

# Linux 内核 0.11 详细注释

A Heavy Commented Linux Kernel Source Code

Linux Version 0.11

Zhao Jiong

赵炯

gohigh@sh163.net

[www.plinux.org](http://www.plinux.org)

[www.oldlinux.org](http://www.oldlinux.org)

2002-10-25

---

## 内容简介

本书对 Linux 早期操作系统内核 (v0.95) 全部代码文件进行了详细全面的注释和说明, 旨在使读者能够在尽量短的时间内对 Linux 的工作机理获得全面而深刻的理解, 为进一步学习和研究 Linux 系统打下坚实的基础。虽然所选择的版本较低, 但该内核已能够正常编译运行, 其中已经包括了 LINUX 工作原理的精髓, 通过阅读其源代码能快速地完全理解内核的运作机制。书中首先以 Linux 源代码版本的变迁历史为主线, 详细介绍了 Linux 系统的发展历史, 着重说明了各个内核版本之间的重要区别和改进方面, 给出了选择 0.11 (0.95) 版作为研究的对象的原因。另外介绍了内核源代码的组织结构及相互关系, 同时还说明了编译和运行该版本内核的方法。然后本书依据内核源代码的组织结构对所有内核程序和文件进行了注释和详细说明。每章的安排基本上分为具体研究对象的概述、每个文件的功能介绍、代码内注释、代码中难点及相关资料介绍、与当前版本的主要区别等部分。最后一章内容总结性地介绍了继续研究 Linux 系统的方法和着手点。

## 版权说明

作者保留本电子书籍的修改和正式出版的所有权利。读者可以自由传播本书全部和部分章节的内容, 但需要注明出处。由于目前本书尚为草稿阶段, 因此存在许多错误和不足之处, 希望读者能踊跃给予批评指正或建议。可以通过电子邮件给我发信息: gohigh@sh163.net, 或直接来信至: 上海同济大学 机械电子工程研究所 (上海四平路 1239 号, 邮编: 200092)。

© 2002, 2003 by Zhao Jiong

© 2002, 2003 赵炯 版权所有。

---

“RTFSC - Read The F\*\*king Source Code ☺!”

- Linus Benedict Torvalds

---

# 目录

|                              |    |
|------------------------------|----|
| 文件列表索引 .....                 | VI |
| 序言 .....                     | 1  |
| 本书的主要目标.....                 | 1  |
| 现有书籍不足之处.....                | 1  |
| 阅读早期内核其它的好处? .....           | 2  |
| 阅读完整源代码的重要性和必要性 .....        | 2  |
| 如何选择要阅读的内核代码版本.....          | 3  |
| 阅读本书需具备的基础知识.....            | 3  |
| 使用早期版本是否过时? .....            | 4  |
| Ext2 文件系统与 Minix 文件系统? ..... | 4  |
| 第 1 章 概述 .....               | 1  |
| 1.1 Linux 的诞生和发展 .....       | 1  |
| 1.2 内容综述.....                | 6  |
| 1.3 Makefile 文件 .....        | 9  |
| 1.4 linux/Makefile 文件 .....  | 9  |
| 1.5 其它信息.....                | 13 |
| 第 2 章 引导启动程序 (boot) .....    | 17 |
| 2.1 概述.....                  | 18 |
| 2.2 bootsect.s 文件 .....      | 18 |
| 2.3 setup.s 程序 .....         | 25 |
| 2.4 head.s 文件 .....          | 35 |
| 2.5 小结.....                  | 42 |
| 第 3 章 初始化程序 (init) .....     | 43 |
| 3.1 概述.....                  | 43 |
| 3.2 main.c 程序 .....          | 43 |
| 3.3 其它信息.....                | 48 |
| 3.4 小结.....                  | 48 |
| 第 4 章 内核代码 (kernel) .....    | 49 |
| 4.1 概述.....                  | 49 |
| 4.2 Makefile 文件 .....        | 50 |
| 4.3 asm.s 文件 .....           | 52 |
| 4.4 system_call.s 程序 .....   | 57 |
| 4.5 traps.c 程序 .....         | 65 |
| 4.6 mktime.c 文件 .....        | 71 |
| 4.7 sched.c 文件 .....         | 73 |
| 4.8 signal.c 文件 .....        | 84 |
| 4.9 exit.c 文件 .....          | 88 |
| 4.10 fork.c 程序 .....         | 93 |
| 4.11 sys.c 程序 .....          | 98 |

---

|                                       |            |
|---------------------------------------|------------|
| 4.12 vsprintf.c 程序 .....              | 104        |
| 4.13 printk.c 程序 .....                | 112        |
| 4.14 panic.c 程序 .....                 | 113        |
| <b>第5章 块设备驱动程序(block driver).....</b> | <b>115</b> |
| 5.1 概述.....                           | 115        |
| 5.2 Makefile 文件 .....                 | 115        |
| 5.3 blk.h 文件 .....                    | 117        |
| 5.4 hd.c 文件 .....                     | 121        |
| 5.5 ll_rw_blk.c 文件 .....              | 133        |
| 5.6 ramdisk.c 文件 .....                | 138        |
| 5.7 floppy.c 文件 .....                 | 141        |
| <b>第6章 字符设备驱动程序(char driver).....</b> | <b>157</b> |
| 6.1 概述.....                           | 157        |
| 6.2 Makefile 文件.....                  | 157        |
| 6.3 serial.c 文件.....                  | 159        |
| 6.4 rs_io.s 文件.....                   | 163        |
| 6.5 tty_io.c 文件.....                  | 166        |
| 6.6 tty_ioctl.c 文件.....               | 176        |
| 6.7 keyboard.S 文件.....                | 182        |
| 6.8 console.c 文件.....                 | 200        |
| <b>第7章 数学协处理器(math) .....</b>         | <b>225</b> |
| 7.1 概述.....                           | 225        |
| 7.2 Makefile 文件 .....                 | 225        |
| 7.3 math-emulation.c 程序 .....         | 227        |
| <b>第8章 文件系统(fs) .....</b>             | <b>229</b> |
| 8.1 概述.....                           | 229        |
| 8.2 Makefile 文件 .....                 | 230        |
| 8.3 bitmap.c 程序 .....                 | 232        |
| 8.4 inode.c 程序 .....                  | 237        |
| 8.5 buffer.c 文件 .....                 | 246        |
| 8.6 block_dev.c 文件 .....              | 256        |
| 8.7 file_dev.c 文件 .....               | 258        |
| 8.8 file_table.c 文件 .....             | 261        |
| 8.9 fcntl.c 文件 .....                  | 262        |
| 8.10 stat.c 文件 .....                  | 264        |
| 8.11 pipe.c 文件 .....                  | 265        |
| 8.12 read_write.c 文件 .....            | 269        |
| 8.13 truncate.c 文件 .....              | 272        |
| 8.14 super.c 文件 .....                 | 274        |
| 8.15 open.c 文件 .....                  | 282        |
| 8.16 exec.c 程序 .....                  | 288        |
| 8.17 namei.c 文件 .....                 | 303        |
| 8.18 char_dev.c 文件 .....              | 325        |
| 8.19 ioctl.c 文件 .....                 | 327        |

---

|                             |            |
|-----------------------------|------------|
| <b>第 9 章 内存管理(mm)</b>       | <b>331</b> |
| 9.1 概述                      | 331        |
| 9.2 Makefile 文件             | 331        |
| 9.3 memory.c 程序             | 332        |
| 9.4 page.s 程序               | 344        |
| <b>第 10 章 包含文件(include)</b> | <b>347</b> |
| 10.1 概述                     | 347        |
| 10.2 include/目录下的文件         | 347        |
| 10.3 a.out.h 文件             | 347        |
| 10.4 const.h 文件             | 358        |
| 10.5 ctype.h 文件             | 358        |
| 10.6 errno.h 文件             | 359        |
| 10.7 fcntl.h 文件             | 361        |
| 10.8 signal.h 文件            | 363        |
| 10.9 stdarg.h 文件            | 365        |
| 10.10 stddef.h 文件           | 366        |
| 10.11 string.h 文件           | 366        |
| 10.12 termios.h 文件          | 376        |
| 10.13 time.h 文件             | 383        |
| 10.14 unistd.h 文件           | 384        |
| 10.15 utime.h 文件            | 390        |
| 10.16 include/asm/目录下的文件    | 391        |
| 10.17 io.h 文件               | 391        |
| 10.18 memory.h 文件           | 392        |
| 10.19 segment.h 文件          | 393        |
| 10.20 system.h 文件           | 395        |
| 10.21 include/linux/目录下的文件  | 398        |
| 10.22 config.h 文件           | 398        |
| 10.23 fdreg.h 头文件           | 400        |
| 10.24 fs.h 文件               | 403        |
| 10.25 hdreg.h 文件            | 409        |
| 10.26 head.h 文件             | 411        |
| 10.27 kernel.h 文件           | 412        |
| 10.28 mm.h 文件               | 413        |
| 10.29 sched.h 文件            | 413        |
| 10.30 sys.h 文件              | 423        |
| 10.31 tty.h 文件              | 425        |
| 10.32 include/sys/目录中的文件    | 428        |
| 10.33 stat.h 文件             | 428        |
| 10.34 times.h 文件            | 430        |
| 10.35 types.h 文件            | 430        |
| 10.36 utsname.h 文件          | 431        |
| 10.37 wait.h 文件             | 432        |
| <b>第 11 章 库文件(lib)</b>      | <b>435</b> |

---

|                                 |            |
|---------------------------------|------------|
| 11.1 概述.....                    | 435        |
| 11.2 Makefile 文件 .....          | 436        |
| 11.3 _exit.c 程序 .....           | 438        |
| 11.4 close.c 程序 .....           | 438        |
| 11.5 ctype.c 程序 .....           | 439        |
| 11.6 dup.c 程序 .....             | 440        |
| 11.7 errno.c 程序 .....           | 441        |
| 11.8 execve.c 程序 .....          | 441        |
| 11.9 malloc.c 程序 .....          | 442        |
| 11.10 open.c 程序 .....           | 450        |
| 11.11 setsid.c 程序 .....         | 451        |
| 11.12 string.c 程序 .....         | 452        |
| 11.13 wait.c 程序 .....           | 452        |
| 11.14 write.c 程序 .....          | 453        |
| <b>第 12 章 建造工具(tools) .....</b> | <b>455</b> |
| 12.1 概述.....                    | 455        |
| 12.2 build.c 程序 .....           | 455        |
| <b>参考文献.....</b>                | <b>463</b> |
| <b>附录.....</b>                  | <b>464</b> |
| 附录 1 80x86 保护运行模式.....          | 464        |
| <b>索引.....</b>                  | <b>473</b> |

# 文件列表索引

|  |     |
|--|-----|
| 列表 1.1 Linux/目录 .....                        | 7   |
| 列表 1.2 linux/Makefile 文件 .....               | 9   |
| 列表 2.1 linux/boot/目录 .....                   | 18  |
| 列表 2.2 linux/boot/bootsect.s 程序 .....        | 18  |
| 列表 2.3 linux/boot/setup.s 程序 .....           | 26  |
| 列表 2.4 linux/boot/head.s 程序 .....            | 35  |
| 列表 3.1 linux/init/main.c 程序 .....            | 43  |
| 列表 4.1 linux/kernel/目录 .....                 | 49  |
| 列表 4.2 linux/kernel/Makefile 文件 .....        | 50  |
| 列表 4.3 linux/kernel/asm.s 程序 .....           | 53  |
| 列表 4.4 linux/kernel/system_call.s 程序 .....   | 58  |
| 列表 4.5 linux/kernel/traps.c 程序 .....         | 65  |
| 列表 4.6 linux/kernel/mktime.c 程序 .....        | 71  |
| 列表 4.7 linux/kernel/sched.c 程序 .....         | 73  |
| 列表 4.8 linux/kernel/signal.c 程序 .....        | 84  |
| 列表 linux/kernel/exit.c 程序 .....              | 88  |
| 列表 linux/kernel/fork.c 程序 .....              | 93  |
| 列表 linux/kernel/sys.c 程序 .....               | 98  |
| 列表 linux/kernel/vsprintf.c 程序 .....          | 104 |
| 列表 linux/kernel/printk.c 程序 .....            | 112 |
| 列表 linux/kernel/panic.c 程序 .....             | 113 |
| 列表 linux/kernel/blk_drv/目录 .....             | 115 |
| 列表 linux/kernel/blk_drv/Makefile 文件 .....    | 115 |
| 列表 linux/kernel/blk_drv/blk.h 文件 .....       | 117 |
| 列表 linux/kernel/blk_drv/hd.c 程序 .....        | 121 |
| 列表 linux/kernel/blk_drv/ll_rw_blk.c 程序 ..... | 133 |
| 列表 linux/kernel/blk_drv/ramdisk.c 程序 .....   | 138 |
| 列表 linux/kernel/blk_drv/floppy.c 程序 .....    | 141 |
| 列表 linux/kernel/chr_drv/目录 .....             | 157 |
| 列表 linux/kernel/chr_drv/Makefile 文件 .....    | 157 |
| 列表 linux/kernel/chr_drv/serial.c 程序 .....    | 159 |
| 列表 linux/kernel/chr_drv/rs_io.s 程序 .....     | 163 |
| 列表 linux/kernel/chr_drv/tty_io.c 程序 .....    | 166 |
| 列表 linux/kernel/chr_drv/tty_ioctl.c 程序 ..... | 176 |
| 列表 linux/kernel/chr_drv/keyboard.S 文件 .....  | 182 |
| 列表 linux/kernel/chr_drv/console.c 程序 .....   | 200 |
| 列表 linux/kernel/math/目录 .....                | 225 |
| 列表 linux/kernel/math/Makefile 文件 .....       | 225 |
| 列表 linux/kernel/math/math_emulate.c 程序 ..... | 227 |
| 列表 linux/fs/Makefile 文件 .....                | 230 |



---

|  |     |
|--|-----|
| 列表 linux/fs/bitmap.c 程序 .....            | 233 |
| 列表 linux/fs/inode.c 程序 .....             | 237 |
| 列表 linux/fs/buffer.c 程序 .....            | 246 |
| 列表 linux/fs/block_dev.c 程序 .....         | 256 |
| 列表 linux/fs/file_dev.c 程序 .....          | 258 |
| 列表 linux/fs/file_table.c 程序 .....        | 261 |
| 列表 linux/fs/fcntl.c 程序 .....             | 262 |
| 列表 linux/fs/stat.c 程序 .....              | 264 |
| 列表 linux/fs/pipe.c 程序 .....              | 266 |
| 列表 linux/fs/read_write.c 程序 .....        | 269 |
| 列表 linux/fs/truncate.c 程序 .....          | 272 |
| 列表 linux/fs/super.c 程序 .....             | 274 |
| 列表 linux/fs/open.c 程序 .....              | 282 |
| 列表 linux/fs/exec.c 程序 .....              | 289 |
| 列表 linux/fs/namei.c 程序 .....             | 303 |
| 列表 linux/fs/char_dev.c 程序 .....          | 325 |
| 列表 linux/fs/ioctl.c 程序 .....             | 328 |
| 列表 linux/mm/Makefile 文件 .....            | 331 |
| 列表 linux/mm/memory.c 程序 .....            | 333 |
| 列表 linux/mm/page.s 程序 .....              | 344 |
| 列表 linux/include/目录下的文件 .....            | 347 |
| 列表 linux/include/a.out.h 文件 .....        | 348 |
| 列表 linux/include/const.h 文件 .....        | 358 |
| 列表 linux/include/ctype.h 文件 .....        | 358 |
| 列表 linux/include/errno.h 文件 .....        | 359 |
| 列表 linux/include/fcntl.h 文件 .....        | 361 |
| 列表 linux/include/signal.h 文件 .....       | 363 |
| 列表 linux/include/stdarg.h 文件 .....       | 365 |
| 列表 linux/include/stddef.h 文件 .....       | 366 |
| 列表 linux/include/string.h 文件 .....       | 366 |
| 列表 linux/include/termios.h 文件 .....      | 376 |
| 列表 linux/include/time.h 文件 .....         | 383 |
| 列表 linux/include/unistd.h 文件 .....       | 384 |
| 列表 linux/include/utime.h 文件 .....        | 390 |
| 列表 linux/include/asm/目录下的文件 .....        | 391 |
| 列表 linux/include/asm/io.h 文件 .....       | 391 |
| 列表 linux/include/asm/memory.h 文件 .....   | 392 |
| 列表 linux/include/asm/segment.h 文件 .....  | 393 |
| 列表 linux/include/asm/system.h 文件 .....   | 395 |
| 列表 linux/include/linux/目录下的文件 .....      | 398 |
| 列表 linux/include/linux/config.h 文件 ..... | 398 |
| 列表 linux/include/linux/fdreg.h 文件 .....  | 401 |
| 列表 linux/include/linux/fs.h 文件 .....     | 404 |
| 列表 linux/include/linux/hdreg.h 文件 .....  | 409 |

---

|   |     |
|---|-----|
| 列表 linux/include/linux/head.h 文件 .....  | 411 |
| 列表 linux/include/linux/kernel.h .....   | 412 |
| 列表 linux/include/linux/mm.h 文件 .....    | 413 |
| 列表 linux/include/linux/sched.h 文件 ..... | 413 |
| 列表 linux/include/linux/sys.h 文件 .....   | 423 |
| 列表 linux/include/linux/tty.h 文件 .....   | 425 |
| 列表 linux/include/sys/目录下的文件 .....       | 428 |
| 列表 linux/include/sys/stat.h 文件 .....    | 428 |
| 列表 linux/include/sys/times.h 文件 .....   | 430 |
| 列表 linux/include/sys/types.h 文件 .....   | 430 |
| 列表 linux/include/sys/utsname.h 文件 ..... | 431 |
| 列表 linux/include/sys/wait.h 文件 .....    | 432 |
| 列表 linux/lib/Makefile 文件 .....          | 436 |
| 列表 linux/lib/_exit.c 程序 .....           | 438 |
| 列表 linux/lib/close.c 程序 .....           | 439 |
| 列表 linux/lib/ctype.c 程序 .....           | 439 |
| 列表 linux/lib/dup.c 程序 .....             | 440 |
| 列表 linux/lib/errno.c 程序 .....           | 441 |
| 列表 linux/lib/execve.c 程序 .....          | 441 |
| 列表 linux/lib/malloc.c 程序 .....          | 443 |
| 列表 linux/lib/open.c 程序 .....            | 450 |
| 列表 linux/lib/setsid.c 程序 .....          | 451 |
| 列表 linux/lib/string.c 程序 .....          | 452 |
| 列表 linux/lib/wait.c 程序 .....            | 453 |
| 列表 linux/lib/write.c 程序 .....           | 453 |
| 列表 linux/tools/build.c 程序 .....         | 456 |

# 序言

本书是一本有关 Linux 操作系统内核基本工作原理的入门读物。

## 本书的主要目标

本书的主要目标是使用尽量少的篇幅或在有限的篇幅内，对完整的 Linux 内核源代码进行解剖，以期对操作系统的基本功能和实际实现方式获得全方位的理解。做到对 linux 内核有一个完整而深刻的理解，对 linux 操作系统的基本工作原理真正理解和入门。

本书读者群的定位是一些知晓 Linux 系统的一般使用方法或具有一定的编程基础，但比较缺乏阅读目前最新内核源代码的基础知识，又急切希望能够进一步理解 UNIX 类操作系统内核工作原理和实际代码实现的爱好者。这部分读者的水平应该介于初级与中级水平之间。目前，这部分读者人数在 Linux 爱好者中所占的比例是很高的，而面向这部分读者以比较易懂和有效的手段讲解内核的书籍资料不多。

## 现有书籍不足之处

目前已有的描述 Linux 内核的书籍，均尽量选用最新 Linux 内核版本（例如 Redhat 7.0 使用的 2.2.16 稳定版等）进行描述，但由于目前 Linux 内核整个源代码的大小已经非常得大（例如 2.2.20 版具有 268 万行代码！），因此这些书籍仅能对 Linux 内核源代码进行选择性地或原理性地说明，许多系统实现细节被忽略。因此并不能给予读者对实际 Linux 内核有清晰而完整的理解。

目前已有的描述 Linux 内核的书籍，均尽量选用最新 Linux 内核版本（例如 Redhat 7.0 使用的 2.2.16 稳定版等）进行描述，但由于目前 Linux 内核整个源代码的大小已经非常得大（例如 2.2.20 版具有 268 万行代码！），因此这些书籍仅能对 Linux 内核源代码进行选择性地或原理性地说明，许多系统实现细节被忽略。因此并不能给予读者对实际 Linux 内核有清晰而完整的理解。

在使用这些书籍进行学习时会有一种“盲人摸象”的感觉，不能真正理解 Linux 内核系统具体实现的整体概念，尤其是对那些 Linux 系统初学者或刚学会如何使用 Linux 系统的人在使用那些书学习内核原理时，内核的整体运作结构并不能清晰地脑海中形成。这在本人多年的 Linux 内核学习过程中也深有体会，也是写作本书的动机之一。

陆丽娜等翻译的 Scott Maxwell 著的一书《Linux 内核源代码分析》基本上是针对 Linux 中级水平的读者，需要较为全面的基础知识才能完全理解。而且可能是由于篇幅所限，该书并没有对所有 Linux 内核代码进行注释，略去了很多内核实现细节，例如其中内核中使用的各个头文件 (\*.h)、生成内核代码映像文件的工具程序、各个 make 文件的作用和实现等均没有涉及。因此对于处于初中级水平之间的读者来说阅读该书有些困难。

去年初浙江大学出版的《Linux 内核源代码情景分析》一书，也基本有这些不足之处。甚至对于一些具有较高 Linux 系统应用水平的计算机本科高年级学生，由于该书篇幅问题以及仅仅选择性地讲解内核源代码，也不能真正吃透内核的实际实现方式，因而往往刚开始阅读就放弃了。这在本人教学的学生中基本都会出现这个问题。该书刚面市时，本人曾极力劝说学生购之阅读，并在二个月后调查阅读学习情况，基本都存在看不下去或不能理解等问题，大多数人都放弃了。

John Lions 著的《莱昂氏 UNIX 源代码分析》一书虽然是一本学习 UNIX 类操作系统内核源代码很好的书籍，但是由于其采用的是 UNIX V6 版，其中系统调用等部分代码是用早已过时的 PDP-11 系列机的汇编

语言编制的，因此在阅读与硬件部分相关的源代码时就会遇到较大的困难。

A. S. Tanenbaum 的书《操作系统：设计与实现》是一本有关操作系统内核实现很好的入门书籍，但该书所叙述的 minix 系统是一种基于消息传递的内核实现机制，与 Linux 内核的实现有所区别。因此在学习该书之后，并不能很顺利地即刻着手进一步学习较新的 Linux 内核源代码实现。

在使用这些书籍进行学习时会有一种“盲人摸象”的感觉，不能真正理解 Linux 内核系统具体实现的整体概念，尤其是对那些 Linux 系统初学者或刚学会如何使用 Linux 系统的人在使用那些书学习内核原理时，内核的整体运作结构并不能清晰地脑海中形成。这在本人多年的 Linux 内核学习过程中也深有体会，也是写作本书的动机之一。

为了填补这个空缺，本书的主要目标是使用尽量少的篇幅或在有限的篇幅内，对完整的 Linux 内核源代码进行全面解剖，以期对操作系统的基本功能和实际实现方式获得全方位的理解。做到对 linux 内核有一个完整而深刻的理解，对 linux 操作系统的基本工作原理真正理解和入门。

## 阅读早期内核其它的好处？

目前，已经出现不少基于 Linux 早期内核而开发的专门用于嵌入式系统的内核版本，如 DJJ 的 x86 操作系统、Uclinux 等（在 [www.linux.org](http://www.linux.org) 上有专门目录），世界上也有许多人认识到通过早期 Linux 内核源代码学习的好处，目前国内也已经有人正在组织人力注释出版类似本文的书籍。因此，通过阅读 Linux 早期内核版本的源代码，的确是学习 Linux 系统的一种行之有效的途径，并且对研究和应用 Linux 嵌入式系统也有很大的帮助。

在对早期内核源代码的注释过程中，作者已经发现，早期内核源代码几乎就是目前所使用的较新内核的一个精简版本。其中已经包括了目前新版本中几乎所有的基本功能原理的内容。正如《系统软件：系统编程导论》一书的作者 Leland L. Bach 在介绍系统程序以及操作系统设计时，引入了一种极其简化的简单指令计算机(SIC)系统来说明所有系统程序的设计和实现原理，从而既避免了实际计算机系统的复杂性，又能透彻地说明问题。这里选择 Linux 的早期内核版本作为学习对象，其指导思想与 Leland 的一致。这对 Linux 内核学习的入门者来说，无非是一个最理想的选择之一。能够在尽可能短的时间内深入理解 Linux 内核的基本工作原理。

当然，使用早期内核作为学习的对象也有不足之处，所选用的 Linux 早期内核版本不包含对虚拟文件系统 VFS 的支持、对网络系统的支持、仅支持 a.out 执行文件和对其它一些现有内核中复杂子系统的说明。但由于本书是作为 Linux 内核工作机制实现的入门教材，因此这也正是选择早期内核版本的优点之一。通过学习本书，可以为进一步学习这些高级内容打下扎实的基础。

## 阅读完整源代码的重要性和必要性

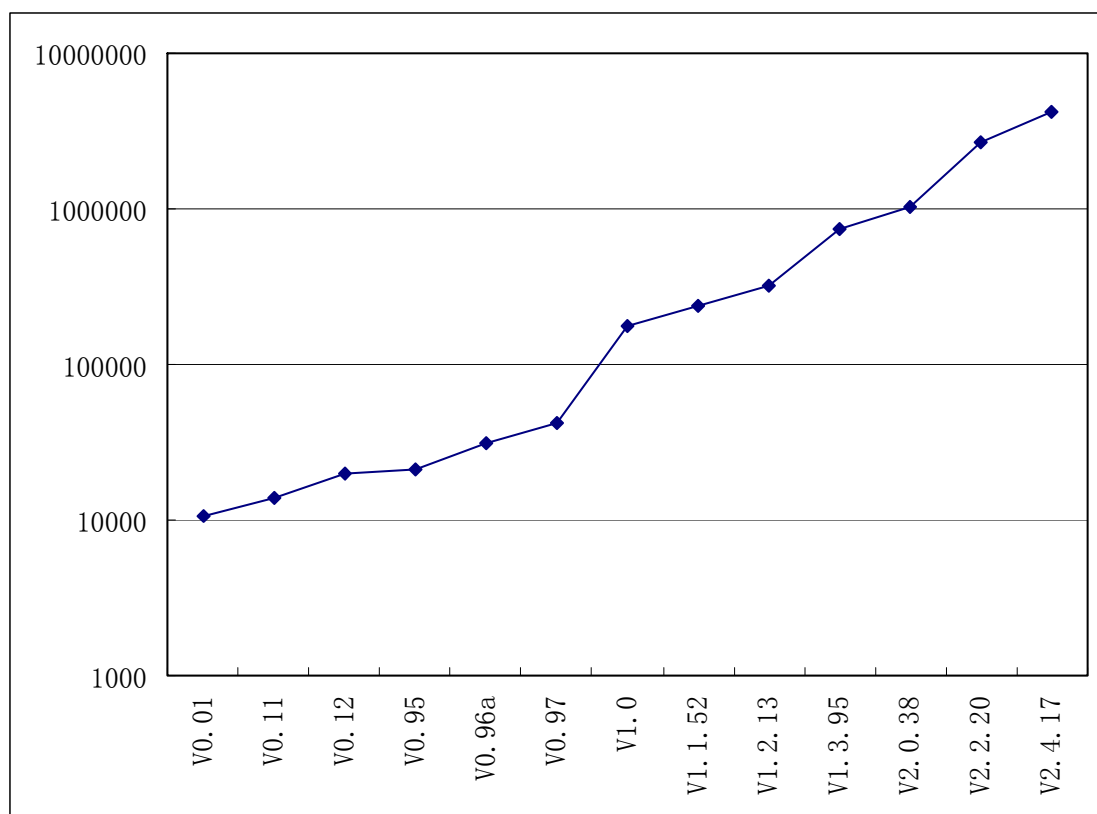
正如 Linux 系统的创始人的一篇新闻组投稿上所说的，要理解一个软件系统的真正运行机制，一定要阅读其源代码（RTFSC - Read The Fucking Source Code）。系统本身是一个完整的整体，具有很多看似不重要的细节存在，但是若忽略这些细节，就会对整个系统的理解带来困难，并且不能真正了解一个实际系统的实现方法和手段。

虽然通过阅读一些操作系统原理经典书籍（例如 M. J. Bach 的《UNIX 操作系统设计》）能够对 UNIX 类操作系统的工作原理有一些理论上的指导作用，但实际上对操作系统的真正组成和内部关系实现的理解仍

不是很清晰。正如 AST 所说的，“许多操作系统教材都是重理论而轻实践”，“多数书籍和课程为调度算法耗费大量的时间和篇幅而完全忽略 I/O，其实，前者通常不足一页代码，而后者往往要占到整个系统三分之一的代码总量。”内核中大量的重要细节均未提到。因此并不能让读者理解一个真正的操作系统实现的奥妙所在。只有在详细阅读过完整的内核源代码之后，才会对系统有一种忽然开朗的感觉，对整个系统的运作过程有深刻的理解。以后再选择最新的或较新内核源代码进行学习时，也不会碰到大问题，基本上都能顺利地理解新代码的内容。

## 如何选择要阅读的内核代码版本

那么，如何选择既能达到上述要求，又不被太多的内容而搞乱头脑，选择一个适合的 Linux 内核版本进行学习，提高学习的效率呢？作者通过对大量内核版本进行比较和选择后，最终选择了与目前 Linux 内核基本功能较为相近，又非常短小的 0.11 版内核作为入门学习的最佳版本。下图是对一些主要 Linux 内核版本行数的统计。



目前的 Linux 内核源代码量都在几百万行的数量上，极其庞大，对这些版本进行完全注释和说明是不可能的，而 0.11 版内核不超过 2 万行代码量，因此完全可以在一本书中解释和注释清楚。麻雀虽小，五脏俱全。

另外，使用该版本可以避免使用现有较新内核版本中已经变得越来越复杂得各子系统部分的研究（如虚拟文件系统 VFS、ext2 或 ext3 文件系统、网络子系统、新的复杂的内存管理机制等）。

## 阅读本书需具备的基础知识

在阅读本书时，作者希望读者具有以下一些基础知识或有相关的参考书籍在手边。其一是有关 80x86 处理器结构和编程的知识或资料。例如可以从网上下载的 80x86 编程手册（INTEL 80386 Programmer's

Reference Manual); 其二是有关 80x86 硬件体系结构和接口编程的知识或资料。有关这方面的资料很多; 其三还应具备初级使用 Linux 系统的简单技能。

另外, 由于 Linux 系统内核实现, 最早是根据 M. J. Bach 的《UNIX 操作系统设计》一书的基本原理开发的, 源代码中许多变量或函数的名称都来自该书, 因此在阅读本书时, 若能适当参考该书, 更易于对源代码的理解。

Linux 在最初开发 Linux 操作系统时, 参照了 Minix 操作系统。例如, 最初的 Linux 内核版本完全照搬了 Minix 1.0 文件系统。因此, 在阅读本书时, A. S. Tanenbaum 的书《操作系统: 设计与实现》也具有一定的参考价值。但 Tanenbaum 的书描述的是一种基于消息传递在内核各模块之间进行通信(信息交换), 这与 Linux 内核的工作机制不一样。因此可以仅参考其中有关一般操作系统工作原理章节和文件系统实现的内容。

## 使用早期版本是否过时?

表面看来本书对 Linux 早期内核版本注释的内容犹如 Linux 操作系统刚公布时 Tanenbaum 就认为其已经过时的(Linux is obsolete)想法一样, 但通过学习本书内容, 你就会发现, 利用本书学习 Linux 内核, 由于内核源代码量短小而精干, 因此会有极高的学习效率, 能够做到事半功倍, 快速入门。并且对继续进一步选择新内核部分源代码的学习打下坚实的基础。在学习完本书之后, 你将对系统的运作原理有一个非常完整而实际的概念, 这种完整概念能使人很容易地进一步选择和学习新内核源代码中的任何部分, 而不需要再去啃读代码量巨大的新内核中的完整的源代码。

## Ext2 文件系统与 Minix 文件系统?

目前 Linux 系统上所使用的 Ext2 (或最新的 Ext3) 文件系统, 是在内核 1. x 之后开发的文件系统, 其功能详尽并且其性能也非常完整和稳固, 是目前 Linux 操作系统上默认的标准文件系统。但是, 作为对 Linux 操作系统完整工作原理入门学习所使用的部分, 原则上是越精简越好。为了达到对一个操作系统有完整的理解, 并且能不被其中各子系统复杂性和过多的细节喧宾夺主, 在选择学习剖析用的内核版本时, 只要系统的部分代码内容能说明实际工作原理, 就越简单越好。

Linux 内核 0. 11 版上当时仅包含最为简单的 minix 1. 0 文件系统, 对于理解一个操作系统中文件系统的实际组成和工作原理已经足够。这也是选择 Linux 早期内核版本进行学习的主要原因之一。

