

程序 12-6 linux/fs/namei.c

```
1  /*
2  *  linux/fs/namei.c
3  *
4  *  (C) 1991  Linus Torvalds
5  */
6
7  /*
8  *  Some corrections by tytso.
9  */
10
11 #include <linux/sched.h> // 调度程序头文件, 定义了任务结构 task_struct、任务 0 的数据等。
12 #include <linux/kernel.h> // 内核头文件。含有一些内核常用函数的原形定义。
13 #include <asm/segment.h> // 段操作头文件。定义了有关段寄存器操作的嵌入式汇编函数。
14
15 #include <string.h> // 字符串头文件。主要定义了一些有关字符串操作的嵌入函数。
16 #include <fcntl.h> // 文件控制头文件。文件及其描述符的操作控制常数符号的定义。
17 #include <errno.h> // 错误号头文件。包含系统中各种出错号。
18 #include <const.h> // 常数符号头文件。目前仅定义 i 节点中 i_mode 字段的各标志位。
19 #include <sys/stat.h> // 文件状态头文件。含有文件或文件系统状态结构 stat{} 和常量。
20
21 // 由文件名查找对应 i 节点的内部函数。
22 static struct m_inode * namei(const char * filename, struct m_inode * base,
23     int follow_links);
24
25 // 下面宏中右侧表达式是访问数组的一种特殊使用方法。它基于这样的一个事实，即用数组名和
26 // 数组下标所表示的数组项（例如 a[b]）的值等同于使用数组首指针（地址）加上该项偏移地址
27 // 的形式的值 *(a + b)，同时可知项 a[b] 也可以表示成 b[a] 的形式。因此对于字符数组项形式
28 // 为 "LoveYou"[2]（或者 2["LoveYou"]）就等同于*("LoveYou" + 2)。另外，字符串 "LoveYou"
29 // 在内存中被存储的位置就是其地址，因此数组项 "LoveYou"[2] 的值就是该字符串中索引值为 2
30 // 的字符 "v" 所对应的 ASCII 码值 0x76，或用八进制表示就是 0166。在 C 语言中，字符也可以用
31 // 其 ASCII 码值来表示，方法是在字符的 ASCII 码值前面加一个反斜杠。例如字符 "v" 可以表示
32 // 成 "\x76" 或者 "\166"。因此对于不可显示的字符（例如 ASCII 码值为 0x00–0x1f 的控制字符）
33 // 就可用其 ASCII 码值来表示。
34
35 // 下面是访问模式宏。x 是头文件 include/fcntl.h 中第 7 行开始定义的文件访问（打开）标志。
36 // 这个宏根据文件访问标志 x 的值来索引双引号中对应的数值。双引号中有 4 个八进制数值（实
37 // 际表示 4 个控制字符）："\004\002\006\377"，分别表示读、写和执行的权限为：r、w、rw
38 // 和 wxrwxrwx，并且分别对应 x 的索引值 0–3。例如，如果 x 为 2，则该宏返回八进制值 006，
39 // 表示可读可写（rw）。另外，其中 0_ACCMODE = 00003，是索引值 x 的屏蔽码。
40 #define ACC_MODE(x) ("|004|002|006|377"[(x)&0_ACCMODE])
41
42 /*
43 * comment out this line if you want names > NAME_LEN chars to be
44 * truncated. Else they will be disallowed.
45 */
46
47 * 如果想让文件名长度 > NAME_LEN 个的字符被截掉，就将下面定义注释掉。
48 */
49
50 /* #define NO_TRUNCATE */
51
52 #define MAY_EXEC 1 // 可执行(可进入)。
53 #define MAY_WRITE 2 // 可写。
```

```

34 #define MAY_READ 4           // 可读。
35
36 /*
37  *      permission()
38  *
39  * is used to check for read/write/execute permissions on a file.
40  * I don't know if we should look at just the euid or both euid and
41  * uid, but that should be easily changed.
42 */
43 /* permission()
44 */
45  * 该函数用于检测一个文件的读/写/执行权限。我不知道是否只需检查 euid,
46  * 还是需要检查 euid 和 uid 两者，不过这很容易修改。
47 */
48 // 检测文件访问许可权限。
49 // 参数: inode - 文件的 i 节点指针; mask - 访问属性屏蔽码。
50 // 返回: 访问许可返回 1, 否则返回 0。
51
52 static int permission(struct m_inode * inode, int mask)
53 {
54     int mode = inode->i_mode;                                // 文件访问属性。
55
56     /* special case: not even root can read/write a deleted file */
57     /* 特殊情况: 即使是超级用户 (root) 也不能读/写一个已被删除的文件 */
58     /* 如果 i 节点有对应的设备, 但该 i 节点的链接计数值等于 0, 表示该文件已被删除, 则返回。 */
59     /* 否则, 如果进程的有效用户 id (euid) 与 i 节点的用户 id 相同, 则取文件宿主的访问权限。 */
60     /* 否则, 如果进程的有效组 id (egid) 与 i 节点的组 id 相同, 则取组用户的访问权限。 */
61     if (inode->i_dev && !inode->i_nlinks)
62         return 0;
63     else if (current->euid==inode->i_uid)
64         mode >= 6;
65     else if (in_group_p(inode->i_gid))
66         mode >= 3;
67
68     // 最后判断如果所取的的访问权限与屏蔽码相同, 或者是超级用户, 则返回 1, 否则返回 0。
69     if (((mode & mask & 0007) == mask) || suser())
70         return 1;
71     return 0;
72 }
73
74 /*
75  * ok, we cannot use strncmp, as the name is not in our data space.
76  * Thus we'll have to use match. No big problem. Match also makes
77  * some sanity tests.
78  *
79  * NOTE! unlike strncmp, match returns 1 for success, 0 for failure.
80 */
81
82 /*
83  * ok, 我们不能使用 strncmp 字符串比较函数, 因为名称不在我们的数据空间
84  * (不在内核空间)。因而我们只能使用 match()。问题不大, match() 同样
85  * 也处理一些完整的测试。
86  *
87  * 注意! 与 strncmp 不同的是 match() 成功时返回 1, 失败时返回 0。
88 */

```

```

//// 指定长度字符串比较函数。
// 参数: len - 比较的字符串长度; name - 文件名指针; de - 目录项结构。
// 返回: 相同返回 1, 不同返回 0。
// 第 68 行上定义了一个局部寄存器变量 same。该变量将被保存在 eax 寄存器中, 以便于高效
// 访问。
66 static int match(int len, const char * name, struct dir_entry * de)
67 {
68     register int same asm ("ax");
69
// 首先判断函数参数的有效性。如果目录项指针空, 或者目录项 i 节点等于 0, 或者要比较的
// 字符串长度超过文件名长度, 则返回 0 (不匹配)。如果比较的长度 len 等于 0 并且目录项
// 中文件名的第一个字符是 '.', 并且只有这么一个字符, 那么我们就认为是相同的, 因此返
// 回 1 (匹配)。如果要比较的长度 len 小于 NAME_LEN, 但是目录项中文件名长度超过 len,
// 则也返回 0 (不匹配)。
// 第 75 行上对目录项中文件名长度是否超过 len 的判断方法是检测 name[len] 是否为 NULL。
// 若长度超过 len, 则 name[len] 处就是一个不是 NULL 的普通字符。而对于长度为 len 的字符
// 串 name, 字符 name[len] 就应该是 NULL。
70     if (!de || !de->inode || len > NAME_LEN)
71         return 0;
72     /* "" means "." --> so paths like "/usr/lib//libc.a" work */
73     /* "" 当作 "." 来看待 --> 这样就能处理象 "/usr/lib//libc.a" 那样的路径名 */
74     if (!len && (de->name[0]=='.') && (de->name[1]=='\0'))
75         return 1;
76     if (len < NAME_LEN && de->name[len])
77         return 0;
// 然后使用嵌入汇编语句进行快速比较操作。它会在用户数据空间 (fs 段) 执行字符串的比较
// 操作。%0 - eax (比较结果 same); %1 - eax (eax 初值 0); %2 - esi (名字指针);
// %3 - edi (目录项名指针); %4 - ecx (比较的字节长度值 len)。
78     asm ("cld|n|t"           // 清方向标志位。
79             "fs ; repe ; cmpsb|n|t" // 用户空间执行循环比较[esi++]和[edi++]操作,
80             "setz %%al"           // 若比较结果一样 (zf=0) 则置 al=1 (same=eax)。
81             :"=a" (same)
82             :"0", "S" ((long) name), "D" ((long) de->name), "c" (len)
83             :"cx", "di", "si");
84     return same;           // 返回比较结果。
85 }
86 /*
87 *     find_entry()
88 *
89 * finds an entry in the specified directory with the wanted name. It
90 * returns the cache buffer in which the entry was found, and the entry
91 * itself (as a parameter - res_dir). It does NOT read the inode of the
92 * entry - you'll have to do that yourself if you want to.
93 *
94 * This also takes care of the few special cases due to '..'-traversal
95 * over a pseudo-root and a mount point.
96 */
/* *
*     find_entry()
*
* 在指定目录中寻找一个与名字匹配的目录项。返回一个含有找到目录项的高速
* 缓冲块以及目录项本身 (作为一个参数 - res_dir)。该函数并不读取目录项

