos,

2361                                 ppos, count, ocount);
```

调用该函数进行写入操作。

```
2362            if (written < 0 || written == count)

2363                goto out;
```

```
2364            /*
2365             * direct-io write to a hole: fall through to buffered I/O
2366             * for completing the rest of the request.
2367             */
2368            pos += written;
2369            count -= written;
2370            written_buffered = generic_file_buffered_write(iocb, iov,
2371                                    nr_segs, pos, ppos, count,
2372                                    written);
2373            /*
2374             * If generic_file_buffered_write() retuned a synchronous error
2375             * then we want to return the number of bytes which were
2376             * direct-written, or the error code if that was zero.   Note
2377             * that this differs from normal direct-io semantics, which
2378             * will return -EFOO even if some bytes were written.
2379             */
2380            if (written_buffered < 0) {
2381                    err = written_buffered;
2382                    goto out;
2383            }
2384
2385            /*
2386             * We need to ensure that the page cache pages are written to
2387             * disk and invalidated to preserve the expected O_DIRECT
2388             * semantics.
2389             */
2390            endbyte = pos + written_buffered - written - 1;
2391            err = filemap_write_and_wait_range(file->f_mapping, pos,
endbyte);
2392            if (err == 0) {
2393                    written = written_buffered;
2394                    invalidate_mapping_pages(mapping,
```

```
2395                          pos >> PAGE_CACHE_SHIFT,
2396                          endbyte >> PAGE_CACHE_SHIFT);
2397        } else {
2398            /*
2399             * We don't know how much we wrote, so just return
2400             * the number of bytes which were direct-written
2401             */
2402        }
2403    } else {
2404        written = generic_file_buffered_write(iocb, iov, nr_segs,
2405                    pos, ppos, count, written);
2406    }
2407 out:
2408    current->backing_dev_info = NULL;
2409    return written ? written : err;
2410 }
2411 EXPORT_SYMBOL(__generic_file_aio_write);
```

### 5.4.4.5 generic_file_direct_write

```
2083 ssize_t
2084 generic_file_direct_write(struct kiocb *iocb, const struct iovec *iov,
2085        unsigned long *nr_segs, loff_t pos, loff_t *ppos,
2086        size_t count, size_t ocount)
2087 {
2088    struct file *file = iocb->ki_filp;
2089    struct address_space *mapping = file->f_mapping;
2090    struct inode    *inode = mapping->host;
2091    ssize_t      written;
2092    size_t       write_len;
2093    pgoff_t      end;
```

```c
2094
2095     if (count != ocount)
2096         *nr_segs = iov_shorten((struct iovec *)iov, *nr_segs, count);
2097
2098     write_len = iov_length(iov, *nr_segs);
2099     end = (pos + write_len - 1) >> PAGE_CACHE_SHIFT;
2100
2101     written = filemap_write_and_wait_range(mapping, pos, pos +
write_len - 1);
2102     if (written)
2103         goto out;
2104
2105     /*
2106      * After a write we want buffered reads to be sure to go to disk to get
2107      * the new data.   We invalidate clean cached page from the region we're
2108      * about to write.   We do this *before* the write so that we can return
2109      * without clobbering -EIOCBQUEUED from ->direct_IO().
2110      */
2111     if (mapping->nrpages) {
2112         written = invalidate_inode_pages2_range(mapping,
2113                         pos >> PAGE_CACHE_SHIFT, end);
2114         /*
2115          * If a page can not be invalidated, return 0 to fall back
2116          * to buffered write.
2117          */
2118         if (written) {
2119             if (written == -EBUSY)
2120                 return 0;
2121             goto out;
2122         }
2123     }
```

```
2124
2125        written = mapping->a_ops->direct_IO(WRITE, iocb, iov, pos,
*nr_segs);
2126
2127        /*
2128         * Finally, try again to invalidate clean pages which might have been
2129         * cached by non-direct readahead, or faulted in by get_user_pages()
2130         * if the source of the write was an mmap'ed region of the file
2131         * we're writing.   Either one is a pretty crazy thing to do,
2132         * so we don't support it 100%.   If this invalidation
2133         * fails, tough, the write still worked...
2134         */
2135        if (mapping->nrpages) {
2136            invalidate_inode_pages2_range(mapping,
2137                                 pos >> PAGE_CACHE_SHIFT, end);
2138        }
2139
2140        if (written > 0) {
2141            loff_t end = pos + written;
2142            if (end > i_size_read(inode) && !S_ISBLK(inode->i_mode)) {
2143                i_size_write(inode,   end);
2144                mark_inode_dirty(inode);
2145            }
2146            *ppos = end;
2147        }
2148 out:
2149        return written;
2150 }
2151 EXPORT_SYMBOL(generic_file_direct_write);
```

## 5.4.4.6 filemap_write_and_wait_range