在完整阅读完本书之后, 相信您定会发出这样的感叹: “对于 Linux 内核系统, 我现在终于入门了! ”。此时, 您应该有十分的把握去进一步学习最新内核中各部分的工作原理和过程了。

同济大学  
赵炯 博士  
2002. 12

# 第 1 章 概述

本章首先回顾了 Linux 操作系统的诞生、开发和成长过程，由此理解本书为什么会选择 Linux 系统早期版本作为学习对象的一些原因。然后具体说明了选择早期 Linux 内核版本进行学习的优点和不足之处以及如何开始进一步的学习。接下来就直接切入正题，从内核源文件 `linux/` 目录下的第一个文件 `Makefile` 开始，对每一行代码进行详细注释说明。

## 1.1 Linux 的诞生和发展

Linux 操作系统是 UNIX 操作系统的一种克隆系统。它诞生于 1991 年的 10 月 5 日（这是第一次正式向外公布的时间）。以后借助于 Internet 网络，并经过全世界各地计算机爱好者的共同努力下，现已成为今天世界上使用最多的一种 UNIX 类操作系统，并且使用人数还在迅猛增长。

Linux 操作系统的诞生、发展和成长过程始终依赖着以下五个重要支柱：UNIX 操作系统、MINIX 操作系统、GNU 计划、POSIX 标准和 Internet 网络。

下面主要根据这五个基本线索来追寻一下 Linux 的开发历程，它的酝酿过程，最初的发展经历。首先分别介绍其中的四个基本要素（UNIX、MINIX、GNU 和 POSIX，Internet 的重要性显而易见，所以不用对其罗嗦），然后根据 Linux 的创始人 Linus Torvalds 从对计算机感兴趣而自学计算机知识，到心里开始酝酿编制一个自己的操作系统，到最初 Linux 内核 0.01 版公布，以及从此如何艰难地一步一个脚印地在全世界 hacker 的帮助下最后推出比较完善的 1.0 版本这段时间的发展经过，也即对 Linux 的早期发展历史进行详细介绍。

当然，目前 Linux 内核版本已经开发到了 2.5.52 版。而大多数 Linux 系统中所用到的内核是稳定的 2.4.20 版内核。（其中第 2 个数字奇数表示是正在开发的版本，不能保证系统的稳定性）对于 Linux 的一般发展史，许多文章和书籍都有介绍，这里就不重复。

### 1.1.1 UNIX 操作系统的诞生

Linux 操作系统是 UNIX 操作系统的一个克隆版本。UNIX 操作系统是美国贝尔实验室的 Ken. Thompson 和 Dennis Ritchie 于 1969 年夏在 DEC PDP-7 小型计算机上开发的一个分时操作系统。

当时 Ken Thompson 为了能在闲置不用的 PDP-7 计算机上运行他非常喜欢的星际旅行（Space travel）游戏，在 1969 年夏天乘他夫人回家乡加利福尼亚渡假期间，在一个月内开发出了 unix 操作系统的原型。当时使用的是 BCPL 语言（基本组合编程语言），后经 Dennis Ritchie 于 1972 年用移植性很强的 C 语言进行了改写，使得 UNIX 系统在大专院校得到了推广。

### 1.1.2 MINIX 操作系统

MINIX 系统是由 Andrew S. Tanenbaum（AST）开发的。AST 是在荷兰 Amsterdam 的 Vrije 大学数学与计算机科学系统工作，是 ACM 和 IEEE 的资深会员（全世界也只有很少人是两会的资深会员）。共发表了 100 多篇文章，5 本计算机书籍。

AST 虽出生在美国纽约，但是是荷兰侨民（1914 年他的祖辈来到美国）。他在纽约上的中学、M. I. T 上的大学、加州大学 Berkeley 分校念的博士学位。由于读博士后的缘故，他来到了家乡荷兰。从此就与家乡一直有来往。后来就在 Vrije 大学开始教书、带研究生了。荷兰首都 Amsterdam 是个常年阴雨绵绵的城市，而对于 AST 来说，这最好不过了，因为这样他就可以待在家里摆弄他的计算机了。

MINIX 是他 1987 年编制的，主要用于学生学习操作系统原理。到 91 年时版本是 1.5。目前主要有两

个版本在使用：1.5 版和 2.0 版，当时该操作系统在大学使用是免费的，但其它用途不是，当然目前都已经是免费的，可以从许多 FTP 上下载。

对于 Linux 系统，他表示对其开发者 Linus 的称赞。但他认为 Linux 的发展有很大原因是因为他为了保持 minix 的小型化，能让学生在一个学期内就能学完，而没有接纳全世界许多人对 Minix 的扩展要求。因此这激发了 Linus 编写 Linux。Linus 正好抓住了这个好时机。

作为一个操作系统，MINIX 并不是优秀者，但它同时提供了用 C 语言和汇编语言写的系统源代码。这是第一次使得有抱负的程序员或 hacker 能够阅读操作系统的源代码，在当时这种源代码是软件商一直小心地守护着的。

### 1.1.3 GNU 计划

GNU 计划和自由软件基金会(the Free Software Foundation - FSF)是由 Richard M. Stallman 于 1984 年一手创办的。旨在开发一个类似 Unix、并且是自由软件的完整操作系统：GNU 系统。(GNU 是“GNU’s Not Unix”的递归缩写，它的发音为“guh-NEW”。)各种使用 linux 作为核心的 GNU 操作系统正在被广泛的使用。虽然这些系统通常被称作“Linux”，但是严格地说，它们应该被称为 GNU/Linux 系统。

到上世纪 90 年代初，GNU 项目已经开发出许多高质量的免费软件，其中包括有名的 emacs 编辑系统、bash shell 程序、gcc 系列编译程序、gdb 调试程序等等。这些软件为 Linux 操作系统的开发创造了一个合适的环境，是 Linux 能够诞生的基础之一。以至于目前许多人都将 Linux 操作系统称为“GNU/Linux”操作系统。

### 1.1.4 POSIX 标准

POSIX(Portable Operating System Interface for Computing Systems)是由 IEEE 和 ISO/IEC 开发的一簇标准。该标准是基于现有的 UNIX 实践和经验，描述了操作系统的调用服务接口，用于保证编制的应用程序可以在源代码一级上在多种操作系统上移植运行。它是在 1980 年早期一个 UNIX 用户组(usr/group)的早期工作的基础上取得的。该 UNIX 用户组原来试图将 AT&T 的系统 V 和 Berkeley CSRG 的 BSD 系统的调用接口之间的区别重新调和集成，从而于 1984 年产生了/usr/group 标准。1985 年，IEEE 操作系统技术委员会标准小组委员会(TCOS-SS)开始在 ANSI 的支持下责成 IEEE 标准委员会制定有关程序源代码可移植性操作系统服务接口正式标准。到了 1986 年 4 月，IEEE 就制定出了试用标准。第一个正式标准是在 1988 年 9 月份批准的(IEEE 1003.1-1988)，也既以后经常提到的 POSIX.1 标准。

1989 年 POSIX 的工作被转移至 ISO/IEC 社团，并由 15 工作组继续将其制定成 ISO 标准。到 1990 年，POSIX.1 与已经通过的 C 语言标准联合，正式批准为 IEEE 1003.1-1990(也是 ANSI 标准)和 ISO/IEC 9945-1:1990 标准。

POSIX.1 仅规定了系统服务应用程序编程接口(API)，仅概括了基本的系统服务标准，因此期望对系统的其它功能也制定出标准。这样 IEEE POSIX 的工作就开始展开了。在 1990 年，刚开始有十个批准的计划在进进行，有近 300 多人参加每季度为期一周的会议。着手的工作有命令与工具标准(POSIX.2)、测试方法标准(POSIX.3)、实时 API(POSIX.4)等。到了 1990 年上半年已经有 25 个计划在进进行，并且有 16 个工作组参与了进来。与此同时，还有一些组织也在制定类似的标准，如 X/Open, AT&T, OSF 等。

在 90 年代初，POSIX 标准的制定正处在最后投票敲定的时候，那是 1991-1993 年间。此时正是 Linux 刚刚起步的时候，这个 UNIX 标准为 Linux 提供了极为重要的信息，使得 Linux 的能够在标准的指导下进行开发，能够与绝大多数 UNIX 系统兼容。在最初的 Linux 内核代码中(0.01 版、0.11 版)就已经为 Linux 与 POSIX 标准的兼容做好了准备工作。在 0.01 版的内核/include/unistd.h 文件中就已经定义了几个有关 POSIX 标准要求的常数符号，并且在注释中就写到“ok，这也许是个玩笑，但我正在着手研究它呢”。

1991 年 7 月 3 日在 comp.os.minix 上发布的 post 上就已经提到了正在搜集 POSIX 的资料。(当然此时还不存在 Linux 这个名称，当时 Linus 的脑子里想的可能是 FREAX ©，FREAX 的英文含义是怪诞的、怪物、异想天开等)。其中透露了他正在进行 Linux 系统的开发，并且在 Linux 最初的时候已经想到要实现与 POSIX



(UNIX 的国际标准) 的兼容问题了。

### 1.1.5 Linux 操作系统的诞生

1981 年 IBM 公司推出享誉全球的微型计算机 IBM PC。在 1981-1991 年间, MS-DOS 操作系统一直是微型计算机上操作系统的主宰。此时计算机硬件价格虽然逐年下降, 但软件价格仍然是居高不下。当时 Apple 的 MACs 操作系统可以说是性能最好的, 但是其天价没人能够轻易靠近。

当时的另一个计算机技术阵营是 Unix 世界。但是 Unix 操作系统就不仅是价格贵的问题了。为了寻求高利率, Unix 经销商将价格抬得极高, PC 小用户就根本不能靠近它。曾经一度受到 Bell Labs 的许可而可以在大学中用于教学的 UNIX 源代码一直被小心地守卫着不需公开。对于广大的 PC 用户, 软件行业的大型供应商始终没有给出有效的解决该问题的手段。

正在此时, 出现了 MINIX 操作系统, 并有一本详细的书本描述它的设计实现原理。由于 AST 的书写得非常详细, 并且叙述有条有理, 几乎全世界的计算机爱好者都在看这本书以理解操作系统的工作原理。其中也包括 Linux 系统的创始者 Linus Benedict Torvalds。

当时(1991 年), Linus Benedict Torvalds 是赫尔辛基大学计算机科学系的二年级学生, 也是一个自学 hacker。这个 21 岁的芬兰年轻人喜欢鼓捣计算机, 测试计算机的能力和限制。但当时缺乏的是一个专业级的操作系统。MINIX 虽然很好, 但只是一个用于教学目的简单操作系统, 而不是一个强有力的实用操作系统。

到 1991 年, GNU 计划已经开发出了许多工具软件。最受期盼的 Gnu C 编译器已经出现, 但还没有开发出免费的 GNU 操作系统。即使是 MINIX 也开始有了版权, 需要购买才能得到源代码。而 GNU 的操作系统 HURD 一直在开发之中, 但并不能在几年内完成。

对于 Linus 来说, 已经不能等待了。从 1991 年 4 月份起, 他开始酝酿并着手编制自己的操作系统。刚开始, 他的目的很简单, 只是为了学习 Intel 386 体系结构保护模式运行方式下的编程技术。但后来 Linux 的发展却完全改变了初衷。

1991 年初, Linux 开始在一台 386sx 兼容微机上学习 minix 操作系统。通过学习, 他逐渐不能满足 minix 系统的现有性能, 并开始酝酿开发一个新的免费操作系统。根据 Linus 在 comp.os.minix 新闻组上发布的消息, 我们可以知道他逐步从学习 minix 系统到开发自己的 Linux 的过程。

Linus 第 1 次向 comp.os.minix 投递消息是在 1991 年 3 月 29 日。题目是“gcc on minix-386 doesn't optimize”, 是有关 gcc 编译器在 minix-386 上运行的优化问题, 由此可知, Linus 在 1991 年的初期已经开始深入研究了 minix 系统, 并在这段时间有了改进 minix 操作系统的思想, 而且在进一步学习 minix 系统中, 逐步演变为想自己重新设计一个基于 Intel 80386 体系结构的新操作系统。

他在回答有人提出 minix 上的一个问题时, 所说的第一句话是“阅读源代码”(“RTFSC (Read the F\*\*ing Source Code :-)” )。他认为答案就在源程序中。这也说明了对于学习系统软件来说, 你不光需要懂得系统的工作基本原理, 还需要结合实际系统, 学习实际系统的实现方法。因为理论毕竟是理论, 其中省略了许多枝节, 而这些枝节问题虽然没有太多的理论含量, 但却是一个系统必要的组成部分, 就向麻雀身上的一根羽毛。

从 1991 年的 4 月份开始, Linus 几乎花了全部时间研究 386-minix 系统(hack the kernel), 并且尝试着移植 GNU 的软件到该系统上(GNU gcc、bash、gdb 等)。并于 4 月 13 日在 comp.os.minix 上发布说自己已经成功地将 bash 移植到了 minix 上, 而且已经爱不释手、不能离开这个 shell 软件了。

第一个与 Linux 有关的消息是在 1991 年 7 月 3 日在 comp.os.minix 上发布的(当然此时还不存在 Linux 这个名称, 当时 Linus 的脑子里想的可能是 FREAX ☺, FREAX 的英文含义是怪诞的、怪物、异想天开等)。其中透露了他正在进行 Linux 系统的开发, 并且在 Linux 最初的时候已经想到要实现与 POSIX (UNIX 的国际标准) 的兼容问题了。

在 Linus 的下一发布的消息中(1991 年 8 月 25 日 comp.os.minix), 他向所有 minix 用户询问“*What would you like to see in minix?*”(“你最想在 minix 中见到什么?”), 在该消息中他首次透露出正

在开发一个(免费的)386(486)操作系统,并且说只是兴趣而已,代码不会很大,也不会象 GNU 的那样专业。开发免费操作系统这个想法从 4 月份就开始酝酿了,希望大家反馈一些对于 minix 系统中喜欢那些特色不喜欢什么等信息,由于实际的和其它一些原因,新开发的系统刚开始与 minix 很象(并且使用了 minix 的文件系统)。并且已经成功地将 bash(1.08 版)和 gcc(1.40 版)移植到了新系统上,而且在过几个月就可以实用了。

最后, Linus 申明他开发的操作系统没有使用一行 minix 的源代码;而且由于使用了 386 的任务切换特性,所以该操作系统不好移植(没有可移植性),并且只能使用 AT 硬盘。对于 Linux 的移植性问题, Linus 当时并没有考虑。但是目前 Linux 几乎可以运行在任何一种硬件体系结构上。

到了 1991 年的 10 月 5 日, Linus 在 comp.os.minix 新闻组上发布消息,正式向外宣布 Linux 内核系统的诞生(Free minix-like kernel sources for 386-AT)。这段消息可以称为 Linux 的诞生宣言,并且一直广为流传。因此 10 月 5 日对 Linux 社区来说是一个特殊的日子,许多后来 Linux 的新版本发布时都选择了这个日子。所以 RedHat 公司选择这个日子发布它的新系统也不是偶然的。

### 1.1.6 Linux 操作系统版本的变迁

0.00 (1991.2-4?) 两个进程分别显示 AAA BBB  
 0.01 (1991.9?) 第一个正式向外公布的 Linux 内核版本。  
 0.02 (1991.10.5) 该版本以及 0.03 版是内部版本,目前已经无法找到。  
 0.03 (1991.10.5)  
 0.10 (1991.10) 由 Ted Ts'o 发布的 Linux 内核版本。  
 0.11 (1991.12.8) 基本可以正常运行的内核版本(也是本书着重注释的版本)。  
 0.12 (1992.1.15) 主要加入对数学协处理器的软件模拟程序。  
 2.0 (1992.3.8) 开始加入虚拟文件系统思想的内核版本。  
 0.96 (1992.5.12) 开始加入网络支持和虚拟文件系统 VFS。  
 0.97 (1992.8.1)  
 0.98 (1992.9.29)  
 0.99 (1992.12.13)  
 2.4 (1994.3.14)  
 1.20 (1995.3.7)  
 2.0 (1996.2.9)  
 2.20 (1999.1.26)  
 2.40 (2001.1.4)  
 到现在为止,最新版是

| 内核版本号  | 发布日期       |  | 代码总字节数     |
|--------|------------|--|------------|
| 2.4.20 | 2002.11.29 |  | 26,200,000 |
| 2.5.52 | 2002.12.16 |  | 30,000,000 |

将 Linux 系统 0.13 版内核直接改称 0.95 版, Linus 的意思是让大家不要觉得离 1.0 版还很遥远。同时,从 0.95 版开始,对内核的许多改进之处(补丁程序的提供)均以其他人为主了,而 Linus 的主要任务开始变成对内核的维护和决定是否采用某个补丁程序。

### 1.1.7 Linux 名称的来由

Linux 操作系统刚开始时并没有被称作 Linux, Linus 给他的操作系统取名为 FREAX,其英文含义是怪诞的、怪物、异想天开等意思。在他将新的操作系统上载到 ftp.funet.fi 服务器上时,管理员 Ari Lemke

很不喜欢这个名称。他认为既然是 Linus 的操作系统就取其谐音 Linux 作为该操作系统的目录吧，于是 Linux 这个名称就开始流传下来。

在 Linus 的自传《Just for Fun》一书中，Linus 解释说：

“坦白地说，我从来没有想到过要用 Linux 这个名称发布这个操作系统，因为这个名字有些太自负了。而我为最终发布版准备的是什么呢？Freax。实际上，内核代码中某些早期的 Makefile - 用于描述如何编译源代码的文件 - 文件中就已经包含有“Freax”这个名字了，大约存在了半年左右。但其实这也没什么关系，在当时还不需要一个名字，因为我还没有向任何人发布过内核代码。”

“而 Ari Lemke，他坚持要用自己的方式将内核代码放到 ftp 站点上，并且非常不喜欢 Freax 这个名字。他坚持要用现在这个名字(Linux)，我承认当时我并没有跟他多争论。但这都是他取的名字。所以我可以光明正大地说我并不自负，或者部分坦白地说我并没有本位主义思想。但我想好吧，这也是个好名字，而且以后为这事我总能说服别人，就象我现在做的这样。”

-- Linus Torvalds 《Just for fun》第 84-88 页。

### 1.1.8 早期 Linux 系统开发的主要贡献者

从 Linux 的早期源代码中可知，Linux 系统的早期主要开发人员除了 Linus 本人以外，最著名的人员之一就是 Theodore Ts'o (Ted Ts'o)。他 1990 年毕业于 MIT 计算机专业。在大学时代他就积极参加学校中举办的各种学生活动。他喜欢烹饪、骑自行车，当然还有就是 hacking on Linux，后来他开始喜欢起业余无线电报运动。目前他在 IBM 工作从事系统编程及其它重要事务。他还是国际网络设计、操作、销售和研究者开放团体 IETF 成员。

Linux 在世界范围内的流行也有他很大的功劳。早在 Linux 操作系统刚问世时，他就怀着极大的热情为 linux 的发展提供了 maillist，几乎是在 Linux 刚开始发布起(1991 年开始)就一直为 Linux 做出贡献的人，也是最早向 Linux 内核添加程序的人(Linux 内核 0.10 版中的虚拟盘驱动程序 ramdisk.c 和内核内存分配程序 kmalloc.c)。直到目前仍然从事着与 Linux 有关的工作。他当时在北美洲地区最早设立了 linux 的 ftp 站点(tsx-11.mit.edu)，而且至今仍然为广大 linux 用户提供服务。他对 linux 作出的最大贡献之一是提出并实现了 ext2 文件系统。该文件系统已成为 linux 世界中事实上的文件系统标准。最近他又推出了 ext3 文件系统，大大提高了文件系统的稳定性和访问效率。作为对他的推崇，第 97 期(2002 年 5 月)的 linuxjournal 期刊将他作为封面人物，并对他进行了采访。目前，他为 IBM linux 技术中心工作，并从事着有关 LSB(Linux Standard Base)等方面的工作。

Linux 社区中另一位著名人物是 Alan Cox。他原工作于英国威尔士斯旺西大学(Swansea University College)。刚开始他特别喜欢玩电脑游戏，尤其是 MUD (Multi-User Dungeon or Dimension，多用户网络游戏)。在 90 年代早期 games.mud 新闻组的 posts 中你可以找到他发表的大量 posts。他甚至为此还写了一篇 MUD 的发展史(rec.games.mud 新闻组，1992 年 3 月 9 日，A history of MUD)。由于 MUD 游戏与网络密切相关，慢慢的他对计算机网络开始感兴趣。为了玩游戏并提高电脑运行游戏的速度以及网络传输的速度，他开始接触各种类型的操作系统，为他的游戏选择一个最为满意的平台。由于没钱，即使 Minix 他都买不起，当 Linux 0.11 和 386BSD 发布时，他考虑良久总算买了一台 386SX 电脑。由于 386BSD 需要数学协处理器的支持，而 386SX 中是不带的，所以他安装了 Linux 系统。于是他开始学习带有免费源代码的 Linux。开始对 Linux 产生了兴趣，尤其是有关网络方面的实现。在关于 Linux 的单用户运行模式问题的讨论中，他甚至赞叹 Linux 实现的巧妙(beautifully)。

Linux 0.95 版发布之后，他开始为 Linux 系统编写补丁程序(修改程序)(记得他最早的两个补丁程序，都没有被 Linus 采纳)，成为 Linux 上 TCP/IP 网络代码的最早使用人之一。后逐渐加入 Linux 的开发队伍，并开始成为维护 Linux 内核源代码的主要负责人之一，也可以说成为 Linux 社团中在 Linus 之后最为重要的人物。以后 Microsoft 公司曾经邀请他加盟，但他却干脆地拒绝了。从 2001 年开始他负责维护 Linux 内核 2.4.x 的代码(Linus 主要负责开发最新开发版内核的研制(奇数版，比如 2.5.x 版)。

《内核骇客手册》一书的作者 Michael K. Johnson 也是最早接触 Linux 操作系统的人之一(从 0.97 版)。

他还是著名 Linux 文档计划 (Linux Document Project - LDP) 的发起者之一。曾经在 Linux Journal 工作, 现在 RedHat 公司工作。

Linux 系统并不是仅有这些中坚力量就能发展成今天这个样子的, 还有许多计算机高手对 Linux 做出了极大的贡献, 这里就不一一列举了。主要贡献者的具体名单可参见 Linux 内核中的 CREDITS 文件, 其中以字母顺序列出了对 Linux 做出较大贡献的近 400 人的名单列表, 包括他们的 email 地址和通信地址、主页以及主要贡献事迹等信息。

通过上述说明, 我们可以对上述 Linux 的五大支柱归纳如下:

- UNIX 操作系统 -- UNIX 于 1969 年诞生在 Bell 实验室。Linux 就是 UNIX 的一种克隆系统。UNIX 的重要性就不用多说了。
- MINIX 操作系统 -- Minix 操作系统也是 UNIX 的一种克隆系统, 它于 1987 年由著名计算机教授 Andrew S. Tanenbaum 开发完成。由于 MINIX 系统的出现并且提供源代码(只能免费用于大学内)在全世界的大学中刮起了学习 UNIX 系统旋风。Linux 刚开始就是参照 Minix 系统于 1991 年才开始开发。
- GNU 计划 -- 开发 Linux 操作系统, 以及 Linux 上所用大多数软件基本上都出自 GNU 计划。Linux 只是操作系统的一个内核, 没有 GNU 软件环境(比如说 bash shell), 则 Linux 将寸步难行。
- POSIX 标准 -- 该标准在推动 Linux 操作系统以后朝着正规路上发展起着重要的作用。是 Linux 前进的灯塔。
- INTERNET -- 如果没有 Internet 网, 没有遍布全世界的无数计算机骇客的无私奉献, 那么 Linux 最多只能发展到 0.13(0.95)版的水平。

## 1.2 内容综述

本文将主要对 Linux 的早期内核 0.11 版进行详细描述和注释。Linux-0.11 版本是在 1991 年 12 月 8 日发布的。在发布时包括以下文件:

bootimage.Z      - 具有美国键盘代码的压缩启动映像文件;  
rootimage.Z      - 以 1200kB 压缩的根文件系统映像文件;  
linux-0.11.tar.Z - 内核源代码文件;  
as86.tar.Z       - linux bruce evans' 二进制执行文件;  
                  是 16 位的汇编程序和装入程序;  
INSTALL-0.11    - 更新过的安装信息文件。

目前除了原来的 rootimage.Z 文件, 其它四个文件均能找到。

本文主要详细分析 linux-0.11 内核中的所有源代码程序, 对每个源程序文件都进行了详细注释, 包括对 Makefile 文件的注释。分析过程主要是按照计算机启动过程进行的。因此分析的连贯性到初始化结束内核开始调用 shell 程序为止。其余的各个程序均针对其自身进行分析, 没有连贯性, 因此可以根据自己的需要进行阅读。但在分析时还是提供了一些应用实例。

所有的程序在分析过程中如果遇到作者认为是较难理解的语句时, 将给出相关知识的详细介绍。比如, 在阅读代码头一次遇到 C 语言内嵌汇编码时, 将对 gnu C 语言的内嵌汇编语言进行较为详细的介绍; 在遇到对中断控制器进行输入/输出操作时, 将对 Intel 中断控制器 (8259A) 芯片给出详细的说明, 并列出的命令和方法。这样做有助于加深对代码的理解, 又能更好的了解所用硬件的使用方法, 作者认为这种










解读方法要比单独列出一章内容来总体介绍硬件或其它知识要效率高得多。

拿 Linux 0.11 版内核来“开刀”是为了提高我们认识 Linux 运行机理的效率。Linux-0.11 版整个内核源代码只有 325K 字节左右，其中包括的内容基本上都是 Linux 的精髓。而目前最新的 2.5.XX 版内核非常大，将近有 188 兆字节，即使你花一生的经历来阅读也未必能全部都看完。也许你要问“既然要从简入手，为什么不分析更小的 Linux 0.01 版内核源代码呢？它只有 240K 字节左右”主要原因是因为 0.01 版的内核代码有太多的不足之处，甚至还没有包括对软盘的驱动程序，也没有很好地涉及数学协处理器的使用以及对登陆程序的说明。并且其引导启动程序的结构也与目前的版本不太一样，而 0.11 版的引导启动程序结构则与现在的基本上是一样的。另外一个原因是可以找到 0.11 版早期的已经编译制作好的内核映像文件 (bootimage)，可以用来进行引导演示。如果再配上简单的根文件系统映像文件 (rootimage)，那么它就可以进行正常的运行了。

拿 Linux 0.11 版进行学习也有不足之处，比如该内核版本中尚不包括有关进程等待队列、TCP/IP 网络等方面的一些当前非常重要的代码，但好在 Linux 中的网络代码基本上是自成一体的，与内核机制关系不是非常大，因此可以在了解了 Linux 工作的基本原理之后再去分析这些代码。

本文对 Linux 内核中所有的代码都进行了说明。为了保持结构的完整性，对代码的说明是以内核中源代码的组成结构来进行的，基本上是以每个源代码中的目录为一章内容进行介绍。介绍的源程序文件的次序可参见前面的文件列表索引。整个 Linux 内核源代码的目录结构如下列表 1.1 所示。所有目录结构均是以 linux 为当前目录。

列表 1.1 Linux/目录

| Name   | Size       | Last modified (GMT) | Description |
|--|------------|---------------------|-------------|
|  <a href="#">boot/</a>    |            | 1991-12-05 22:48:49 |             |
|  <a href="#">fs/</a>      |            | 1991-12-08 14:08:27 |             |
|  <a href="#">include/</a> |            | 1991-09-22 19:58:04 |             |
|  <a href="#">init/</a>    |            | 1991-12-05 19:59:04 |             |
|  <a href="#">kernel/</a>  |            | 1991-12-08 14:08:00 |             |
|  <a href="#">lib/</a>     |            | 1991-12-08 14:10:00 |             |
|  <a href="#">mm/</a>      |            | 1991-12-08 14:08:21 |             |
|  <a href="#">tools/</a>   |            | 1991-12-04 13:11:56 |             |
|  <a href="#">Makefile</a> | 2887 bytes | 1991-12-06 03:12:46 |             |

本章主要说明内核程序包中 Linux/目录下的所看到的第一个文件，也即内核代码的总体 Makefile 文

件的内容。

第 2 章将详细注释 boot/目录下的三个汇编程序,其中包括磁盘引导程序 bootsect.s、32 位运行启动代码程序 head.s 和获取 BIOS 中参数的 setup.s 汇编程序。

第 3 章主要介绍 init/目录中内核系统的初始化程序 main.c。它是内核完成所有初始化工作并进入正常运行的关键地方。在介绍该程序时将需要查看其所调用的其它程序,因此对后续章节的阅读可以按照这里调用的顺序进行。当你能真正看懂直到 main.c 程序为止的所有程序时,你应该已经对 Linux 内核有了一定的了解,可以说已经有一半入门了☺,但你还需要对文件系统、系统调用、各种驱动程序等进行更进一步的阅读。

第 4 章主要介绍 kernel/目录中的所有程序。此时你应该已经对其中的一些重要程序有所了解。

第 5 章介绍内核源代码 fs/目录中的文件系统程序,在看这章内容时建议你能够暂停一下而去阅读 Andrew S. Tanenbaum 的《操作系统设计与实现》一书中有关 minix 文件系统的章节,因为最初的 Linux 系统是只支持 minix 一种文件系统, Linux 0.11 版也不例外。

第 6 章解说 mm/目录中的内存管理程序。要透彻地理解这方面的内容,需要对 Intel 80X86 微处理器的保护模式运行方式有足够的理解,因此本章在适当的地方包含有较为完整的有关 80X86 保护模式运行方式的说明,这些知识基本上是从 Intel 80386 程序员编程手册(Intel 80386 Programmer's Reference Manual)中摘录下来的,但在此是以源代码中的运用实例为对象进行解说的,因此应该能更好地理解它的原理。

第 7 章是对 include/目录中的所有头文件进行详细说明,基本上对每一个定义、每一个常量或数据结构都进行了详细注释。为了便于在阅读时参考查阅,本书在附录中还对一些经常要用到的重要的数据结构和变量进行了归纳注释,但这些内容其实都能在这一章中找到。

第 8 章介绍 Linux 0.11 版内核源代码中 lib/目录中的所有文件。由于这个版本较低,所以这里的内容并不是很多,可以很快地看完。这也是我们为什么选择 0.11 版的原因之一。

第 9 章介绍 tools/目录下的 build.c 程序。这个程序并不会包括在编译生成的内核映像(image)文件中,它仅用于将内核中的磁盘引导程序块与其它主要内核模块连接成一个完整了内核映像(kernel image)文件。

为了便于查阅,在本书的附录中还单独列出了内核中要用到的有关 PC 机硬件方面的信息。在参考文献中,我们仅给出了在阅读源代码时可以参考的书籍、文章等信息,并没有包罗万象地给出一大堆的繁杂凌乱的文献列表。比如在引用 Linux 文档项目 LDP(Linux Document Project)中的文件时,我们会明确地列出具体需要参考哪一篇 HOWTO 文章,而并不是仅仅给出 LDP 的网站地址了事。

在对每个程序进行解说时,我们首先简单说明程序的主要用途和目的、输入输出参数以及与其它程序的关系,然后列出程序的完整代码并在其中对代码进行详细注释,注释时对原程序代码或文字不作任何方面的改动或删除,因为 C 语言是一种英语类语言,程序中原有的少量英文注释对常数符号、变量名等也提供了不少有用的信息。在代码之后是对程序更为深入的解剖,并对代码中出现的一些语言或硬件方面的相关知识进行说明。如果在看完这些信息后回头再浏览一遍程序,你会有更深一层的体会。

最后要说明的是当你已经完全理解了本文说解说的一切时,并不代表你已经成为一个 Linux 行家了,你只是刚刚踏上 Linux 的征途,具有了一定的成为一个 Linux GURU 的初步知识。这时你应该去阅读更多的源代码,最好是循序渐进地从 1.0 版本开始直到最新的正在开发中的奇数编号的版本。在撰写这本书时最新的 Linux 内核是 2.5.44 版。当你能快速理解这些开发中的最新版本甚至能提出自己的建议和补丁(patch)程序时,我也甘拜下风了☺。

最后祝你征途愉快!