```

```

* 的 i 节点 - 如果需要的话则自己操作。
*
* 由于有'..'目录项，因此在操作期间也会对几种特殊情况分别处理 - 比如横越
* 一个伪根目录以及安装点。
*/
//// 查找指定目录和文件名的目录项。
// 参数: *dir - 指定目录 i 节点的指针; name - 文件名; namelen - 文件名长度;
// 该函数在指定目录的数据(文件)中搜索指定文件名的目录项。并对指定文件名是'..'的情况根据当前进行的相关设置进行特殊处理。关于函数参数传递指针的指针的作用, 请参见 linux/sched.c 第 151 行前的注释。
// 返回: 成功则函数高速缓冲区指针, 并在*res_dir 处返回的目录项结构指针。失败则返回
// 空指针 NULL。
97 static struct buffer_head * find_entry(struct m_inode ** dir,
98         const char * name, int namelen, struct dir_entry ** res_dir)
99 {
100     int entries;
101     int block, i;
102     struct buffer_head * bh;
103     struct dir_entry * de;
104     struct super_block * sb;
105
// 同样, 本函数一上来也需要对函数参数的有效性进行判断和验证。如果我们在前面第 30 行
// 定义了符号常数 NO_TRUNCATE, 那么如果文件名长度超过最大长度 NAME_LEN, 则不予处理。
// 如果没有定义过 NO_TRUNCATE, 那么在文件名长度超过最大长度 NAME_LEN 时截短之。
106 #ifdef NO_TRUNCATE
107     if (namelen > NAME_LEN)
108         return NULL;
109 #else
110     if (namelen > NAME_LEN)
111         namelen = NAME_LEN;
112 #endif
// 首先计算本目录中目录项项数 entries。目录 i 节点 i_size 字段中含有本目录包含的数据
// 长度, 因此其除以一个目录项的长度(16 字节)即可得到该目录中目录项数。然后置空返回
// 目录项结构指针。
113     entries = (*dir)->i_size / (sizeof (struct dir_entry));
114     *res_dir = NULL;
// 接下来我们对目录项文件名是'..'的情况进行特殊处理。如果当前进程指定的根 i 节点就是
// 函数参数指定的目录, 则说明对于本进程来说, 这个目录就是它的伪根目录, 即进程只能访问
// 该目录中的项而不能后退到其父目录中去。也即对于该进程本目录就如同是文件系统的根
// 目录。因此我们需要将文件名修改为'.'。
// 否则, 如果该目录的 i 节点号等于 ROOT_INO(1 号)的话, 说明确实是文件系统的根 i 节点。
// 则取文件系统的超级块。如果被安装到的 i 节点存在, 则先放回原 i 节点, 然后对被安装到
// 的 i 节点进行处理。于是我们让*dir 指向该被安装到的 i 节点; 并且该 i 节点的引用数加 1。
// 即针对这种情况, 我们悄悄地进行了“偷梁换柱”工程:
115 /* check for '..', as we might have to do some "magic" for it */
/* 检查目录项'..', 因为我们可能需要对其进行特殊处理 */
116     if (namelen==2 && get_fs_byte(name)=='..' && get_fs_byte(name+1)=='..') {
117 /* '..' in a pseudo-root results in a faked '.' (just change namelen) */
/* 伪根中的'..'如同一个假'.'(只需改变名字长度) */
118         if ((*dir) == current->root)
119             namelen=1;
120         else if ((*dir)->i_num == ROOT_INO) {

```

```

121 /* '...' over a mount-point results in 'dir' being exchanged for the mounted
122 directory-inode. NOTE! We set mounted, so that we can input the new dir */
123 /* 在一个安装点上的 '...' 将导致目录交换到被安装文件系统的目录 i 节点上。注意！
124 由于我们设置了 mounted 标志，因而我们能够放回该新目录 */
125     sb=get_super((*dir)->i_dev);
126     if (sb->s_imount) {
127         iput(*dir);
128         (*dir)=sb->s_imount;
129         (*dir)->i_count++;
130     }

```

// 现在我们开始正常操作，查找指定文件名的目录项在什么地方。因此我们需要读取目录的数据，即取出目录 i 节点对应块设备数据区中的数据块（逻辑块）信息。这些逻辑块的块号保存在 i 节点结构的 i_zone[9]数组中。我们先取其中第 1 个块号。如果目录 i 节点指向的第一个直接磁盘块号为 0，则说明该目录竟然不含数据，这不正常。于是返回 NULL 退出。否则我们就从节点所在设备读取指定的目录项数据块。当然，如果不成功，则也返回 NULL 退出。

```

131     if (!block = (*dir)->i_zone[0]))
132         return NULL;
133     if (!bh = bread((*dir)->i_dev, block))
134         return NULL;

```

// 此时我们就在这个读取的目录 i 节点数据块中搜索匹配指定文件名的目录项。首先让 de 指向缓冲块中的数据块部分，并在不超过目录中目录项数的条件下，循环执行搜索。其中 i 是目录中的目录项索引号，在循环开始时初始化为 0。

```

135     i = 0;
136     de = (struct dir_entry *) bh->b_data;
137     while (i < entries) {
138         // 如果当前目录项数据块已经搜索完，还没有找到匹配的目录项，则释放当前目录项数据块。
139         // 再读入目录的下一个逻辑块。若这块为空，则只要还没有搜索完目录中的所有目录项，就
140         // 跳过该块，继续读目录的下一逻辑块。若该块不空，就让 de 指向该数据块，然后在其中
141         // 继续搜索。其中 141 行上 i/DIR_ENTRIES_PER_BLOCK 可得到当前搜索的目录项所在目录文
142         // 件中的块号，而 bmap() 函数 (inode.c, 第 142 行) 则可计算出在设备上对应的逻辑块号。
143         if ((char *)de >= BLOCK_SIZE+bh->b_data) {
144             brelse(bh);
145             bh = NULL;
146             if (!block = bmap(*dir, i/DIR_ENTRIES_PER_BLOCK)) ||
147                 !bh = bread((*dir)->i_dev, block)) {
148                 i += DIR_ENTRIES_PER_BLOCK;
149                 continue;
150             }
151             de = (struct dir_entry *) bh->b_data;
152         }

```

// 如果找到匹配的目录项的话，则返回该目录项结构指针 de 和该目录项 i 节点指针*dir 以及该目录项数据块指针 bh，并退出函数。否则继续在目录项数据块中比较下一个目录项。

```

153         if (match(namelen, name, de)) {
154             *res_dir = de;
155             return bh;
156         }
157         de++;
158         i++;
159     }

```

```

// 如果指定目录中的所有目录项都搜索完后, 还没有找到相应的目录项, 则释放目录的数据
// 块, 最后返回 NULL (失败)。
155     brelse(bh);
156     return NULL;
157 }
158
159 /*
160 *     add_entry()
161 *
162 *     adds a file entry to the specified directory, using the same
163 *     semantics as find_entry(). It returns NULL if it failed.
164 *
165 *     NOTE!! The inode part of 'de' is left at 0 - which means you
166 *     may not sleep between calling this and putting something into
167 *     the entry, as someone else might have used it while you slept.
168 */
169 /*
170 *     add_entry()
171 * 使用与 find_entry()同样的方法, 往指定目录中添加一指定文件名的目
172 * 录项。如果失败则返回 NULL。
173 *
174 * 注意! ! 'de' (指定目录项结构指针) 的 i 节点部分被设置为 0 - 这表
175 * 示在调用该函数和往目录项中添加信息之间不能去睡眠。因为如果睡眠,
176 * 那么其他人(进程)可能会使用了该目录项。
177 */
178 ///// 根据指定的目录和文件名添加目录项。
179 // 参数: dir - 指定目录的 i 节点; name - 文件名; namelen - 文件名长度;
180 // 返回: 高速缓冲区指针; res_dir - 返回的目录项结构指针;
181 static struct buffer_head * add_entry(struct m_inode * dir,
182                                     const char * name, int namelen, struct dir_entry ** res_dir)
183 {
184     int block, i;
185     struct buffer_head * bh;
186     struct dir_entry * de;
187
188     // 同样, 本函数一上来也需要对函数参数的有效性进行判断和验证。如果我们在前面第 30 行
189     // 定义了符号常数 NO_TRUNCATE, 那么如果文件名长度超过最大长度 NAME_LEN, 则不予处理。
190     // 如果没有定义过 NO_TRUNCATE, 那么在文件名长度超过最大长度 NAME_LEN 时截短之。
191     *res_dir = NULL;                                // 用于返回目录项结构指针。
192 #ifdef NO_TRUNCATE
193     if (namelen > NAME_LEN)
194         return NULL;
195 #else
196     if (namelen > NAME_LEN)
197         namelen = NAME_LEN;
198 #endif
199     // 现在我们开始操作, 向指定目录中添加一个指定文件名的目录项。因此我们需要先读取目录
200     // 的数据, 即取出目录 i 节点对应块设备数据区中的数据块(逻辑块)信息。这些逻辑块的块
201     // 号保存在 i 节点结构的 i_zone[9]数组中。我们先取其中第 1 个块号。如果目录 i 节点指向
202     // 的第一个直接磁盘块号为 0, 则说明该目录竟然不含数据, 这不正常。于是返回 NULL 退出。
203     // 否则我们就从节点所在设备读取指定的目录项数据块。当然, 如果不成功, 则也返回 NULL
204     // 退出。另外, 如果参数提供的文件名长度等于 0, 则也返回 NULL 退出。
205     if (!namelen)

