

程序 12-10 linux/fs/pipe.c

```

1  /*
2  *  linux/fs/pipe.c
3  *
4  *  (C) 1991 Linus Torvalds
5  */
6
7  #include <signal.h>      // 信号头文件。定义信号符号常量，信号结构及操作函数原型。
8  #include <errno.h>      // 错误号头文件。包含系统中各种出错号。
9  #include <termios.h>    // 终端输入输出函数头文件。主要定义控制异步通信口的终端接口。
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 数据等。
12 #include <linux/mm.h>   /* for get_free_page */ /* 使用其中的 get_free_page */
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 #include <linux/kernel.h> // 内核头文件。含有一些内核常用函数的原形定义。
15
16  // 管道读操作函数。
17  // 参数 inode 是管道对应的 i 节点，buf 是用户数据缓冲区指针，count 是读取的字节数。
18  int read_pipe(struct m_inode * inode, char * buf, int count)
19  {
20      int chars, size, read = 0;
21
22      // 如果需要读取的字节计数 count 大于 0，我们就循环执行以下操作。在循环读操作过程中，
23      // 若当前管道中没有数据 (size=0)，则唤醒等待该节点的进程，这通常是写管道进程。如果
24      // 已没有写管道者，即 i 节点引用计数值小于 2，则返回已读字节数退出。如果当前收到有非
25      // 阻塞信号，则立刻返回已读取字节数退出；若还没有收到任何数据，则返回重新启动系统
26      // 调用号退出。否则就让进程在该管道上睡眠，用以等待信息的到来。宏 PIPE_SIZE 定义在
27      // include/linux/fs.h 中。关于“重新启动系统调用”，请参见 kernel/signal.c 程序。
28      while (count>0) {
29          while (!(size=PIPE_SIZE(*inode))) { // 取管道中数据长度值。
30              wake_up(& PIPE_WRITE_WAIT(*inode));
31              if (inode->i_count != 2) /* are there any writers? */
32                  return read;
33              if (current->signal & ~current->blocked)
34                  return read?read:-ERESTARTSYS;
35              interruptible_sleep_on(& PIPE_READ_WAIT(*inode));
36          }
37          // 此时说明管道（缓冲区）中有数据。于是我们取管道尾指针到缓冲区末端的字节数 chars。
38          // 如果其大于还需要读取的字节数 count，则令其等于 count。如果 chars 大于当前管道中含
39          // 有数据的长度 size，则令其等于 size。然后把需读字节数 count 减去此次可读的字节数
40          // chars，并累加已读字节数 read。
41          chars = PAGE_SIZE-PIPE_TAIL(*inode);
42          if (chars > count)
43              chars = count;
44          if (chars > size)
45              chars = size;
46          count -= chars;
47          read += chars;
48          // 再令 size 指向管道尾指针处，并调整当前管道尾指针（前移 chars 字节）。若尾指针超过
49          // 管道末端则绕回。然后将管道中的数据复制到用户缓冲区中。对于管道 i 节点，其 i_size
50          // 字段中是管道缓冲块指针。
51          size = PIPE_TAIL(*inode);
52          PIPE_TAIL(*inode) += chars;

```

```

38         PIPE\_TAIL(*inode) &= (PAGE\_SIZE-1);
39         while (chars-->0)
40             put\_fs\_byte((char *)inode->i_size)[size++],buf++);
41     }
    // 当此次读管道操作结束，则唤醒等待该管道的进程，并返回读取的字节数。
42     wake\_up(& PIPE\_WRITE\_WAIT(*inode));
43     return read;
44 }
45
    // 管道写操作函数。
    // 参数 inode 是管道对应的 i 节点，buf 是数据缓冲区指针，count 是将写入管道的字节数。
46 int write\_pipe(struct m\_inode * inode, char * buf, int count)
47 {
48     int chars, size, written = 0;
49
    // 如果要写入的字节数 count 还大于 0，那么我们就循环执行以下操作。在循环操作过程中，
    // 如果当前管道中已经满了（空闲空间 size = 0），则唤醒等待该管道的进程，通常唤醒
    // 的是读管道进程。如果已没有读管道者，即 i 节点引用计数值小于 2，则向当前进程发送
    // SIGPIPE 信号，并返回已写入的字节数退出；若写入 0 字节，则返回 -1。否则让当前进程
    // 在该管道上睡眠，以等待读管道进程来读取数据，从而让管道腾出空间。宏 PIPE\_SIZE()、
    // PIPE\_HEAD() 等定义在文件 include/linux/fs.h 中。
50     while (count>0) {
51         while (!(size=(PAGE\_SIZE-1)-PIPE\_SIZE(*inode))) {
52             wake\_up(& PIPE\_READ\_WAIT(*inode));
53             if (inode->i_count != 2) { /* no readers */
54                 current->signal |= (1<<(SIGPIPE-1));
55                 return written?written:-1;
56             }
57             sleep\_on(& PIPE\_WRITE\_WAIT(*inode));
58         }
    // 程序执行到这里表示管道缓冲区中有可写空间 size。于是我们取管道头指针到缓冲区末端空
    // 间字节数 chars。写管道操作是从管道头指针处开始写的。如果 chars 大于还需要写入的字节
    // 数 count，则令其等于 count。如果 chars 大于当前管道中空闲空间长度 size，则令其等于
    // size。然后把需要写入字节数 count 减去此次可写入的字节数 chars，并把写入字节数累加到
    // written 中。
59         chars = PAGE\_SIZE-PIPE\_HEAD(*inode);
60         if (chars > count)
61             chars = count;
62         if (chars > size)
63             chars = size;
64         count -= chars;
65         written += chars;
    // 再令 size 指向管道数据头指针处，并调整当前管道数据头部指针（前移 chars 字节）。若头
    // 指针超过管道末端则绕回。然后从用户缓冲区复制 chars 个字节到管道头指针开始处。对于
    // 管道 i 节点，其 i_size 字段中是管道缓冲块指针。
66         size = PIPE\_HEAD(*inode);
67         PIPE\_HEAD(*inode) += chars;
68         PIPE\_HEAD(*inode) &= (PAGE\_SIZE-1);
69         while (chars-->0)
70             ((char *)inode->i_size)[size++]=get\_fs\_byte(buf++);
71     }
    // 当此次写管道操作结束，则唤醒等待管道的进程，返回已写入的字节数，退出。
72     wake\_up(& PIPE\_READ\_WAIT(*inode));