## 1.3 Makefile 文件

Makefile 文件相当于程序编译过程中的批处理文件。是工具程序 make 运行时的输入数据文件。只要在含有 Makefile 的当前目录中键入 make 命令，它就会依据 Makefile 文件中的设置对源程序或目标代码文件进行编译、链接或进行安装等活动。

make 工具程序能自动地确定一个大程序系统中那些程序文件需要被重新编译，并发出命令对这些程序文件进行编译。在使用 make 之前，需要编写 Makefile 信息文件，该文件描述了整个程序包中各程序之间的关系，并针对每个需要更新的文件给出具体的控制命令。通常，执行程序是根据其目标文件进行更新的，而这些目标文件则是由编译程序创建的。一旦编写好一个合适的 Makefile 文件，那么在你每次修改过程序系统中的某些源代码文件后，执行 make 命令就能进行所有必要的重新编译工作。make 程序是使用 Makefile 数据文件和代码文件的最后修改时间(last-modification time)来确定那些文件需要进行更新，对于每一个需要更新的文件它会根据 Makefile 中的信息发出相应的命令。在 Makefile 文件中，开头为'#'的行是注释行。文件开头部分的 '=' 赋值语句定义了一些参数或命令的缩写。

这个 Makefile 文件的主要作用是指示 make 程序最终使用独立编译连接成的 tools/目录中的 build 执行程序将所有内核编译代码连接和合并成一个可运行的内核映像文件 image。具体是对 boot/ 中的 bootsect.s、setup.s 使用 8086 汇编器进行编译，分别生成各自的执行模块。再对源代码中的其它所有程序使用 GNU 的编译器 gcc/gas 进行编译，并连接成模块 system。再用 build 工具将这三块组合成一个内核映像文件 image。基本编译连接/组合结构如下图 1.1 所示。

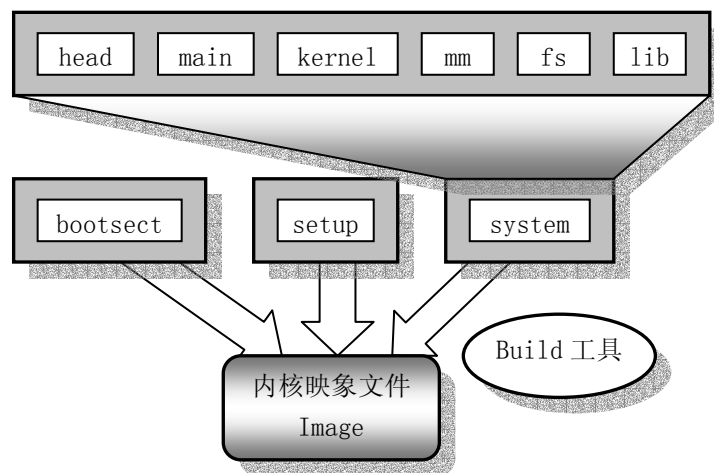


图 1.1 内核编译连接/组合结构

## 1.4 linux/Makefile 文件

列表 1.2 linux/Makefile 文件

|  |  |
|--|--|
| <p><u>1</u> #</p> <p><u>2</u> # if you want the ram-disk device, define this to be the</p> <p><u>3</u> # size in blocks.</p> <p><u>4</u> #</p> | <p># 如果你要使用 RAM 盘设备的话，就</p> <p># 定义块的大小。</p> |
|--|--|



```

5 RAMDISK = #-DRAMDISK=512
6
7 AS86      =as86 -O -a      # 8086 汇编编译器和连接器，见列表后的介绍。后带的参数含义分别
8 LD86      =ld86 -O        # 是：-O 生成 8086 目标程序；-a 生成与 gas 和 gld 部分兼容的代码。
9
10 AS        =gas            # GNU 汇编编译器和连接器，见列表后的介绍。
11 LD        =gld
12 LDFLAGS   =-s -x -M      # GNU 连接器 gld 运行时用到的选项。含义是：-s 输出文件中省略所
# 有的符号信息；-x 删除所有局部符号；-M 表示需要在标准输出设备
# (显示器)上打印连接映象(link map)，是指由连接程序产生的一种
# 内存地址映象，其中列出了程序段装入到内存中的位置信息。具体
# 来讲有如下信息：
# • 目标文件及符号信息映射到内存中的位置；
# • 公共符号如何放置；
# • 连接中包含的所有文件成员及其引用的符号。
13 CC        =gcc $(RAMDISK) # gcc 是 GNU C 程序编译器。对于 UNIX 类的脚本(script)程序而言，
# 在引用定义的标识符时，需在前面加上$符号并用括号括住标识符。
14 CFLAGS    =-Wall -O -fstrength-reduce -fomit-frame-pointer \
15 -fcombine-regs -mstring-insns # gcc 的选项。前一行最后的'\ '符号表示下一行是续行。
# 选项含义为：-Wall 打印所有警告信息；-O 对代码进行优化；
# -fstrength-reduce 优化循环语句；-mstring-insns 是
# Linus 自己为 gcc 增加的选项，可以去掉。
16 CPP       =cpp -nostdinc -Iinclude # cpp 是 gcc 的前(预)处理程序。-nostdinc -Iinclude 的含
# 义是不要搜索标准的头文件目录中的文件，而是使用-I
# 选项指定的目录或者是在当前目录里搜索头文件。
17
18 #
19 # ROOT_DEV specifies the default root-device when making the image.
20 # This can be either FLOPPY, /dev/xxxx or empty, in which case the
21 # default of /dev/hd6 is used by 'build'.
22 #
23 ROOT_DEV=/dev/hd6 # ROOT_DEV 指定在创建内核映像(image)文件时所使用的默认根文件系统所
# 在的设备，这可以是软盘(FLOPPY)、/dev/xxxx 或者干脆空着，空着时
# build 程序(在 tools/目录中)就使用默认值/dev/hd6。
24
25 ARCHIVES=kernel/kernel.o mm/mm.o fs/fs.o # kernel 目录、mm 目录和 fs 目录所产生的目标代
# 码文件。为了方便引用在这里将它们用
# ARCHIVES (归档文件)标识符表示。
26 DRIVERS =kernel/blk_drv/blk_drv.a kernel/chr_drv/chr_drv.a # 块和字符设备库文件。.a 表
# 示该文件是个归档文件，也即包含有许多可执行二进制代码子程
# 序集合的库文件，通常是用 GNU 的 ar 程序生成。ar 是 GNU 的二进制
# 文件处理程序，用于创建、修改以及从归档文件中抽取文件。
27 MATH      =kernel/math/math.a # 数学运算库文件。
28 LIBS      =lib/lib.a # 由 lib/目录中的文件所编译生成的通用库文件。
29
30 .c.s:      # 该行指示 make 利用下面的命令将所有的.c 文件编译生成.s 汇
# 编程序。':'表示下面是该规则的命令。
31 $(CC) $(CFLAGS) \
32 -nostdinc -Iinclude -S -o $*.s $< # 指使 gcc 采用前面 CFLAGS 所指定的选项以及
# 仅使用 include/目录中的头文件，在适当地编译后不进行汇编就
# 停止(-S)，从而产生与输入的各个 C 文件对应的汇编语言形式的
# 代码文件。默认情况下所产生的汇编程序文件是原 C 文件名去掉.c
# 而加上.s 后缀。-o 表示其后是输出文件的形式。其中$*.s (或$@)

```



```

# 是自动目标变量，$<代表第一个先决条件，这里即是符合条件
# *.c 的文件。
33 .s.o:      # 表示将所有.s 汇编程序文件编译成.o 目标文件。下一条是实
              # 现该操作的具体命令。
34 $(AS) -c -o $*.o $< # 使用 gas 编译器将汇编程序编译成.o 目标文件。-c 表示只编译
              # 或汇编，但不进行连接操作。
35 .c.o:      # 类似上面，*.c 文件→*.o 目标文件。
36 $(CC) $(CFLAGS) \
37 -nostdinc -Iinclude -c -o $*.o $< # 使用 gcc 将 C 语言文件编译成目标文件但不连接。
38
39 all:       Image      # all 表示创建 Makefile 所知的最顶层的目标。这里即是 image 文件。
40
41 Image: boot/bootsect boot/setup tools/system tools/build # 说明目标 (Image 文件) 是由
              # 分号后面的 4 个元素产生，分别是 boot/目录中的 bootsect 和
              # setup 文件、tools/目录中的 system 和 build 文件。
42 tools/build boot/bootsect boot/setup tools/system $(ROOT_DEV) > Image
43 sync      # 这两行是执行的命令。第一行表示使用 tools 目录下的 build 工具
              # 程序 (下面会说明如何生成) 将 bootsect、setup 和 system 文件
              # 以$(ROOT_DEV)为根文件系统设备组装成内核映像文件 Image。
              # 第二行的 sync 同步命令是迫使缓冲块数据立即写盘并更新超级块。

44
45 disk: Image      # 表示 disk 这个目标要由 Image 产生。
46 dd bs=8192 if=Image of=/dev/PS0 # dd 为 UNIX 标准命令：复制一个文件，根据选项
              # 进行转换和格式化。bs=表示一次读/写的字节数。
              # if=表示输入的文件，of=表示输出到的文件。
              # 这里/dev/PS0 是指第一个软盘驱动器 (设备文件)。

47
48 tools/build: tools/build.c      # 由 tools 目录下的 build.c 程序生成执行程序 build。
49 $(CC) $(CFLAGS) \
50 -o tools/build tools/build.c # 编译生成执行程序 build 的命令。
51
52 boot/head.o: boot/head.s      # 利用上面给出的.s.o 规则生成 head.o 目标文件。
53
54 tools/system: boot/head.o init/main.o \
55 $(ARCHIVES) $(DRIVERS) $(MATH) $(LIBS) # 表示 tools 目录中的 system 文件
              # 要由分号右边所列的元素生成。

56 $(LD) $(LDFLAGS) boot/head.o init/main.o \
57 $(ARCHIVES) \
58 $(DRIVERS) \
59 $(MATH) \
60 $(LIBS) \
61 -o tools/system > System.map # 生成 system 的命令。最后的 > System.map 表示
              # gld 需要将连接映像重定向存放在 System.map 文件中。
              # 关于 System.map 文件的用途参见注释后的说明。

62
63 kernel/math/math.a:      # 数学协处理函数文件 math.a 由下一行上的命令实现。
64 (cd kernel/math; make)   # 进入 kernel/math/目录；运行 make 工具程序。
              # 下面从 66--82 行的含义与此处的类似。

65
66 kernel/blk_drv/blk_drv.a: # 块设备函数文件 blk_drv.a
67 (cd kernel/blk_drv; make)
68
69 kernel/chr_drv/chr_drv.a: # 字符设备函数文件 chr_drv.a

```

```

70      (cd kernel/chr_drv; make)
71
72 kernel/kernel.o:                # 内核目标模块 kernel.o
73      (cd kernel; make)
74
75 mm/mm.o:                        # 内存管理模块 mm.o
76      (cd mm; make)
77
78 fs/fs.o:                        # 文件系统目标模块 fs.o
79      (cd fs; make)
80
81 lib/lib.a:                      # 库函数 lib.a
82      (cd lib; make)
83
84 boot/setup: boot/setup.s        # 这里开始的三行是使用 8086 汇编和连接器
85      $(AS86) -o boot/setup.o boot/setup.s # 对 setup.s 文件进行编译生成 setup 文件。
86      $(LD86) -s -o boot/setup boot/setup.o # -s 选项表示要去除目标文件中的符号信息。
87
88 boot/bootsect: boot/bootsect.s  # 同上。生成 bootsect.o 磁盘引导块。
89      $(AS86) -o boot/bootsect.o boot/bootsect.s
90      $(LD86) -s -o boot/bootsect boot/bootsect.o
91
92 tmp.s: boot/bootsect.s tools/system # 从 92--95 这四行的作用是在 bootsect.s 程序开头添加
# 一行有关 system 文件长度信息。方法是首先生成含有 "SYSSIZE = system 文件实际长度"
# 一行信息的 tmp.s 文件，然后将 bootsect.s 文件添加在其后。取得 system 长度的方法是：
# 首先利用命令 ls 对 system 文件进行长列表显示，用 grep 命令取得列表行上文件字节数字段
# 信息，并定向保存在 tmp.s 临时文件中。cut 命令用于剪切字符串，tr 用于去除行尾的回车符。
# 其中：(实际长度 + 15)/16 用于获得用 '节' 表示的长度信息。1 节 = 16 字节。
93      (echo -n "SYSSIZE = (" ; ls -l tools/system | grep system \
94      | cut -c25-31 | tr '\012' ' '; echo "+ 15 ) / 16") > tmp.s
95      cat boot/bootsect.s >> tmp.s
96
97 clean: # 当执行 'make clean' 时，就会执行 98--103 行上的命令，去除所有编译连接生成的文件。
# 'rm' 是文件删除命令，选项 -f 含义是忽略不存在的文件，并且不显示删除信息。
98      rm -f Image System.map tmp_make core boot/bootsect boot/setup
99      rm -f init/*.o tools/system tools/build boot/*.o
100     (cd mm; make clean) # 进入 mm/目录；执行该目录 Makefile 文件中的 clean 规则。
101     (cd fs; make clean)
102     (cd kernel; make clean)
103     (cd lib; make clean)
104
105 backup: clean # 该规则将首先执行上面的 clean 规则，然后对 linux/目录进行压缩，生成
# backup.Z 压缩文件。'cd ..' 表示退到 linux/的上一级（父）目录；
# 'tar cf - linux' 表示对 linux/目录执行 tar 归档程序。-cf 表示需要创建
# 新的归档文件 '| compress -' 表示将 tar 程序的执行通过管道操作('|')
# 传递给压缩程序 compress，并将压缩程序的输出存成 backup.Z 文件。
106     (cd .. ; tar cf - linux | compress - > backup.Z)
107     sync # 迫使缓冲块数据立即写盘并更新磁盘超级块。
108
109 dep: # 该目标或规则用于各文件之间的依赖关系。方法如下：
# 使用字符串编辑程序 sed 对 Makefile 文件（这里即是自己）进行处理，输出为删除 Makefile
# 文件中 '### Dependencies' 行后面的所有行（下面从 118 开始的行），并生成 tmp_make
# 临时文件（也即 110 行的作用）。然后对 init/目录下的每一个 C 文件（其实只有一个文件

```

```
# main.c) 执行 gcc 预处理操作, -M 标志告诉预处理程序输出描述每个目标文件相关性的规则,
# 并且这些规则符合 make 语法。对于每一个源文件, 预处理程序输出一个 make 规则, 其结果
# 形式是相应源程序文件的目标文件名加上其依赖关系--该源文件中包含的所有头文件列表。
# 把预处理结果都添加到临时文件 tmp_make 中, 然后将该临时文件复制成新的 Makefile 文件。
```

```
110 sed '/\#\#\# Dependencies/q' < Makefile > tmp_make
111 (for i in init/*.c;do echo -n "init/";$(CPP) -M $$i;done) >> tmp_make
112 cp tmp_make Makefile
113 (cd fs; make dep)          # 对 fs/目录下的 Makefile 文件也作同样的处理。
114 (cd kernel; make dep)
115 (cd mm; make dep)
116
117 ### Dependencies:
118 init/main.o : init/main.c include/unistd.h include/sys/stat.h \
119 include/sys/types.h include/sys/times.h include/sys/utsname.h \
120 include/utime.h include/time.h include/linux/tty.h include/termios.h \
121 include/linux/sched.h include/linux/head.h include/linux/fs.h \
122 include/linux/mm.h include/signal.h include/asm/system.h include/asm/io.h \
123 include/stddef.h include/stdarg.h include/fcntl.h
```

## 1.5 其它信息

### 1.5.1 Makefile 简介

你需要一个称为 makefile 的文件来告诉 make 要做什么。通常, makefile 文件会告诉 make 如何编译和连接一个文件。这里我们将讨论如何编译和连接一个含有 8 个 C 程序文件和 3 个头文件(head file)的字编辑程序(text editor)。当明确指出时, makefile 还可以告诉 make 运行各种命令(例如, 作为清理操作而删除某些文件)。

当 make 重新编译 editor 时, 每个修改过的 C 代码文件必须被重新编译。如果一个头文件被修改过了, 那么为了确保正确, 每一个包含该头文件的 C 代码程序都将被重新编译。每次编译操作都产生一个与源程序对应的目标文件(object file)。最终, 如果任何源代码文件被编译过了, 那么所有的目标文件不管是刚编译完的还是以前就编译好的必须连接在一起以生成新的可执行文件 editor。

简单的 makefile 文件含有一些规则, 这些规则具有如下的形式:

```
目标(target)... : 先决条件(prerequisites)...
                  命令(command)
                  ...
                  ...
```

目标通常是程序生成的一个文件的名称; 例如可执行文件或目标文件。目标也可以是所要采取活动的名字, 比如'清除'('clean')。先决条件是一个文件, 用作产生目标的输入。通常一个目标依赖几个文件。命令是 make 需要执行的操作。一个规则可以有多个命令, 每一个命令自成一。请注意, 你需要在每个命令之前键入一个制表符! 这是粗心者常常忽略的地方。

通常命令是属于一个具有先决条件的规则, 并在任何先决条件改变时用于生成一个目标(target)文件。然而, 为目标而指定命令的规则并不一定要有先决条件的, 例如, 与目标'clean'相关的含有删除

(delete) 命令的规则并不需要有先决条件。此时，一个规则说明了如何以及何时来重新制作某些文件，而这些文件是特定规则的目标。make 根据先决条件来执行命令以创建或更新目标。一个规则也可以说明如何及何时执行一个操作。

一个 makefile 文件也可以含有除规则以外的其它文字，但一个简单的 makefile 文件只需要含有适当的规则。规则可能看上去要比上面示出的模板复杂得多，但基本上都是符合的。

### 1.5.2 as86,ld86 简介

as86 和 ld86 是由 Bruce Evans 编写的 Intel 8086 汇编编译程序和连接程序。它完全是一个 8086 的汇编编译器，但却可以为 386 处理器编制 32 位的代码。Linux 使用它仅仅是为了创建 16 位的启动扇区 (bootsector) 代码和 setup 二进制执行代码。该编译器的语法与 GNU 的汇编编译器的语法是不兼容的，但近似于 Intel 的汇编语言语法（如操作数的次序相反等）。

Bruce Evans 是 minix 操作系统 32 位版本的主要编制者，他与 Linux 的创始人 Linus Torvalds 是很好的朋友。Linus 本人也从 Bruce Evans 那里学到了不少有关 UNIX 类操作系统的知识，minix 操作系统的不足之处也是两个好朋友互相探讨得出的结果，这激发了 Linus 在 Intel 386 体系结构上开发一个全新概念的操作系统，因此 Linux 操作系统的诞生与 Bruce Evans 也有着密切的关系。

有关这个编译器和连接器的源代码可以从 FTP 服务器 ftp.funet.fi 上或从我的网站 (www.oldlinux.org) 上下载。

这两个程序的使用方法和选项如下：

as 的使用方法和选项：

```
as [-O3agjuw] [-b [bin]] [-lm [list]] [-n name] [-o obj] [-s sym] src
```

默认设置（除了以下默认值以外，其它选项默认为关闭或无；若没有明确说明 a 标志，则不会有输出）：

```
-O3      32 位输出；
list     在标准输出上显示；
name     源文件的基本名称（也即不包括“.”后的扩展名）；
```

选项含义：

```
-O      从 16 比特代码段开始；
-3      从 32 比特代码段开始；
-a      开启与 as、ld 的部分兼容性选项；
-b      产生二进制文件，后面可以跟文件名；
-g      在目标文件中仅存入全局符号；
-j      使所有跳转语句均为长跳转；
-l      产生列表文件，后面可以跟随列表文件名；
-m      在列表中扩展宏定义；
-n      后面跟随模块名称（取代源文件名称放入目标文件中）；
-o      产生目标文件，后跟目标文件名；
-s      产生符号文件，后跟符号文件名；
-u      将未定义符号作为输入的未指定段的符号；
-w      不显示警告信息；
```

ld 连接器的使用语法和选项：

对于生成 Minix a.out 格式的版本：

```
ld [-O3Mims[-]] [-T textaddr] [-llib_extension] [-o outfile] infile...
```

对于生成 GNU-Minix 的 a.out 格式的版本:

```
ld [-03Mimrs[-]] [-T textaddr] [-llib_extension] [-o outfile] infile...
```

默认设置(除了以下默认值以外, 其它选项默认为关闭或无):

-03 32 位输出;

outfile a.out 格式输出;

- 0 产生具有 16 比特魔数的头结构, 并且对-lx 选项使用 i86 子目录;
- 3 产生具有 32 比特魔数的头结构, 并且对-lx 选项使用 i386 子目录;
- M 在标准输出设备上显示已链接的符号;
- T 后面跟随文本基地址 (使用适合于 strtoul 的格式);
- i 分离的指令与数据段 (I&D) 输出;
- lx 将库/local/lib/subdir/libx.a 加入链接的文件列表中;
- m 在标准输出设备上显示已链接的模块;
- o 指定输出文件名, 后跟输出文件名;
- r 产生适合于进一步重定位的输出;
- s 在目标文件中删除所有符号。

### 1.5.3 System.map 文件

System.map 文件用于存放内核符号表信息。符号表是所有符号及其对应地址的一个列表。随着每次内核的编译, 就会产生一个新的对应 System.map 文件。当内核运行出错时, 通过 System.map 文件中的符号表解析, 就可以查到一个地址值对应的变量名, 或反之。

利用 System.map 符号表文件, 在内核或相关程序出错时, 就可以获得我们比较容易识别的信息。符号表的样例如下所示:

```
c03441a0 B dmi_broken
c03441a4 B is_sony_vaio_laptop
c03441c0 b dmi_ident
c0344200 b pci_bios_present
c0344204 b pirq_table
```

可以看出名称为 dmi\_broken 的变量位于内核地址 c03441a0 处。

System.map 位于使用它的软件(例如内核日志记录后台程序 klogd)能够寻找到的地方。在系统启动时, 如果没有以一个参数的形式为 klogd 给出 System.map 的位置, 则 klogd 将会在三个地方搜寻 System.map。依次为:

1. /boot/System.map
2. /System.map
3. /usr/src/linux/System.map

尽管内核本身实际上不使用 System.map, 但其它程序, 象 klogd, lsof, ps 以及其它许多软件, 象 dosemu, 都需要有一个正确的 System.map 文件。



## 第2章 引导启动程序 (boot)

这里先总的说明一下 Linux 操作系统启动部分的主要执行流程。当 PC 的电源打开后, 80x86 结构的 CPU 将自动进入实模式, 并从地址 0xFFFF0 开始自动执行程序代码, 这个地址通常是 ROM-BIOS 中的地址。PC 机的 BIOS 将执行某些系统的检测, 并在物理地址 0 处开始初始化中断向量。此后, 它将可启动设备的第一个扇区 (磁盘引导扇区, 512 字节) 读入内存地址 0x7C00 处, 并跳转到这个地方。启动设备通常是软驱或是硬盘。这里的叙述是非常简单的, 但这已经足够理解内核初始化的工作过程了。

Linux 的最前面部分是用 8086 汇编语言编写的 (boot/bootsect.s), 它将由 BIOS 读入到内存 0x7C00 处, 当它被执行时就会把自己移到绝对地址 0x90000 处, 并将启动设备 (boot/setup.s) 的下 2kB 字节的代码读入内存 0x90200 处, 而内核的其它部分 (system 模块) 则被读入到从地址 0x10000 开始处, 因为当时 system 模块的长度不会超过 0x80000 (即 512KB), 所以它不会覆盖在 0x90000 处开始的 bootsect 和 setup 模块。图 2.1 清晰地显示出 Linux 系统启动时这几个程序或模块在内存中的动态位置。其中, 每一竖条框代表某一时刻内存中各程序的映像位置图。在系统加载期间将显示信息 "Loading..."。然后控制权将传递给 boot/setup.s 中的代码, 这是另一个实模式汇编语言程序。

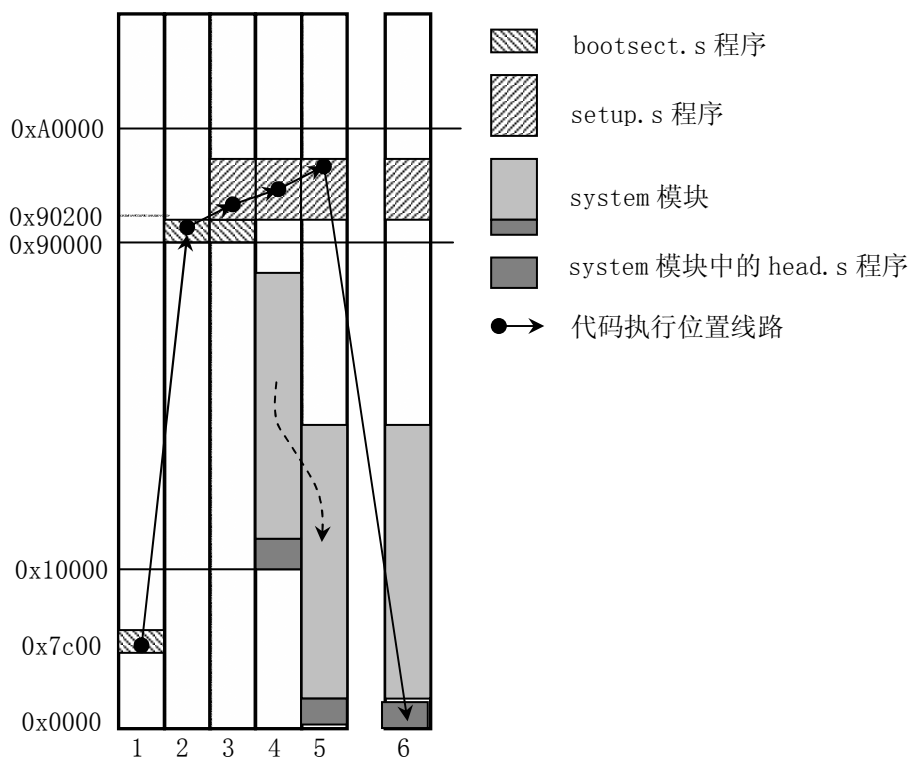


图 2.1 启动引导时内核在内存中的位置和移动后的位置情况




启动部分识别主机的某些特性以及 vga 卡的类型。如果需要, 它会要求用户为控制台选择显示模式。然后将整个系统从地址 0x10000 移至 0x0000 处, 进入保护模式并跳转至系统的余下部分 (在 0x0000 处)。此时所有 32 位运行方式的设置启动被完成: IDT、GDT 以及 LDT 被加载, 处理器和协处理器也已确认, 分页工作也设置好了; 最终调用 init/main.c 中的 main() 程序。上述操作的源代码是在 boot/head.S 中的, 这可能是整个内核中最有诀窍的代码了。注意如果在前述任何一步中出了错, 计算机就会死锁。在操作系

统还没有完全运转之前是处理不了出错的。

## 2.1 概述

本章主要描述 boot/目录中的三个汇编代码文件，见列表 2.1 所示。正如在前一章中提到的，这三个文件虽然都是汇编程序，但却使用了两种语法格式。Bootsect.s 和 setup.s 采用近似于 Intel 的汇编语言语法，需要使用 Intel 8086 汇编编译器和连接器 as86 和 ld86，而 head.s 则使用 GNU 的汇编程序格式，需要用 GNU 的 as 进行编译。这是一种 AT&T 语法的汇编语言程序。

列表 2.1 linux/boot/目录

| 文件名  | 长度(字节)     | 最后修改时间 (GMT)        | 说明 |
|--|------------|---------------------|----|
|  <a href="#">bootsect.s</a> | 5052 bytes | 1991-12-05 22:47:58 |    |
|  <a href="#">head.s</a>     | 5938 bytes | 1991-11-18 15:05:09 |    |
|  <a href="#">setup.s</a>    | 5364 bytes | 1991-12-05 22:48:10 |    |

阅读这些代码除了你需要知道一些一般 8086 汇编语言的知识以外，还要对采用 Intel 80X86 微处理器的 PC 机的体系结构以及 80386 32 位保护模式下的编程原理有些了解。所以在开始阅读源代码之前可以先大概浏览一下附录中有关 PC 机硬件接口控制编程和 80386 32 位保护模式的编程方法，在阅读代码时再就事论事地针对具体问题参考附录中的详细说明。

## 2.2 bootsect.s 文件

### 2.2.1 功能描述

bootsect.s 代码是磁盘引导块程序，驻留在磁盘的第一个扇区中（引导扇区，0 磁道（柱面），0 磁头，第 1 个扇区）。在 PC 机加电 ROM BIOS 自检后，引导扇区由 BIOS 加载到内存 0x7C00 处，然后将自己移动到内存 0x90000 处。该程序的主要作用是首先将 setup 模块（由 setup.s 编译成）从磁盘加载到内存紧接着 bootsect 的后面位置（0x90200），然后利用 BIOS 中断 0x13 取磁盘参数表中当前启动引导盘的参数，接着在屏幕上显示“Loading system...”字符串。再者将 system 模块从磁盘上加载到内存 0x10000 开始的地方。随后确定根文件系统的设备号，若没有指定，则根据所保存的引导盘的每磁道扇区数判别出盘的类型和种类（是 1.44M A 盘？）并保存其设备号于 root\_dev（引导块的 0x508 地址处），最后长跳转到 setup 程序的开始处（0x90200）执行 setup 程序。

### 2.2.2 代码注释

列表 2.2 linux/boot/bootsect.s 程序

1 !

2 ! SYS\_SIZE is the number of clicks (16 bytes) to be loaded.



```

3 ! 0x3000 is 0x30000 bytes = 196kB, more than enough for current
4 ! versions of linux ! SYS_SIZE 是要加载的节数(16字节为1节)。0x3000 共为
! 0x30000 字节 = 196 kB, 对于当前版本的 linux 空间已足够了。
5 !
6 SYSSIZE = 0x3000 ! 指编译连接后 system 模块的大小。参见列表 1.2 中第 92 的说明。
! 这里给出了一个最大默认值。
7 !
8 ! bootsect.s (C) 1991 Linus Torvalds
9 !
10 ! bootsect.s is loaded at 0x7c00 by the bios-startup routines, and moves
11 ! itself out of the way to address 0x90000, and jumps there.
12 !
13 ! It then loads 'setup' directly after itself (0x90200), and the system
14 ! at 0x10000, using BIOS interrupts.
15 !
16 ! NOTE! currently system is at most 8*65536 bytes long. This should be no
17 ! problem, even in the future. I want to keep it simple. This 512 kB
18 ! kernel size should be enough, especially as this doesn't contain the
19 ! buffer cache as in minix
20 !
21 ! The loader has been made as simple as possible, and continuous
22 ! read errors will result in a unbreakable loop. Reboot by hand. It
23 ! loads pretty fast by getting whole sectors at a time whenever possible.
!
! 以下是前面这些文字的翻译:
! bootsect.s (C) 1991 Linus Torvalds 版权所有
!
! bootsect.s 被 bios-启动子程序加载至 0x7c00 (31k)处, 并将自己
! 移到了地址 0x90000 (576k)处, 并跳转至那里。
!
! 它然后使用 BIOS 中断将'setup'直接加载到自己的后面(0x90200) (576.5k),
! 并将 system 加载到地址 0x10000 处。
!
! 注意! 目前的内核系统最大长度限制为(8*65536) (512k)字节, 即使是在
! 将来这也应该没有问题的。我想让它保持简单明了。这样 512k 的最大内核长度应该
! 足够了, 尤其是这里没有象 minix 中一样包含缓冲区高速缓冲。
!
! 加载程序已经做的够简单了, 所以持续的读出错将导致死循环。只能手工重启。
! 只要可能, 通过一次取取所有的扇区, 加载过程可以做的很快的。
24
25 .globl begtext, begdata, begbss, endtext, enddata, endbss ! 定义了6个全局标识符;
26 .text ! 文本段;
27 begtext:
28 .data ! 数据段;
29 begdata:
30 .bss ! 堆栈段;
31 begbss:
32 .text ! 文本段;
33
34 SETUPLEN = 4 ! nr of setup-sectors
! setup 程序的扇区数(setup-sectors)值;
35 BOOTSEG = 0x07c0 ! original address of boot-sector
! bootsect 的原始地址(是段地址, 以下同);

```

```

36 INITSEG = 0x9000          ! we move boot here - out of the way
                             ! 将 bootsect 移到这里 -- 避开;
37 SETUPSEG = 0x9020        ! setup starts here
                             ! setup 程序从这里开始;
38 SYSSEG = 0x1000          ! system loaded at 0x10000 (65536).
                             ! system 模块加载到 0x10000 (64 kB) 处;
39 ENDSEG = SYSSEG + SYSSIZE ! where to stop loading
                             ! 停止加载的段地址;

40
41 ! ROOT_DEV:      0x000 - same type of floppy as boot.
!                      根文件系统设备使用与引导时同样的软驱设备;
42 !              0x301 - first partition on first drive etc
!                      根文件系统设备在第一个硬盘的第一个分区上, 等等;
43 ROOT_DEV = 0x306 ! 指定根文件系统设备是第 2 个硬盘的第 1 个分区。这是 Linux 老式的硬盘命名
! 方式, 具体值的含义如下:
! 设备号=主设备号*256 + 次设备号 (也即 dev_no = (major<<8) + minor )
! (主设备号: 1-内存, 2-磁盘, 3-硬盘, 4-ttyx, 5-tty, 6-并行口, 7-非命名管道)
! 0x300 - /dev/hd0 - 代表整个第 1 个硬盘;
! 0x301 - /dev/hd1 - 第 1 个盘的第 1 个分区;
! ...
! 0x304 - /dev/hd4 - 第 1 个盘的第 4 个分区;
! 0x305 - /dev/hd5 - 代表整个第 2 个硬盘盘;
! 0x306 - /dev/hd6 - 第 2 个盘的第 1 个分区;
! ...
! 0x309 - /dev/hd9 - 第 2 个盘的第 4 个分区;
! 从 linux 内核 0.95 版后已经使用与现在相同的命名方法了。

44
45 entry start      ! 告知连接程序, 程序从 start 标号开始执行。
46 start:          ! 47--56 行作用是将自身 (bootsect) 从目前段位置 0x07c0 (31k)
                  ! 移动到 0x9000 (576k) 处, 共 256 字 (512 字节), 然后跳转到
                  ! 移动后代码的 go 标号处, 也即本程序的下一语句处。

47      mov      ax, #BOOTSEG ! 将 ds 段寄存器置为 0x7C0;
48      mov      ds, ax
49      mov      ax, #INITSEG ! 将 es 段寄存器置为 0x9000;
50      mov      es, ax
51      mov      cx, #256     ! 移动计数值=256 字;
52      sub      si, si       ! 源地址 ds:si = 0x07C0:0x0000
53      sub      di, di       ! 目的地址 es:di = 0x9000:0x0000
54      rep      ! 重复执行, 直到 cx = 0
55      movw     ! 移动 1 个字;
56      jmp     go, INITSEG   ! 间接跳转。这里 INITSEG 指出跳转到的段地址。
57 go:      mov      ax, cs    ! 将 ds、es 和 ss 都置成移动后代码所在的段处 (0x9000)。
58      mov      ds, ax
59      mov      es, ax
60 ! put stack at 0x9ff00.    ! 将堆栈指针 sp 指向 0x9ff00 (即 0x9000:0xff00) 处
61      mov      ss, ax
62      mov      sp, #0xFF00  ! arbitrary value >>512
                             ! sp 只要指向远大于 512 偏移 (即地址 0x90400) 处
                             ! 都可以。因为从 0x90200 地址开始要放置 setup 程序,
                             ! 而此时 setup 程序大约为 4 个扇区, 因此 sp 要指向大
                             ! 于 (200 + 200 * 4 + 堆栈大小) 处。

63
64 ! load the setup-sectors directly after the bootblock.