```

```

185         return NULL;
186     if (! (block = dir->i_zone[0]))
187         return NULL;
188     if (! (bh = bread(dir->i_dev, block)))
189         return NULL;
// 此时我们就在这个目录 i 节点数据块中循环查找最后未使用的空目录项。首先让目录项结构
// 指针 de 指向缓冲块中的数据块部分，即第一个目录项处。其中 i 是目录中的目录项索引号，
// 在循环开始时初始化为 0。
190     i = 0;
191     de = (struct dir_entry *) bh->b_data;
192     while (1) {
// 如果当前目录项数据块已经搜索完毕，但还没有找到需要的空目录项，则释放当前目录项数
// 据块，再读入目录的下一个逻辑块。如果对应的逻辑块不存在就创建一块。若读取或创建操
// 作失败则返回空。如果此次读取的磁盘逻辑块数据返回的缓冲块指针为空，说明这块逻辑块
// 可能是因为不存在而新创建的空块，则把目录项索引值加上一块逻辑块所能容纳的目录项数
// DIR_ENTRIES_PER_BLOCK，用以跳过该块并继续搜索。否则说明新读入的块上有目录项数据，
// 于是让目录项结构指针 de 指向该块的缓冲块数据部分，然后在其中继续搜索。其中 196 行
// 上的 i/DIR_ENTRIES_PER_BLOCK 可计算得到当前搜索的目录项 i 所在目录文件中的块号，
// 而 create_block() 函数 (inode.c, 第 147 行) 则可读取或创建出在设备上对应的逻辑块。
193         if ((char *)de >= BLOCK_SIZE+bh->b_data) {
194             brelse(bh);
195             bh = NULL;
196             block = create_block(dir, i/DIR_ENTRIES_PER_BLOCK);
197             if (!block)
198                 return NULL;
199             if (! (bh = bread(dir->i_dev, block))) { // 若空则跳过该块继续。
200                 i += DIR_ENTRIES_PER_BLOCK;
201                 continue;
202             }
203             de = (struct dir_entry *) bh->b_data;
204         }
// 如果当前所操作的目录项序号 i 乘上目录结构大小所的长度值已经超过了该目录 i 节点信息
// 所指出的目录数据长度值 i_size，则说明整个目录文件数据中没有由于删除文件留下的空
// 目录项，因此我们只能把需要添加的新目录项附加到目录文件数据的末端处。于是对该处目
// 录项进行设置（置该目录项的 i 节点指针为空），并更新该目录文件的长度值（加上一个目
// 录项的长度），然后设置目录的 i 节点已修改标志，再更新该目录的改变时间为当前时间。
205         if (i*sizeof(struct dir_entry) >= dir->i_size) {
206             de->inode=0;
207             dir->i_size = (i+1)*sizeof(struct dir_entry);
208             dir->i_dirt = 1;
209             dir->i_ctime = CURRENT_TIME;
210         }
// 若当前搜索的目录项 de 的 i 节点为空，则表示找到一个还未使用的空闲目录项或是添加的
// 新目录项。于是更新目录的修改时间为当前时间，并从用户数据区复制文件名到该目录项的
// 文件名字段，置含有本目录项的相应高速缓冲块已修改标志。返回该目录项的指针以及该高
// 速缓冲块的指针，退出。
211         if (!de->inode) {
212             dir->i_mtime = CURRENT_TIME;
213             for (i=0; i < NAME_LEN ; i++)
214                 de->name[i]=(i<namelen)?get_fs_byte(name+i):0;
215             bh->b_dirt = 1;
216             *res_dir = de;
217             return bh;

```

```

218         }
219         de++;           // 如果该目录项已经被使用，则继续检测下一个目录项。
220         i++;
221     }
222     // 本函数执行不到这里。这也许是 Linus 在写这段代码时，先复制了上面 find_entry() 函数
223     // 的代码，而后修改成本函数的②。
224     brelse(bh);
225     return NULL;
226 }
227
228 ///// 查找符号链接的 i 节点。
229 // 参数: dir - 目录 i 节点; inode - 目录项 i 节点。
230 // 返回: 返回符号链接到文件的 i 节点指针。出错返回 NULL。
231 static struct m_inode * follow_link(struct m_inode * dir, struct m_inode * inode)
232 {
233     unsigned short fs;           // 用于临时保存 fs 段寄存器值。
234     struct buffer_head * bh;
235
236     // 首先判断函数参数的有效性。如果没有给出目录 i 节点，我们就使用进程任务结构中设置的
237     // 根 i 节点，并把链接数增 1。如果没有给出目录项 i 节点，则放回目录 i 节点后返回 NULL。
238     // 如果指定目录项不是一个符号链接，就直接返回目录项对应的 i 节点 inode。
239     if (!dir) {
240         dir = current->root;
241         dir->i_count++;
242     }
243     if (!inode) {
244         iput(dir);
245         return NULL;
246     }
247     if (!S_ISLNK(inode->i_mode)) {
248         iput(dir);
249         return inode;
250     }
251
252     // 然后取 fs 段寄存器值。fs 通常保存着指向任务数据段的选择符 0x17。如果 fs 没有指向用户
253     // 数据段，或者给出的目录项 i 节点第 1 个直接块块号等于 0，或者是读取第 1 个直接块出错，
254     // 则放回 dir 和 inode 两个 i 节点并返回 NULL 退出。否则说明现在 fs 正指向用户数据段、并且
255     // 我们已经成功地读取了这个符号链接目录项的文件内容，并且文件内容已经在 bh 指向的缓冲
256     // 块数据区中。实际上，这个缓冲块数据区中仅包含一个链接指向的文件路径名字符串。
257     __asm__ ("mov %%fs, %0": "=r" (fs));
258     if (fs != 0x17 || !inode->i_zone[0] ||
259         !(bh = bread(inode->i_dev, inode->i_zone[0]))) {
260         iput(dir);
261         iput(inode);
262         return NULL;
263     }
264
265     // 此时我们已经不需要符号链接目录项的 i 节点了，于是把它放回。现在碰到一个问题，那就
266     // 是内核函数处理的用户数据都是存放在用户数据空间中的，并使用了 fs 段寄存器来从用户
267     // 空间传递数据到内核空间中。而这里需要处理的数据却在内核空间中。因此为了正确地处理
268     // 位于内核中的用户数据，我们需要让 fs 段寄存器临时指向内核空间，即让 fs =0x10。并在
269     // 调用函数处理完后再恢复原 fs 的值。最后释放相应缓冲块，并返回 _namei() 解析得到的符
270     // 号链接指向的文件 i 节点。
271     iput(inode);
272     __asm__ ("mov %0, %%fs": "r" ((unsigned short) 0x10));

```

```

252     inode = namei(bh->b_data, dir, 0);
253     __asm__ ( "mov %0, %%fs": "r" (fs));
254     brelse(bh);
255     return inode;
256 }
257
258 /*
259 *      get_dir()
260 *
261 * Getdir traverses the pathname until it hits the topmost directory.
262 * It returns NULL on failure.
263 */
264
265 /*
266 *      get_dir()
267 *
268 * 该函数根据给出的路径名进行搜索，直到达到最顶端的目录。
269 * 如果失败则返回 NULL。
270 */
271
272 // 从指定目录开始搜寻指定路径名的目录（或文件名）的 i 节点。
273 // 参数： pathname - 路径名； inode - 指定起始目录的 i 节点。
274 // 返回： 目录或文件的 i 节点指针。失败时返回 NULL。
275
276 static struct m_inode * get_dir(const char * pathname, struct m_inode * inode)
277 {
278     char c;
279     const char * thisname;
280     struct buffer_head * bh;
281     int namelen, inr;
282     struct dir_entry * de;
283     struct m_inode * dir;
284
285     // 首先判断参数有效性。如果给出的指定目录的 i 节点指针 inode 为空，则使用当前进程的当
286     // 前工作目录 i 节点。如果用户指定路径名的第一个字符是'/'，则说明路径名是绝对路径名。
287     // 则应该从当前进程任务结构中设置的根（或伪根）i 节点开始操作。于是我们需要先放回参
288     // 数指定的或者设定的目录 i 节点，并取得进程使用的根 i 节点。然后把该 i 节点的引用计数
289     // 加 1，并删除路径名的第一个字符 '/'. 这样就可以保证当前进程只能以其设定的根 i 节点
290     // 作为搜索的起点。
291
292     if (!inode) {
293         inode = current->pwd;           // 进程的当前工作目录 i 节点。
294         inode->i_count++;
295     }
296
297     if ((c=get_fs_byte(pathname))== '/') {
298         iput(inode);                // 放回原 i 节点。
299         inode = current->root;       // 为进程指定的根 i 节点。
300         pathname++;
301         inode->i_count++;
302     }
303
304     // 然后针对路径名中的各个目录名部分和文件名进行循环处理。在循环处理过程中，我们先要
305     // 对当前正在处理的目录名部分的 i 节点进行有效性判断，并且把变量 thisname 指向当前正
306     // 在处理的目录名部分。如果该 i 节点表明当前处理的目录名部分不是目录类型，或者没有可
307     // 进入该目录的访问许可，则放回该 i 节点，并返回 NULL 退出。当然在刚进入循环时，当前
308     // 目录的 i 节点 inode 就是进程根 i 节点或者是当前工作目录的 i 节点，或者是参数指定的某
309     // 个搜索起始目录的 i 节点。
310
311     while (1) {