```
366 int filemap_write_and_wait_range(struct address_space *mapping,
367                     loff_t lstart, loff_t lend)
368 {
369     int err = 0;
370
371     if (mapping->nrpages) {
372         err = __filemap_fdatawrite_range(mapping, lstart, lend,
373                         WB_SYNC_ALL);
374         /* See comment of filemap_write_and_wait() */
375         if (err != -EIO) {
376             int err2 = filemap_fdatawait_range(mapping,
377                             lstart, lend);
378             if (!err)
379                 err = err2;
380         }
381     }
382     return err;
383 }
384 EXPORT_SYMBOL(filemap_write_and_wait_range);
```

## 5.4.4.7 __filemap_fdatawrite_range

```
212 int __filemap_fdatawrite_range(struct address_space *mapping, loff_t start,
213                     loff_t end, int sync_mode)
214 {
215     int ret;
216     struct writeback_control wbc = {
217         .sync_mode = sync_mode,
218         .nr_to_write = LONG_MAX,
```

| | |
|---|---|
| 219 | .range_start = start, |
| 220 | .range_end = end, |
| 221 | }; |

分配一个 writeback_control 结构.

| | |
|---|---|
| 222 | |
| 223 | if (!mapping_cap_writeback_dirty(mapping)) |
| 224 | return 0; |
| 225 | |
| 226 | ret = do_writepages(mapping, &wbc); |

调用 do_writepages。

| | |
|---|---|
| 227 | return ret; |
| 228 | } |

## 5.4.4.8 mm/page-writeback.c

1012 int do_writepages(struct address_space *mapping, struct writeback_control *wbc)

1013 {

1014    int ret;

1015

1016    if (wbc->nr_to_write <= 0)

1017        return 0;

1018    if (mapping->a_ops->writepages)

1019        ret = mapping->a_ops->writepages(mapping, wbc);

调用 ext2_writepages 函数。

1020    else

1021        ret = generic_writepages(mapping, wbc);

1022    return ret;

1023 }

### 5.4.4.9 **ext2_writepages**

```
810 static int
811 ext2_writepages(struct address_space *mapping, struct
writeback_control *wbc)
812 {
813      return mpage_writepages(mapping, wbc, ext2_get_block);
```

传递 ext2_get_block 参数。

```
814 }
```

### 5.4.4.10 **mpage_writepages**

```
680 int
681 mpage_writepages(struct address_space *mapping,
682          struct writeback_control *wbc, get_block_t get_block)
683 {
684      int ret;
685
686      if (!get_block)
687          ret = generic_writepages(mapping, wbc);
688      else {
689          struct mpage_data mpd = {
690              .bio = NULL,
691              .last_block_in_bio = 0,
692              .get_block = get_block,
693              .use_writepage = 1,
694          };
```

初始化一个 mpage_data 结构

```
695
696          ret = write_cache_pages(mapping, wbc, __mpage_writepage,
&mpd);
```

传入产生__mpage_writepage 调用 write_cache_pages 函数。

| | |
|---|---|
| 697 | if (mpd.bio) |
| 698 | mpage_bio_submit(WRITE, mpd.bio); |

bio 请求准备好后，则调用 mpage_bio_submit 进行写入操作。

| | |
|---|---|
| 699 | } |
| 700 | return ret; |
| 701 } | |
| 702 EXPORT_SYMBOL(mpage_writepages); | |

## 5.4.4.11  write_cache_pages

| | |
|---|---|
| 820 int write_cache_pages(struct address_space *mapping, | |
| 821 | struct writeback_control *wbc, writepage_t writepage, |
| 822 | void *data) |
| 823 { | |
| 824 | int ret = 0; |
| 825 | int done = 0; |
| 826 | struct pagevec pvec; |
| 827 | int nr_pages; |
| 828 | pgoff_t uninitialized_var(writeback_index); |
| 829 | pgoff_t index; |
| 830 | pgoff_t end;          /* Inclusive */ |
| 831 | pgoff_t done_index; |
| 832 | int cycled; |
| 833 | int range_whole = 0; |
| 834 | long nr_to_write = wbc->nr_to_write; |
| 835 | |
| 836 | pagevec_init(&pvec, 0); |
| 837 | if (wbc->range_cyclic) { |
| 838 | writeback_index = mapping->writeback_index; /* prev offset */ |
| 839 | index = writeback_index; |
| 840 | if (index == 0) |