```

```

65 ! Note that 'es' is already set up.
! 在 bootsect 程序块后紧跟着加载 setup 模块的代码数据。
! 注意 es 已经设置好了。(在移动代码时 es 已经指向目的段地址处 0x9000)。

66
67 load_setup:
! 68--77 行的用途是利用 BIOS 中断 INT 0x13 将 setup 模块从磁盘第 2 个扇区
! 开始读到 0x90200 开始处, 共读 4 个扇区。如果读出错, 则复位驱动器, 并
! 重试, 没有退路。INT 0x13 的使用方法如下:
! 读扇区:
! ah = 0x02 - 读磁盘扇区到内存; al = 需要读出的扇区数量;
! ch = 磁道(柱面)号的低 8 位;    cl = 开始扇区(0-5 位), 磁道号高 2 位(6-7);
! dh = 磁头号;                    dl = 驱动器号(如果是硬盘则要置位 7);
! es:bx → 指向数据缓冲区; 如果出错则 CF 标志置位。
68     mov     dx, #0x0000          ! drive 0, head 0
69     mov     cx, #0x0002          ! sector 2, track 0
70     mov     bx, #0x0200          ! address = 512, in INITSEG
71     mov     ax, #0x0200+SETUPLEN ! service 2, nr of sectors
72     int     0x13                 ! read it
73     jnc     ok_load_setup        ! ok - continue
74     mov     dx, #0x0000
75     mov     ax, #0x0000          ! reset the diskette
76     int     0x13
77     j       load_setup
78
79 ok_load_setup:
80
81 ! Get disk drive parameters, specifically nr of sectors/track
! 取磁盘驱动器的参数, 特别是每道的扇区数量。
! 取磁盘驱动器参数 INT 0x13 调用格式和返回信息如下:
! ah = 0x08    dl = 驱动器号(如果是硬盘则要置位 7 为 1)。
! 返回信息:
! 如果出错则 CF 置位, 并且 ah = 状态码。
! ah = 0,  al = 0,      bl = 驱动器类型(AT/PS2)
! ch = 最大磁道号的低 8 位, cl = 每磁道最大扇区数(位 0-5), 最大磁道号高 2 位(位 6-7)
! dh = 最大磁头号,      dl = 驱动器数量,
! es:di → 软驱磁盘参数表。
82
83     mov     dl, #0x00
84     mov     ax, #0x0800          ! AH=8 is get drive parameters
85     int     0x13
86     mov     ch, #0x00
87     seg     cs
88     mov     sectors, cx          ! 保存每磁道扇区数。
89     mov     ax, #INITSEG
90     mov     es, ax              ! 因为上面取磁盘参数中断改掉了 es 的值, 这里重新改回。
91
92 ! Print some inane message ! 在显示一些信息('Loading system ...' 回车换行, 共 24 个字符)。
93
94     mov     ah, #0x03            ! read cursor pos
95     xor     bh, bh               ! 读光标位置。
96     int     0x10
97
98     mov     cx, #24              ! 共 24 个字符。

```

```

99      mov     bx, #0x0007          ! page 0, attribute 7 (normal)
100     mov     bp, #msg1            ! 指向要显示的字符串。
101     mov     ax, #0x1301          ! write string, move cursor
102     int     0x10                 ! 写字符串并移动光标。
103
104 ! ok, we've written the message, now
105 ! we want to load the system (at 0x10000) ! 现在开始将 system 模块加载到 0x10000 (64k) 处。
106
107     mov     ax, #SYSSEG
108     mov     es, ax                ! segment of 0x010000 ! es = 存放 system 的段地址。
109     call    read_it              ! 读磁盘上 system 模块, es 为输入参数。
110     call    kill_motor          ! 关闭驱动器马达, 这样就可以知道驱动器的状态了。
111
112 ! After that we check which root-device to use. If the device is
113 ! defined (!= 0), nothing is done and the given device is used.
114 ! Otherwise, either /dev/PS0 (2, 28) or /dev/at0 (2, 8), depending
115 ! on the number of sectors that the BIOS reports currently.
! 此后, 我们检查要使用哪个根文件系统设备 (简称根设备)。如果已经指定了设备 (!=0)
! 就直接使用给定的设备。否则就需要根据 BIOS 报告的每磁道扇区数来
! 确定到底使用 /dev/PS0 (2, 28) 还是 /dev/at0 (2, 8)。
! 上面一行中两个设备文件的含义:
! 在 Linux 中软驱的主设备号是 2 (参见第 43 行的注释), 次设备号 = type*4 + nr, 其中
! nr 为 0-3 分别对应软驱 A、B、C 或 D; type 是软驱的类型 (2→1.2M 或 7→1.44M 等)。
! 因为 7*4 + 0 = 28, 所以 /dev/PS0 (2, 28) 指的是 1.44M A 驱动器, 其设备号是 0x021c
! 同理 /dev/at0 (2, 8) 指的是 1.2M A 驱动器, 其设备号是 0x0208。
116
117     seg cs
118     mov     ax, root_dev          ! 将根设备号
119     cmp     ax, #0
120     jne     root_defined
121     seg cs
122     mov     bx, sectors          ! 取上面第 88 行保存的每磁道扇区数。如果 sectors=15
! 则说明是 1.2Mb 的驱动器; 如果 sectors=18, 则说明是
! 1.44Mb 软驱。因为是可引导的驱动器, 所以肯定是 A 驱。
123     mov     ax, #0x0208          ! /dev/ps0 - 1.2Mb
124     cmp     bx, #15              ! 判断每磁道扇区数是否=15
125     je      root_defined        ! 如果等于, 则 ax 中就是引导驱动器的设备号。
126     mov     ax, #0x021c          ! /dev/PS0 - 1.44Mb
127     cmp     bx, #18
128     je      root_defined
129 undef_root:                      ! 如果都不一样, 则死循环 (死机)。
130     jmp     undef_root
131 root_defined:
132     seg cs
133     mov     root_dev, ax         ! 将检查过的设备号保存起来。
134
135 ! after that (everything loaded), we jump to
136 ! the setup-routine loaded directly after
137 ! the bootblock:
! 到此, 所有程序都加载完毕, 我们就跳转到被
! 加载在 bootsect 后面的 setup 程序去。
138
139     jmp     0, SETUPSEG          ! 跳转到 0x9020:0000 (setup.s 程序的开始处)。

```

!!!! 本程序到此就结束了。!!!!

! 下面是两个子程序。

```

140
141 ! This routine loads the system at address 0x10000, making sure
142 ! no 64kB boundaries are crossed. We try to load it as fast as
143 ! possible, loading whole tracks whenever we can.
144 !
145 ! in:  es - starting address segment (normally 0x1000)
146 !
! 该子程序将系统模块加载到内存地址 0x10000 处，并确定没有跨越 64KB 的内存边界。我们试图尽快
! 地进行加载，只要可能，就每次加载整条磁道的数据。
! 输入: es - 开始内存地址段值 (通常是 0x1000)
147 sread: .word 1+SETUPLEN          ! sectors read of current track
! 当前磁道中已读的扇区数。开始时已经读进 1 扇区的引导扇区
! bootsect 和 setup 程序所占的扇区数 SETUPLEN。
148 head:  .word 0                  ! current head    !当前磁头号。
149 track:  .word 0                  ! current track  !当前磁道号。
150
151 read_it:
! 测试输入的段值。必须位于内存地址 64KB 边界处，否则进入死循环。清 bx 寄存器，用于表示当前段内
! 存放数据的开始位置。
152     mov ax, es
153     test ax, #0x0fff
154 die:   jne die                   ! es must be at 64kB boundary ! es 值必须位于 64KB 地址边界!
155     xor bx, bx                   ! bx is starting address within segment ! bx 为段内偏移位置。
156 rp_read:
! 判断是否已经读入全部数据。比较当前所读段是否就是系统数据末端所处的段 (#ENDSEG)，如果不是就
! 跳转至下面 ok1_read 标号处继续读数据。否则退出子程序返回。
157     mov ax, es
158     cmp ax, #ENDSEG             ! have we loaded all yet? ! 是否已经加载了全部数据?
159     jb ok1_read
160     ret
161 ok1_read:
! 计算和验证当前磁道需要读取的扇区数，放在 ax 寄存器中。
! 根据当前磁道还未读取的扇区数以及段内数据字节开始偏移位置，计算如果全部读取这些未读扇区，所
! 读总字节数是否会超过 64KB 段长度的限制。若会超过，则根据此次最多能读入的字节数 (64KB - 段内
! 偏移位置)，反算出此次需要读取的扇区数。
162     seg cs
163     mov ax, sectors             ! 取每磁道扇区数。
164     sub ax, sread               ! 减去当前磁道已读扇区数。
165     mov cx, ax                  ! cx = ax = 当前磁道未读扇区数。
166     shl cx, #9                  ! cx = cx * 512 字节。
167     add cx, bx                  ! cx = cx + 段内当前偏移值 (bx)
! = 此次读操作后，段内共读入的字节数。
168     jnc ok2_read               ! 若没有超过 64KB 字节，则跳转至 ok2_read 处执行。
169     je ok2_read
170     xor ax, ax                  ! 若加上此次将读磁道上所有未读扇区时会超过 64KB，则计算
171     sub ax, bx                  ! 此时最多能读入的字节数 (64KB - 段内读偏移位置)，再转换
172     shr ax, #9                  ! 成需要读取的扇区数。
173 ok2_read:
174     call read_track
175     mov cx, ax                  ! cx = 该次操作已读取的扇区数。
176     add ax, sread               ! 当前磁道上已经读取的扇区数。

```

```

177      seg cs
178      cmp ax, sectors      ! 如果当前磁道上的还有扇区未读, 则跳转到 ok3_read 处。
179      jne ok3_read
! 读该磁道的下一磁头面(1 号磁头)上的数据。如果已经完成, 则去读下一磁道。
180      mov ax, #1
181      sub ax, head        ! 判断当前磁头号。
182      jne ok4_read        ! 如果是 0 磁头, 则再去读 1 磁头面上的扇区数据。
183      inc track           ! 否则去读下一磁道。
184 ok4_read:
185      mov head, ax        ! 保存当前磁头号。
186      xor ax, ax          ! 清当前磁道已读扇区数。
187 ok3_read:
188      mov sread, ax       ! 保存当前磁道已读扇区数。
189      shl cx, #9          ! 上次已读扇区数*512 字节。
190      add bx, cx          ! 调整当前段内数据开始位置。
191      jnc rp_read         ! 若小于 64KB 边界值, 则跳转到 rp_read(156 行)处, 继续读数据。
! 否则调整当前段, 为读下一段数据作准备。
192      mov ax, es
193      add ax, #0x1000     ! 将段基址调整为指向下一个 64KB 段内存。
194      mov es, ax
195      xor bx, bx          ! 清段内数据开始偏移值。
196      jmp rp_read        ! 跳转至 rp_read(156 行)处, 继续读数据。
197
! 读当前磁道上指定开始扇区和需读扇区数的数据到 es:bx 开始处。参见第 67 行下对 BIOS 磁盘读中断
! int 0x13, ah=2 的说明。
! al - 需读扇区数; es:bx - 缓冲区开始位置。
198 read_track:
199      push ax
200      push bx
201      push cx
202      push dx
203      mov dx, track       ! 取当前磁道号。
204      mov cx, sread       ! 取当前磁道上已读扇区数。
205      inc cx              ! cl = 开始读扇区。
206      mov ch, dl          ! ch = 当前磁道号。
207      mov dx, head        ! 取当前磁头号。
208      mov dh, dl          ! dh = 磁头号。
209      mov dl, #0          ! dl = 驱动器号(为 0 表示当前驱动器)。
210      and dx, #0x0100     ! 磁头号不大于 1。
211      mov ah, #2          ! ah = 2, 读磁盘扇区功能号。
212      int 0x13
213      jc bad_rt           ! 若出错, 则跳转至 bad_rt。
214      pop dx
215      pop cx
216      pop bx
217      pop ax
218      ret
! 执行驱动器复位操作(磁盘中断功能号 0), 再跳转到 read_track 处重试。
219 bad_rt: mov ax, #0
220      mov dx, #0
221      int 0x13
222      pop dx
223      pop cx

```

```

224         pop bx
225         pop ax
226         jmp read_track
227
228 /*
229  * This procedure turns off the floppy drive motor, so
230  * that we enter the kernel in a known state, and
231  * don't have to worry about it later.
232  */
    ! 这个子程序用于关闭软驱的马达，这样我们进入内核后它处于已知状态，以后也就无须担心它了。
233 kill_motor:
234         push dx
235         mov dx, #0x3f2      ! 软驱控制卡的驱动端口，只写。
236         mov al, #0         ! A 驱动器，关闭 FDC，禁止 DMA 和中断请求，关闭马达。
237         outb               ! 将 al 中的内容输出到 dx 指定的端口去。
238         pop dx
239         ret
240
241 sectors:
242         .word 0            ! 存放当前启动软盘每磁道的扇区数。
243
244 msg1:
245         .byte 13,10        ! 回车、换行的 ASCII 码。
246         .ascii "Loading system ..."
247         .byte 13,10,13,10 ! 共 24 个 ASCII 码字符。
248
249 .org 508                  ! 表示下面语句从地址 508(0x1FC) 开始，所以 root_dev
                           ! 在启动扇区的第 508 开始的 2 个字节中。
250 root_dev:
251         .word ROOT_DEV    ! 这里存放根文件系统所在的设备号(init/main.c 中会用)。
252 boot_flag:
253         .word 0xAA55      ! 硬盘有效标识。
254
255 .text
256 endtext:
257 .data
258 enddata:
259 .bss
260 endbss:

```

## 2.3 setup.s 程序

### 2.3.1 功能描述

setup 程序的作用主要是利用 ROM BIOS 中断读取机器系统数据，并将这些数据保存到 0x90000 开始的位置（覆盖掉了 bootsect 程序所在的地方），所取得的参数和保留的内存位置见下表 2.1 所示。这些参数将被内核中相关程序使用，例如 ttyio.c 程序等。

表 2.1 setup 程序读取并保留的参数



| 内存地址    | 长度(字节) | 名称    | 描述   |
|---------|--------|-------|--|
| 0x90000 | 2      | 光标位置  | 列号 (0x00-最左端), 行号 (0x00-最顶端)                     |
| 0x90002 | 2      | 扩展内存数 | 系统从 1M 开始的扩展内存数值 (KB)。                           |
| 0x90004 | 2      | 显示页面  | 当前显示页面   |
| 0x90006 | 1      | 显示模式  |  |
| 0x90007 | 1      | 字符列数  |  |
| 0x90008 | 2      | ??    |  |
| 0x9000A | 1      | 显示内存  | 显示内存 (0x00-64k, 0x01-128k, 0x02-192k, 0x03=256k) |
| 0x9000B | 1      | 显示状态  | 0x00-彩色, I/O=0x3dX; 0x11-单色, I/O=0x3bX           |
| 0x9000C | 2      | 特性参数  | 显示卡特性参数  |
| ...     |        |       |  |
| 0x90080 | 16     | 硬盘参数表 | 第 1 个硬盘的参数表                                      |
| 0x90090 | 16     | 硬盘参数表 | 第 2 个硬盘的参数表 (如果没有, 则清零)                          |
| 0x901FC | 2      | 根设备号  | 根文件系统所在的设备号 (bootsec.s 中设置)                      |

然后 setup 程序将 system 模块从 0x10000-0x8ffff (当时认为内核系统模块 system 的长度不会超过此值: 512KB) 整块向下移动到内存绝对地址 0x00000 处。接下来加载中断描述符表寄存器(idtr)和全局描述符表寄存器(gdtr), 开启 A20 地址线, 重新设置两个中断控制芯片 8259A, 将硬件中断号重新设置为 0x20 - 0x2f。最后设置 CPU 的控制寄存器 CR0 (也称机器状态字), 从而进入 32 位保护模式运行, 并跳转到位于 system 模块最前面部分的 head.s 程序继续运行。

## 2.3.2 代码注释

列表 2.3 linux/boot/setup.s 程序

```

1 !
2 !      setup.s      (C) 1991 Linus Torvalds
3 !
4 ! setup.s is responsible for getting the system data from the BIOS,
5 ! and putting them into the appropriate places in system memory.
6 ! both setup.s and system has been loaded by the bootblock.
7 !
8 ! This code asks the bios for memory/disk/other parameters, and
9 ! puts them in a "safe" place: 0x90000-0x901FF, ie where the
10 ! boot-block used to be. It is then up to the protected mode
11 ! system to read them from there before the area is overwritten
12 ! for buffer-blocks.
13 !
14 ! setup.s 负责从 BIOS 中获取系统数据, 并将这些数据放到系统内存的适当地方。
15 ! 此时 setup.s 和 system 已经由 bootsect 引导块加载到内存中。
16 !
17 ! 这段代码询问 bios 有关内存/磁盘/其它参数, 并将这些参数放到一个
18 ! "安全的" 地方: 0x90000-0x901FF, 也即原来 bootsect 代码块曾经在
19 ! 的地方, 然后在被缓冲块覆盖掉之前由保护模式的 system 读取。
20 !
21 ! NOTE! These had better be the same as in bootsect.s!
22 ! 以下这些参数最好和 bootsect.s 中的相同!
23 !
24 !

```



```

17 INITSEG = 0x9000      ! we move boot here - out of the way ! 原来 bootsect 所处的段。
18 SYSSEG  = 0x1000      ! system loaded at 0x10000 (65536). ! system 在 0x10000(64k)处。
19 SETUPSEG = 0x9020      ! this is the current segment ! 本程序所在的段地址。
20
21 .globl begtext, begdata, begbss, endtext, enddata, endbss
22 .text
23 begtext:
24 .data
25 begdata:
26 .bss
27 begbss:
28 .text
29
30 entry start
31 start:
32
33 ! ok, the read went well so we get current cursor position and save it for
34 ! posterity.
35 ! ok, 整个读磁盘过程都正常, 现在将光标位置保存以备今后使用。
36
37     mov     ax, #INITSEG      ! this is done in bootsect already, but...
38                                     ! 将 ds 置成 #INITSEG (0x9000)。这已经在 bootsect 程序中
39                                     ! 设置过, 但是现在是 setup 程序, Linus 觉得需要再重新
40                                     ! 设置一下。
41
42     mov     ds, ax
43     mov     ah, #0x03        ! read cursor pos
44                                     ! BIOS 中断 0x10 的读光标功能号 ah = 0x03
45                                     ! 输入: bh = 页号
46                                     ! 返回: ch = 扫描开始线, cl = 扫描结束线,
47                                     ! dh = 行号 (0x00 是顶端), dl = 列号 (0x00 是左边)。
48
49     xor     bh, bh
50
51     int     0x10             ! save it in known place, con_init fetches
52                                     ! it from 0x90000.
53     mov     [0], dx
54                                     ! 上两句是说将光标位置信息存放在 0x90000 处, 控制台
55                                     ! 初始化时会来取。
56
57 ! Get memory size (extended mem, kB) ! 下面 3 句取扩展内存的大小值 (KB)。
58                                     ! 是调用中断 0x15, 功能号 ah = 0x88
59                                     ! 返回: ax = 从 0x100000 (1M) 处开始的扩展内存大小 (KB)。
60                                     ! 若出错则 CF 置位, ax = 出错码。
61
62     mov     ah, #0x88
63     int     0x15
64     mov     [2], ax          ! 将扩展内存数值存在 0x90002 处 (1 个字)。
65
66 ! Get video-card data: ! 下面这段用于取显卡当前显示模式。
67                                     ! 调用 BIOS 中断 0x10, 功能号 ah = 0x0f
68                                     ! 返回: ah = 字符列数, al = 显示模式, bh = 当前显示页。
69                                     ! 0x90004 (1 字) 存放当前页, 0x90006 显示模式, 0x90007 字符列数。
70
71     mov     ah, #0x0f
72     int     0x10
73     mov     [4], bx          ! bh = display page

```

```

54      mov     [6],ax          ! al = video mode, ah = window width
55
56 ! check for EGA/VGA and some config parameters ! 检查显示方式 (EGA/VGA) 并取参数。
      ! 调用 BIOS 中断 0x10, 附加功能选择 - 取方式信息
      ! 功能号: ah = 0x12, bl = 0x10
      ! 返回: bh = 显示状态
      !      (0x00 - 彩色模式, I/O 端口=0x3dX)
      !      (0x01 - 单色模式, I/O 端口=0x3bX)
      ! bl = 安装的显示内存
      ! (0x00 - 64k, 0x01 - 128k, 0x02 - 192k, 0x03 = 256k)
      ! cx = 显卡特性参数(参见程序后的说明)。

57
58      mov     ah,#0x12
59      mov     bl,#0x10
60      int     0x10
61      mov     [8],ax          ! 0x90008 = ??
62      mov     [10],bx         ! 0x9000A = 安装的显示内存, 0x9000B = 显示状态(彩色/单色)
63      mov     [12],cx         ! 0x9000C = 显卡特性参数。
64
65 ! Get hd0 data              ! 取第一个硬盘的信息 (复制硬盘参数表)。
      ! 第 1 个硬盘参数表的首地址竟然是中断向量 0x41 的向量值! 而第 2 个硬盘
      ! 参数表紧接第 1 个表的后面, 中断向量 0x46 的向量值也指向这第 2 个硬盘
      ! 的参数表首址。表的长度是 16 个字节(0x10)。
      ! 下面两段程序分别复制 BIOS 有关两个硬盘的参数表, 0x90080 处存放第 1 个
      ! 硬盘的表, 0x90090 处存放第 2 个硬盘的表。

66
67      mov     ax,#0x0000
68      mov     ds,ax
69      lds     si,[4*0x41]     ! 取中断向量 0x41 的值, 也即 hd0 参数表的地址→ds:si
70      mov     ax,#INITSEG
71      mov     es,ax
72      mov     di,#0x0080      ! 传输的目的地址: 0x9000:0x0080 → es:di
73      mov     cx,#0x10        ! 共传输 0x10 字节。
74      rep
75      movsb
76
77 ! Get hd1 data
78
79      mov     ax,#0x0000
80      mov     ds,ax
81      lds     si,[4*0x46]     ! 取中断向量 0x46 的值, 也即 hd1 参数表的地址→ds:si
82      mov     ax,#INITSEG
83      mov     es,ax
84      mov     di,#0x0090      ! 传输的目的地址: 0x9000:0x0090 → es:di
85      mov     cx,#0x10
86      rep
87      movsb
88
89 ! Check that there IS a hd1 :-) ! 检查系统是否存在第 2 个硬盘, 如果不存在则第 2 个表清零。
      ! 利用 BIOS 中断调用 0x13 的取盘类型功能。
      ! 功能号 ah = 0x15;
      ! 输入: dl = 驱动器号 (0x8X 是硬盘: 0x80 指第 1 个硬盘, 0x81 第 2 个硬盘)
      ! 输出: ah = 类型码; 00 --没有这个盘, CF 置位; 01 --是软驱, 没有 change-line 支持;

```

! 02 --是软驱(或其它可移动设备), 有 change-line 支持; 03 --是硬盘。

```

90
91      mov     ax, #0x01500
92      mov     dl, #0x81
93      int     0x13
94      jc      no_disk1
95      cmp     ah, #3          ! 是硬盘吗? (类型 = 3 ? )。
96      je      is_disk1
97 no_disk1:
98      mov     ax, #INITSEG    ! 第2个硬盘不存在, 则对第2个硬盘表清零。
99      mov     es, ax
100     mov     di, #0x0090
101     mov     cx, #0x10
102     mov     ax, #0x00
103     rep
104     stosb
105 is_disk1:
106
107 ! now we want to move to protected mode ... ! 从这里开始我们要保护模式方面的工作了。
108
109     cli          ! no interrupts allowed !! 此时不允许中断。
110
111 ! first we move the system to it's rightful place
! 首先我们将 system 模块移到正确的位置。
! bootsect 引导程序是将 system 模块读入到从 0x10000 (64k) 开始的位置。由于当时假设
! system 模块最大长度不会超过 0x80000 (512k), 也即其末端不会超过内存地址 0x90000,
! 所以 bootsect 会将自己移动到 0x90000 开始的地方, 并把 setup 加载到它的后面。
! 下面这段程序的用途是再把整个 system 模块移动到 0x00000 位置, 即把从 0x10000 到 0x8ffff
! 的内存数据块(512k), 整块地向内存低端移动了 0x10000 (64k) 的位置。
112
113     mov     ax, #0x0000
114     cld          ! 'direction'=0, movs moves forward
115 do_move:
116     mov     es, ax          ! destination segment ! es:di→目的地址(初始为 0x0000:0x0)
117     add     ax, #0x1000
118     cmp     ax, #0x9000    ! 已经把从 0x8000 段开始的 64k 代码移动完?
119     jz      end_move
120     mov     ds, ax          ! source segment ! ds:si→源地址(初始为 0x1000:0x0)
121     sub     di, di
122     sub     si, si
123     mov     cx, #0x8000    ! 移动 0x8000 字 (64k 字节)。
124     rep
125     movsw
126     jmp     do_move
127
128 ! then we load the segment descriptors
! 此后, 我们加载段描述符。
! 从这里开始会遇到 32 位保护模式的操作, 因此需要 Intel 32 位保护模式编程方面的知识了,
! 有关这方面的信息请查阅列表后的简单介绍或附录中的详细说明。这里仅作概要说明。
!
! lidt 指令用于加载中断描述符表(idt)寄存器, 它的操作数是 6 个字节, 0-1 字节是描述符表的
! 长度值(字节); 2-5 字节是描述符表的 32 位线性基地址(首地址), 其形式参见下面
! 219-220 行和 223-224 行的说明。中断描述符表中的每一个表项(8 字节)指出发生中断时

```

! 需要调用的代码的信息, 与中断向量有些相似, 但要包含更多的信息。

!

! lgdt 指令用于加载全局描述符表(gdt)寄存器, 其操作数格式与 lidt 指令的相同。全局描述符

! 表中的每个描述符项(8 字节)描述了保护模式下数据和代码段(块)的信息。其中包括段的

! 最大长度限制(16 位)、段的线性基址(32 位)、段的特权级、段是否在内存、读写许可以及

! 其它一些保护模式运行的标志。参见后面 205-216 行。

!

129

130 end\_move:

131 mov ax, #SETUPSEG ! right, forgot this at first. didn't work :-)

132 mov ds, ax ! ds 指向本程序(setup)段。

133 lidt idt\_48 ! load idt with 0,0

! 加载中断描述符表(idt)寄存器, idt\_48 是 6 字节操作数的位置

! (见 218 行)。前 2 字节表示 idt 表的限长, 后 4 字节表示 idt 表

! 所处的基地址。

134 lgdt gdt\_48 ! load gdt with whatever appropriate

! 加载全局描述符表(gdt)寄存器, gdt\_48 是 6 字节操作数的位置

! (见 222 行)。

135

136 ! that was painless, now we enable A20

! 以上的操作很简单, 现在我们开启 A20 地址线。参见程序列表后有关 A20 信号线的说明。

137

138 call empty\_8042 ! 等待输入缓冲器空。

! 只有当输入缓冲器为空时才可以对其进行写命令。

139 mov al, #0xD1 ! command write ! 0xD1 命令码-表示要写数据到

140 out #0x64, al ! 8042 的 P2 端口。P2 端口的位 1 用于 A20 线的选通。

! 数据要写到 0x60 口。

141 call empty\_8042 ! 等待输入缓冲器空, 看命令是否被接受。

142 mov al, #0xDF ! A20 on ! 选通 A20 地址线的参数。

143 out #0x60, al

144 call empty\_8042 ! 输入缓冲器为空, 则表示 A20 线已经选通。

145

146 ! well, that went ok, I hope. Now we have to reprogram the interrupts :-)

147 ! we put them right after the intel-reserved hardware interrupts, at

148 ! int 0x20-0x2F. There they won't mess up anything. Sadly IBM really

149 ! messed this up with the original PC, and they haven't been able to

150 ! rectify it afterwards. Thus the bios puts interrupts at 0x08-0x0f,

151 ! which is used for the internal hardware interrupts as well. We just

152 ! have to reprogram the 8259's, and it isn't fun.

!! 希望以上一切正常。现在我们必须重新对中断进行编程⊗

!! 我们将它们放在正好处于 intel 保留的硬件中断后面, 在 int 0x20-0x2F。

!! 在那里它们不会引起冲突。不幸的是 IBM 在原 PC 机中搞糟了, 以后也没有纠正过来。

!! PC 机的 bios 将中断放在了 0x08-0x0f, 这些中断也被用于内部硬件中断。

!! 所以我们就必须重新对 8259 中断控制器进行编程, 这一点都没劲。

153

154 mov al, #0x11 ! initialization sequence

! 0x11 表示初始化命令开始,

! 是 ICW1 命令字, 表示边沿触发、多片 8259

! 级连、最后要发送 ICW4 命令字。

155 out #0x20, al ! send it to 8259A-1 ! 发送到 8259A 主芯片。

156 .word 0x00eb, 0x00eb ! jmp \$+2, jmp \$+2 ! \$ 表示当前指令的地址,

! 两条跳转指令, 跳到下一条指令, 起延时作用。

157 out #0xA0, al ! and to 8259A-2 ! 再发送到 8259A 从芯片。

```

158      .word    0x00eb, 0x00eb
159      mov      al, #0x20                ! start of hardware int's (0x20)
                                           ! 送主芯片 ICW2 命令字, 起始中断号, 要送奇地址。
160      out      #0x21, al
161      .word    0x00eb, 0x00eb
162      mov      al, #0x28                ! start of hardware int's 2 (0x28)
                                           ! 送从芯片 ICW2 命令字, 起始中断号。
163      out      #0xA1, al
164      .word    0x00eb, 0x00eb
165      mov      al, #0x04                ! 8259-1 is master
                                           ! 送主芯片 ICW3 命令字, 主芯片的 IR2 连从芯片 INT。
166      out      #0x21, al
167      .word    0x00eb, 0x00eb
168      mov      al, #0x02                ! 8259-2 is slave
                                           ! 送从芯片 ICW3 命令字, 表示从芯片的 INT 连到主芯
                                           ! 片的 IR2 引脚上。
169      out      #0xA1, al
170      .word    0x00eb, 0x00eb
171      mov      al, #0x01                ! 8086 mode for both
                                           ! 送两个芯片的 ICW4 命令字。8086 模式; 普通 EOI
                                           ! 方式, 需发送指令来复位。初始化结束, 芯片就绪。
172      out      #0x21, al
173      .word    0x00eb, 0x00eb
174      out      #0xA1, al
175      .word    0x00eb, 0x00eb
176      mov      al, #0xFF                ! mask off all interrupts for now
                                           ! 当前屏蔽所有中断请求。
177      out      #0x21, al
178      .word    0x00eb, 0x00eb
179      out      #0xA1, al
180
181 ! well, that certainly wasn't fun :-(. Hopefully it works, and we don't
182 ! need no steenking BIOS anyway (except for the initial loading :-).
183 ! The BIOS-routine wants lots of unnecessary data, and it's less
184 ! "interesting" anyway. This is how REAL programmers do it.
185 !
186 ! Well, now's the time to actually move into protected mode. To make
187 ! things as simple as possible, we do no register set-up or anything,
188 ! we let the gnu-compiled 32-bit programs do that. We just jump to
189 ! absolute address 0x00000, in 32-bit protected mode.
!! 哼, 上面这段当然没劲☹, 希望这样能工作, 而且我们也不再需要乏味的 BIOS 了(除了
!! 初始的加载☺。BIOS 子程序要求很多不必要的数据, 而且它一点都没趣。那是“真正”的
!! 程序员所做的事。
190
! 这里设置进入 32 位保护模式运行。首先加载机器状态字(lmsw - Load Machine Status Word), 也称
! 控制寄存器 CRO, 其比特位 0 置 1 将导致 CPU 工作在保护模式。
191      mov      ax, #0x0001                ! protected mode (PE) bit ! 保护模式比特位(PE)。
192      lmsw     ax                        ! This is it!! 就这样加载机器状态字!
193      jmp      0, 8                      ! jmp offset 0 of segment 8 (cs) ! 跳转至 cs 段 8, 偏移 0 处。
! 我们已经将 system 模块移动到 0x00000 开始的地方, 所以这里的偏移地址是 0。这里的段
! 值的 8 已经是保护模式下的段选择符了, 用于选择描述符表和描述符表项以及所要求的
! 特权级。
! 段选择符长度为 16 位 (2 字节); 位 0-1 表示请求的特权级 0-3, linux 操作系统只

```

！ 用到两级：0 级（系统级）和 3 级（用户级）；位 2 用于选择全局描述符表(0)还是局部描述符表(1)；位 3-15 是描述符表项的索引，指出选择第几项描述符。所以段选择符 8(0b0000, 0000, 0000, 1000)表示请求特权级 0、使用全局描述符表中的第 1 项，该项指出代码的基地址是 0（参见 209 行），因此这里的跳转指令就会去执行 system 中的代码。

194

195 ! This routine checks that the keyboard command queue is empty

196 ! No timeout is used - if this hangs there is something wrong with

197 ! the machine, and we probably couldn't proceed anyway.