```

```

284         thisname = pathname;
285         if (!S_ISDIR(inode->i_mode) || !permission(inode, MAY_EXEC)) {
286             iput(inode);
287             return NULL;
288         }
// 每次循环我们处理路径名中一个目录名（或文件名）部分。因此在每次循环中我们都要从路
// 径名字符串中分离出一个目录名（或文件名）。方法是从当前路径名指针 pathname 开始处
// 搜索检测字符，直到字符是一个结尾符（NULL）或者是一个'/'字符。此时变量 namelen 正
// 好是当前处理目录名部分的长度，而变量 thisname 正指向该目录名部分的开始处。此时如
// 果字符是结尾符 NULL，则表明已经搜索到路径名末尾，并已到达最后指定目录名或文件名，
// 则返回该 i 节点指针退出。
// 注意！如果路径名中最后一个名称也是一个目录名，但其后面没有加上 '/' 字符，则函数不
// 会返回该最后目录名的 i 节点！例如：对于路径名/usr/src/linux，该函数将只返回 src/目
// 录名的 i 节点。
289         for (namelen=0; (c=get_fs_byte(pathname++))&&(c!= '/') ;namelen++)
290             /* nothing */;
291         if (!c)
292             return inode;
// 在得到当前目录名部分（或文件名）后，我们调用查找目录项函数 find_entry() 在当前处
// 理的目录中寻找指定名称的目录项。如果没有找到，则放回该 i 节点，并返回 NULL 退出。
// 然后在找到的目录项中取出其 i 节点号 inr 和设备号 idev，释放包含该目录项的高速缓冲
// 块并放回该 i 节点。然后取节点号 inr 的 i 节点 inode，并以该目录项为当前目录继续循
// 环处理路径名中的下一目录名部分（或文件名）。如果当前处理的目录项是一个符号链接
// 名，则使用 follow_link() 就可以得到其指向的目录项名的 i 节点。
293         if (! (bh = find_entry(&inode, thisname, namelen, &de))) {
294             iput(inode);
295             return NULL;
296         }
297         inr = de->inode;                                // 当前目录名部分的 i 节点号。
298         brelse(bh);
299         dir = inode;
300         if (! (inode = iget(dir->i_dev, inr))) {      // 取 i 节点内容。
301             iput(dir);
302             return NULL;
303         }
304         if (! (inode = follow_link(dir, inode)))
305             return NULL;
306     }
307 }
308
309 */
310 *      dir_namei()
311 *
312 * dir_namei() returns the inode of the directory of the
313 * specified name, and the name within that directory.
314 */
315 */
316 *      dir_namei()
317 *
318 * dir_namei() 函数返回指定目录名的 i 节点指针，以及在最顶层
319 * 目录的名称。
320 */
321 // 参数： pathname - 目录路径名； namelen - 路径名长度； name - 返回的最顶层目录名。

```

```

// base - 搜索起始目录的 i 节点。
// 返回: 指定目录名最顶层目录的 i 节点指针和最顶层目录名称及长度。出错时返回 NULL。
// 注意! ! 这里“最顶层目录”是指路径名中最靠近末端的目录。
315 static struct m_inode * dir_namei(const char * pathname,
316     int * namelen, const char ** name, struct m_inode * base)
317 {
318     char c;
319     const char * basename;
320     struct m_inode * dir;
321
322     // 首先取得指定路径名最顶层目录的 i 节点。然后对路径名 pathname 进行搜索检测，查出最后
323     // 一个'/'字符后面的名字字符串，计算其长度，并且返回最顶层目录的 i 节点指针。注意！如
324     // 果路径名最后一个字符是斜杠字符'/'，那么返回的目录名为空，并且长度为 0。但返回的 i
325     // 节点指针仍然指向最后一个'/'字符前目录名的 i 节点。参见第 289 行上的“注意”说明。
326     if (!(dir = get_dir(pathname, base))) // base 是指定的起始目录 i 节点。
327         return NULL;
328     basename = pathname;
329     while (c=get_fs_byte(pathname++))
330         if (c=='/')
331             basename=pathname;
332     *namelen = pathname-basename-1;
333     *name = basename;
334     return dir;
335 }
336
337 // // 取指定路径名的 i 节点内部函数。
338 // 参数: pathname - 路径名; base - 搜索起点目录 i 节点; follow_links - 是否跟随
339 // 符号链接的标志，1 - 需要，0 不需要。
340 // 返回: 对应的 i 节点。
341 struct m_inode * namei(const char * pathname, struct m_inode * base,
342     int follow_links)
343 {
344     const char * basename;
345     int inr, namelen;
346     struct m_inode * inode;
347     struct buffer_head * bh;
348     struct dir_entry * de;
349
350     // 首先查找指定路径名中最顶层目录的目录名并得到其 i 节点。若不存在，则返回 NULL 退出。
351     // 如果返回的最顶层名字的长度是 0，则表示该路径名以一个目录名为最后一项。因此说明我
352     // 们已经找到对应目录的 i 节点，可以直接返回该 i 节点退出。如果返回的名字长度不是 0，
353     // 则我们以指定的起始目录 base，再次调用 dir_namei() 函数来搜索顶层目录名，并根据返回
354     // 的信息作类似判断。
355     if (!(dir = dir_namei(pathname, &namelen, &basename)))
356         return NULL;
357     if (!namelen) /* special case: '/usr/' etc */
358         return dir; /* 对应于 '/usr/' 等情况 */
359     if (!(base = dir_namei(pathname, &namelen, &basename, base)))
360         return NULL;
361     if (!namelen) /* special case: '/usr/' etc */
362         return base;
363
364     // 然后在返回的顶层目录中寻找指定文件名目录项的 i 节点。注意！因为如果最后也是一个目
365     // 录名，但其后没有加'/'，则不会返回该最后目录的 i 节点！例如：/usr/src/linux，将只

```

```

// 返回 src/目录名的 i 节点。因为函数 dir_namei() 将不以'/'结束的最后一个名字当作一个
// 文件名来看待，因此这里需要单独对这种情况使用寻找目录项 i 节点函数 find_entry() 进行
// 处理。此时 de 中含有寻找到的目录项指针，而 dir 是包含该目录项的目录的 i 节点指针。
346     bh = find_entry(&base, basename, namelen, &de);
347     if (!bh) {
348         iput(base);
349         return NULL;
350     }
// 接着取该目录项的 i 节点号，并释放包含该目录项的高速缓冲块并放回目录 i 节点。然后取
// 对应节点号的 i 节点，修改其被访问时间为当前时间，并置已修改标志。最后返回该 i 节点
// 指针 inode。如果当前处理的目录项是一个符号链接名，则使用 follow_link() 得到其指向的
// 目录项名的 i 节点。
351     inr = de->inode;
352     brelse(bh);
353     if (!(inode = iget(base->i_dev, inr))) {
354         iput(base);
355         return NULL;
356     }
357     if (follow_links)
358         inode = follow_link(base, inode);
359     else
360         iput(base);
361     inode->i_atime=CURRENT_TIME;
362     inode->i_dirt=1;
363     return inode;
364 }
365
//// 取指定路径名的 i 节点，不跟随符号链接。
// 参数: pathname - 路径名。
// 返回: 对应的 i 节点。
366 struct m_inode * lnamei(const char * pathname)
367 {
368     return namei(pathname, NULL, 0);
369 }
370
371 /*
372 *      namei()
373 *
374 * is used by most simple commands to get the inode of a specified name.
375 * Open, link etc use their own routines, but this is enough for things
376 * like 'chmod' etc.
377 */
378 /*
379 *      namei()
380 *
381 * 该函数被许多简单命令用于取得指定路径名称的 i 节点。open、link 等则使用它们
382 * 自己的相应函数。但对于象修改模式'chmod'等这样的命令，该函数已足够用了。
383 */
384
//// 取指定路径名的 i 节点，跟随符号链接。
// 参数: pathname - 路径名。
// 返回: 对应的 i 节点。
385 struct m_inode * namei(const char * pathname)
386 {