```

```

73         return written;
74     }
75
76     // 创建管道系统调用。
77     // 在 fildes 所指的数组中创建一对文件句柄（描述符）。这对文件句柄指向一管道 i 节点。
78     // 参数：fildes - 文件句柄数组。fildes[0] 用于读管道数据，fildes[1] 向管道写入数据。
79     // 成功时返回 0，出错时返回 -1。
80     int sys_pipe(unsigned long * fildes)
81     {
82         struct m_inode * inode;
83         struct file * f[2];           // 文件结构数组。
84         int fd[2];                   // 文件句柄数组。
85         int i, j;
86
87         // 首先从系统文件表中取两个空闲项（引用计数字段为 0 的项），并分别设置引用计数为 1。
88         // 若只有 1 个空闲项，则释放该项（引用计数复位）。若没有找到两个空闲项，则返回 -1。
89         j=0;
90         for(i=0; j<2 && i<NR_FILE; i++)
91             if (!file_table[i].f_count)
92                 (f[j++] = i + file_table) ->f_count++;
93         if (j==1)
94             f[0]->f_count=0;
95         if (j<2)
96             return -1;
97
98         // 针对上面取得的两个文件表结构项，分别分配一文件句柄号，并使进程文件结构指针数组的
99         // 两项分别指向这两个文件结构。而文件句柄即是该数组的索引号。类似地，如果只有一个空
100        // 闲文件句柄，则释放该句柄（置空相应数组项）。如果没有找到两个空闲句柄，则释放上面
101        // 获取的两个文件结构项（复位引用计数值），并返回 -1。
102         j=0;
103         for(i=0; j<2 && i<NR_OPEN; i++)
104             if (!current->filp[i]) {
105                 current->filp[fd[j]=i] = f[j];
106                 j++;
107             }
108         if (j==1)
109             current->filp[fd[0]]=NULL;
110         if (j<2) {
111             f[0]->f_count=f[1]->f_count=0;
112             return -1;
113         }
114
115         // 然后利用函数 get_pipe_inode() 申请一个管道使用的 i 节点，并为管道分配一页内存作为缓
116         // 冲区。如果不成功，则相应释放两个文件句柄和文件结构项，并返回 -1。
117         if (!(inode=get_pipe_inode())) { // fs/inode.c, 第 231 行开始处。
118             current->filp[fd[0]] =
119                 current->filp[fd[1]] = NULL;
120             f[0]->f_count = f[1]->f_count = 0;
121             return -1;
122         }
123
124         // 如果管道 i 节点申请成功，则对两个文件结构进行初始化操作，让它们都指向同一个管道 i 节
125         // 点，并把读写指针都置零。第 1 个文件结构的文件模式置为读，第 2 个文件结构的文件模式置
126         // 为写。最后将文件句柄数组复制到对应的用户空间数组中，成功返回 0，退出。
127         f[0]->f_inode = f[1]->f_inode = inode;
128         f[0]->f_pos = f[1]->f_pos = 0;

```

```
111         f[0]->f_mode = 1;                /* read */
112         f[1]->f_mode = 2;                /* write */
113         put fs long(fd[0],0+fildes);
114         put fs long(fd[1],1+fildes);
115         return 0;
116     }
117
118     ///// 管道 io 控制函数。
119     // 参数: pino - 管道 i 节点指针; cmd - 控制命令; arg - 参数。
120     // 函数返回 0 表示执行成功, 否则返回出错码。
121     int pipe\_ioctl(struct m\_inode *pino, int cmd, int arg)
122     {
123         // 如果命令是取管道中当前可读数据长度, 则把管道数据长度值添入用户参数指定的位置处,
124         // 并返回 0。否则返回无效命令错误码。
125         switch (cmd) {
126             case FIONREAD:
127                 verify area((void *) arg,4);
128                 put fs long(PIPE\_SIZE(*pino), (unsigned long *) arg);
129                 return 0;
130             default:
131                 return -EINVAL;
132         }
133     }
134 }
```