！ 下面这个子程序检查键盘命令队列是否为空。这里不使用超时方法 - 如果这里死机，

！ 则说明 PC 机有问题，我们就没有办法再处理下去了。

！ 只有当输入缓冲器为空时（状态寄存器位 2 = 0）才可以对其进行写命令。

198 empty\_8042:

199 .word 0x00eb, 0x00eb ! 这是两个跳转指令的机器码(跳转到下一句)，相当于延时空操作。

200 in al, #0x64 ! 8042 status port ! 读 AT 键盘控制器状态寄存器。

201 test al, #2 ! is input buffer full? ! 测试位 2，输入缓冲器满？

202 jnz empty\_8042 ! yes - loop

203 ret

204

205 gdt: ! 全局描述符表开始处。描述符表由多个 8 字节长的描述符项组成。

！ 这里给出了 3 个描述符项。第 1 项无用（206 行），但须存在。第 2 项是系统代码段

！ 描述符（208-211 行），第 3 项是系统数据段描述符（213-216 行）。每个描述符的具体

！ 含义参见列表后说明。

206 .word 0, 0, 0, 0 ! dummy ! 第 1 个描述符，不用。

207 ! 这里在 gdt 表中的偏移量为 0x08，当加载代码段寄存器(段选择符)时，使用的是这个偏移值。

208 .word 0x07FF ! 8Mb - limit=2047 (2048\*4096=8Mb)

209 .word 0x0000 ! base address=0

210 .word 0x9A00 ! code read/exec

211 .word 0x00C0 ! granularity=4096, 386

212 ! 这里在 gdt 表中的偏移量是 0x10，当加载数据段寄存器(如 ds 等)时，使用的是这个偏移值。

213 .word 0x07FF ! 8Mb - limit=2047 (2048\*4096=8Mb)

214 .word 0x0000 ! base address=0

215 .word 0x9200 ! data read/write

216 .word 0x00C0 ! granularity=4096, 386

217

218 idt\_48:

219 .word 0 ! idt limit=0

220 .word 0, 0 ! idt base=0L

221

222 gdt\_48:

223 .word 0x800 ! gdt limit=2048, 256 GDT entries

！ 全局表长度为 2k 字节，因为每 8 字节组成一个段描述符项

！ 所以表中共可有 256 项。

224 .word 512+gdt, 0x9 ! gdt base = 0x9xxxx

！ 4 个字节构成的内存线性地址：0x0009<<16 + 0x0200+gdt

！ 也即 0x90200 + gdt(即在本程序段中的偏移地址，205 行)。

225

226 .text

227 endtext:

228 .data

229 enddata:

230 .bss

231 endbss:

## 2.3.3 其它信息

### 2.3.3.1 BIOS 视频中断 0x10

这里说明上面程序中用到的 ROM BIOS 中视频中断调用的几个子功能。

A. 获取显示卡信息（其它辅助功能选择）：

功能号： ah = 0x12, bh = 0x10

| 输入/返回信息 | 寄存器 | 内容说明   |
|---------|-----|--|
| 输入信息    | ah  | 功能号=0x12，获取显示卡信息   |
|         | bh  | 子功能号=0x10。   |
| 返回信息    | bh  | 视频状态：<br>0x00 - 彩色模式（此时视频硬件 I/O 端口基地址为 0x3DX）；<br>0x01 - 单色模式（此时视频硬件 I/O 端口基地址为 0x3BX）；<br>注：其中端口地址中的 X 值可为 0 - f。   |
|         | bl  | 已安装的显示内存大小：<br>00 = 64K, 01 = 128K, 02 = 192K, 03 = 256K   |
|         | ch  | 特性连接器比特位信息：<br>比特位 说明<br>0 特性线 1，状态 2；<br>1 特性线 0，状态 2；<br>2 特性线 1，状态 1；<br>3 特性线 0，状态 1；<br>4-7 未使用(为 0)  |
|         | cl  | 视频开关设置信息：<br>比特位 说明<br>0 开关 1 关闭；<br>1 开关 2 关闭；<br>2 开关 3 关闭；<br>3 开关 4 关闭；<br>4-7 未使用。<br>原始 EGA/VGA 开关设置值：<br>0x00 MDA/HGC；<br>0x01-0x03 MDA/HGC；<br>0x04 CGA 40x25；<br>0x05 CGA 80x25；<br>0x06 EGA+ 40x25；<br>0x07-0x09 EGA+ 80x25；<br>0x0A EGA+ 80x25 单色；<br>0x0B EGA+ 80x25 单色。 |

### 2.3.3.2 硬盘基本参数表（“INT 0x41”）

中断向量表中，int 0x41 的中断向量位置（4 \* 0x41 = 0x0000:0x0104）存放的并不是中断程序的地址，而是第一个硬盘的基本参数表。对于 100%兼容的 BIOS 来说，这里存放着硬盘参数表阵列的首地址 F000h:E401h。第二个硬盘的基本参数表入口地址存于 int 0x46 中断向量位置处。

表 硬盘基本参数信息表

| 位移   | 大小 | 说明  |
|------|----|---|
| 0x00 | 字  | 柱面数   |
| 0x02 | 字节 | 磁头数   |
| 0x03 | 字  | 开始减小写电流的柱面(仅 PC XT 使用, 其它为 0)   |
| 0x05 | 字  | 开始写前预补偿柱面号(乘 4)   |
| 0x07 | 字节 | 最大 ECC 猝发长度(仅 XT 使用, 其它为 0)   |
| 0x08 | 字节 | 控制字节(驱动器步进选择)<br>位 0 未用<br>位 1 保留(0)(关闭 IRQ)<br>位 2 允许复位<br>位 3 若磁头数大于 8 则置 1<br>位 4 未用(0)<br>位 5 若在柱面数+1 处有生产商的坏区图, 则置 1<br>位 6 禁止 ECC 重试<br>位 7 禁止访问重试。 |
| 0x09 | 字节 | 标准超时值(仅 XT 使用, 其它为 0)   |
| 0x0A | 字节 | 格式化超时值(仅 XT 使用, 其它为 0)  |
| 0x0B | 字节 | 检测驱动器超时值(仅 XT 使用, 其它为 0)  |
| 0x0C | 字  | 磁头着陆(停止)柱面号   |
| 0x0E | 字节 | 每磁道扇区数  |
| 0x0F | 字节 | 保留。   |

### 2.3.3.3 A20 地址线问题

1981 年 8 月, IBM 公司最初推出的个人计算机 IBM PC 使用的 CPU 是 Intel 8088。在该微机中地址线只有 20 根(A0 - A19)。在当时内存 RAM 只有几百 KB 或不到 1MB 时, 20 根地址线已足够用来寻址这些内存。其所能寻址的最高地址是 0xffff:0xffff, 也即 0x10ffef。对于超出 0x100000(1MB)的寻址地址将默认地环绕到 0x0ffef。当 IBM 公司于 1985 年引入 AT 机时, 使用的是 Intel 80286 CPU, 具有 24 根地址线, 最高可寻址 16MB, 并且有一个与 8088 完全兼容的实模式运行方式。然而, 在寻址值超过 1MB 时它却不能象 8088 那样实现地址寻址的环绕。但是当时已经有一些程序是利用这种地址环绕机制进行工作的。为了实现完全的兼容性, IBM 公司发明了使用一个开关来开启或禁止 0x100000 地址比特位。由于在当时的 8042 键盘控制器上恰好有空闲的端口引脚(输出端口 P2, 引脚 P21), 于是便使用了该引脚来作为与门控制这个地址比特位。该信号即被称为 A20。如果它为零, 则比特 20 及以上地址都被清除。从而实现了兼容性。

由于在机器启动时, 默认条件下, A20 地址线是禁止的, 所以操作系统必须使用适当的方法来开启它。但是由于各种兼容机所使用的芯片集不同, 要做到这一点却是非常的麻烦。因此通常要在几种控制方法中选择。

对 A20 信号线进行控制的常用方法是通过设置键盘控制器的端口值。这里的 setup.s 程序(138-144 行)即使用了这种典型的控制方式。对于其它一些兼容微机还可以使用其它方式来做对 A20 线的控制。

有些操作系统将 A20 的开启和禁止作为实模式与保护运行模式之间进行转换的标准过程中的一部分。由于键盘的控制器速度很慢, 因此就不能使用键盘控制器对 A20 线来进行操作。为此引进了一个 A20 快速门选项(Fast Gate A20), 它使用 I/O 端口 0x92 来处理 A20 信号线, 避免了使用慢速的键盘控制器操作方式。对于不含键盘控制器的系统就只能使用 0x92 端口来控制, 但是该端口也有可能被其它兼容微机上的设备(如显示芯片)所使用, 从而造成系统错误的操作。

还有一种方式是通过读 0xee 端口来开启 A20 信号线, 写该端口则会禁止 A20 信号线。



### 2.3.3.4 Intel CPU 32 位保护运行模式

Intel CPU 一般可以在两种模式下运行，即实地址模式和保护模式。早期的 Intel CPU (8088/8086) 只能工作在实模式下，某一时刻只能运行单个任务。对于 Intel 80386 以上的芯片则还可以运行在 32 位保护模式下。在保护模式下运行可以支持多任务；支持 4G 的物理内存；支持虚拟内存；支持内存的页式管理和段式管理；支持特权级。

虽然对保护模式下的运行机制是理解 Linux 内核的重要基础，但由于篇幅所限，对其的简单介绍可以参考书后的附录。但仍然建议初学者能够使用书后列出相关书籍，首先作一番仔细研究。为真正理解 setup.s 程序和下面 head.s 程序的作用，起码要先明白段选择符、段描述符和 80x86 的页表寻址机制。

段选择符

逻辑地址，线性地址，物理地址

段选择符的格式

## 2.4 head.s 文件

### 2.4.1 功能描述

head.s 程序在被编译后，会被连接成 system 模块的最前面开始部分，这也就是为什么称其为头部(head)程序的原因。从这里开始，内核完全都是在保护模式下运行了。heads.s 汇编程序与前面的语法格式不同，它采用的是 AT&T 的汇编语言格式，并且需要使用 GNU 的 gas 和 gld 进行编译连接。因此请注意代码中付值的方向是从左到右。

这段程序实际上处于内存绝对地址 0 处开始的地方。这个程序的功能比较单一。首先是加载各个数据段寄存器，重新设置中断描述符表，共 256 项，并使各个表项均指向一个只报错误的哑中断程序。然后重新设置全局描述符表。接着使用物理地址 0 与 1M 开始处的内容相比较的方法，检测 A20 地址线是否真的开启（如果没有开启，则在访问高于 1Mb 物理内存地址时 CPU 实际只会访问 (IP MOD 1Mb) 地址处的内容），如果检测下来发现没有开启则进入死循环。然后测试 PC 机是否含有数学协处理器芯片，并在控制寄存器 CR0 中置相应的标志位。接着设置分页处理机制，将页目录表放在绝对物理地址 0 开始处，紧随后面放置 4 个页表(可寻址 16MB 内存)，并分别设置它们的表项。最后利用返回指令将预先放置在堆栈中 /init/main.c 程序的地址弹出，去运行 main() 程序。

### 2.4.2 代码注释

列表 2.4 linux/boot/head.s 程序

```

1 /*
2  * linux/boot/head.s
3  *
4  * (C) 1991 Linus Torvalds
5  */
6
7 /*
8  * head.s contains the 32-bit startup code.
9  *
10 * NOTE!!! Startup happens at absolute address 0x00000000, which is also where
11 * the page directory will exist. The startup code will be overwritten by

```

```

12 * the page directory.
13 */
/*
   * head.s 含有 32 位启动代码。
   * 注意!!! 32 位启动代码是从绝对地址 0x00000000 开始的, 这里也同样是页目录将存在的地方,
   * 因此这里的启动代码将被页目录覆盖掉。
   */
14 .text
15 .globl _idt, _gdt, _pg_dir, _tmp_floppy_area
16 _pg_dir:    # 页目录将会存放在这里。
17 startup_32:    # 18-22 行设置各个数据段寄存器。
18     movl $0x10, %eax    # 对于 GNU 汇编来说, 每个直接数要以 '$' 开始, 否则是表示地址。
                           # 每个寄存器名都要以 '%' 开头, eax 表示是 32 位的 ax 寄存器。

# 再此注意!!! 这里已经处于 32 位运行模式, 因此这里的 $0x10 并不是把地址 0x10 装入各个
# 段寄存器, 它现在其实是全局段描述符表中的偏移值, 或者更正确地说是一个描述符表项
# 的选择符。有关选择符的说明请参见 setup.s 中 193 行下的说明。这里 $0x10 的含义是请求
# 特权级 0 (位 0-1=0)、选择全局描述符表 (位 2=0)、选择表中第 1 项 (位 3-15=1)。它正好指
# 向表中的数据段描述符项。(描述符的具体数值参见前面 setup.s 中 212, 213 行)
19     mov %ax, %ds
20     mov %ax, %es
21     mov %ax, %fs
22     mov %ax, %gs
23     lss _stack_start, %esp    # 表示 _stack_start → ss:esp。_stack_start 应是编译程序
                               # 为程序自动生成的存有堆栈信息地方。
24     call setup_idt            # 调用设置中断描述符表子程序。
25     call setup_gdt            # 调用设置全局描述符表子程序。
26     movl $0x10, %eax          # reload all the segment registers
                               # 重新加载所有的段寄存器。
                               # 在改变过全局描述符表 gdt 之后, CS 代码段寄存器已经
                               # 在 setup_gdt 中重新加载过了。
27     mov %ax, %ds              # after changing gdt. CS was already
28     mov %ax, %es              # reloaded in 'setup_gdt'
29     mov %ax, %fs
30     mov %ax, %gs
31     lss _stack_start, %esp
# 32-36 行用于测试 A20 地址线是否已经开启。采用的方法是向内存地址 0x000000 处写入任意
# 一个数值, 然后看内存地址 0x100000 (1M) 处是否也是这个数值。如果一直相同的话, 就一直
# 比较下去, 也即死循环、死机。表示地址 A20 线没有选通, 结果不能对 1M 以上内存寻址。
32     xorl %eax, %eax
33 1:    incl %eax                # check that A20 really IS enabled
34     movl %eax, 0x000000        # loop forever if it isn't
35     cmpl %eax, 0x100000
36     je 1b                     # '1b' 表示向后 (backward) 跳转到标号 1 去 (33 行)。
                               # 若是 '5f' 则表示向前 (forward) 跳转到标号 5 去。
37 /*
38 * NOTE! 486 should set bit 16, to check for write-protect in supervisor
39 * mode. Then it would be unnecessary with the "verify_area()" -calls.
40 * 486 users probably want to set the NE (#5) bit also, so as to use
41 * int 16 for math errors.
42 */
/*
   * 注意! 在下面这段程序中, 486 应该将位 16 置位, 以检查在超级用户模式下的写保护,

```

```

* 此后“verify_area()”调用中就不需要了。486 的用户通常也会想将 NE(#5)置位，以便
* 对数学协处理器的出错使用 int 16。
*/
# 下面这段程序（43-65）用于检查数学协处理器芯片是否存在。方法是修改控制寄存器 CR0，在
# 假设存在协处理器的情况下执行一个协处理器指令，如果出错的话则说明协处理器芯片不存在，
# 需要设置 CR0 中的协处理器仿真位 EM（位 2），并复位协处理器存在标志 MP（位 1）。
43      movl %cr0,%eax      # check math chip
44      andl $0x80000011,%eax # Save PG,PE,ET
45 /* "orl $0x10020,%eax" here for 486 might be good */
46      orl $2,%eax        # set MP
47      movl %eax,%cr0
48      call check_x87
49      jmp after_page_tables # 跳转到 135 行。
50
51 /*
52 * We depend on ET to be correct. This checks for 287/387.
53 */
54 /*
55 * 我们依赖于 ET 标志的正确性来检测 287/387 存在与否。
56 */
57 check_x87:
58      fninit
59      fstsw %ax
60      cmpb $0,%al
61      je 1f              /* no coprocessor: have to set bits */
62      movl %cr0,%eax      # 如果存在的则向前跳转到标号 1 处，否则改写 cr0。
63      xorl $6,%eax        /* reset MP, set EM */
64      movl %eax,%cr0
65      ret
66
67 .align 2 # 这里“.align 2”的含义是指存储边界对齐调整。“2”表示调整到地址最后 2 位为零，
68 # 即按 4 字节方式对齐内存地址。
69 1:      .byte 0xDB,0xE4    /* fsetpm for 287, ignored by 387 */ # 287 协处理器码。
70      ret
71
72 /*
73 * setup_idt
74 *
75 * sets up a idt with 256 entries pointing to
76 * ignore_int, interrupt gates. It then loads
77 * idt. Everything that wants to install itself
78 * in the idt-table may do so themselves. Interrupts
79 * are enabled elsewhere, when we can be relatively
80 * sure everything is ok. This routine will be over-
81 * written by the page tables.
82 */
83 /*
84 * 下面这段是设置中断描述符表子程序 setup_idt
85 *
86 * 将中断描述符表 idt 设置成具有 256 个项，并都指向 ignore_int 中断门。然后加载中断
87 * 描述符表寄存器(用 lidt 指令)。真正实用的中断门以后再安装。当我们在其它地方认为一切
88 * 都正常时再开启中断。该子程序将会被页表覆盖掉。
89 */
90 # 中断描述符表中的项虽然也是 8 字节组成，但其格式与全局表中的不同，被称为门描述符

```

# (Gate Descriptor)。它的 0-1, 6-7 字节是偏移量, 2-3 字节是选择符, 4-5 字节是一些标志。

```

78 setup_idt:
79     lea ignore_int,%edx      # 将 ignore_int 的有效地址 (偏移值) 值→edx 寄存器
80     movl $0x00080000,%eax    # 将选择符 0x0008 置入 eax 的高 16 位中。
81     movw %dx,%ax            /* selector = 0x0008 = cs */
                                # 偏移值的低 16 位置入 eax 的低 16 位中。此时 eax 含有
                                # 门描述符低 4 字节的值。
82     movw $0x8E00,%dx        /* interrupt gate - dpl=0, present */
83                                # 此时 edx 含有门描述符高 4 字节的值。
84     lea _idt,%edi           # _idt 是中断描述符表的地址。
85     mov $256,%ecx
86 rp_sidt:
87     movl %eax, (%edi)        # 将哑中断门描述符存入表中。
88     movl %edx, 4(%edi)
89     addl $8,%edi            # edi 指向表中下一项。
90     dec %ecx
91     jne rp_sidt
92     lidt idt_descr          # 加载中断描述符表寄存器值。
93     ret
94
95 /*
96 *  setup_gdt
97 *
98 *  This routines sets up a new gdt and loads it.
99 *  Only two entries are currently built, the same
100 *  ones that were built in init.s. The routine
101 *  is VERY complicated at two whole lines, so this
102 *  rather long comment is certainly needed :-).
103 *  This routine will beoverwritten by the page tables.
104 */
/*
 * 设置全局描述符表项 setup_gdt
 * 这个子程序设置一个新的全局描述符表 gdt, 并加载。此时仅创建了两个表项, 与前
 * 面的一样。该子程序只有两行, “非常的”复杂, 所以当然需要这么长的注释了©。
105 setup_gdt:
106     lgdt gdt_descr          # 加载全局描述符表寄存器(内容已设置好, 见 232-238 行)。
107     ret
108
109 /*
110 *  I put the kernel page tables right after the page directory,
111 *  using 4 of them to span 16 Mb of physical memory. People with
112 *  more than 16MB will have to expand this.
113 */
/* Linus 将内核的内存页表直接放在页目录之后, 使用了 4 个表来寻址 16 Mb 的物理内存。
 * 如果你有多于 16 Mb 的内存, 就需要在这里进行扩充修改。
 */
# 每个页表长为 4 Kb 字节, 而每个页表项需要 4 个字节, 因此一个页表共可以存放 1000 个表项,
# 如果一个表项寻址 4 Kb 的地址空间, 则一个页表就可以寻址 4 Mb 的物理内存。
# 页表项的格式为: 项的前 0-11 位存放一些标志, 如是否在内存中(P 位 0)、读写许可(R/W 位 1)、
# 普通用户还是超级用户使用(U/S 位 2)、是否修改过(是否脏了)(D 位 6)等; 表项的位 12-31 是
# 页框地址, 用于指出一页内存的物理起始地址。
114 .org 0x1000                # 从偏移 0x1000 处开始是第 1 个页表 (偏移 0 开始处将存放页表目录)。
115 pg0:

```

```

116
117 .org 0x2000
118 pg1:
119
120 .org 0x3000
121 pg2:
122
123 .org 0x4000
124 pg3:
125
126 .org 0x5000      # 定义下面的内存数据块从偏移 0x5000 处开始。
127 /*
128  * tmp_floppy_area is used by the floppy-driver when DMA cannot
129  * reach to a buffer-block. It needs to be aligned, so that it isn't
130  * on a 64kB border.
131  */
132 /* 当 DMA（直接存储器访问）不能访问缓冲块时，下面的 tmp_floppy_area 内存块
133  * 就可供软盘驱动程序使用。其地址需要对齐调整，这样就不会跨越 64kB 边界。
134  */
135 _tmp_floppy_area:
136     .fill 1024, 1, 0      # 共保留 1024 项，每项 1 字节，填充数值 0。
137
138 # 下面这几个入栈操作 (pushl) 用于为调用 /init/main.c 程序和返回作准备。
139 # 前面 3 个入栈指令不知道作什么用的，也许是 Linus 用于在调试时能看清机器码用的 ☺。
140 # 139 行的入栈操作是模拟调用 main.c 程序时首先将返回地址入栈的操作，所以如果
141 # main.c 程序真的退出时，就会返回到这里的标号 L6 处继续执行下去，也即死循环。
142 # 140 行将 main.c 的地址压入堆栈，这样，在设置分页处理 (setup_paging) 结束后
143 # 执行 'ret' 返回指令时就会将 main.c 程序的地址弹出堆栈，并去执行 main.c 程序去了。
144
145 after_page_tables:
146     pushl $0              # These are the parameters to main :-)
147     pushl $0              # 这些是调用 main 程序的参数（指 init/main.c）。
148     pushl $0
149     pushl $L6             # return address for main, if it decides to.
150     pushl $_main          # '_main' 是编译程序对 main 的内部表示方法。
151     jmp setup_paging      # 跳转至第 198 行。
152
153 L6:
154     jmp L6                # main should never return here, but
155                          # just in case, we know what happens.
156
157 /* This is the default interrupt "handler" :-) */
158 /* 下面是默认的中断“向量句柄” ☺ */
159
160 int_msg:
161     .asciz "Unknown interrupt\n\r"      # 定义字符串“未知中断(回车换行)”。
162     .align 2                            # 按 4 字节方式对齐内存地址。
163
164 ignore_int:
165     pushl %eax
166     pushl %ecx
167     pushl %edx
168     push %ds      # 这里请注意！！ds, es, fs, gs 等虽然是 16 位的寄存器，但入栈后
169                  # 仍然会以 32 位的形式入栈，也即需要占用 4 个字节的堆栈空间。
170
171     push %es
172     push %fs
173     movl $0x10, %eax  # 置段选择符（使 ds, es, fs 指向 gdt 表中的数据段）。

```

```

158     mov %ax,%ds
159     mov %ax,%es
160     mov %ax,%fs
161     pushl $int_msg    # 把调用 printk 函数的参数指针（地址）入栈。
162     call _printk      # 该函数在/kernel/printk.c 中。
                        # '_printk' 是 printk 编译后模块中的内部表示法。
163     popl %eax
164     pop %fs
165     pop %es
166     pop %ds
167     popl %edx
168     popl %ecx
169     popl %eax
170     iret              # 中断返回（把中断调用时压入栈的 CPU 标志寄存器（32 位）值也弹出）。
171
172
173 /*
174  * Setup_paging
175  *
176  * This routine sets up paging by setting the page bit
177  * in cr0. The page tables are set up, identity-mapping
178  * the first 16MB. The pager assumes that no illegal
179  * addresses are produced (ie >4Mb on a 4Mb machine).
180  *
181  * NOTE! Although all physical memory should be identity
182  * mapped by this routine, only the kernel page functions
183  * use the >1Mb addresses directly. All "normal" functions
184  * use just the lower 1Mb, or the local data space, which
185  * will be mapped to some other place - mm keeps track of
186  * that.
187  *
188  * For those with more memory than 16 Mb - tough luck. I've
189  * not got it, why should you :-). The source is here. Change
190  * it. (Seriously - it shouldn't be too difficult. Mostly
191  * change some constants etc. I left it at 16Mb, as my machine
192  * even cannot be extended past that (ok, but it was cheap :-).
193  * I've tried to show which constants to change by having
194  * some kind of marker at them (search for "16Mb"), but I
195  * won't guarantee that's all :-( )
196  */
/*
 * 这个子程序通过设置控制寄存器 cr0 的标志（PG 位 31）来启动对内存的分页处理功能，
 * 并设置各个页表项的内容，以恒等映射前 16 MB 的物理内存。分页器假定不会产生非法的
 * 地址映射（也即在只有 4Mb 的机器上设置出大于 4Mb 的内存地址）。
 * 注意！尽管所有的物理地址都应该由这个子程序进行恒等映射，但只有内核函数直接使用
 * >1Mb 的地址。所有“一般”函数仅使用低于 1Mb 的地址空间，或者是使用局部数据空间，
 * 地址空间将被映射到其它一些地方去 -- mm(内存管理程序)会管理这些事的。
 * 对于那些有多于 16Mb 内存的家伙 - 太幸运了，我还没有，为什么你会有☺。代码就在这里，
 * 对它进行修改吧。（实际上，这并不太困难的。通常只需修改一些常数等。我把它设置为
 * 16Mb，因为我的机器再怎么扩充甚至不能超过这个界限（当然，我的机器很便宜的☺）。
 * 我已经通过设置某类标志来给出需要改动的地方（搜索“16Mb”），但我不能保证作这些
 * 改动就行了☺）。
 */

```



```

197 .align 2          # 按4字节方式对齐内存地址边界。
198 setup_paging:     # 首先对5页内存(1页目录 + 4页页表)清零
199     movl $1024*5,%ecx          /* 5 pages - pg_dir+4 page tables */
200     xorl %eax,%eax
201     xorl %edi,%edi            /* pg_dir is at 0x000 */
                                # 页目录从0x000地址开始。
202     cld;rep;stosl
# 下面4句设置页目录中的项,我们共有4个页表所以只需设置4项。
# 页目录项的结构与页表中项的结构一样,4个字节为1项。参见上面113行下的说明。
# "$pg0+7"表示: 0x00001007,是页目录表中的第1项。
# 则第1个页表所在的地址 = 0x00001007 & 0xfffff000 = 0x1000;
# 第1个页表的属性标志 = 0x00001007 & 0x00000fff = 0x07,表示该页存在、用户可读写。
203     movl $pg0+7,_pg_dir      /* set present bit/user r/w */
204     movl $pg1+7,_pg_dir+4    /* ----- " " ----- */
205     movl $pg2+7,_pg_dir+8    /* ----- " " ----- */
206     movl $pg3+7,_pg_dir+12   /* ----- " " ----- */
# 下面6行填写4个页表中所有项的内容,共有: 4(页表)*1024(项/页表)=4096项(0 - 0xfff),
# 也即能映射物理内存 4096*4Kb = 16Mb。
# 每项的内容是: 当前项所映射的物理内存地址 + 该页的标志(这里均为7)。
# 使用的方法是从最后一个页表的最后一项开始按倒退顺序填写。一个页表的最后一项在页表中的
# 位置是 1023*4 = 4092。因此最后一页的最后一项的位置就是$pg3+4092。
207     movl $pg3+4092,%edi      # edi→最后一页的最后一项。
208     movl $0xffff007,%eax     /* 16Mb - 4096 + 7 (r/w user,p) */
                                # 最后1项对应物理内存页面的地址是0xffff000,
                                # 加上属性标志7,即为0xffff007。
209     std                # 方向位置位,edi值递减(4字节)。
210 1:     stosl           /* fill pages backwards - more efficient :- ) */
211     subl $0x1000,%eax     # 每填写好一项,物理地址值减0x1000。
212     jge 1b              # 如果小于0则说明全添写好了。
# 设置页目录基址寄存器cr3的值,指向页目录表。
213     xorl %eax,%eax       /* pg_dir is at 0x0000 */ # 页目录表在0x0000处。
214     movl %eax,%cr3       /* cr3 - page directory start */
# 设置启动使用分页处理(cr0的PG标志,位31)
215     movl %cr0,%eax
216     orl $0x80000000,%eax  # 添上PG标志。
217     movl %eax,%cr0       /* set paging (PG) bit */
218     ret                 /* this also flushes prefetch-queue */
# 在改变分页处理标志后要求使用转移指令刷新预取指令队列,这里用的是返回指令ret。
# 该返回指令的另一个作用是将堆栈中的main程序的地址弹出,并开始运行/init/main.c程序。
# 本程序到此真正结束了。
219
220 .align 2          # 按4字节方式对齐内存地址边界。
221 .word 0
222 idt_descr:        #下面两行是lidt指令的6字节操作数: 长度,基址。
223     .word 256*8-1    # idt contains 256 entries
224     .long _idt
225 .align 2
226 .word 0
227 gdt_descr:        # 下面两行是lgdt指令的6字节操作数: 长度,基址。
228     .word 256*8-1    # so does gdt (not that that's any
229     .long _gdt       # magic number, but it works for me :)
230
231 .align 3          # 按8字节方式对齐内存地址边界。

```

---

```
232 _idt:  .fill 256,8,0          # idt is uninitialized  # 256 项，每项 8 字节，填 0。
233
# 全局表。前 4 项分别是空项（不用）、代码段描述符、数据段描述符、系统段描述符，其中
# 系统段描述符 linux 没有派用处。后面还预留了 252 项的空间，用于放置所创建任务的
# 局部描述符(LDT)和对应的任务状态段 TSS 的描述符。
# (0-nul, 1-cs, 2-ds, 3-sys, 4-TSS0, 5-LDT0, 6-TSS1, 7-LDT1, 8-TSS2 etc...)
234 _gdt:  .quad 0x0000000000000000      /* NULL descriptor */
235        .quad 0x00c09a0000000fff      /* 16Mb */ # 代码段最大长度 16M。
236        .quad 0x00c0920000000fff      /* 16Mb */ # 数据段最大长度 16M。
237        .quad 0x0000000000000000      /* TEMPORARY - don't use */
238        .fill 252,8,0                /* space for LDT's and TSS's etc */
```

---

### 2.4.3 其它信息

## 2.5 小结



## 第3章 初始化程序(init)

### 3.1 概述

在内核源代码的 init/ 目录中只有一个文件 main.c。系统在执行完 boot/ 目录中的 head.s 程序后就将执行权交给了 main.c。该程序虽然不长，但却包括了内核初始化的所有工作。因此在阅读该程序的代码时需要参照很多其它程序中的初始化部分。如果能完全理解这里所调用的所有程序，那么看完这章内容后你应该对 Linux 内核有了大致的了解。

从这一章开始，我们将接触大量的 C 程序代码，因此读者最好具有一定的 C 语言知识。最好的一本参考书还是 Brian W. Kernighan 和 Dennis M. Ritchie 编著的《C 程序设计语言》，对该书第五章关于指针和数组的理解，可以说是弄懂 C 语言的关键。

在注释 C 语言代码时，为了与程序中原有的注释相区别，我们使用 ‘//’ 作为注释语句的开始，有关原有注释的翻译则采用与其一样的注释标志。对于程序中包含的头文件 (\*.h)，仅作概要含义的解释，具体内容将在注释相应头文件时给出。

### 3.2 main.c 程序

列表 3.1 linux/init/main.c 程序

```

1 /*
2  * linux/init/main.c
3  *
4  * (C) 1991 Linus Torvalds
5  */
6
7 #define LIBRARY // 定义该变量是为了包括定义在 unistd.h 中的内嵌汇编代码等信息。
8 #include <unistd.h> // *.h 头文件所在的默认目录是 include/，则在代码中就不用明确指明位置。
                      // 如果不是 UNIX 的标准头文件，则需要指明所在的目录，并用双引号括住。
                      // 标准符号常数与类型文件。定义了各种符号常数和类型，并声明了各种函数。
                      // 如果定义了 __LIBRARY__，则还包括系统调用号和内嵌汇编代码 _syscall0() 等。
9 #include <time.h> // 时间类型头文件。其中最主要定义了 tm 结构和一些有关时间的函数原形。
10
11 /*
12  * we need this inline - forking from kernel space will result
13  * in NO COPY ON WRITE (!!!), until an execve is executed. This
14  * is no problem, but for the stack. This is handled by not letting
15  * main() use the stack at all after fork(). Thus, no function
16  * calls - which means inline code for fork too, as otherwise we
17  * would use the stack upon exit from 'fork()'.
18  *
19  * Actually only pause and fork are needed inline, so that there
20  * won't be any messing with the stack from main(), but we define
21  * some others too.