```
841                 cycled = 1;
842             else
843                 cycled = 0;
844             end = -1;
845         } else {
846             index = wbc->range_start >> PAGE_CACHE_SHIFT;
847             end = wbc->range_end >> PAGE_CACHE_SHIFT;
848             if (wbc->range_start == 0 && wbc->range_end == LLONG_MAX)
849                 range_whole = 1;
850             cycled = 1; /* ignore range_cyclic tests */
851         }
852 retry:
853     done_index = index;
854     while (!done && (index <= end)) {
855         int i;
856
857         nr_pages = pagevec_lookup_tag(&pvec, mapping, &index,
858                 PAGECACHE_TAG_DIRTY,
859                 min(end - index, (pgoff_t)PAGEVEC_SIZE-1) + 1);
860         if (nr_pages == 0)
861             break;
862
863         for (i = 0; i < nr_pages; i++) {
864             struct page *page = pvec.pages[i];
865
866             /*
867              * At this point, the page may be truncated or
868              * invalidated (changing page->mapping to NULL), or
869              * even swizzled back from swapper_space to tmpfs file
870              * mapping. However, page->index will not change
871              * because we have a reference on the page.
```

```
872                    */
873                 if (page->index > end) {
874                    /*
875                      * can't be range_cyclic (1st pass) because
876                      * end == -1 in that case.
877                      */
878                    done = 1;
879                    break;
880                 }
881
882              done_index = page->index + 1;
883
884              lock_page(page);
885
886              /*
887               * Page truncated or invalidated. We can freely skip it
888               * then, even for data integrity operations: the page
889               * has disappeared concurrently, so there could be no
890               * real expectation of this data interity operation
891               * even if there is now a new, dirty page at the same
892               * pagecache address.
893              */
894              if (unlikely(page->mapping != mapping)) {
895 continue_unlock:
896                    unlock_page(page);
897                    continue;
898              }
899
900              if (!PageDirty(page)) {
901                    /* someone wrote it for us */
902                    goto continue_unlock;
```

```
903                 }
904
905                 if (PageWriteback(page)) {
906                         if (wbc->sync_mode != WB_SYNC_NONE)
907                                 wait_on_page_writeback(page);
908                         else
909                                 goto continue_unlock;
910                 }
911
912                 BUG_ON(PageWriteback(page));
913                 if (!clear_page_dirty_for_io(page))
914                         goto continue_unlock;
915
916                 ret = (*writepage)(page, wbc, data);
917                 if (unlikely(ret)) {
918                         if (ret == AOP_WRITEPAGE_ACTIVATE) {
919                                 unlock_page(page);
920                                 ret = 0;
921                         } else {
922                                 /*
923                                  * done_index is set past this page,
924                                  * so media errors will not choke
925                                  * background writeout for the entire
926                                  * file. This has consequences for
927                                  * range_cyclic semantics (ie. it may
928                                  * not be suitable for data integrity
929                                  * writeout).
930                                  */
931                                 done = 1;
932                                 break;
933                         }
```

```
934                    }
935
936                if (nr_to_write > 0) {
937                    nr_to_write--;
938                    if (nr_to_write == 0 &&
939                        wbc->sync_mode == WB_SYNC_NONE) {
940                        /*
941                         * We stop writing back only if we are
942                         * not doing integrity sync. In case of
943                         * integrity sync we have to keep going
944                         * because someone may be concurrently
945                         * dirtying pages, and we might have
946                         * synced a lot of newly appeared dirty
947                         * pages, but have not synced all of the
948                         * old dirty pages.
949                         */
950                        done = 1;
951                        break;
952                    }
953                }
954            }
955        pagevec_release(&pvec);
956        cond_resched();
957    }
958    if (!cycled && !done) {
959        /*
960         * range_cyclic:
961         * We hit the last page and there is more work to be done: wrap
962         * back to the start of the file
963         */
964        cycled = 1;
```