```

```

380         return namei(pathname, NULL, 1);
381     }
382
383     /*
384     *      open_namei()
385     *
386     * namei for open - this is in fact almost the whole open-routine.
387     */
388
389     /*
390     *      open()
391     *
392     * open() 函数使用的 namei 函数 - 这其实几乎是完整的打开文件程序。
393     */
394
395     ///// 文件打开 namei 函数。
396
397     // 参数 filename 是文件路径名, flag 是打开文件标志, 可取值 O_RDONLY (只读)、O_WRONLY
398     // (只写) 或 O_RDWR (读写), 以及 O_CREAT (创建)、O_EXCL (被创建文件必须不存在)、
399     // O_APPEND (在文件尾添加数据) 等其他一些标志的组合。如果本调用创建了一个新文件, 则
400     // mode 就用于指定文件的许可属性。这些属性有 S_IRWXU (文件宿主具有读、写和执行权限)、
401     // S_IRUSR (用户具有读文件权限)、S_IWUGR (组成员具有读、写和执行权限) 等等。对于新
402     // 创建的文件, 这些属性只应用于将来对文件的访问, 创建了只读文件的打开调用也将返回一
403     // 个可读写的文件句柄。参见包含文件 sys/stat.h、fcntl.h。
404
405     // 返回: 成功返回 0, 否则返回出错码; res_inode - 返回对应文件路径名的 i 节点指针。
406
407     int open_namei(const char * pathname, int flag, int mode,
408                 struct m_inode ** res_inode)
409     {
410         const char * basename;
411         int inr, dev, namelen;
412         struct m_inode * dir, *inode;
413         struct buffer_head * bh;
414         struct dir_entry * de;
415
416         // 首先对函数参数进行合理的处理。如果文件访问模式标志是只读 (0), 但是文件截零标志
417         // O_TRUNC 却置位了, 则在文件打开标志中添加只写标志 O_WRONLY。这样做的原因是由于截零
418         // 标志 O_TRUNC 必须在文件可写情况下才有效。然后使用当前进程的文件访问许可屏蔽码, 屏
419         // 蔽掉给定模式中的相应位, 并添上普通文件标志 I_REGULAR。该标志将用于打开的文件不存
420         // 在而需要创建文件时, 作为新文件的默认属性。参见下面 411 行上的注释。
421
422         if ((flag & O_TRUNC) && !(flag & O_ACCMODE))
423             flag |= O_WRONLY;
424         mode &= 0777 & ~current->umask;
425         mode |= I_REGULAR; // 常规文件标志。见参见 include/const.h 文件)。
426
427         // 然后根据指定的路径名寻找到对应的 i 节点, 以及最顶端目录名及其长度。此时如果最顶端
428         // 目录名长度为 0 (例如 '/usr/' 这种路径名的情况), 那么若操作不是读写、创建和文件长
429         // 度截 0, 则表示是在打开一个目录名文件操作。于是直接返回该目录的 i 节点并返回 0 退出。
430         // 否则说明进程操作非法, 于是放回该 i 节点, 返回出错码。
431
432         if (!(dir = dir_namei(pathname, &namelen, &basename, NULL)))
433             return ENOENT;
434         if (!namelen) { /* special case: '/usr/' etc */
435             if (!(flag & (O_ACCMODE|O_CREAT|O_TRUNC))) {
436                 *res_inode=dir;
437                 return 0;
438             }
439             iput(dir);
440             return EISDIR;

```

```

410         }
411         // 接着根据上面得到的最顶层目录名的 i 节点 dir, 在其中查找取得路径名字符串中最后的文
412         // 件名对应的目录项结构 de, 并同时得到该目录项所在的高速缓冲区指针。如果该高速缓冲
413         // 指针为 NULL, 则表示没有找到对应文件名的目录项, 因此只可能是创建文件操作。此时如
414         // 果不是创建文件, 则放回该目录的 i 节点, 返回出错号退出。如果用户在该目录没有写的权
415         // 力, 则放回该目录的 i 节点, 返回出错号退出。
416         bh = find\_entry(&dir, basename, namelen, &de);
417         if (!bh) {
418             if (!(flag & O\_CREAT)) {
419                 iuput(dir);
420                 return -ENOENT;
421             }
422             if (!permission(dir, MAY\_WRITE)) {
423                 iuput(dir);
424                 return -EACCES;
425             }
426             // 现在我们确定了是创建操作并且有写操作许可。因此我们就在目录 i 节点对应设备上申请
427             // 一个新的 i 节点给路径名上指定的文件使用。若失败则放回目录的 i 节点, 并返回没有空
428             // 间出错码。否则使用该新 i 节点, 对其进行初始设置: 置节点的用户 id; 对应节点访问模
429             // 式; 置已修改标志。然后并在指定目录 dir 中添加一个新目录项。
430             inode = new\_inode(dir->i_dev);
431             if (!inode) {
432                 iuput(dir);
433                 return -ENOSPC;
434             }
435             inode->i_uid = current->euid;
436             inode->i_mode = mode;
437             inode->i_dirt = 1;
438             bh = add\_entry(dir, basename, namelen, &de);
439             // 如果返回的应该含有新目录项的高速缓冲区指针为 NULL, 则表示添加目录项操作失败。于是
440             // 将该新 i 节点的引用连接计数减 1, 放回该 i 节点与目录的 i 节点并返回出错码退出。否则
441             // 说明添加目录项操作成功。于是我们来设置该新目录项的一些初始值: 置 i 节点号为新申请
442             // 到的 i 节点的号码; 并置高速缓冲区已修改标志。然后释放该高速缓冲区, 放回目录的 i 节
443             // 点。返回新目录项的 i 节点指针, 并成功退出。
444             if (!bh) {
445                 inode->i_nlinks--;
446                 iuput(inode);
447                 iuput(dir);
448                 return -ENOSPC;
449             }
450             de->inode = inode->i_num;
451             bh->b_dirt = 1;
452             brelse(bh);
453             iuput(dir);
454             *res_inode = inode;
455             return 0;
456         }
457         // 若上面 (411 行) 在目录中取文件名对应目录项结构的操作成功 (即 bh 不为 NULL), 则说
458         // 明指定打开的文件已经存在。于是取出该目录项的 i 节点号和其所在设备号, 并释放该高速
459         // 缓冲区以及放回目录的 i 节点。如果此时独占操作标志 O\_EXCL 置位, 但现在文件已经存在,
460         // 则返回文件已存在出错码退出。
461         inr = de->inode;
462         dev = dir->i_dev;

```

```

445     brelse (bh);
446     if (flag & O_EXCL) {
447         iput(dir);
448         return -EXIST;
449     }
450 // 然后我们读取该目录项的 i 节点内容。若该 i 节点是一个目录的 i 节点并且访问模式是只
451 // 写或读写，或者没有访问的许可权限，则放回该 i 节点，返回访问权限出错码退出。
452     if (!(inode = follow_link(dir, iget(dev, inr))))
453         return -EACCES;
454     if ((S_ISDIR(inode->i_mode) && (flag & O_ACCMODE)) ||
455         !permission(inode, ACC_MODE(flag))) {
456         iput(inode);
457         return -EPERM;
458     }
459 // 接着我们更新该 i 节点的访问时间字段值为当前时间。如果设立了截 0 标志，则将该 i 节
460 // 点的文件长度截为 0。最后返回该目录项 i 节点的指针，并返回 0（成功）。
461     inode->i_atime = CURRENT_TIME;
462     if (flag & O_TRUNC)
463         truncate(inode);
464     *res_inode = inode;
465     return 0;
466 }
467 // //// 创建一个设备特殊文件或普通文件节点（node）。
468 // 该函数创建名称为 filename，由 mode 和 dev 指定的文件系统节点（普通文件、设备特殊文
469 // 件或命名管道）。
470 // 参数：filename - 路径名；mode - 指定使用许可以及所创建节点的类型；dev - 设备号。
471 // 返回：成功则返回 0，否则返回出错码。
472 int sys_mknod(const char * filename, int mode, int dev)
473 {
474     const char * basename;
475     int namelen;
476     struct m_inode * dir, * inode;
477     struct buffer_head * bh;
478     struct dir_entry * de;
479
480 // 首先检查操作许可和参数的有效性并取路径名中顶层目录的 i 节点。如果不是超级用户，则
481 // 返回访问许可出错码。如果找不到对应路径名中顶层目录的 i 节点，则返回出错码。如果最
482 // 顶端的文件名长度为 0，则说明给出的路径名最后没有指定文件名，放回该目录 i 节点，返
483 // 回出错码退出。如果在该目录中没有写的权限，则放回该目录的 i 节点，返回访问许可出错
484 // 码退出。如果不是超级用户，则返回访问许可出错码。
485     if (!suser())
486         return -EPERM;
487     if (!(dir = dir_namei(filename, &namelen, &basename, NULL)))
488         return -ENOENT;
489     if (!namelen) {
490         iput(dir);
491         return -ENOENT;
492     }
493     if (!permission(dir, MAY_WRITE)) {
494         iput(dir);
495         return -EPERM;
496     }
497 }