```

```

22  */
   /*
   * 我们需要下面这些内嵌语句 - 从内核空间创建进程(forking)将导致没有写时复制(COPY ON WRITE)!!!
   * 直到一个执行 execve 调用。这对堆栈可能带来问题。处理的方法是在 fork() 调用之后不让 main() 使用
   * 任何堆栈。因此就不能有函数调用 - 这意味着 fork 也要使用内嵌的代码, 否则我们在从 fork() 退出
   * 时就要使用堆栈了。
   * 实际上只有 pause 和 fork 需要使用内嵌方式, 以保证从 main() 中不会弄乱堆栈, 但是我们同时还
   * 定义了其它一些函数。
   */
23 static inline \_syscall0(int, fork) // 是 unistd.h 中的宏代码。以嵌入汇编的形式进行调用
   // Linux 的系统调用中断 0x80。该中断是所有系统调用的
   // 入口。该条语句实际上是 int fork() 创建进程系统调用。
   // _syscall0 名中最后的数 0 表示无参数, 1 表示 1 个参数。
24 static inline \_syscall0(int, pause) // int pause() 系统调用: 暂停进程的执行, 直到
   // 收到一个信号。
25 static inline \_syscall1(int, setup, void *, BIOS) // int setup(void * BIOS) 系统调用, 仅用于
   // linux 初始化 (仅在这个程序中被调用)。
26 static inline \_syscall0(int, sync) // int sync() 系统调用: 更新文件系统。
27
28 #include <linux/tty.h> // tty 头文件, 定义了有关 tty_io, 串行通信方面的参数、常数。
29 #include <linux/sched.h> // 调度程序头文件, 定义了任务结构 task_struct、第 1 个初始任务
   // 的数据。还有一些以宏的形式定义的有关描述符参数设置和获取的
   // 嵌入式汇编函数程序。
30 #include <linux/head.h> // head 头文件, 其中定义了与 head.s 程序描述符表 (GDT, LDT)
   // 有关的 C 程序。
31 #include <asm/system.h> // 系统头文件。以宏的形式定义了许多有关设置或修改
   // 描述符/中断门等的嵌入式汇编子程序。
32 #include <asm/io.h> // io 头文件。以宏的嵌入式汇编程序形式定义对 io 端口操作的函数。
33
34 #include <stddef.h> // 标准定义头文件。定义了 NULL, offsetof(TYPE, MEMBER)。
35 #include <stdarg.h> // 标准参数头文件。以宏的形式定义变量参数列表。主要说明了一个
   // 类型(va_list)和三个宏(va_start, va_arg 和 va_end), vsprintf、
   // vprintf、vfprintf。
36 #include <unistd.h>
37 #include <fcntl.h> // 文件控制头文件。用于文件及其描述符的操作控制常数符号的定义。
38 #include <sys/types.h> // 类型头文件。定义了基本的系统数据类型。
39
40 #include <linux/fs.h> // 文件系统头文件。定义文件表结构 (file, buffer_head, m_inode 等)。
41
42 static char printbuf[1024]; // 静态字符串数组。
43
44 extern int vsprintf(); // 送格式化输出到一字符串中 (在 kernel/vsprintf.c, 92 行)。
45 extern void init(void); // 函数原形, 初始化 (在 168 行)。
46 extern void blk\_dev\_init(void); // 块设备初始化子程序 (kernel/blk_drv/ll_rw_blk.c, 157 行)
47 extern void chr\_dev\_init(void); // 字符设备初始化 (kernel/chr_drv/tty_io.c, 347 行)
48 extern void hd\_init(void); // 硬盘初始化程序 (kernel/blk_drv/hd.c, 343 行)
49 extern void floppy\_init(void); // 软驱初始化程序 (kernel/blk_drv/floppy.c, 457 行)
50 extern void mem\_init(long start, long end); // 内存管理初始化 (mm/memory.c, 399 行)
51 extern long rd\_init(long mem_start, int length); // 虚拟盘初始化 (kernel/blk_drv/ramdisk.c, 52)
52 extern long kernel\_mktime(struct tm * tm); // 建立内核时间 (秒)。
53 extern long startup\_time; // 内核启动时间 (开机时间) (秒)。
54
55 /*

```

```

56 * This is set up by the setup-routine at boot-time
57 */
58 /*
59  * 以下这些数据是由 setup.s 程序在引导时间设置的（参见第2章2.3.1节中的表2.1）。
60  */
61 #define EXT_MEM_K (*(unsigned short *)0x90002) // 1M以后的扩展内存大小（KB）。
62 #define DRIVE_INFO (*(struct drive_info *)0x90080) // 硬盘参数表基址。
63 #define ORIG_ROOT_DEV (*(unsigned short *)0x901FC) // 根文件系统所在设备号。
64 /*
65  * Yeah, yeah, it's ugly, but I cannot find how to do this correctly
66  * and this seems to work. I anybody has more info on the real-time
67  * clock I'd be interested. Most of this was trial and error, and some
68  * bios-listing reading. Urghh.
69  */
70 /*
71  * 是啊，是啊，下面这段程序很差劲，但我不知道如何正确地实现，而且好象它还能运行。如果有
72  * 关于实时时钟更多的资料，那我很感兴趣。这些都是试探出来的，以及看了一些 bios 程序，呵！
73  */
74 #define CMOS_READ(addr) ({ \ // 这段宏读取 CMOS 实时时钟信息。
75 outb_p(0x80|addr, 0x70); \ // 0x70 是写端口号，0x80|addr 是要读取的 CMOS 内存地址。
76 inb_p(0x71); \ // 0x71 是读端口号。
77 })
78 #define BCD_TO_BIN(val) ((val)=((val)&15) + ((val)>>4)*10) // 将 BCD 码转换成数字。
79 static void time_init(void) // 该子程序取 CMOS 时钟，并设置开机时间→startup_time(秒)。
80 {
81     struct tm time;
82     do {
83         time.tm_sec = CMOS_READ(0); // 参见后面 CMOS 内存列表。
84         time.tm_min = CMOS_READ(2);
85         time.tm_hour = CMOS_READ(4);
86         time.tm_mday = CMOS_READ(7);
87         time.tm_mon = CMOS_READ(8);
88         time.tm_year = CMOS_READ(9);
89     } while (time.tm_sec != CMOS_READ(0));
90     BCD_TO_BIN(time.tm_sec);
91     BCD_TO_BIN(time.tm_min);
92     BCD_TO_BIN(time.tm_hour);
93     BCD_TO_BIN(time.tm_mday);
94     BCD_TO_BIN(time.tm_mon);
95     BCD_TO_BIN(time.tm_year);
96     time.tm_mon--;
97     startup_time = kernel_mktime(&time);
98 }
99 static long memory_end = 0; // 机器具有的内存（字节数）。
100 static long buffer_memory_end = 0; // 高速缓冲区末端地址。
101 static long main_memory_start = 0; // 主内存（将用于分页）开始的位置。

```

```

102 struct drive_info { char dummy[32]; } drive_info; // 用于存放硬盘参数表信息。
103
104 void main(void)          /* This really IS void, no error here. */
105 {                        /* The startup routine assumes (well, ...) this */
106 /*
107  * Interrupts are still disabled. Do necessary setups, then
108  * enable them
109  */
110 /*
111  * 此时中断仍被禁止着，做完必要的设置后就将其开启。
112  */
113 // 下面这段代码用于保存：
114 // 根设备号 → ROOT_DEV; 高速缓存末端地址 → buffer_memory_end;
115 // 机器内存数 → memory_end; 主内存开始地址 → main_memory_start;
116
117     ROOT_DEV = ORIG_ROOT_DEV;
118     drive_info = DRIVE_INFO;
119     memory_end = (1<<20) + (EXT_MEM_K<<10); // 内存大小=1Mb 字节+扩展内存(k)*1024 字节。
120     memory_end &= 0xfffff000; // 忽略不到 4Kb (1 页) 的内存数。
121     if (memory_end > 16*1024*1024) // 如果内存超过 16Mb, 则按 16Mb 计。
122         memory_end = 16*1024*1024;
123     if (memory_end > 12*1024*1024) // 如果内存>12Mb, 则设置缓冲区末端=4Mb
124         buffer_memory_end = 4*1024*1024;
125     else if (memory_end > 6*1024*1024) // 否则如果内存>6Mb, 则设置缓冲区末端=2Mb
126         buffer_memory_end = 2*1024*1024;
127     else
128         buffer_memory_end = 1*1024*1024; // 否则则设置缓冲区末端=1Mb
129     main_memory_start = buffer_memory_end; // 主内存起始位置=缓冲区末端;
130 #ifdef RAMDISK // 如果定义了虚拟盘, 则主内存将减少。
131     main_memory_start += rd_init(main_memory_start, RAMDISK*1024);
132 #endif
133 // 以下是内核进行所有方面的初始化工作。阅读时最好跟着调用的程序深入进去看, 实在看
134 // 不下去了, 就先放一放, 看下一个初始化调用 -- 这是经验之谈☺。
135
136     mem_init(main_memory_start, memory_end);
137     trap_init(); // 陷阱门初始化。
138     blk_dev_init(); // 块设备初始化。
139     chr_dev_init(); // 字符设备初始化。
140     tty_init(); // tty 初始化。
141     time_init(); // 设置开机启动时间 → startup_time (见 76 行)。
142     sched_init(); // 调度程序初始化(加载了任务 0 的 tr, ldtr)
143     buffer_init(buffer_memory_end); // 内存管理初始化 (建立内存链表等)。
144     hd_init(); // 硬盘初始化。
145     floppy_init(); // 软驱初始化。
146     sti(); // 所有初始化工作都做完了, 开启中断。
147     move_to_user_mode(); // 移到用户模式。
148     if (!fork()) { /* we count on this going ok */
149         init();
150     }
151 /*
152  * NOTE!! For any other task 'pause()' would mean we have to get a
153  * signal to awaken, but task0 is the sole exception (see 'schedule()')
154  * as task 0 gets activated at every idle moment (when no other tasks
155  * can run). For task0 'pause()' just means we go check if some other
156  * task can run, and if not we return here.

```

```

147 */
/* 注意!! 对于任何其它的任务, 'pause()' 将意味着我们必须等待收到一个信号才会返
* 回就绪运行态, 但任务 0 (task0) 是唯一的意外情况 (参见 'schedule()'), 因为任务 0 在
* 任何空闲时间里都会被激活 (当没有其它任务在运行时), 因此对于任务 0 'pause()' 仅意味着
* 我们返回来查看是否有其它任务可以运行, 如果没有的话我们就回到这里, 一直循环执行 'pause()'。
*/
148     for(;;) pause();
149 }
150
151 static int printf(const char *fmt, ...)
// 产生格式化信息并输出到标准输出设备 stdout(1), 这里是指屏幕上显示。参数 '*fmt' 指定输出将
// 采用的格式, 参见各种标准 C 语言书籍。该子程序正好是 vsprintf 如何使用的一个例子。
// 该程序使用 vsprintf() 将格式化的字符串放入 printf 缓冲区, 然后用 write() 将缓冲区的内容
// 输出到标准设备 (1--stdout)。
152 {
153     va_list args;
154     int i;
155
156     va_start(args, fmt);
157     write(1, printfbuf, i=vsprintf(printfbuf, fmt, args));
158     va_end(args);
159     return i;
160 }
161
162 static char * argv_rc[] = { "/bin/sh", NULL }; // 调用执行程序时参数的字符串数组。
163 static char * envp_rc[] = { "HOME=/", NULL }; // 调用执行程序时的环境字符串数组。
164
165 static char * argv[] = { "-/bin/sh", NULL }; // 同上。
166 static char * envp[] = { "HOME=/usr/root", NULL };
167
168 void init(void)
169 {
170     int pid, i;
171
172     setup((void *) &drive_info); // 读取硬盘参数包括分区表信息并建立虚拟盘和
// 安装根文件系统设备。(kernel/blk_drv/hd.c, 71)
173     (void) open("/dev/tty0", O_RDWR, 0); // 用读写访问方式打开设备 "/dev/tty0",
// 这里对应终端控制台。
// 返回的句柄号 0 -- stdin 标准输入设备。
174     (void) dup(0); // 复制句柄, 产生句柄 1 号 -- stdout 标准输出设备。
175     (void) dup(0); // 复制句柄, 产生句柄 2 号 -- stderr 标准出错输出设备。
176     printf("%d buffers = %d bytes buffer space\n\r", NR_BUFFERS,
177           NR_BUFFERS*BLOCK_SIZE); // 打印缓冲区块数和总字节数, 每块 1024 字节。
178     printf("Free mem: %d bytes\n\r", memory_end-main memory start); // 空闲内存字节数。
// 下面 fork() 用于创建一个子进程(子任务)。对于被创建的子进程, fork() 将返回 0 值,
// 对于原(父进程)将返回子进程的进程号。所以 180-184 句是子进程执行的内容。该子进程
// 关闭了句柄 0(stdin), 以只读方式打开/etc/rc 文件, 并执行/bin/sh 程序, 所带参数和
// 环境变量分别由 argv_rc 和 envp_rc 数组给出。
179     if (!(pid=fork())) {
180         close(0);
181         if (open("/etc/rc", O_RDONLY, 0))
182             _exit(1); // 如果打开文件失败, 则退出 (/lib/_exit.c, 10)。
183         execve("/bin/sh", argv_rc, envp_rc); // 装入/bin/sh 程序并执行。

```

```

184         _exit(2);
185     }
    // 下面是父进程执行的语句。wait() 是等待子进程停止或终止，其返回值应是子进程的进程号(pid)。
    // 这三句的作用是父进程等待子进程的结束。&i 是存放返回状态信息的位置。如果 wait() 返回值不
    // 等于子进程号，则继续等待。
186     if (pid>0)
187         while (pid != wait(&i))
188             /* nothing */;
    // 如果执行到这里，说明刚创建的子进程的执行已停止或终止了。下面循环中首先再创建一个子进程，
    // 如果出错，则显示“初始化程序创建子进程失败”的信息并继续执行。对于所创建的子进程关闭所有
    // 以前还遗留的句柄(stdin, stdout, stderr)，新建一个会话并设置进程组号，然后重新打开
    // /dev/tty0 作为 stdin，并复制成 stdout 和 stderr。再次执行系统解释程序/bin/sh。但这次执行所
    // 选用的参数和环境数组另选了一套（见上面 165-167 行）。然后父进程再次运行 wait() 等待。如果
    // 子进程又停止了执行，则在标准输出上显示出错误信息“子进程 pid 停止了运行，返回码是 i”，然后
    // 继续重试下去…，形成“大”死循环。
189     while (1) {
190         if ((pid=fork())<0) {
191             printf("Fork failed in init\r\n");
192             continue;
193         }
194         if (!pid) {
195             close(0);close(1);close(2);
196             setsid();
197             (void) open("/dev/tty0",O_RDWR,0);
198             (void) dup(0);
199             (void) dup(0);
200             _exit(execve("/bin/sh",argv,envp));
201         }
202         while (1)
203             if (pid == wait(&i))
204                 break;
205         printf("\n\rchild %d died with code %04x\r\n",pid,i);
206         sync();
207     }
208     _exit(0);        /* NOTE! _exit, not exit() */
209 }
210

```

### 3.3 其它信息






### 3.4 小结

对于 0.11 版内核，通过上面代码分析可知只要根文件系统是一个 MINIX 文件系统，并且其中只要包含文件/etc/rc、/bin/sh、/dev/\* 以及一些目录/etc/、/dev/、/bin/、/home/、/home/root/ 就可以构成一个最简单的根文件系统，让 Linux 运行起来。

## 第 4 章 内核代码(kernel)

### 4.1 概述

列表 4.1 linux/kernel/目录

| Name  | Size       | Last modified (GMT) | Description |
|---|------------|---------------------|-------------|
|  <a href="#">blk_drv/</a>        |            | 1991-12-08 14:09:29 |             |
|  <a href="#">chr_drv/</a>        |            | 1991-12-08 18:36:09 |             |
|  <a href="#">math/</a>           |            | 1991-12-08 14:09:58 |             |
|  <a href="#">Makefile</a>        | 3309 bytes | 1991-12-02 03:21:37 | m           |
|  <a href="#">asm.s</a>         | 2335 bytes | 1991-11-18 00:30:28 | m           |
|  <a href="#">exit.c</a>        | 4175 bytes | 1991-12-07 15:47:55 | m           |
|  <a href="#">fork.c</a>        | 3693 bytes | 1991-11-25 15:11:09 | m           |
|  <a href="#">mktime.c</a>      | 1461 bytes | 1991-10-02 14:16:29 | m           |
|  <a href="#">panic.c</a>       | 448 bytes  | 1991-10-17 14:22:02 | m           |
|  <a href="#">printk.c</a>      | 734 bytes  | 1991-10-02 14:16:29 | m           |
|  <a href="#">sched.c</a>       | 8242 bytes | 1991-12-04 19:55:28 | m           |
|  <a href="#">signal.c</a>      | 2651 bytes | 1991-12-07 15:47:55 | m           |
|  <a href="#">sys.c</a>         | 3706 bytes | 1991-11-25 19:31:13 | m           |
|  <a href="#">system_call.s</a> | 5265 bytes | 1991-12-04 13:56:34 | m           |
|  <a href="#">traps.c</a>       | 4951 bytes | 1991-10-30 20:20:40 | m           |
|  <a href="#">vsprintf.c</a>    | 4800 bytes | 1991-10-02 14:16:29 | m           |



其中三个子目录中的代码注释将放在下一章。

## 4.2 Makefile 文件

### 4.2.1 功能简介

该文件的组成格式与第一章中列表 1.2 的基本相同，在阅读时可以参考列表 1.2 中的有关注释。

### 4.2.2 文件注释

列表 4.2 linux/kernel/Makefile 文件

```

1 #
2 # Makefile for the FREAX-kernel.
3 #
4 # Note! Dependencies are done automagically by 'make dep', which also
5 # removes any old dependencies. DON'T put your own dependencies here
6 # unless it's something special (ie not a .c file).
7 #
8 # FREAX 内核的 Makefile 文件。
9 #
10 # 注意！依赖关系是由'make dep'自动进行的，它也会自动去除原来的依赖信息。不要把你自己的
11 # 依赖关系信息放在这里，除非是特别文件的（也即不是一个.c文件的信息）。
12 # (Linux最初的名字叫FREAX，后来被ftp.funet.fi的管理员改成Linux这个名字)
13
14 AR      =gar      # GNU 的二进制文件处理程序，用于创建、修改以及从归档文件中抽取文件。
15 AS      =gas      # GNU 的汇编程序。
16 LD      =gld      # GNU 的连接程序。
17 LDFLAGS =-s -x    # 连接程序所有的参数，-s 输出文件中省略所有符号信息。-x 删除所有局部符号。
18 CC      =gcc      # GNU C 语言编译器。
19 CFLAGS  =-Wall -O -fstrength-reduce -fomit-frame-pointer -fcombine-regs \
20         -finline-functions -mstring-insns -nostdinc -I../include
21 # C 编译程序选项。-Wall 显示所有的警告信息；-O 优化选项，优化代码长度和执行时间；
22 # -fstrength-reduce 优化循环执行代码，排除重复变量；-fomit-frame-pointer 省略保存不必要
23 # 的框架指针；-fcombine-regs 合并寄存器，减少寄存器类的使用；-finline-functions 将所有简
24 # 单短小的函数代码嵌入调用程序中；-mstring-insns Linux 自己添加的优化选项，以后不再使用；
25 # -nostdinc -I../include 不使用默认路径中的包含文件，而使用这里指定目录中的(../include)。
26
27 CPP      =gcc -E -nostdinc -I../include
28 # C 前处理选项。-E 只运行 C 前处理，对所有指定的 C 程序进行预处理并将处理结果输出到标准输
29 # 出设备或指定的输出文件中；-nostdinc -I../include 同前。
30
31 # 下面的规则指示 make 利用下面的命令将所有的.c 文件编译生成.s 汇编程序。该规则的命令
32 # 指使 gcc 采用 CFLAGS 所指定的选项对 C 代码编译后不进行汇编就停止（-S），从而产生与
33 # 输入的各个 C 文件对应的汇编代码文件。默认情况下所产生的汇编程序文件名是原 C 文件名
34 # 去掉.c 而加上.s 后缀。-o 表示其后是输出文件的名称。其中$.s（或$@）是自动目标变量，
35 # $<代表第一个先决条件，这里即是符合条件*.c 的文件。
36
37 .c.s:
38     $(CC) $(CFLAGS) \
39     -S -o $.s $<

```



# 下面规则表示将所有.s 汇编程序文件编译成.o 目标文件。22 行是实现该操作的具体命令。

```

21 .s.o:
22     $(AS) -c -o $*.o $<
23 .c.o:          # 类似上面, *.c 文件→*.o 目标文件。不进行连接。
24     $(CC) $(CFLAGS) \
25     -c -o $*.o $<
26
27 OBJS = sched.o system_call.o traps.o asm.o fork.o \    # 定义目标文件变量 OBJS。
28     panic.o printk.o vsprintf.o sys.o exit.o \
29     signal.o mktime.o
30
31 kernel.o: $(OBJS)          # 在有了先决条件 OBJS 后使用下面的命令连接成目标 kernel.o
32     $(LD) -r -o kernel.o $(OBJS)
33     sync
34
35 # 下面的规则用于清理工作。当执行'make clean'时, 就会执行 36--40 行上的命令, 去除所有编译
36 # 连接生成的文件。'rm' 是文件删除命令, 选项-f 含义是忽略不存在的文件, 并且不显示删除信息。
37 clean:
38     rm -f core *.o *.a tmp_make keyboard.s
39     for i in *.c;do rm -f `basename $$i .c`.s;done
40     (cd chr_drv; make clean)    # 进入 chr_drv/目录; 执行该目录 Makefile 中的 clean 规则。
41     (cd blk_drv; make clean)
42     (cd math; make clean)
43
44 # 下面得目标或规则用于检查各文件之间的依赖关系。方法如下:
45 # 使用字符串编辑程序 sed 对 Makefile 文件 (这里即是自己) 进行处理, 输出为删除 Makefile
46 # 文件中'### Dependencies' 行后面的所有行 (下面从 51 开始的行), 并生成 tmp_make
47 # 临时文件 (43 行的作用)。然后对 kernel/目录下的每一个 C 文件执行 gcc 预处理操作。
48 # -M 标志告诉预处理程序输出描述每个目标文件相关性的规则, 并且这些规则符合 make 语法。
49 # 对于每一个源文件, 预处理程序输出一个 make 规则, 其结果形式是相应源程序文件的目标
50 # 文件名加上其依赖关系--该源文件中包含的所有头文件列表。把预处理结果都添加到临时
51 # 文件 tmp_make 中, 然后将该临时文件复制成新的 Makefile 文件。
52 dep:
53     sed '/\#\#\# Dependencies/q' < Makefile > tmp_make
54     (for i in *.c;do echo -n `echo $$i | sed 's,\.\,s,`"; \
55     $(CPP) -M $$i;done) >> tmp_make
56     cp tmp_make Makefile
57     (cd chr_drv; make dep)    # 对 chr_drv/目录下的 Makefile 文件也作同样的处理。
58     (cd blk_drv; make dep)
59
60 ### Dependencies:
61 exit.s exit.o : exit.c ../include/errno.h ../include/signal.h \
62     ../include/sys/types.h ../include/sys/wait.h ../include/linux/sched.h \
63     ../include/linux/head.h ../include/linux/fs.h ../include/linux/mm.h \
64     ../include/linux/kernel.h ../include/linux/tty.h ../include/termios.h \
65     ../include/asm/segment.h
66 fork.s fork.o : fork.c ../include/errno.h ../include/linux/sched.h \
67     ../include/linux/head.h ../include/linux/fs.h ../include/sys/types.h \
68     ../include/linux/mm.h ../include/signal.h ../include/linux/kernel.h \
69     ../include/asm/segment.h ../include/asm/system.h
70 mktime.s mktime.o : mktime.c ../include/time.h
71 panic.s panic.o : panic.c ../include/linux/kernel.h ../include/linux/sched.h \
72     ../include/linux/head.h ../include/linux/fs.h ../include/sys/types.h \

```

```

63  ../include/linux/mm.h ../include/signal.h
64  printk.s printk.o : printk.c ../include/stdarg.h ../include/stddef.h \
65  ../include/linux/kernel.h
66  sched.s sched.o : sched.c ../include/linux/sched.h ../include/linux/head.h \
67  ../include/linux/fs.h ../include/sys/types.h ../include/linux/mm.h \
68  ../include/signal.h ../include/linux/kernel.h ../include/linux/sys.h \
69  ../include/linux/fdreg.h ../include/asm/system.h ../include/asm/io.h \
70  ../include/asm/segment.h
71  signal.s signal.o : signal.c ../include/linux/sched.h ../include/linux/head.h \
72  ../include/linux/fs.h ../include/sys/types.h ../include/linux/mm.h \
73  ../include/signal.h ../include/linux/kernel.h ../include/asm/segment.h
74  sys.s sys.o : sys.c ../include/errno.h ../include/linux/sched.h \
75  ../include/linux/head.h ../include/linux/fs.h ../include/sys/types.h \
76  ../include/linux/mm.h ../include/signal.h ../include/linux/tty.h \
77  ../include/termios.h ../include/linux/kernel.h ../include/asm/segment.h \
78  ../include/sys/times.h ../include/sys/utsname.h
79  traps.s traps.o : traps.c ../include/string.h ../include/linux/head.h \
80  ../include/linux/sched.h ../include/linux/fs.h ../include/sys/types.h \
81  ../include/linux/mm.h ../include/signal.h ../include/linux/kernel.h \
82  ../include/asm/system.h ../include/asm/segment.h ../include/asm/io.h
83  vsprintf.s vsprintf.o : vsprintf.c ../include/stdarg.h ../include/string.h

```

## 4.3 asm.s 文件

### 4.3.1 功能描述

asm.s 汇编程序中包括大部分 CPU 探测到的异常故障处理的底层次代码，也包括数学协处理器（FPU）的异常处理。该程序与 kernel/traps.c 程序有着密切的关系。该程序的主要处理方式是在中断处理程序中调用相应的 C 函数程序，显示出错位置和出错号，然后退出中断。

中断信号通常可以分为两类：硬件中断和软件中断(异常)。每个中断是由 0-255 之间的一个数字来标识。中断 int0—int31(0x00—0x1f) 每个的功能是由 Intel 固定设定或保留用的，属于软件中断，但 Intel 称之为异常。是由 CPU 执行指令时探测到异常时引起的。通常分为故障(Fault)和陷阱(traps)。中断 int32—int255 (0x20-0xff)可以由用户自己设定。Linux 中则将 int32-int47(0x20-0x2f)对应于 8259A 中断控制芯片发出的硬件中断请求信号 IRQ0-IRQ15；并将程序编程发出的系统调用(system\_call)中断设置为 int128(0x80)。

本代码文件主要涉及对 Intel 保留中断 int0—int16 的处理(其余保留的 int17-int31 由 Intel 留作今后扩充使用)。对应于中断控制器芯片各 IRQ 发出的 int32-int47 的处理程序将分别在各种硬件（如时钟、键盘、软盘、数学协处理器、硬盘等）初始化程序中处理。Linux 系统调用中断 int128(0x80)的处理将在 kernel/system\_call.s 中给出。各个中断的具体定义见代码注释后其它信息一节中的说明。

在阅读这段代码时参照下面堆栈变化示意图将是很有帮助的（图中每个行代表 4 个字节）。在开始执行程序之前，堆栈指针 esp 指在中断返回地址一栏(图中 esp0 处)。当把将要调用的 C 函数 do\_divide\_error 或其它 C 函数地址入栈后，指针位置是 esp1 处，此时通过交换指令，该函数的地址被放入 eax 寄存器中，

而原来 `eax` 的值被保存到堆栈上。在把一些寄存器入栈后，堆栈指针位置在 `esp2` 处。当正式调用 `do_divide_error()` 之前，程序将开始执行时的 `esp0` 堆栈指针值压入堆栈，放到了 `esp3` 处，并在中断返回弹出入栈的寄存器之前指针通过加上 8 又回到 `esp2` 处。

正式调用 `do_divide_error()` 之前把出错代码以及 `esp0` 入栈的原因是为了作为调用 C 函数 `do_divide_error()` 的参数。在 `traps.c` 中该函数的原形为：

```
void do_divide_error(long esp, long error_code)。
```

因此在这个 C 函数中就可以打印出出错的位置和错误号。程序中其余异常出错的处理过程与这里描述的过程基本类似。

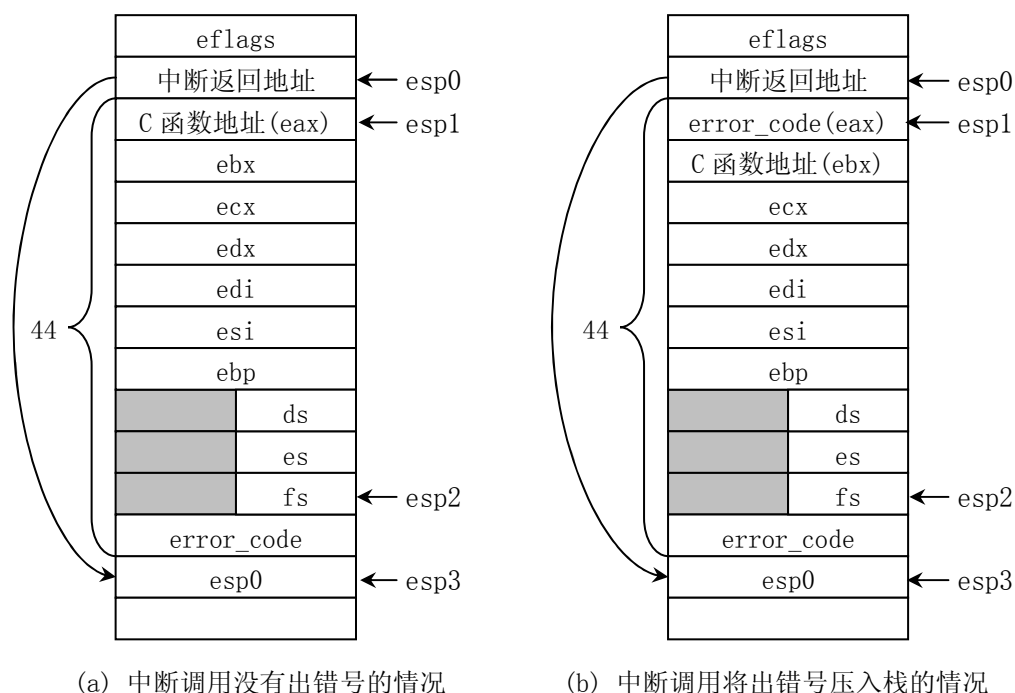


图 4.1 出错处理堆栈变化示意图

## 4.3.2 代码注释

列表 4.3 linux/kernel/asm.s 程序

```
1 /*
2  * linux/kernel/asm.s
3  *
4  * (C) 1991 Linus Torvalds
5  */
6
7 /*
8  * asm.s contains the low-level code for most hardware faults.
9  * page_exception is handled by the mm, so that isn't here. This
10 * file also handles (hopefully) fpu-exceptions due to TS-bit, as
11 * the fpu must be properly saved/resored. This hasn't been tested.
12 */
```

```

/*
 * asm.s 程序中包括大部分的硬件故障（或出错）处理的底层次代码。页异常是由内存管理程序
 * mm 处理的，所以不在此处。此程序还处理（希望是这样）由于 TS-位而造成的 fpu 异常，
 * 因为 fpu 必须正确地进行保存/恢复处理，这些还没有测试过。
 */
13 # 本代码文件主要涉及对 Intel 保留的中断 int0--int16 的处理（int17--int31 留作今后使用）。
   # 以下是一些全局函数名的声明，其原形在 traps.c 中说明。
14 .globl _divide_error, _debug, _nmi, _int3, _overflow, _bounds, _invalid_op
15 .globl _double_fault, _coprocessor_segment_overrun
16 .globl _invalid_TSS, _segment_not_present, _stack_segment
17 .globl _general_protection, _coprocessor_error, _irq13, _reserved
18
   # int0 -- （下面这段代码的含义参见图 4.1(a)）。
   # 下面是被零除出错(divide_error)处理代码。标号'_divide_error'实际上是 C 语言函
   # 数 divide_error() 编译后所生成模块中对应的名称。'_do_divide_error' 函数在 traps.c 中。
19 _divide_error:
20     pushl $_do_divide_error # 首先把将要调用的函数地址入栈。
21 no_error_code:             # 这段程序的出错号为 0。
22     xchgl %eax, (%esp)      # _do_divide_error 的地址 → eax, eax 被交换入栈。
23     pushl %ebx
24     pushl %ecx
25     pushl %edx
26     pushl %edi
27     pushl %esi
28     pushl %ebp
29     push %ds                # !! 16 位的段寄存器入栈后也要占用 4 个字节。
30     push %es
31     push %fs
32     pushl $0                # "error code" # 将出错码入栈。
33     lea 44(%esp), %edx      # 取原调用返回地址处堆栈指针位置，并压入堆栈。
34     pushl %edx
35     movl $0x10, %edx        # 内核代码数据段选择符。
36     mov %dx, %ds
37     mov %dx, %es
38     mov %dx, %fs
39     call *%eax              # 调用 C 函数 do_divide_error()。
40     addl $8, %esp           # 让堆栈指针重新指向寄存器 fs 入栈处。
41     pop %fs
42     pop %es
43     pop %ds
44     popl %ebp
45     popl %esi
46     popl %edi
47     popl %edx
48     popl %ecx
49     popl %ebx
50     popl %eax              # 弹出原来 eax 中的内容。
51     iret
52
   # int1 -- debug 调试中断入口点。处理过程同上。
53 _debug:
54     pushl $_do_int3         # _do_debug