```
965            index = 0;
966            end = writeback_index - 1;
967            goto retry;
968        }
969        if (!wbc->no_nrwrite_index_update) {
970            if (wbc->range_cyclic || (range_whole && nr_to_write > 0))
971                mapping->writeback_index = done_index;
972            wbc->nr_to_write = nr_to_write;
973        }
974
975        return ret;
976 }
977 EXPORT_SYMBOL(write_cache_pages);
```

## 5.4.4.12 __mpage_writepage

```
451 static int __mpage_writepage(struct page *page, struct writeback_control *wbc,
452                        void *data)
453 {
454     struct mpage_data *mpd = data;
455     struct bio *bio = mpd->bio;
456     struct address_space *mapping = page->mapping;
457     struct inode *inode = page->mapping->host;
458     const unsigned blkbits = inode->i_blkbits;
459     unsigned long end_index;
460     const unsigned blocks_per_page = PAGE_CACHE_SIZE >> blkbits;
461     sector_t last_block;
462     sector_t block_in_file;
463     sector_t blocks[MAX_BUF_PER_PAGE];
464     unsigned page_block;
465     unsigned first_unmapped = blocks_per_page;
```

```
466        struct block_device *bdev = NULL;

467        int boundary = 0;

468        sector_t boundary_block = 0;

469        struct block_device *boundary_bdev = NULL;

470        int length;

471        struct buffer_head map_bh;

472        loff_t i_size = i_size_read(inode);

473        int ret = 0;

474

475        if (page_has_buffers(page)) {

476            struct buffer_head *head = page_buffers(page);

477            struct buffer_head *bh = head;

479            /* If they're all mapped and dirty, do it */

480            page_block = 0;

481            do {

482                BUG_ON(buffer_locked(bh));

483                if (!buffer_mapped(bh)) {

484                    /*

485                     * unmapped dirty buffers are created by

486                     * __set_page_dirty_buffers -> mmapped data

487                     */

488                    if (buffer_dirty(bh))

489                        goto confused;

490                    if (first_unmapped == blocks_per_page)

491                        first_unmapped = page_block;

492                    continue;

493                }

494

495                if (first_unmapped != blocks_per_page)

496                    goto confused;    /* hole -> non-hole */

497
```

```
498                    if (!buffer_dirty(bh) || !buffer_uptodate(bh))
499                        goto confused;
500                    if (page_block) {
501                        if (bh->b_blocknr != blocks[page_block-1] + 1)
502                            goto confused;
503                    }
504                    blocks[page_block++] = bh->b_blocknr;
505                    boundary = buffer_boundary(bh);
506                    if (boundary) {
507                        boundary_block = bh->b_blocknr;
508                        boundary_bdev = bh->b_bdev;
509                    }
510                    bdev = bh->b_bdev;
511                } while ((bh = bh->b_this_page) != head);
512
513            if (first_unmapped)
514                goto page_is_mapped;
515
516            /*
517             * Page has buffers, but they are all unmapped. The page was
518             * created by pagein or read over a hole which was handled by
519             * block_read_full_page().   If this address_space is also
520             * using mpage_readpages then this can rarely happen.
521             */
522            goto confused;
523        }
```

475-523 如果 page 已经有对应的块缓存，则进行处理。

```
524
525        /*
526         * The page has no buffers: map it to disk
527         */
```

| 528 | BUG_ON(!PageUptodate(page)); |
| 529 | block_in_file = (sector_t)page->index << (PAGE_CACHE_SHIFT - blkbits); |
| 530 | last_block = (i_size - 1) >> blkbits; |
| 531 | map_bh.b_page = page; |

文件中写入的块号范围计算[block_in_file，last_block]。

| 532 | for (page_block = 0; page_block < blocks_per_page; ) { |

for 循环依次写入每个块。

| 533 | |
| 534 | map_bh.b_state = 0; |
| 535 | map_bh.b_size = 1 << blkbits; |
| 536 | if (mpd->get_block(inode, block_in_file, &map_bh, 1)) |

将文件中的逻辑块号转换为磁盘中的逻辑块号。

| 537 | goto confused; |
| 538 | if (buffer_new(&map_bh)) |
| 539 | unmap_underlying_metadata(map_bh.b_bdev, |
| 540 | map_bh.b_blocknr); |
| 541 | if (buffer_boundary(&map_bh)) { |
| 542 | boundary_block = map_bh.b_blocknr; |
| 543 | boundary_bdev = map_bh.b_bdev; |
| 544 | } |
| 545 | if (page_block) { |
| 546 | if (map_bh.b_blocknr != blocks[page_block-1] + 1) |
| 547 | goto confused; |
| 548 | } |
| 549 | blocks[page_block++] = map_bh.b_blocknr; |
| 550 | boundary = buffer_boundary(&map_bh); |
| 551 | bdev = map_bh.b_bdev; |
| 552 | if (block_in_file == last_block) |
| 553 | break; |
| 554 | block_in_file++; |