```

```

// 然后我们搜索一下路径名指定的文件是否已经存在。若已经存在则不能创建同名文件节点。
// 如果对应路径名上最后的文件名的目录项已经存在，则释放包含该目录项的缓冲区块并放回
// 目录的 i 节点，返回文件已经存在的出错码退出。
484     bh = find_entry(&dir, basename, namelen, &de);
485     if (bh) {
486         brelse(bh);
487         iput(dir);
488         return -EXIST;
489     }
// 否则我们就申请一个新的 i 节点，并设置该 i 节点的属性模式。如果要创建的是块设备文件
// 或者是字符设备文件，则令 i 节点的直接逻辑块指针 0 等于设备号。即对于设备文件来说，
// 其 i 节点的 i_zone[0] 中存放的是该设备文件所定义设备的设备号。然后设置该 i 节点的修
// 改时间、访问时间为当前时间，并设置 i 节点已修改标志。
490     inode = new_inode(dir->i_dev);
491     if (!inode) { // 若不成功则放回目录 i 节点，返回无空间出错码退出。
492         iput(dir);
493         return -ENOSPC;
494     }
495     inode->i_mode = mode;
496     if (S_ISBLK(mode) || S_ISCHR(mode))
497         inode->i_zone[0] = dev;
498     inode->i_mtime = inode->i_atime = CURRENT_TIME;
499     inode->i_dirt = 1;
// 接着为这个新的 i 节点在目录中新添加一个目录项。如果失败（包含该目录项的高速缓冲
// 块指针为 NULL），则放回目录的 i 节点；把所申请的 i 节点引用连接计数复位，并放回该
// i 节点，返回出错码退出。
500     bh = add_entry(dir, basename, namelen, &de);
501     if (!bh) {
502         iput(dir);
503         inode->i_nlinks=0;
504         iput(inode);
505         return -ENOSPC;
506     }
// 现在添加目录项操作也成功了，于是我们来设置这个目录项内容。令该目录项的 i 节点字
// 段等于新 i 节点号，并置高速缓冲区已修改标志，放回目录和新的 i 节点，释放高速缓冲
// 区，最后返回 0(成功)。
507     de->inode = inode->i_num;
508     bh->b_dirt = 1;
509     iput(dir);
510     iput(inode);
511     brelse(bh);
512     return 0;
513 }
514
//// 创建一个目录。
// 参数: pathname - 路径名; mode - 目录使用的权限属性。
// 返回: 成功则返回 0, 否则返回出错码。
515 int sys_mkdir(const char * pathname, int mode)
516 {
517     const char * basename;
518     int namelen;
519     struct m_inode * dir, * inode;
520     struct buffer_head * bh, *dir_block;

```

```

521     struct dir_entry * de;
522
// 首先检查参数的有效性并取路径名中顶层目录的 i 节点。如果找不到对应路径名中顶层目录
// 的 i 节点，则返回出错码。如果最顶端的文件名长度为 0，则说明给出的路径名最后没有指
// 定文件名，放回该目录 i 节点，返回出错码退出。如果在该目录中没有写的权限，则放回该
// 目录的 i 节点，返回访问许可出错码退出。如果不是超级用户，则返回访问许可出错码。
523     if (!(dir = dir_namei(pathname, &namelen, &basename, NULL)))
524         return ENOENT;
525     if (!namelen) {
526         iput(dir);
527         return ENOENT;
528     }
529     if (!permission(dir, MAY_WRITE)) {
530         iput(dir);
531         return EPERM;
532     }
// 然后我们搜索一下路径名指定的目录名是否已经存在。若已经存在则不能创建同名目录节点。
// 如果对应路径名上最后的目录项已经存在，则释放包含该目录项的缓冲区块并放回
// 目录的 i 节点，返回文件已经存在的出错码退出。否则我们就申请一个新的 i 节点，并设置
// 该 i 节点的属性模式：置该新 i 节点对应的文件长度为 32 字节（2 个目录项的大小）、置
// 节点已修改标志，以及节点的修改时间和访问时间。2 个目录项分别用于'.' 和'..' 目录。
533     bh = find_entry(&dir, basename, namelen, &de);
534     if (bh) {
535         brelse(bh);
536         iput(dir);
537         return EXIST;
538     }
539     inode = new_inode(dir->i_dev);
540     if (!inode) { // 若不成功则放回目录的 i 节点，返回无空间出错码。
541         iput(dir);
542         return ENOSPC;
543     }
544     inode->i_size = 32;
545     inode->i_dirt = 1;
546     inode->i_mtime = inode->i_atime = CURRENT_TIME;
// 接着为该新 i 节点申请一用于保存目录项数据的磁盘块，并令 i 节点的第一个直接块指针等
// 于该块号。如果申请失败则放回对应目录的 i 节点；复位新申请的 i 节点连接计数；放回该
// 新的 i 节点，返回没有空间出错码退出。否则置该新的 i 节点已修改标志。
547     if (!(inode->i_zone[0] = new_block(inode->i_dev))) {
548         iput(dir);
549         inode->i_nlinks--;
550         iput(inode);
551         return ENOSPC;
552     }
553     inode->i_dirt = 1;
// 从设备上读取新申请的磁盘块（目的是把对应块放到高速缓冲区中）。若出错，则放回对应
// 目录的 i 节点；释放申请的磁盘块；复位新申请的 i 节点连接计数；放回该新的 i 节点，返
// 回没有空间出错码退出。
554     if (!(dir_block = bread(inode->i_dev, inode->i_zone[0]))) {
555         iput(dir);
556         inode->i_nlinks--;
557         iput(inode);
558         return ERROR;

```

```

559         }
560         // 然后我们在缓冲块中建立起所创建目录文件中的 2 个默认的新目录项 ('.' 和'..') 结构数
561         // 据。首先令 de 指向存放目录项的数据块，然后置该目录项的 i 节点号字段等于新申请的 i
562         // 节点号，名字字段等于".."。然后 de 指向下一个目录项结构，并在该结构中存放上级目录
563         // 的 i 节点号和名字".."。然后设置该高速缓冲块已修改标志，并释放该缓冲块。再初始化
564         // 设置新 i 节点的模式字段，并置该 i 节点已修改标志。
565         de = (struct dir_entry *) dir_block->b_data;
566         de->inode=inode->i_num;                                // 设置'.'目录项。
567         strcpy(de->name, ".");
568         de++;
569         de->inode = dir->i_num;                                // 设置'..'目录项。
570         strcpy(de->name, "..");
571         inode->i_nlinks = 2;
572         dir_block->b_dirt = 1;
573         brelse(dir_block);
574         inode->i_mode = I_DIRECTORY | (mode & 0777 & ~current->umask);
575         inode->i_dirt = 1;
576         // 现在我们在指定目录中新添加一个目录项，用于存放新建目录的 i 节点和目录名。如果失
577         // 败（包含该目录项的高速缓冲区指针为 NULL），则放回目录的 i 节点；所申请的 i 节点引
578         // 用连接计数复位，并放回该 i 节点。返回出错码退出。
579         bh = add_entry(dir, basename, namelen, &de);
580         if (!bh) {
581             iput(dir);
582             inode->i_nlinks=0;
583             iput(inode);
584             return -ENOSPC;
585         }
586         // 最后令该新目录项的 i 节点字段等于新 i 节点号，并置高速缓冲块已修改标志，放回目录
587         // 和新的 i 节点，释放高速缓冲区，最后返回 0（成功）。
588         de->inode = inode->i_num;
589         bh->b_dirt = 1;
590         dir->i_nlinks++;
591         dir->i_dirt = 1;
592         iput(dir);
593         iput(inode);
594         brelse(bh);
595         return 0;
596     }
597
598 /*
599  * routine to check that the specified directory is empty (for rmdir)
600  */
601 /*
602  * 用于检查指定的目录是否为空的子程序（用于 rmdir 系统调用）。
603  */
604 /**
605  * /////////////////////////////////////////////////////////////////// 检查指定目录是否空。
606  * // 参数: inode - 指定目录的 i 节点指针。
607  * // 返回: 1 - 目录中是空的; 0 - 不空。
608  */
609 static int empty_dir(struct m_inode * inode)
610 {
611     int nr,block;
612     int len;
613     struct buffer_head * bh;