```

```

55         jmp no_error_code
56
57     # int2 -- 非屏蔽中断调用入口点。
58     _nmi:
59         pushl $_do_nmi
60         jmp no_error_code
61
62     # int3 -- 同_debug。
63     _int3:
64         pushl $_do_int3
65         jmp no_error_code
66
67     # int4 -- 溢出出错处理中断入口点。
68     _overflow:
69         pushl $_do_overflow
70         jmp no_error_code
71
72     # int5 -- 边界检查出错中断入口点。
73     _bounds:
74         pushl $_do_bounds
75         jmp no_error_code
76
77     # int6 -- 无效操作指令出错中断入口点。
78     _invalid_op:
79         pushl $_do_invalid_op
80         jmp no_error_code
81
82     # int9 -- 协处理器段超出出错中断入口点。
83     _coprocessor_segment_overrun:
84         pushl $_do_coprocessor_segment_overrun
85         jmp no_error_code
86
87     # int15 --
88     _reserved:
89         pushl $_do_reserved
90         jmp no_error_code
91
92     # int45 -- (32+13) 数学协处理器发出的中断。
93     _irq13:
94         pushl %eax
95         xorb %al,%al
96         outb %al,$0xF0
97         movb $0x20,%al
98         outb %al,$0x20
99         jmp 1f                                # 这两个跳转指令起延时作用。
100    1:     jmp 1f
101    1:     outb %al,$0xA0
102         popl %eax
103         jmp _coprocessor_error # _coprocessor_error 原来在本文件中，现在已经放到
                                # (kernel/system_call.s, 131)
104
105     # 以下中断在调用时会在中断返回地址之后将出错号压入堆栈，因此返回时也需要将出错号弹出。
106     # int8 -- 双出错故障。（下面这段代码的含义参见图 4.1(b)）。

```

```

97 _double_fault:
98     pushl $_do_double_fault      # C 函数地址入栈。
99 error_code:
100     xchgl %eax, 4(%esp)          # error code <-> %eax
101     xchgl %ebx, (%esp)          # &function <-> %ebx
102     pushl %ecx
103     pushl %edx
104     pushl %edi
105     pushl %esi
106     pushl %ebp
107     push %ds
108     push %es
109     push %fs
110     pushl %eax                  # error code    # 出错号入栈。
111     lea 44(%esp), %eax          # offset      # 程序返回地址处堆栈指针位置值入栈。
112     pushl %eax
113     movl $0x10, %eax           # 置内核数据段选择符。
114     mov %ax, %ds
115     mov %ax, %es
116     mov %ax, %fs
117     call *%ebx                 # 调用相应的 C 函数，其参数已入栈。
118     addl $8, %esp              # 堆栈指针重新指向栈中放置 fs 内容的位置。
119     pop %fs
120     pop %es
121     pop %ds
122     popl %ebp
123     popl %esi
124     popl %edi
125     popl %edx
126     popl %ecx
127     popl %ebx                  # 弹出 C 函数地址。
128     popl %eax                  # 弹出出错号。
129     iret
130
    # int10 -- 无效的任务状态段(TSS)。
131 _invalid_TSS:
132     pushl $_do_invalid_TSS
133     jmp error_code
134
    # int11 -- 段不存在。
135 _segment_not_present:
136     pushl $_do_segment_not_present
137     jmp error_code
138
    # int12 -- 堆栈段错误。
139 _stack_segment:
140     pushl $_do_stack_segment
141     jmp error_code
142
    # int13 -- 一般保护性出错。
143 _general_protection:
144     pushl $_do_general_protection
145     jmp error_code

```

146

```
# int7 -- 设备不存在(_device_not_available)在(kernel/system_call.s,148)
# int14 -- 页错误(_page_fault)在(mm/page.s,14)
# int16 -- 协处理器错误(_coprocessor_error)在(kernel/system_call.s,131)
# 时钟中断 int 0x20 (_timer_interrupt)在(kernel/system_call.s,176)
# 系统调用 int 0x80 (_system_call)在(kernel/system_call.s,80)
```

### 4.3.3 其它信息

#### 4.3.3.1 Intel 保留中断向量的定义

这里给出了 Intel 保留中断向量具体含义的说明，见表 4.2 所示。

表 4.2 Intel 保留的中断号含义

| 中断号 | 名称                          | 类型       | 信号      | 说明   |
|-----|-----------------------------|----------|---------|--|
| 0   | Device error                | 故障       | SIGFPE  | 当进行除以零的操作时产生。  |
| 1   | Debug                       | 陷阱<br>故障 | SIGTRAP | 当进行程序单步跟踪调试时，设置了标志寄存器 eflags 的 T 标志时产生这个中断。  |
| 2   | nmi                         | 硬件       |         | 由不可屏蔽中断 NMI 产生。  |
| 3   | Breakpoint                  | 陷阱       | SIGTRAP | 由断点指令 int3 产生，与 debug 处理相同。  |
| 4   | Overflow                    | 陷阱       | SIGSEGV | eflags 的溢出标志 OF 引起。  |
| 5   | Bounds check                | 故障       | SIGSEGV | 寻址到有效地址以外时引起。  |
| 6   | Invalid Opcode              | 故障       | SIGILL  | CPU 执行时发现一个无效的指令操作码。   |
| 7   | Device not available        | 故障       | SIGSEGV | 设备不存在，指协处理器。在两种情况下会产生该中断：(a)CPU 遇到一个转意指令并且 EM 置位时。在这种情况下处理程序应该模拟导致异常的指令。(b)MP 和 TS 都在置位状态时，CPU 遇到 WAIT 或一个转意指令。在这种情况下，处理程序在必要时应该更新协处理器的状态。 |
| 8   | Double fault                | 故障       | SIGSEGV | 双故障出错。   |
| 9   | Coprocessor segment overrun | 故障       | SIGFPE  | 协处理器段超出。   |
| 10  | Invalid TSS                 | 故障       | SIGSEGV | CPU 切换时发觉 TSS 无效。  |
| 11  | Segment not present         | 故障       | SIGBUS  | 描述符所指的段不存在。  |
| 12  | Stack segment               | 故障       | SIGBUS  | 堆栈段不存在或寻址越出堆栈段。  |
| 13  | General protection          | 故障       | SIGSEGV | 没有符合 80386 保护机制（特权级）的操作引起。   |
| 14  | Page fault                  | 故障       | SIGSEGV | 页不在内存。   |
| 15  | Reserved                    |          |         |  |
| 16  | Coprocessor error           | 故障       | SIGFPE  | 协处理器发出的发出的出错信号引起。  |

## 4.4 system\_call.s 程序

### 4.4.1 功能描述

本程序主要实现系统调用(system\_call)中断 int 0x80 的入口处理过程以及信号检测处理，同时给出

了两个系统功能的底层接口, 分别是 `sys_execve` 和 `sys_fork`。还列出了处理过程类似的协处理器出错(int 16)、设备不存在(int7)、时钟中断(int32)、硬盘中断(int46)、软盘中断(int38)的中断处理程序。

对于软中断(system\_call、coprocessor\_error、device\_not\_available), 处理过程基本上是首先为调用相应 C 函数处理程序作准备, 将一些参数压入堆栈, 然后调用 C 函数进行相应功能的处理, 处理返回后再去检测当前任务的信号位图, 对值最小的一个信号进行处理并复位位图中的该信号。

对于硬件中断请求信号 IRQ 发来的中断, 其处理过程首先是向中断控制芯片 8259A 发送结束硬件中断控制字指令 EOI, 然后调用相应的 C 函数处理程序。对于时钟中断也要对当前任务的信号位图进行检测处理。

## 4.4.2 代码注释

列表 4.4 linux/kernel/system\_call.s 程序

---

```

1  /*
2  *  linux/kernel/system_call.s
3  *
4  *  (C) 1991  Linus Torvalds
5  */
6
7  /*
8  *  system_call.s  contains the system-call low-level handling routines.
9  *  This also contains the timer-interrupt handler, as some of the code is
10 *  the same. The hd- and floppy-interrupts are also here.
11 *
12 *  NOTE: This code handles signal-recognition, which happens every time
13 *  after a timer-interrupt and after each system call. Ordinary interrupts
14 *  don't handle signal-recognition, as that would clutter them up totally
15 *  unnecessarily.
16 *
17 *  Stack layout in 'ret_from_system_call':
18 *
19 *      0(%esp) - %eax
20 *      4(%esp) - %ebx
21 *      8(%esp) - %ecx
22 *      C(%esp) - %edx
23 *      10(%esp) - %fs
24 *      14(%esp) - %es
25 *      18(%esp) - %ds
26 *      1C(%esp) - %eip
27 *      20(%esp) - %cs
28 *      24(%esp) - %eflags
29 *      28(%esp) - %oldesp
30 *      2C(%esp) - %oldss
31 */
/*
 * system_call.s 文件包含系统调用(system-call)底层处理子程序。由于有些代码比较类似, 所以
 * 同时也包括时钟中断处理(timer-interrupt)句柄。硬盘和软盘的中断处理程序也在这里。
 *
 * 注意: 这段代码处理信号(signal)识别, 在每次时钟中断和系统调用之后都会进行识别。一般

```



```

* 中断信号并不处理信号识别，因为会给系统造成混乱。
*
* 从系统调用返回 ('ret_from_system_call') 时堆栈的内容见上面 19-30 行。
*/

32
33 SIG_CHLD          = 17          # 定义一些位移常数符号。
34
35 EAX               = 0x00
36 EBX               = 0x04
37 ECX               = 0x08
38 EDX               = 0x0C
39 FS                = 0x10
40 ES                = 0x14
41 DS                = 0x18
42 EIP               = 0x1C
43 CS                = 0x20
44 EFLAGS            = 0x24
45 OLDESP            = 0x28
46 OLDSS             = 0x2C
47
# 以下这些是任务结构(task_struct)中变量的偏移值。
48 state = 0          # these are offsets into the task-struct. # 进程状态码
49 counter = 4        #
50 priority = 8
51 signal = 12
52 sigaction = 16      # MUST be 16 (=len of sigaction)
53 blocked = (33*16)   # 阻塞的信号位图偏移量。
54
# 以下定义在 sigaction 结构中的偏移量。
55 # offsets within sigaction
56 sa_handler = 0
57 sa_mask = 4
58 sa_flags = 8
59 sa_restorer = 12
60
61 nr_system_calls = 72 # Linux 0.11 版内核中的系统调用总数。
62
63 /*
64 * Ok, I get parallel printer interrupts while using the floppy for some
65 * strange reason. Urgel. Now I just ignore them.
66 */
/*
* 好了，在使用软驱时我收到了并行打印机中断，很奇怪。呵，现在不管它。
*/
# 定义入口点。
67 .globl _system_call, _sys_fork, _timer_interrupt, _sys_execve
68 .globl _hd_interrupt, _floppy_interrupt, _parallel_interrupt
69 .globl _device_not_available, _coprocessor_error
70
# 错误的系统调用号。
71 .align 2          # 内存 4 字节对齐。
72 bad_sys_call:
73     movl $-1, %eax    # eax 中置-1，退出中断。

```

```

74         iret
# 重新执行调度程序入口。调度程序 schedule 在(kernel/sched.c, 104)。
75 .align 2
76 reschedule:
77     pushl $ret_from_sys_call # 将 ret_from_sys_call 的地址入栈。
78     jmp _schedule
#### int 0x80 --linux 系统调用入口点(调用中断 int 0x80, 并在 eax 中放入调用号)。
79 .align 2
80 _system_call:
81     cmpl $nr_system_calls-1,%eax # 调用号如果超出范围的话就在 eax 中置-1 并退出。
82     ja bad_sys_call
83     push %ds                    # 保存原段寄存器值。
84     push %es
85     push %fs
86     pushl %edx                 # ebx, ecx, edx 中放着系统调用相应的 C 语言函数的调用参数。
87     pushl %ecx                 # push %ebx, %ecx, %edx as parameters
88     pushl %ebx                 # to the system call
89     movl $0x10, %edx           # set up ds, es to kernel space
90     mov %dx, %ds               # ds, es 指向内核数据段(全局描述符表中数据段描述符)。
91     mov %dx, %es
92     movl $0x17, %edx           # fs points to local data space
93     mov %dx, %fs               # fs 指向局部数据段(局部描述符表中数据段描述符)。
# 下面这句操作数的含义是: 调用地址 = _sys_call_table + %eax * 4。参见列表后的说明。
# 对应的 C 程序中的 sys_call_table 在 include/linux/sys.h 中, 其中定义了一个包括 72 个
# 系统调用 C 处理函数的地址数组表。
94     call _sys_call_table(, %eax, 4)
95     pushl %eax                 # 把系统调用号入栈。
96     movl _current, %eax        # 取当前任务(进程)数据结构地址→eax。
# 下面 97-100 行查看当前任务的运行状态。如果不在就绪状态(state 不等于 0)就去执行调度程序。
# 如果该任务在就绪状态但 counter[??]值等于 0, 则也去执行调度程序。
97     cmpl $0, state(%eax)       # state
98     jne reschedule
99     cmpl $0, counter(%eax)     # counter
100    je reschedule
# 以下这段代码执行从系统调用 C 函数返回后, 对信号量进行识别处理。
101 ret_from_sys_call:
# 首先判别当前任务是否是初始任务 task0, 如果是则不必对其进行信号量方面的处理, 直接返回。
# 103 行上的 _task 对应 C 程序中的 task[] 数组, 直接引用 task 相当于引用 task[0]。
102    movl _current, %eax         # task[0] cannot have signals
103    cmpl _task, %eax
104    je 3f                      # 向前(forward)跳转到标号 3。
# 通过对原调用程序代码选择符的检查来判断调用程序是否是超级用户。如果是超级用户就直接
# 退出中断, 否则需进行信号量的处理。这里比较选择符是否为普通用户代码段的选择符 0x000f
# (RPL=3, 局部表, 第 1 个段(代码段)), 如果不是则跳转退出中断程序。
105    cmpw $0x0f, CS(%esp)        # was old code segment supervisor ?
106    jne 3f
# 如果原堆栈段选择符不为 0x17 (也即堆栈不在用户数据段中), 则也退出。
107    cmpw $0x17, OLDSS(%esp)     # was stack segment = 0x17 ?
108    jne 3f
# 下面这段代码(109-120)的用途是首先取当前任务结构中的信号位图(32 位, 每位代表 1 种信号),
# 然后用任务结构中的信号阻塞(屏蔽)码, 阻塞不允许的信号位, 取得数值最小的信号值, 再把
# 原信号位图中该信号对应的位复位(置 0), 最后将该信号值作为参数之一调用 do_signal()。
# do_signal() 在(kernel/signal.c, 82)中, 其参数包括 13 个入栈的信息。

```

```

109      movl signal(%eax),%ebx      # 取信号位图→ebx, 每1位代表1种信号, 共32个信号。
110      movl blocked(%eax),%ecx    # 取阻塞(屏蔽)信号位图→ecx。
111      notl %ecx                  # 每位取反。
112      andl %ebx,%ecx             # 获得许可的信号位图。
113      bsfl %ecx,%ecx             # 从低位(位0)开始扫描位图, 看是否有1的位,
                                   # 若有, 则ecx保留该位的偏移值(即第几位0-31)。
114      je 3f                      # 如果没有信号则向前跳转退出。
115      btrl %ecx,%ebx             # 复位该信号(ebx含有原signal位图)。
116      movl %ebx,signal(%eax)     # 重新保存signal位图信息→current->signal。
117      incl %ecx                  # 将信号调整为从1开始的数(1-32)。
118      pushl %ecx                 # 信号值入栈作为调用do_signal的参数之一。
119      call _do_signal            # 调用C函数信号处理程序(kernel/signal.c, 82)
120      popl %eax                  # 弹出信号值。
121 3:      popl %eax
122      popl %ebx
123      popl %ecx
124      popl %edx
125      pop %fs
126      pop %es
127      pop %ds
128      iret
129
#### int16 -- 下面这段代码处理协处理器发出的出错信号。跳转执行C函数math_error()
# (kernel/math/math_emulate.c, 82), 返回后将跳转到ret_from_sys_call处继续执行。
130 .align 2
131 _coprocessor_error:
132     push %ds
133     push %es
134     push %fs
135     pushl %edx
136     pushl %ecx
137     pushl %ebx
138     pushl %eax
139     movl $0x10,%eax             # ds, es 置为指向内核数据段。
140     mov %ax,%ds
141     mov %ax,%es
142     movl $0x17,%eax             # fs 置为指向局部数据段(出错程序的数据段)。
143     mov %ax,%fs
144     pushl $ret_from_sys_call    # 把下面调用返回的地址入栈。
145     jmp _math_error             # 执行C函数math_error() (kernel/math/math_emulate.c, 37)
146
#### int7 -- 设备不存在或协处理器不存在(Coprocessor not available)。
# 如果控制寄存器CR0的EM标志置位, 则当CPU执行一个ESC转义指令时就会引发该中断, 这样就
# 可以有机会让这个中断处理程序模拟ESC转义指令(169行)。
# CR0的TS标志是在CPU执行任务转换时设置的。TS可以用来确定什么时候协处理器中的内容(上下文)
# 与CPU正在执行的任务不匹配了。当CPU在运行一个转义指令时发现TS置位了, 就会引发该中断。
# 此时就应该恢复新任务的协处理器执行状态(165行)。参见(kernel/sched.c, 77)中的说明。
# 该中断最后将转移到标号ret_from_sys_call处执行下去(检测并处理信号)。
147 .align 2
148 _device_not_available:
149     push %ds
150     push %es
151     push %fs

```

```

152     pushl %edx
153     pushl %ecx
154     pushl %ebx
155     pushl %eax
156     movl $0x10,%eax      # ds, es 置为指向内核数据段。
157     mov %ax,%ds
158     mov %ax,%es
159     movl $0x17,%eax      # fs 置为指向局部数据段（出错程序的数据段）。
160     mov %ax,%fs
161     pushl $ret_from_sys_call # 把下面跳转或调用的返回地址入栈。
162     clls                  # clear TS so that we can use math
163     movl %cr0,%eax
164     testl $0x4,%eax      # EM (math emulation bit)
                            # 如果不是 EM 引起的中断，则恢复新任务协处理器状态，
165     je _math_state_restore # 执行 C 函数 math_state_restore() (kernel/sched.c, 77)。
166     pushl %ebp
167     pushl %esi
168     pushl %edi
169     call _math_emulate    # 调用 C 函数 math_emulate(kernel/math/math_emulate.c, 18)。
170     popl %edi
171     popl %esi
172     popl %ebp
173     ret
174
##### int32 -- (int 0x20) 时钟中断处理程序。中断频率被设置为 100Hz(include/linux/sched.h, 5),
# 定时芯片 8253/8254 是在(kernel/sched.c, 406)处初始化的。因此这里 jiffies 每 10 毫秒加 1。
# 这段代码将 jiffies 增 1，发送结束中断指令给 8259 控制器，然后用当前特权级作为参数调用
# C 函数 do_timer(long CPL)。当调用返回时转去检测并处理信号。
175 .align 2
176 _timer_interrupt:
177     push %ds              # save ds, es and put kernel data space
178     push %es              # into them. %fs is used by _system_call
179     push %fs
180     pushl %edx            # we save %eax,%ecx,%edx as gcc doesn't
181     pushl %ecx            # save those across function calls. %ebx
182     pushl %ebx            # is saved as we use that in ret_sys_call
183     pushl %eax
184     movl $0x10,%eax      # ds, es 置为指向内核数据段。
185     mov %ax,%ds
186     mov %ax,%es
187     movl $0x17,%eax      # fs 置为指向局部数据段（出错程序的数据段）。
188     mov %ax,%fs
189     incl _jiffies
# 由于初始化中断控制芯片时没有采用自动 EOI，所以这里需要发指令结束该硬件中断。
190     movb $0x20,%al       # EOI to interrupt controller #1
191     outb %al,$0x20        # 操作命令字 OCW2 送 0x20 端口。
# 下面 3 句从选择符中取出当前特权级别(0 或 3)并压入堆栈，作为 do_timer 的参数。
192     movl CS(%esp),%eax
193     andl $3,%eax         # %eax is CPL (0 or 3, 0=supervisor)
194     pushl %eax
195     call _do_timer        # 'do_timer(long CPL)' does everything from
196     addl $4,%esp          # task switching to accounting ...
197     jmp ret_from_sys_call

```

```

198 ##### 这是 sys_execve() 系统调用。取中断调用程序的代码指针作为参数调用 C 函数 do_execve()。
    # do_execve() 在 (fs/exec.c, 182)。
199 .align 2
200 _sys_execve:
201     lea EIP(%esp), %eax
202     pushl %eax
203     call _do_execve
204     addl $4, %esp          # 丢弃调用时压入栈的 EIP 值。
205     ret
206 ##### sys_fork() 调用，用于创建子进程，是 system_call 功能 2。原形在 include/linux/sys.h 中。
    # 首先调用 C 函数 find_empty_process()，取得一个进程号 pid。若返回负数则说明目前任务数组
    # 已满。然后调用 copy_process() 复制进程。
207 .align 2
208 _sys_fork:
209     call _find_empty_process    # 调用 find_empty_process() (kernel/fork.c, 135)。
210     testl %eax, %eax
211     js 1f
212     push %gs
213     pushl %esi
214     pushl %edi
215     pushl %ebp
216     pushl %eax
217     call _copy_process          # 调用 C 函数 copy_process() (kernel/fork.c, 68)。
218     addl $20, %esp             # 丢弃这里所有压栈内容。
219 1:     ret
220 ##### int 46 -- (int 0x2E) 硬盘中断处理程序，响应硬件中断请求 IRQ14。(kernel/blk_drv/hd.c)。
    # 首先向 8259A 中断控制从芯片发送结束硬件中断指令 (EOI)，然后取变量 do_hd 中的函数指针放入 edx
    # 寄存器中，并置 do_hd 为 NULL，接着判断 edx 函数指针是否为空。如果为空，则给 edx 赋值指向
    # unexpected_hd_interrupt()，用于显示出错信息。随后向 8259A 主芯片送 EOI 指令，并调用 edx 中
    # 指针指向的函数：read_intr()、write_intr() 或 unexpected_hd_interrupt()。
221 _hd_interrupt:
222     pushl %eax
223     pushl %ecx
224     pushl %edx
225     push %ds
226     push %es
227     push %fs
228     movl $0x10, %eax           # ds, es 置为内核数据段。
229     mov %ax, %ds
230     mov %ax, %es
231     movl $0x17, %eax           # fs 置为调用程序的局部数据段。
232     mov %ax, %fs
    # 由于初始化中断控制芯片时没有采用自动 EOI，所以这里需要发指令结束该硬件中断。
233     movb $0x20, %al
234     outb %al, $0xA0             # EOI to interrupt controller #1    # 送从 8259A。
235     jmp 1f                     # give port chance to breathe
236 1:     jmp 1f                     # 延时作用。
237 1:     xorl %edx, %edx
238     xchgl _do_hd, %edx          # do_hd 定义为一个函数指针，将被赋值 read_intr() 或
    # write_intr() 函数地址。(kernel/blk_drv/hd.c)

```

```

239         testl %edx,%edx          # 放到 edx 寄存器后就将 do_hd 指针变量置为 NULL。
240         jne 1f                    # 测试函数指针是否为 Null。
241         movl $_unexpected_hd_interrupt,%edx  # (kernel/blk_drv/hdc, 237)。
242 1:      outb %al,$0x20            # 送主 8259A 中断控制器 EOI 指令（结束硬件中断）。
243         call *%edx                # "interesting" way of handling intr.
244         pop %fs                   # 上句调用 do_hd 指向的 C 函数。
245         pop %es
246         pop %ds
247         popl %edx
248         popl %ecx
249         popl %eax
250         iret
251
##### int38 -- (int 0x26) 软盘驱动器中断处理程序，响应硬件中断请求 IRQ6。
# 其处理过程与上面对硬盘的处理基本一样。(kernel/blk_drv/floppy.c)。
# 首先向 8259A 中断控制器主芯片发送 EOI 指令，然后取变量 do_floppy 中的函数指针放入 eax
# 寄存器中，并置 do_floppy 为 NULL，接着判断 eax 函数指针是否为空。如为空，则给 eax 赋值指向
# unexpected_floppy_interrupt ()，用于显示出错信息。随后调用 eax 指向的函数：rw_interrupt,
# seek_interrupt, recal_interrupt, reset_interrupt 或 unexpected_floppy_interrupt。
252 _floppy_interrupt:
253     pushl %eax
254     pushl %ecx
255     pushl %edx
256     push %ds
257     push %es
258     push %fs
259     movl $0x10,%eax              # ds, es 置为内核数据段。
260     mov %ax,%ds
261     mov %ax,%es
262     movl $0x17,%eax              # fs 置为调用程序的局部数据段。
263     mov %ax,%fs
264     movb $0x20,%al               # 送主 8259A 中断控制器 EOI 指令（结束硬件中断）。
265     outb %al,$0x20               # EOI to interrupt controller #1
266     xorl %eax,%eax
267     xchgl _do_floppy,%eax        # do_floppy 为一函数指针，将被赋值实际处理 C 函数程序，
                                   # 放到 eax 寄存器后就将 do_floppy 指针变量置空。
268     testl %eax,%eax              # 测试函数指针是否=NULL?
269     jne 1f                       # 若空，则使指针指向 C 函数 unexpected_floppy_interrupt()。
270     movl $_unexpected_floppy_interrupt,%eax
271 1:      call *%eax                # "interesting" way of handling intr.
272     pop %fs                      # 上句调用 do_floppy 指向的函数。
273     pop %es
274     pop %ds
275     popl %edx
276     popl %ecx
277     popl %eax
278     iret
279
##### int 39 -- (int 0x27) 并行口中断处理程序，对应硬件中断请求信号 IRQ7。
# 本版本内核还未实现。这里只是发送 EOI 指令。
280 _parallel_interrupt:
281     pushl %eax

```





```

6
7 /*
8  * 'Traps.c' handles hardware traps and faults after we have saved some
9  * state in 'asm.s'. Currently mostly a debugging-aid, will be extended
10 * to mainly kill the offending process (probably by giving it a signal,
11 * but possibly by killing it outright if necessary).
12 */
13 /*
14  * 在程序 asm.s 中保存了一些状态后，本程序用来处理硬件陷阱和故障。目前主要用于调试目的，
15  * 以后将扩展用来杀死遭损坏的进程（主要是通过发送一个信号，但如果必要也会直接杀死）。
16 */
17 #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。
18 #include <linux/head.h>      // head 头文件，定义了与 head.s 程序描述符表有关的 C 程序。
19 #include <linux/sched.h>     // 调度程序头文件，定义了任务结构 task_struct、初始任务 0 的数据，
20                               // 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
21 #include <linux/kernel.h>    // 内核头文件。含有一些内核常用函数的原形定义。
22 #include <asm/system.h>      // 系统头文件。定义了设置或修改描述符/中断门等的嵌入式汇编宏。
23 #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。
24 #include <asm/io.h>          // 输入/输出头文件。定义硬件端口输入/输出宏汇编语句。
25
26 // 以下语句定义了三个嵌入式汇编宏语句函数。有关嵌入式汇编的基本语法见列表后或参见附录。
27 // 取段 seg 中地址 addr 处的一个字节。
28 #define get_seg_byte(seg, addr) ({ \
29     register char __res; \
30     __asm__ ("push %%fs;mov %%ax,%%fs;movb %%fs:%2,%%al;pop %%fs" \
31             : "=a" (__res): "" (seg), "m" (*(addr))); \
32     __res;})
33
34 // 取段 seg 中地址 addr 处的一个长字（4 字节）。
35 #define get_seg_long(seg, addr) ({ \
36     register unsigned long __res; \
37     __asm__ ("push %%fs;mov %%ax,%%fs;movl %%fs:%2,%%eax;pop %%fs" \
38             : "=a" (__res): "" (seg), "m" (*(addr))); \
39     __res;})
40
41 // 取 fs 段寄存器的值（选择符）。
42 #define fs() ({ \
43     register unsigned short __res; \
44     __asm__ ("mov %%fs,%%ax": "=a" (__res):); \
45     __res;})
46
47 // 以下定义了一些函数原型。
48 int do_exit(long code);          // (kernel/exit.c, 102)
49
50 void page_exception(void);       // [??]
51
52 // 以下定义了一些中断处理程序原型，代码在 (kernel/asm.s 或 system_call.s) 中。
53 void divide_error(void);         // int0 (kernel/asm.s, 19)。
54 void debug(void);               // int1 (kernel/asm.s, 53)。
55 void nmi(void);                 // int2 (kernel/asm.s, 57)。
56 void int3(void);                // int3 (kernel/asm.s, 61)。
57 void overflow(void);            // int4 (kernel/asm.s, 65)。

```



```

48 void bounds(void); // int5 (kernel/asm.s, 69)。
49 void invalid\_op(void); // int6 (kernel/asm.s, 73)。
50 void device\_not\_available(void); // int7 (kernel/system_call.s, 148)。
51 void double\_fault(void); // int8 (kernel/asm.s, 97)。
52 void coprocessor\_segment\_overrun(void); // int9 (kernel/asm.s, 77)。
53 void invalid\_TSS(void); // int10 (kernel/asm.s, 131)。
54 void segment\_not\_present(void); // int11 (kernel/asm.s, 135)。
55 void stack\_segment(void); // int12 (kernel/asm.s, 139)。
56 void general\_protection(void); // int13 (kernel/asm.s, 143)。
57 void page\_fault(void); // int14 (mm/page.s, 14)。
58 void coprocessor\_error(void); // int16 (kernel/system_call.s, 131)。
59 void reserved(void); // int15 (kernel/asm.s, 81)。
60 void parallel\_interrupt(void); // int39 (kernel/system_call.s, 280)。
61 void irq13(void); // int45 协处理器中断处理(kernel/asm.s, 85)。
62
// 该子程序用来打印出错中断的名称、出错号、调用程序的EIP、EFLAGS、ESP、fs段寄存器值、
// 段的基址、段的长度、进程号pid、任务号、10字节指令码。如果堆栈在用户数据段，则还
// 打印16字节的堆栈内容。
63 static void die(char * str, long esp_ptr, long nr)
64 {
65     long * esp = (long *) esp_ptr;
66     int i;
67
68     printk("s: %04x\n\r", str, nr&0xffff);
69     printk("EIP: |t%04x:%p\nEFLAGS: |t%p\nESP: |t%04x:%p\n",
70         esp[1], esp[0], esp[2], esp[4], esp[3]);
71     printk("fs: %04x\n", fs());
72     printk("base: %p, limit: %p\n", get\_base(current->ldt[1]), get\_limit(0x17));
73     if (esp[4] == 0x17) {
74         printk("Stack: ");
75         for (i=0; i<4; i++)
76             printk("p ", get\_seg\_long(0x17, i+(long *) esp[3]));
77         printk("n");
78     }
79     str(i);
80     printk("Pid: %d, process nr: %d\n\r", current->pid, 0xffff & i);
81     for(i=0; i<10; i++)
82         printk("02x ", 0xff & get\_seg\_byte(esp[1], (i+(char *) esp[0])));
83     printk("n\r");
84     do\_exit(11); // play segment exception
85 }
86
// 以下这些以do开头的函数是对应名称中断处理程序调用的C函数。
87 void do\_double\_fault(long esp, long error_code)
88 {
89     die("double fault", esp, error_code);
90 }
91
92 void do\_general\_protection(long esp, long error_code)
93 {
94     die("general protection", esp, error_code);
95 }
96