```
555         }
556         BUG_ON(page_block == 0);
557
558         first_unmapped = page_block;
559
560 page_is_mapped:
561         end_index = i_size >> PAGE_CACHE_SHIFT;
562         if (page->index >= end_index) {
563             /*
564              * The page straddles i_size.   It must be zeroed out on each
565              * and every writepage invocation because it may be mmapped.
566              * "A file is mapped in multiples of the page size.   For a file
567              * that is not a multiple of the page size, the remaining memory
568              * is zeroed when mapped, and writes to that region are not
569              * written out to the file."
570              */
571             unsigned offset = i_size & (PAGE_CACHE_SIZE - 1);
572
573             if (page->index > end_index || !offset)
574                 goto confused;
575             zero_user_segment(page, offset, PAGE_CACHE_SIZE);
576         }
577
578         /*
579          * This page will go to BIO.   Do we need to send this BIO off first?
580          */
581         if (bio && mpd->last_block_in_bio != blocks[0] - 1)
582             bio = mpage_bio_submit(WRITE, bio);
583
584 alloc_new:
585         if (bio == NULL) {
```

```
586         bio = mpage_alloc(bdev, blocks[0] << (blkbits - 9),
587                 bio_get_nr_vecs(bdev), GFP_NOFS|__GFP_HIGH);
588         if (bio == NULL)
589             goto confused;
590     }
591
592     /*
593      * Must try to add the page before marking the buffer clean or
594      * the confused fail path above (OOM) will be very confused when
595      * it finds all bh marked clean (i.e. it will not write anything)
596      */
597     length = first_unmapped << blkbits;
598     if (bio_add_page(bio, page, length, 0) < length) {
599         bio = mpage_bio_submit(WRITE, bio);
600         goto alloc_new;
601     }
```

创建并设置 bio。

```
602
603     /*
604      * OK, we have our BIO, so we can now mark the buffers clean.    Make
605      * sure to only clean buffers which we know we'll be writing.
606      */
607     if (page_has_buffers(page)) {
608         struct buffer_head *head = page_buffers(page);
609         struct buffer_head *bh = head;
610         unsigned buffer_counter = 0;
611
612         do {
613             if (buffer_counter++ == first_unmapped)
614                 break;
615             clear_buffer_dirty(bh);
```

```
616                bh = bh->b_this_page;
617            } while (bh != head);
618
619            /*
620             * we cannot drop the bh if the page is not uptodate
621             * or a concurrent readpage would fail to serialize with the bh
622             * and it would read from disk before we reach the platter.
623             */
624            if (buffer_heads_over_limit && PageUptodate(page))
625                try_to_free_buffers(page);
626        }
627
628    BUG_ON(PageWriteback(page));
629    set_page_writeback(page);
630    unlock_page(page);
631    if (boundary || (first_unmapped != blocks_per_page)) {
632        bio = mpage_bio_submit(WRITE, bio);
```

提交 bio 写请求。

```
633            if (boundary_block) {
634                write_boundary_block(boundary_bdev,
635                        boundary_block, 1 << blkbits);
636            }
637        } else {
638            mpd->last_block_in_bio = blocks[blocks_per_page - 1];
639        }
640    goto out;
641
642 confused:
643    if (bio)
644        bio = mpage_bio_submit(WRITE, bio);
645
```

```
646      if (mpd->use_writepage) {
647          ret = mapping->a_ops->writepage(page, wbc);
648      } else {
649          ret = -EAGAIN;
650          goto out;
651      }
652      /*
653       * The caller has a ref on the inode, so *mapping is stable
654       */
655      mapping_set_error(mapping, ret);
656 out:
657      mpd->bio = bio;
658      return ret;
659 }
```

# 6. 参考书籍

[1] linux 内核 2.6.34 源码
[2] 深入了解 linux 内核 （第三版）
[3] linxu 内核源代码情景分析 毛德超
[4] 深入 linxu 内核架构
[5] http://blog.csdn.net/yunsongice(网络资源)