```

```

596     struct dir_entry * de;
597
// 首先计算指定目录中现有目录项个数并检查开始两个特定目录项中信息是否正确。一个目录
// 中应该起码有 2 个目录项：即“.”和“..”。如果目录项个数少于 2 个或者该目录 i 节点的第
// 1 个直接块没有指向任何磁盘块号，或者该直接块读不出，则显示警告信息“设备 dev 上目
// 录错”，返回 0(失败)。
598     len = inode->i_size / sizeof (struct dir_entry);           // 目录中目录项个数。
599     if (len<2 || !inode->i_zone[0] ||
600         !(bh=bread(inode->i_dev, inode->i_zone[0]))) {
601         printf("warning - bad directory on dev %04x\n", inode->i_dev);
602         return 0;
603     }
// 此时 bh 所指缓冲块中含有目录项数据。我们让目录项指针 de 指向缓冲块中第 1 个目录项。
// 对于第 1 个目录项 (“.”），它的 i 节点号字段 inode 应该等于当前目录的 i 节点号。对于
// 第 2 个目录项 (“..”），它的 i 节点号字段 inode 应该等于上一层目录的 i 节点号，不会
// 为 0。因此如果第 1 个目录项的 i 节点号字段值不等于该目录的 i 节点号，或者第 2 个目录
// 项的 i 节点号字段为零，或者两个目录项的名字字段不分别等于“.”和“..”，则显示出错警
// 告信息“设备 dev 上目录错”，并返回 0。
604     de = (struct dir_entry *) bh->b_data;
605     if (de[0].inode != inode->i_num || !de[1].inode ||
606         strcmp(".", de[0].name) || strcmp("..", de[1].name)) {
607         printf("warning - bad directory on dev %04x\n", inode->i_dev);
608         return 0;
609     }
// 然后我们令 nr 等于目录项序号（从 0 开始计）；de 指向第三个目录项。并循环检测该目录
// 中其余所有的 (len - 2) 个目录项，看有没有目录项的 i 节点号字段不为 0（被使用）。
610     nr = 2;
611     de += 2;
612     while (nr<len) {
// 如果该块磁盘块中的目录项已经全部检测完毕，则释放该磁盘块的缓冲块，并读取目录数据
// 文件中下一块含有目录项的磁盘块。读取的方法是根据当前检测的目录项序号 nr 计算出对
// 应目录项在目录数据文件中的数据块号 (nr/DIR_ENTRIES_PER_BLOCK)，然后使用 bmap()
// 函数取得对应的盘块号 block，再使用读设备盘块函数 bread() 把相应盘块读入缓冲块中，
// 并返回该缓冲块的指针。若所读取的相应盘块没有使用（或已经不用，如文件已经删除等），
// 则继续读下一块，若读不出，则出错返回 0。否则让 de 指向读出块的首个目录项。
613     if ((void *) de >= (void *) (bh->b_data+BLOCK_SIZE)) {
614         brelse(bh);
615         block=bmap(inode, nr/DIR_ENTRIES_PER_BLOCK);
616         if (!block) {
617             nr += DIR_ENTRIES_PER_BLOCK;
618             continue;
619         }
620         if (!(bh=bread(inode->i_dev, block)))
621             return 0;
622         de = (struct dir_entry *) bh->b_data;
623     }
// 对于 de 指向的当前目录项，如果该目录项的 i 节点号字段不等于 0，则表示该目录项目前正
// 被使用，则释放该高速缓冲区，返回 0 退出。否则，若还没有查询完该目录中的所有目录项，
// 则把目录项序号 nr 增 1、de 指向下一个目录项，继续检测。
624     if (de->inode) {
625         brelse(bh);
626         return 0;
627     }

```

```

628             de++;
629             nr++;
630         }
631         // 执行到这里说明该目录中没有找到已用的目录项(当然除了头两个以外)，则释放缓冲块返回 1。
632         brelse(bh);
633         return 1;
634     }
635     // // // 删除目录。
636     // 参数: name - 目录名 (路径名)。
637     // 返回: 返回 0 表示成功, 否则返回出错号。
638     int sys_rmdir(const char * name)
639     {
640         const char * basename;
641         int namelen;
642         struct m_inode * dir, * inode;
643         struct buffer_head * bh;
644         struct dir_entry * de;
645
646         // 首先检查参数的有效性并取路径名中顶层目录的 i 节点。如果找不到对应路径名中顶层目录
647         // 的 i 节点，则返回出错码。如果最顶端的文件名长度为 0，则说明给出的路径名最后没有指
648         // 定文件名，放回该目录 i 节点，返回出错码退出。如果在该目录中没有写的权限，则放回该
649         // 目录的 i 节点，返回访问许可出错码退出。如果不是超级用户，则返回访问许可出错码。
650         if (!(dir = dir_namei(name, &namelen, &basename, NULL)))
651             return -ENOENT;
652         if (!namelen) {
653             iput(dir);
654             return -ENOENT;
655         }
656         if (!permission(dir, MAY_WRITE)) {
657             iput(dir);
658             return -EPERM;
659         }
660
661         // 然后根据指定目录的 i 节点和目录名利用函数 find_entry() 寻找对应目录项，并返回包含该
662         // 目录项的缓冲块指针 bh、包含该目录项的目录的 i 节点指针 dir 和该目录项指针 de。再根据
663         // 该目录项 de 中的 i 节点号利用 ige() 函数得到对应的 i 节点 inode。如果对应路径名上最
664         // 后目录名的目录项不存在，则释放包含该目录项的高速缓冲区，放回目录的 i 节点，返回文
665         // 件已经存在出错码，并退出。如果取目录项的 i 节点出错，则放回目录的 i 节点，并释放含
666         // 有目录项的高速缓冲区，返回出错号。
667         bh = find_entry(&dir, basename, namelen, &de);
668         if (!bh) {
669             iput(dir);
670             return -ENOENT;
671         }
672         if (!(inode = ige(dir->i_dev, de->inode))) {
673             iput(dir);
674             brelse(bh);
675             return -EPERM;
676         }
677
678         // 此时我们已有包含要被删除目录项的目录 i 节点 dir、要被删除目录项的 i 节点 inode 和要
679         // 被删除目录项指针 de。下面我们通过对这 3 个对象中信息的检查来验证删除操作的可行性。
680
681         // 若该目录设置了受限删除标志并且进程的有效用户 id (euid) 不是 root，并且进程的有效

```

```

// 用户 id (euid) 不等于该 i 节点的用户 id，则表示当前进程没有权限删除该目录，于是放
// 回包含要删除目录名的目录 i 节点和该要删除目录的 i 节点，然后释放高速缓冲区，返回
// 出错码。
663     if ((dir->i_mode & S_ISVTX) && current->euid &&
664         inode->i_uid != current->euid) {
665         iput(dir);
666         iput(inode);
667         brelse(bh);
668         return -EPERM;
669     }
// 如果要被删除的目录项 i 节点的设备号不等于包含该目录项的目录的设备号，或者该被删除
// 目录的引用连接计数大于 1（表示有符号连接等），则不能删除该目录。于是释放包含要删
// 除目录名的目录 i 节点和该要删除目录的 i 节点，释放高速缓冲块，返回出错码。
670     if (inode->i_dev != dir->i_dev || inode->i_count>1) {
671         iput(dir);
672         iput(inode);
673         brelse(bh);
674         return -EPERM;
675     }
// 如果要被删除目录的目录项 i 节点就等于包含该需删除目录的目录 i 节点，则表示试图删除
// “.”目录，这是不允许的。于是放回包含要删除目录名的目录 i 节点和要删除目录的 i 节点，
// 释放高速缓冲块，返回出错码。
676     if (inode == dir) { /* we may not delete ".", but "../dir" is ok */
677         iput(inode);
678         iput(dir);
679         brelse(bh);
680         return -EPERM;
681     }
// 若要被删除目录 i 节点的属性表明这不是一个目录，则本删除操作的前提完全不存在。于是
// 放回包含要删除目录名的目录 i 节点和该要删除目录的 i 节点，释放高速缓冲块，返回出错码。
682     if (!S_ISDIR(inode->i_mode)) {
683         iput(inode);
684         iput(dir);
685         brelse(bh);
686         return -ENOTDIR;
687     }
// 若该需被删除的目录不空，则也不能删除。于是放回包含要删除目录名的目录 i 节点和该要
// 删除目录的 i 节点，释放高速缓冲块，返回出错码。
688     if (!empty_dir(inode)) {
689         iput(inode);
690         iput(dir);
691         brelse(bh);
692         return -ENOTEMPTY;
693     }
// 对于一个空目录，其目录项链接数应该为 2（链接到上层目录和本目录）。若该需被删除目
// 录的 i 节点的连接数不等于 2，则显示警告信息。但删除操作仍然继续执行。于是置该需被
// 删除目录的目录项的 i 节点号字段为 0，表示该目录项不再使用，并置含有该目录项的高速
// 缓冲块已修改标志，并释放该缓冲块。然后再置被删除目录 i 节点的链接数为 0（表示空闲），
// 并置 i 节点已修改标志。
694     if (inode->i_nlinks != 2)
695         printk("empty directory has nlink!=2 (%d)", inode->i_nlinks);
696     de->inode = 0;
697     bh->b_dirt = 1;

```



```

732     if (!(inode = iget(dir->i_dev, de->inode))) {
733         iput(dir);
734         brelse(bh);
735         return -ENOENT;
736     }
// 此时我们已有包含要被删除目录项的目录 i 节点 dir、要被删除目录项的 i 节点 inode 和要
// 被删除目录项指针 de。下面我们通过对这 3 个对象中信息的检查来验证删除操作的可行性。
// 若该目录设置了受限删除标志并且进程的有效用户 id (euid) 不是 root，并且进程的 euid
// 不等于该 i 节点的用户 id，并且进程的 euid 也不等于目录 i 节点的用户 id，则表示当前进
// 程没有权限删除该目录，于是放回包含要删除目录名的目录 i 节点和该要删除目录的 i 节点，
// 然后释放高速缓冲块，返回出错码。
737     if ((dir->i_mode & S_ISVTX) && !suser() &&
738         current->euid != inode->i_uid &&
739         current->euid != dir->i_uid) {
740         iput(dir);
741         iput(inode);
742         brelse(bh);
743         return -EPERM;
744     }
// 如果该指定文件名是一个目录，则也不能删除。放回该目录 i 节点和该文件名目录项的 i 节
// 点，释放包含该目录项的缓冲块，返回出错号。
745     if (S_ISDIR(inode->i_mode)) {
746         iput(inode);
747         iput(dir);
748         brelse(bh);
749         return -EPERM;
750     }
// 如果该 i 节点的链接计数值已经为 0，则显示警告信息，并修正其为 1。
751     if (!inode->i_nlinks) {
752         printk("Deleting nonexistent file (%04x:%d), %d\n",
753                 inode->i_dev, inode->i_num, inode->i_nlinks);
754         inode->i_nlinks=1;
755     }
// 现在我们可以删除文件名对应的目录项了。于是将该文件名目录项中的 i 节点号字段置为 0，
// 表示释放该目录项，并设置包含该目录项的缓冲块已修改标志，释放该高速缓冲块。
756     de->inode = 0;
757     bh->b_dirt = 1;
758     brelse(bh);
// 然后把文件名对应 i 节点的链接数减 1，置已修改标志，更新改变时间为当前时间。最后放
// 回该 i 节点和目录的 i 节点，返回 0（成功）。如果是文件的最后一个链接，即 i 节点链接
// 数减 1 后等于 0，并且此时没有进程正打开该文件，那么在调用 iput() 放回 i 节点时，该文
// 件也将被删除，并释放所占用的设备空间。参见 fs/inode.c，第 183 行。
759     inode->i_nlinks--;
760     inode->i_dirt = 1;
761     inode->i_ctime = CURRENT_TIME;
762     iput(inode);
763     iput(dir);
764     return 0;
765 }
766
// // 建立符号链接。
// 为一个已存在文件创建一个符号链接（也称为软连接 - hard link）。

```

```

// 参数: oldname - 原路径名; newname - 新的路径名。
// 返回: 若成功则返回 0, 否则返回出错号。
767 int sys_symlink(const char * oldname, const char * newname)
768 {
769     struct dir_entry * de;
770     struct m_inode * dir, * inode;
771     struct buffer_head * bh, * name_block;
772     const char * basename;
773     int namelen, i;
774     char c;
775
// 首先查找新路径名的最顶层目录的 i 节点 dir, 并返回最后的文件名及其长度。如果目录的
// i 节点没有找到, 则返回出错号。如果新路径名中不包括文件名, 则放回新路径名目录的 i
// 节点, 返回出错号。另外, 如果用户没有在新目录中写的权限, 则也不能建立连接, 于是放
// 回新路径名目录的 i 节点, 返回出错号。
776     dir = dir_namei(newname, &namelen, &basename, NULL);
777     if (!dir)
778         return -EACCES;
779     if (!namelen) {
780         iput(dir);
781         return -EPERM;
782     }
783     if (!permission(dir, MAY_WRITE)) {
784         iput(dir);
785         return -EACCES;
786     }
// 现在我们在目录指定设备上申请一个新的 i 节点, 并设置该 i 节点模式为符号链接类型以及
// 进程规定的模式屏蔽码。并且设置该 i 节点已修改标志。
787     if (!(inode = new_inode(dir->i_dev))) {
788         iput(dir);
789         return -ENOSPC;
790     }
791     inode->i_mode = S_IFLNK | (0777 & ~current->umask);
792     inode->i_dirt = 1;
// 为了保存符号链接路径名字符串信息, 我们需要为该 i 节点申请一个磁盘块, 并让 i 节点的
// 第 1 个直接块号 i_zone[0] 等于得到的逻辑块号。然后置 i 节点已修改标志。如果申请失败
// 则放回对应目录的 i 节点; 复位新申请的 i 节点链接计数; 放回该新的 i 节点, 返回没有空
// 间出错码退出。
793     if (!(inode->i_zone[0] = new_block(inode->i_dev))) {
794         iput(dir);
795         inode->i_nlinks--;
796         iput(inode);
797         return -ENOSPC;
798     }
799     inode->i_dirt = 1;
// 然后从设备上读取新申请的磁盘块(目的是把对应块放到高速缓冲区中)。若出错, 则放回
// 对应目录的 i 节点; 复位新申请的 i 节点链接计数; 放回该新的 i 节点, 返回没有空间出错
// 码退出。
800     if (!(name_block = bread(inode->i_dev, inode->i_zone[0]))) {
801         iput(dir);
802         inode->i_nlinks--;
803         iput(inode);
804         return -ERROR;

```

```

805         }
806         i = 0;
807         while (i < 1023 && (c=get_fs_byte(oldname++)))
808             name_block->b_data[i++] = c;
809         name_block->b_data[i] = 0;
810         name_block->b_dirt = 1;
811         brelse(name_block);
812         inode->i_size = i;
813         inode->i_dirt = 1;
814         // 然后我们搜索一下路径名指定的符号链接文件名是否已经存在。若已经存在则不能创建同名
815         // 目录项 i 节点。如果对应符号链接文件名已经存在，则释放包含该目录项的缓冲区块，复位
816         // 新申请的 i 节点连接计数，并放回目录的 i 节点，返回文件已经存在的出错码退出。
817         bh = find_entry(&dir, basename, namelen, &de);
818         if (bh) {
819             inode->i_nlinks--;
820             iput(inode);
821             brelse(bh);
822             iput(dir);
823             return -EXIST;
824         }
825         // 现在我们在指定目录中新添加一个目录项，用于存放新建符号链接文件名的 i 节点号和目录
826         // 名。如果失败（包含该目录项的高速缓冲区指针为 NULL），则放回目录的 i 节点；所申请的
827         // i 节点引用连接计数复位，并放回该 i 节点。返回出错码退出。
828         bh = add_entry(dir, basename, namelen, &de);
829         if (!bh) {
830             inode->i_nlinks--;
831             iput(inode);
832             iput(dir);
833             return -ENOSPC;
834         }
835     }
836
837     ///// 为文件建立一个文件名目录项。
838     // 为一个已存在的文件创建一个新链接（也称为硬连接 - hard link）。
839     // 参数: oldname - 原路径名; newname - 新的路径名。
840     // 返回: 若成功则返回 0, 否则返回出错号。
841     int sys_link(const char * oldname, const char * newname)
842     {
843         struct dir_entry * de;
844         struct m_inode * oldinode, * dir;

```

```

841     struct buffer_head * bh;
842     const char * basename;
843     int namelen;
844
// 首先对原文件名进行有效性验证，它应该存在并且不是一个目录名。所以我们先取原文件路
// 径名对应的 i 节点 oldinode。如果为 0，则表示出错，返回出错号。如果原路径名对应的是
// 一个目录名，则放回该 i 节点，也返回出错号。
845     oldinode=namei(oldname);
846     if (!oldinode)
847         return -ENOENT;
848     if (S_ISDIR(oldinode->i_mode)) {
849         iput(oldinode);
850         return -EPERM;
851     }
// 然后查找新路径名的最顶层目录的 i 节点 dir，并返回最后的文件名及其长度。如果目录的
// i 节点没有找到，则放回原路径名的 i 节点，返回出错号。如果新路径名中不包括文件名，
// 则放回原路径名 i 节点和新路径名目录的 i 节点，返回出错号。
852     dir = dir_namei(newname, &namelen, &basename, NULL);
853     if (!dir) {
854         iput(oldinode);
855         return -EACCES;
856     }
857     if (!namelen) {
858         iput(oldinode);
859         iput(dir);
860         return -EPERM;
861     }
// 我们不能跨设备建立硬链接。因此如果新路径名顶层目录的设备号与原路径名的设备号
// 不一样，则放回新路径名目录的 i 节点和原路径名的 i 节点，返回出错号。另外，如果用户
// 没有在新目录中写的权限，则也不能建立连接，于是放回新路径名目录的 i 节点和原路径名
// 的 i 节点，返回出错号。
862     if (dir->i_dev != oldinode->i_dev) {
863         iput(dir);
864         iput(oldinode);
865         return -EXDEV;
866     }
867     if (!permission(dir, MAY_WRITE)) {
868         iput(dir);
869         iput(oldinode);
870         return -EACCES;
871     }
// 现在查询该新路径名是否已经存在，如果存在则也不能建立链接。于是释放包含该已存在目
// 录项的高速缓冲块，放回新路径名目录的 i 节点和原路径名的 i 节点，返回出错号。
872     bh = find_entry(&dir, basename, namelen, &de);
873     if (bh) {
874         brelse(bh);
875         iput(dir);
876         iput(oldinode);
877         return -EXIST;
878     }
// 现在所有条件都满足了，于是我们在新目录中添加一个目录项。若失败则放回该目录的 i 节
// 点和原路径名的 i 节点，返回出错号。否则初始设置该目录项的 i 节点号等于原路径名的 i
// 节点号，并置包含该新添目录项的缓冲块已修改标志，释放该缓冲块，放回目录的 i 节点。

```

```
879     bh = add_entry(dir, basename, namelen, &de);
880     if (!bh) {
881         iput(dir);
882         iput(oldinode);
883         return -ENOSPC;
884     }
885     de->inode = oldinode->i_num;
886     bh->b_dirt = 1;
887     brelse(bh);
888     iput(dir);
// 再将原节点的链接计数加 1，修改其改变时间为当前时间，并设置 i 节点已修改标志。最后
// 放回原路径名的 i 节点，并返回 0（成功）。
889     oldinode->i_nlinks++;
890     oldinode->i_ctime = CURRENT_TIME;
891     oldinode->i_dirt = 1;
892     iput(oldinode);
893     return 0;
894 }
895
```