```

```

97 void do\_divide\_error(long esp, long error_code)
98 {
99     die("divide error", esp, error_code);
100 }
101
102 void do\_int3(long * esp, long error_code,
103             long fs, long es, long ds,
104             long ebp, long esi, long edi,
105             long edx, long ecx, long ebx, long eax)
106 {
107     int tr;
108
109     __asm__ ("str %%ax": "=a" (tr): "" (0)); // 取任务寄存器值→tr。
110     printk("eax|t|tebx|t|tecX|t|tedx|n|r%8x|t%8x|t%8x|t%8x|n|r",
111          eax, ebx, ecx, edx);
112     printk("esi|t|tedi|t|tebp|t|tesp|n|r%8x|t%8x|t%8x|t%8x|n|r",
113          esi, edi, ebp, (long) esp);
114     printk("|n|rds|tes|tfs|ttr|n|r%4x|t%4x|t%4x|t%4x|n|r",
115          ds, es, fs, tr);
116     printk("EIP: %8x   CS: %4x   EFLAGS: %8x|n|r", esp[0], esp[1], esp[2]);
117 }
118
119 void do\_nmi(long esp, long error_code)
120 {
121     die("nmi", esp, error_code);
122 }
123
124 void do\_debug(long esp, long error_code)
125 {
126     die("debug", esp, error_code);
127 }
128
129 void do\_overflow(long esp, long error_code)
130 {
131     die("overflow", esp, error_code);
132 }
133
134 void do\_bounds(long esp, long error_code)
135 {
136     die("bounds", esp, error_code);
137 }
138
139 void do\_invalid\_op(long esp, long error_code)
140 {
141     die("invalid operand", esp, error_code);
142 }
143
144 void do\_device\_not\_available(long esp, long error_code)
145 {
146     die("device not available", esp, error_code);
147 }
148
149 void do\_coprocessor\_segment\_overrun(long esp, long error_code)

```

```

150 {
151     die("coprocessor segment overrun", esp, error_code);
152 }
153
154 void do_invalid_TSS(long esp, long error_code)
155 {
156     die("invalid TSS", esp, error_code);
157 }
158
159 void do_segment_not_present(long esp, long error_code)
160 {
161     die("segment not present", esp, error_code);
162 }
163
164 void do_stack_segment(long esp, long error_code)
165 {
166     die("stack segment", esp, error_code);
167 }
168
169 void do_coprocessor_error(long esp, long error_code)
170 {
171     if (last_task_used_math != current)
172         return;
173     die("coprocessor error", esp, error_code);
174 }
175
176 void do_reserved(long esp, long error_code)
177 {
178     die("reserved (15, 17-47) error", esp, error_code);
179 }
180
181 // 下面是异常（陷阱）中断程序初始化子程序。设置它们的中断调用门（中断向量）。
182 // set_trap_gate() 与 set_system_gate() 的主要区别在于前者设置的特权级为 0，后者是 3。因此
183 // 断点陷阱中断 int3、溢出中断 overflow 和边界出错中断 bounds 可以由任何程序产生。
184 // 这两个函数均是嵌入式汇编宏程序(include/asm/system.h, 36, 39)。
185 void trap_init(void)
186 {
187     int i;
188
189     set_trap_gate(0, &divide_error);
190     set_trap_gate(1, &debug);
191     set_trap_gate(2, &nmi);
192     set_system_gate(3, &int3); /* int3-5 can be called from all */
193     set_system_gate(4, &overflow);
194     set_system_gate(5, &bounds);
195     set_trap_gate(6, &invalid_op);
196     set_trap_gate(7, &device_not_available);
197     set_trap_gate(8, &double_fault);
198     set_trap_gate(9, &coprocessor_segment_overrun);
199     set_trap_gate(10, &invalid_TSS);
200     set_trap_gate(11, &segment_not_present);
201     set_trap_gate(12, &stack_segment);
202     set_trap_gate(13, &general_protection);

```

```

199     set_trap_gate(14, &page_fault);
200     set_trap_gate(15, &reserved);
201     set_trap_gate(16, &coprocessor_error);
202     for (i=17; i<48; i++)
203         set_trap_gate(i, &reserved); // 将 int17-48 的陷阱门均设置为 reserved。
204     set_trap_gate(45, &irq13);
205     outb_p(inb_p(0x21)&0xfb, 0x21); // 允许主 8259A 芯片的 IRQ2 中断请求。
206     outb(inb_p(0xA1)&0xdf, 0xA1); // 允许从 8259A 芯片的 IRQ13 中断请求。
207     set_trap_gate(39, &parallel_interrupt);
208 }
209

```

## 4.5.3 其它信息

### 4.5.3.1 嵌入式汇编的基本格式

具有输入和输出参数的嵌入式汇编的基本格式为：

```

asm(“汇编语句”
    : 输出寄存器
    : 输入寄存器
    : 会被修改的寄存器 );

```

这里给出一个例子，然后进行详细解说。

```

asm("cld\n\t" "rep\n\t" "stol"
    : /* 没有输出寄存器 */
    : "c"(count-1), "a"(fill_value), "D"(dest)
    : "%ecx", "%edi");

```

我们一点一点来看。

```
asm("cld\n\t" "rep\n\t" "stosl"
```

这三句是通常的汇编语句，以清方向位，重复保存值。

```
: /* 没有输出寄存器 */
```

说明这段嵌入汇编程序没有用到输出寄存器。

```
: "c"(count-1), "a"(fill_value), "D"(dest)
```

这里我们将 count-1 的值加载到 ecx 寄存器中，fill\_value 加载到 eax 中，dest 放到 edi 中。为什么要让 gcc 编译程序去做这样的加载而不让我们自己做呢？因为 gcc 在它进行寄存器分配时可以进行某些优化，例如 fill\_value 值可能已经在 eax 中。如果是在一个循环语句中的话，gcc 就可能在整个循环操作中保留 eax，这样就可以在每次循环中少用一个 movl 语句。

```
: "%ecx", "%edi");
```

这句告诉 gcc 这些寄存器中的值已经改变了。

很奇怪吧？不过在 gcc 知道你拿这些寄存器做些什么后，这确实能够对 gcc 的优化操作有所帮助。以下是一些你可能会用到的寄存器加载代码：

```

a    eax,      b    ebx,      c    ecx,      d    edx,      S    esi,      D    edi,
q, r 动态分配的寄存器,  g    eax, ebx, ecx, edx 或内存变量,      A    eax 与 edx 联合(64 位),

```

m 内存, o 内存并可加偏移值, I 常数(0-31), J 常数(0-63),  
 K 常数(0-255) L 常数(0-65535), M 常数(0-3), N 常数(0-255)(1字节常数),  
 O 常数(0-32)。

下面的例子不是让你自己指定哪个变量使用哪个寄存器, 而是让 gcc 为你选择。

```
asm("leal (%1, %1, 4), %0"
    : "=r"(y)
    : "0"(x));
```

这个例子可以非常快地将 x 乘 5。[ leal (r1, r2, 4), r3 语句表示  $r1+r2*4 \rightarrow r3$  ]

其中"%0", "%1"是指 gcc 自动分配的寄存器。编号的顺序是从输出寄存器序列从左到右从上到下以"%0"开始。所以这里"%1"代表输入值 x 要放入的寄存器, "%0"表示输出值寄存器。输出寄存器代码前一定要加等于号。如果输入寄存器的代码是 0 或为空时, 则说明使用与相应输出一样的寄存器。所以如果 gcc 将 r 指定为 eax 的话, 那么上面汇编语句的含义即为:

```
"leal (eax, eax, 4), eax"
```

注意: 在执行时如果汇编语句不希望被 gcc 优化而挪动地方, 就需要在 asm 符号后面添加 volatile 关键词:

```
asm volatile (.....);
```

或者更详细的说明为:

```
__asm__ __volatile__ (.....);
```

下面是几个简单例子。

```
#define disable() __asm__ __volatile__ ("cli");
#define enable() __asm__ __volatile__ ("sti");
#define times3(arg1, arg2) \
__asm__ ( \
    "leal (%0, %0 2), %0" \
    : "=r" (arg2) \
    : "0" (arg1) );
```

## 4.6 mktime.c 文件

### 4.6.1 功能描述

该程序只有一个函数 mktime(), 仅供内核使用。计算从 1970 年 1 月 1 日 0 时起到开机当日经过的秒数, 作为开机时间。

### 4.6.2 代码注释

列表 4.6 linux/kernel/mktime.c 程序

```
1 /*
2  * linux/kernel/mktime.c
3  *
```

```

4  * (C) 1991 Linus Torvalds
5  */
6
7  #include <time.h>          // 时间头文件，定义了标准时间数据结构 tm 和一些处理时间函数原型。
8
9  /*
10 * This isn't the library routine, it is only used in the kernel.
11 * as such, we don't care about years<1970 etc, but assume everything
12 * is ok. Similarly, TZ etc is happily ignored. We just do everything
13 * as easily as possible. Let's find something public for the library
14 * routines (although I think minix times is public).
15 */
16 /*
17 * PS. I hate whoever though up the year 1970 - couldn't they have gotten
18 * a leap-year instead? I also hate Gregorius, pope or no. I'm grumpy.
19 */
20 /*
21 * 这不是库函数，它仅供内核使用。因此我们不关心小于 1970 年的年份等，但假定一切均很正常。
22 * 同样，时间区域 TZ 问题也先忽略。我们只是尽可能简单地处理问题。最好能找到一些公开的库函数
23 * （尽管我认为 minix 的时间函数是公开的）。
24 * 另外，我恨那个设置 1970 年开始的人 - 难道他们就不能选择从一个闰年开始？我也恨罗马教皇、
25 * 主教或者什么也不在乎。我是个脾气暴躁的人。
26 */
27 #define MINUTE 60          // 1 分钟的秒数。
28 #define HOUR (60*MINUTE)  // 1 小时的秒数。
29 #define DAY (24*HOUR)     // 1 天的秒数。
30 #define YEAR (365*DAY)    // 1 年的秒数。
31
32 /* interestingly, we assume leap-years */
33 /* 有趣的是我们考虑进了闰年 */
34 // 下面以年为界限，定义了每个月开始时的秒数时间数组。
35 static int month[12] = {
36     0,
37     DAY*(31),
38     DAY*(31+29),
39     DAY*(31+29+31),
40     DAY*(31+29+31+30),
41     DAY*(31+29+31+30+31),
42     DAY*(31+29+31+30+31+30),
43     DAY*(31+29+31+30+31+30+31),
44     DAY*(31+29+31+30+31+30+31+31),
45     DAY*(31+29+31+30+31+30+31+31+30),
46     DAY*(31+29+31+30+31+30+31+31+30+31),
47     DAY*(31+29+31+30+31+30+31+31+30+31+30)
48 };
49
50 // 该函数计算从 1970 年 1 月 1 日 0 时起到开机当日经过的秒数，作为开机时间。
51 long kernel_mktime(struct tm * tm)
52 {
53     long res;
54     int year;
55
56     year = tm->tm_year - 70;          // 从 70 年到现在经过的年数(2 位表示方式),

```

```

// 因此会有 2000 年问题。
47 /* magic offsets (y+1) needed to get leapyears right. */
/* 为了获得正确的闰年数，这里需要这样一个魔幻偏值(y+1) */
48     res = YEAR*year + DAY*((year+1)/4); // 这些年经过的秒数时间 + 每个闰年时多 1 天
49     res += month[tm->tm_mon];           // 的秒数时间，在加上当年到当月时的秒数。
50 /* and (y+2) here. If it wasn't a leap-year, we have to adjust */
/* 以及(y+2)。如果(y+2)不是闰年，那么我们就必须进行调整(减去一天的秒数时间)。*/
51     if (tm->tm_mon>1 && ((year+2)%4))
52         res -= DAY;
53     res += DAY*(tm->tm_mday-1);           // 再加上本月过去的天数的秒数时间。
54     res += HOUR*tm->tm_hour;             // 再加上当天过去的小时数的秒数时间。
55     res += MINUTE*tm->tm_min;            // 再加上 1 小时内过去的分钟数的秒数时间。
56     res += tm->tm_sec;                   // 再加上 1 分钟内已过的秒数。
57     return res;                          // 即等于从 1970 年以来经过的秒数时间。
58 }
59

```

## 4.6.3 其它信息

### 4.6.3.1 闰年时间的计算方法

闰年的基本计算方法是：

如果 y 能被 4 除尽且不能被 100 除尽，或者能被 400 除尽，则 y 是闰年。

## 4.7 sched.c 文件

### 4.7.1 功能描述

sched.c 是内核中有关任务调度函数的程序。

### 4.7.2 代码注释

列表 4.7 linux/kernel/sched.c 程序

```

1 /*
2  * linux/kernel/sched.c
3  *
4  * (C) 1991 Linus Torvalds
5  */
6
7 /*
8  * 'sched.c' is the main kernel file. It contains scheduling primitives
9  * (sleep_on, wakeup, schedule etc) as well as a number of simple system
10 * call functions (type getpid(), which just extracts a field from
11 * current-task
12 */
13 /*
14  * 'sched.c' 是主要的内核文件。其中包括有关调度的基本函数(sleep_on、wakeup、schedule 等)以及
15  * 一些简单的系统调用函数（比如 getpid()，仅从当前任务中获取一个字段）。
16 */
17 #include <linux/sched.h> // 调度程序头文件。定义了任务结构 task_struct、第 1 个初始任务
                           // 的数据。还有一些以宏的形式定义的有关描述符参数设置和获取的

```

```

// 嵌入式汇编函数程序。
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 #include <linux/sys.h> // 系统调用头文件。含有 72 个系统调用 C 函数处理程序, 以 'sys_' 开头。
16 #include <linux/fdreg.h> // 软驱头文件。含有软盘控制器参数的一些定义。
17 #include <asm/system.h> // 系统头文件。定义了设置或修改描述符/中断门等的嵌入式汇编宏。
18 #include <asm/io.h> // io 头文件。定义硬件端口输入/输出宏汇编语句。
19 #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。
20
21 #include <signal.h> // 信号头文件。定义信号符号常量, sigaction 结构, 操作函数原型。
22
23 #define _S(nr) (1<<((nr)-1)) // 取信号 nr 在信号位图中对应位的二进制数值。信号编号 1-32。
// 比如信号 5 的位图数值 = 1<<(5-1) = 16 = 00010000b。
24 #define BLOCKABLE (~(_S(SIGKILL) | _S(SIGSTOP))) // 除了 SIGKILL 和 SIGSTOP 信号以外其它都是
// 可阻塞的(...1011111111101111111b)。
25
// 显示任务号 nr 的进程号、进程状态和内核堆栈空闲字节数 (大约)。
26 void show_task(int nr, struct task_struct * p)
27 {
28     int i, j = 4096-sizeof(struct task_struct);
29
30     printk("%d: pid=%d, state=%d, ", nr, p->pid, p->state);
31     i=0;
32     while (i<j && !((char *) (p+1))[i]) // 检测指定任务数据结构以后等于 0 的字节数。
33         i++;
34     printk("%d (of %d) chars free in kernel stack\n", i, j);
35 }
36
// 显示所有任务的任务号、进程号、进程状态和内核堆栈空闲字节数 (大约)。
37 void show_stat(void)
38 {
39     int i;
40
41     for (i=0; i<NR_TASKS; i++)
42         if (task[i])
43             show_task(i, task[i]);
44 }
45
// 定义每个时间片的滴答数☺。
46 #define LATCH (1193180/HZ)
47
48 extern void mem_use(void); // [??]
49
50 extern int timer_interrupt(void); // 时钟中断处理程序(kernel/system_call.s, 176)。
51 extern int system_call(void); // 系统调用中断处理程序(kernel/system_call.s, 80)。
52
53 union task_union { // 定义任务联合(任务结构成员和 stack 字符数组成员)。
54     struct task_struct task; // 因为一个任务数据结构与其堆栈放在同一内存页中, 所以
55     char stack[PAGE_SIZE]; // 从堆栈段寄存器 ss 可以获得其数据段选择符。
56 };
57
58 static union task_union init_task = {INIT_TASK,}; // 定义初始任务的数据(sched.h 中)。
59
60 long volatile jiffies=0;

```



```

61 long startup_time=0; // 开机时间。从 1970:0:0:0 开始计时的秒数。
62 struct task_struct *current = &(init_task.task); // 当前任务指针（初始化为初始任务）。
63 struct task_struct *last_task_used_math = NULL; // 使用过协处理器任务的指针。
64
65 struct task_struct * task[NR_TASKS] = {&(init_task.task), }; // 定义任务指针数组。
66
67 long user_stack [ PAGE_SIZE>>2 ] ; // 定义用户堆栈数组（1024 项）。
68
69 struct { // 堆栈开始结构（地址指针，数据段选择符）。
70     long * a;
71     short b;
72     } stack_start = { & user_stack [PAGE_SIZE>>2] , 0x10 };
73 /*
74  * 'math_state_restore()' saves the current math information in the
75  * old math state array, and gets the new ones from the current task
76  */
77 /*
78  * 将当前协处理器内容保存到协处理器状态数组中，并将当前任务的协处理器
79  * 内容加载进协处理器。
80  */
81 // 当任务被调度交换过以后，用以保存原任务的协处理器状态（上下文）并恢复新调度进来的当前
82 // 任务的协处理器执行状态。
83 void math_state_restore()
84 {
85     if (last_task_used_math == current) // 如果任务没变则返回(上一个任务就是当前任务)。
86         return; // 这里所指的“上一个任务”是刚被交换出去的任务。
87     __asm__("fwait"); // 在发送协处理器命令之前要先发 WAIT 指令。
88     if (last_task_used_math) { // 如果上个任务使用了协处理器，则保存其状态。
89         __asm__("fnsave %0"::"m" (last_task_used_math->tss.i387));
90     }
91     last_task_used_math=current; // 现在，last_task_used_math 指向当前任务，
92     // 以备当前任务被交换出去时使用。
93     if (current->used_math) { // 如果当前任务用过协处理器，则恢复其状态。
94         __asm__("frstor %0"::"m" (current->tss.i387));
95     } else { // 否则的话说明是第一次使用，
96         __asm__("fninit"::); // 于是就向协处理器发初始化命令，
97         current->used_math=1; // 并设置使用了协处理器标志。
98     }
99 }
100 /*
101  * 'schedule()' is the scheduler function. This is GOOD CODE! There
102  * probably won't be any reason to change this, as it should work well
103  * in all circumstances (ie gives IO-bound processes good response etc).
104  * The one thing you might take a look at is the signal-handler code here.
105  */
106 /*
107  * NOTE!! Task 0 is the 'idle' task, which gets called when no other
108  * tasks can run. It can not be killed, and it cannot sleep. The 'state'
109  * information in task[0] is never used.
110  */
111 /*
112  * 'schedule()' 是调度函数。这是个很好的代码！没有任何理由对它进行修改，因为它可以在所有的
113  * 环境下工作（比如能够对 IO-边界处理很好的响应等）。只有一件事值得留意，那就是这里的信号

```

```

    * 处理代码。
    * 注意！！任务 0 是个闲置('idle')任务，只有当没有其它任务可以运行时才调用它。它不能被杀
    * 死，也不能睡眠。任务 0 中的状态信息'state'是从来不用的。
    */
104 void schedule(void)
105 {
106     int i,next,c;
107     struct task\_struct ** p;           // 任务结构指针的指针。
108
109     /* check alarm, wake up any interruptible tasks that have got a signal */
    /* 检测 alarm，唤醒得到了一个信号的任何可以中断的任务 */
110
    // 从任务数组中最后一个任务开始检测 alarm。
111     for(p = &LAST\_TASK ; p > &FIRST\_TASK ; --p)
112         if (*p) {
113             // 如果任务的 alarm 时间已经过期(alarm<jiffies), 在信号位图中置 SIGALRM 信号，然后清 alarm。
114             if ((*p)->alarm && (*p)->alarm < jiffies) {
115                 (*p)->signal |= (1<<(SIGALRM-1));
116                 (*p)->alarm = 0;
117             }
118             // 如果信号位图中除去被阻塞的信号外还有其它信号并且任务处于可中断状态，则置任务为就绪状态。
119             // 其中'~(BLOCKABLE & (*p)->blocked)'用于忽略被阻塞的信号，但 SIGKILL 和 SIGSTOP 不能被阻塞。
120             if (((*p)->signal & ~(BLOCKABLE & (*p)->blocked)) &&
121                 (*p)->state==TASK\_INTERRUPTIBLE)
122                 (*p)->state=TASK\_RUNNING;      //置为就绪（可执行）状态。
123         }
124
125     /* this is the scheduler proper: */
    /* 这里是调度程序的主要部分 */
126
127     while (1) {
128         c = -1;
129         next = 0;
130         i = NR\_TASKS;
131         p = &task[NR\_TASKS];
132         // 这段代码比较每个就绪状态任务的 counter 值，哪一个大 next 就指向哪个的任务号。
133         while (--i) {
134             if (!*--p)
135                 continue;
136             if ((*p)->state == TASK\_RUNNING && (*p)->counter > c)
137                 c = (*p)->counter, next = i;
138         }
139         // 如果比较得出有 counter 值大于 0 的结果，则退出循环，执行任务调度切换（141 行）。
140         if (c) break;
141         // 否则就根据每个任务的优先权值，更新每一个任务的 counter 值，然后回到 125 行重新比较。
142         // counter 值的计算方式为 counter = counter /2 + priority。
143         for(p = &LAST\_TASK ; p > &FIRST\_TASK ; --p)
144             if (*p)
145                 (*p)->counter = ((*p)->counter >> 1) +
146                     (*p)->priority;
147     }
148     switch\_to(next);    // 切换到任务号为 next 的任务，并运行之。
149 }

```

```

143 // pause() 系统调用。转换当前任务的状态为可中断的等待状态，并重新调度。
144 int sys_pause(void)
145 {
146     current->state = TASK_INTERRUPTIBLE;
147     schedule();
148     return 0;
149 }
150 // 将当前任务置为不可中断的等待状态，并让睡眠对象的任务指针指向当前任务。
// 只有明确地唤醒时才会返回。该函数提供了进程与中断处理程序之间的同步机制。
151 void sleep_on(struct task_struct **p)
152 {
153     struct task_struct *tmp;
154
155     if (!p)
156         return;
157     if (current == &(init_task.task)) // 如果当前任务是任务 0，则死机(impossible!)。
158         panic("task[0] trying to sleep");
159     tmp = *p; // 让 tmp 指向已经在等待此对象的任务(如果有的话)。
160     *p = current; // 将睡眠对象的等待指针指向当前任务。
161     current->state = TASK_UNINTERRUPTIBLE; // 将当前任务置为不可中断的等待状态。
162     schedule(); // 重新调度。
163     // 只有当这个等待任务被唤醒时，调度程序才又返回到这里，则表示被明确地唤醒。
164     if (tmp) // 若存在早已经等待该对象的任务，则也将其
165         tmp->state=0; // 置为就绪(可运行)状态。
166
167     // 将当前任务置为可中断的等待状态，等待指定 (*p) 任务的完成。
168 void interruptible_sleep_on(struct task_struct **p)
169 {
170     struct task_struct *tmp;
171
172     if (!p)
173         return;
174     if (current == &(init_task.task))
175         panic("task[0] trying to sleep");
176     tmp=*p;
177     *p=current;
178 repeat: current->state = TASK_INTERRUPTIBLE;
179     schedule();
180     if (*p && *p != current) {
181         (**p).state=0;
182         goto repeat;
183     }
184     *p=NULL;
185     if (tmp)
186         tmp->state=0;
187
188 // 唤醒指定任务*p。
189 void wake_up(struct task_struct **p)
190 {

```

```

190     if (p && *p) {
191         (**p).state=0;    // 置为就绪（可运行）状态。
192         *p=NULL;
193     }
194 }
195
196 /*
197  * OK, here are some floppy things that shouldn't be in the kernel
198  * proper. They are here because the floppy needs a timer, and this
199  * was the easiest way of doing it.
200  */
201 /*
202  * 好了，从这里开始是一些有关软盘的子程序，本不应该放在内核的主要部分中的。将它们放在这里
203  * 是因为软驱需要一个时钟，而放在这里是最方便的办法。
204  */
205 static struct task_struct * wait_motor[4] = {NULL, NULL, NULL, NULL};
206 static int mon_timer[4]={0, 0, 0, 0};
207 static int moff_timer[4]={0, 0, 0, 0};
208 unsigned char current_DOR = 0x0C;    // 数字输出寄存器(初值：允许 DMA 和请求中断、启动 FDC)。
209
210 // 指定软盘到正常运转状态所需延迟滴答数（时间）。
211 // nr -- 软驱号(0-3)，返回值为滴答数。
212 int ticks_to_floppy_on(unsigned int nr)
213 {
214     extern unsigned char selected;    // 当前选中的软盘号(kernel/blk_drv/floppy.c, 122)。
215     unsigned char mask = 0x10 << nr;    // 所选软驱对应数字输出寄存器中启动马达比特位。
216
217     if (nr>3)
218         panic("floppy_on: nr>3");    // 最多 4 个软驱。
219     moff_timer[nr]=10000;    /* 100 s = very big :-) */
220     cli();    /* use floppy_off to turn it off */
221     mask |= current_DOR;
222 // 如果不是当前软驱，则首先复位其它软驱的选择位，然后置对应软驱选择位。
223     if (!selected) {
224         mask &= 0xFC;
225         mask |= nr;
226     }
227 // 如果数字输出寄存器的当前值与要求的值不同，则向 FDC 数字输出端口输出新值(mask)。并且如果
228 // 要求启动的马达还没有启动，则置相应软驱的马达启动定时器值(HZ/2 = 0.5 秒或 50 个滴答)。
229 // 此后更新当前数字输出寄存器值 current_DOR。
230     if (mask != current_DOR) {
231         outb(mask, FD_DOR);
232         if ((mask ^ current_DOR) & 0xf0)
233             mon_timer[nr] = HZ/2;
234         else if (mon_timer[nr] < 2)
235             mon_timer[nr] = 2;
236         current_DOR = mask;
237     }
238     sti();
239     return mon_timer[nr];
240 }
241
242 // 等待指定软驱马达启动所需时间。

```

```

232 void floppy_on(unsigned int nr)
233 {
234     cli();    // 关中断。
235     while (ticks_to_floppy_on(nr))    // 如果马达启动定时还没到，就一直把当前进程置
236         sleep_on(nr+wait_motor);    // 为不可中断睡眠状态并放入等待马达运行的队列中。
237     sti();    // 开中断。
238 }
239
// 置关闭相应软驱马达停转定时器（3 秒）。
240 void floppy_off(unsigned int nr)
241 {
242     moff_timer[nr]=3*HZ;
243 }
244
// 软盘定时处理子程序。更新马达启动定时值和马达关闭停转计时值。该子程序是在时钟定时
// 中断中被调用，因此每一个滴答(10ms)被调用一次，更新马达开启或停转定时器的值。如果某
// 一个马达停转定时到，则将数字输出寄存器马达启动位复位。
245 void do_floppy_timer(void)
246 {
247     int i;
248     unsigned char mask = 0x10;
249
250     for (i=0 ; i<4 ; i++,mask <= 1) {
251         if (!(mask & current_DOR))    // 如果不是 DOR 指定的马达则跳过。
252             continue;
253         if (mon_timer[i]) {
254             if (--mon_timer[i])
255                 wake_up(i+wait_motor);    // 如果马达启动定时到则唤醒进程。
256             } else if (!moff_timer[i]) {    // 如果马达停转定时到则
257                 current_DOR &= ~mask;    // 复位相应马达启动位，并
258                 outb(current_DOR, FD_DOR);    // 更新数字输出寄存器。
259             } else
260                 moff_timer[i]--;    // 马达停转计时递减。
261         }
262     }
263
264 #define TIME_REQUESTS 64    // 最多可有 64 个定时器链表（64 个任务）。
265
// 定时器链表结构和定时器数组。
266 static struct timer_list {
267     long jiffies;    // 定时滴答数。
268     void (*fn)();    // 定时处理程序。
269     struct timer_list * next;    // 下一个定时器。
270 } timer_list[TIME_REQUESTS], * next_timer = NULL;
271
// 添加定时器。输入参数为指定的定时值(滴答数)和相应的处理程序指针。
//
272 void add_timer(long jiffies, void (*fn)(void))
273 {
274     struct timer_list * p;
275
// 如果定时处理程序指针为空，则退出。
276     if (!fn)

```

```

277         return;
278         cli();
// 如果定时值<=0, 则立刻调用其处理程序。并且该定时器不加入链表中。
279         if (jiffies <= 0)
280             (fn)();
281         else {
// 从定时器数组中, 找一个空闲项。
282             for (p = timer_list; p < timer_list + TIME_REQUESTS; p++)
283                 if (!p->fn)
284                     break;
// 如果已经用完了定时器数组, 则系统崩溃☹。
285             if (p >= timer_list + TIME_REQUESTS)
286                 panic("No more time requests free");
// 向定时器数据结构填入相应信息。并链入链表头
287             p->fn = fn;
288             p->jiffies = jiffies;
289             p->next = next_timer;
290             next_timer = p;
// 链表项按定时值从小到大排序。在排序时减去排在前面需要的滴答数, 这样在处理定时器时只要
// 查看链表头的第一项的定时是否到期即可。[[?? 这段程序好象没有考虑周全。如果新插入的定时
// 器值<原来头一个定时器值时, 也应该将所有后面的定时值均减去新的第 1 个的定时值。]]
291             while (p->next && p->next->jiffies < p->jiffies) {
292                 p->jiffies -= p->next->jiffies;
293                 fn = p->fn;
294                 p->fn = p->next->fn;
295                 p->next->fn = fn;
296                 jiffies = p->jiffies;
297                 p->jiffies = p->next->jiffies;
298                 p->next->jiffies = jiffies;
299                 p = p->next;
300             }
301         }
302         sti();
303     }
304
// 时钟中断 C 函数处理程序。
305 void do_timer(long cpl)
306 {
307     extern int beepcount; // 扬声器发声时间滴答数(kernel/chr_drv/console.c, 697)
308     extern void sysbeepstop(void); // 关闭扬声器(kernel/chr_drv/console.c, 691)
309
// 如果发声计数次数到, 则关闭发声。(向 0x61 口发送命令, 复位位 0 和 1。位 0 控制 8253
// 计数器 2 的工作, 位 1 控制扬声器)。
310     if (beepcount)
311         if (!--beepcount)
312             sysbeepstop();
313
// 如果当前特权级(cpl)为 0 (最高, 表示是内核程序在工作), 则将超级用户运行时间 stime 递增;
// 如果 cpl > 0, 则表示是一般用户程序在工作, 增加 utime。
314     if (cpl)
315         current->utime++;
316     else
317         current->stime++;

```

```

318 // 如果有用户的定时器存在，则将链表第 1 个定时器的值减 1。如果已等于 0，则调用相应的处理
// 程序，并将该处理程序指针置为空。然后去掉该项定时器。
319 if (next_timer) { // next_timer 是定时器链表的头指针(见 270 行)。
320     next_timer->jiffies--;
321     while (next_timer && next_timer->jiffies <= 0) {
322         void (*fn)(void); // 这里插入了一个函数指针定义!!! ⊗
323
324         fn = next_timer->fn;
325         next_timer->fn = NULL;
326         next_timer = next_timer->next;
327         (fn)(); // 调用处理函数。
328     }
329 }
// 如果当前软盘控制器 FDC 的数字输出寄存器中马达启动位有置位的，则执行软盘定时程序(245 行)。
330 if (current_DOR & 0xf0)
331     do_floppy_timer();
332 if ((--current->counter)>0) return; // 如果进程运行时间还没完，则退出。
333 current->counter=0;
334 if (!cpl) return; // 对于超级用户程序，不依赖 counter 值进行调度。
335 schedule();
336 }
337
// 系统调用功能 - 设置报警定时时间值(秒)。
// 如果已经设置过 alarm 值，则返回旧值，否则返回 0。
338 int sys_alarm(long seconds)
339 {
340     int old = current->alarm;
341
342     if (old)
343         old = (old - jiffies) / HZ;
344     current->alarm = (seconds>0)?(jiffies+HZ*seconds):0;
345     return (old);
346 }
347
// 取当前进程号 pid。
348 int sys_getpid(void)
349 {
350     return current->pid;
351 }
352
// 取父进程号 ppid。
353 int sys_getppid(void)
354 {
355     return current->father;
356 }
357
// 取用户号 uid。
358 int sys_getuid(void)
359 {
360     return current->uid;
361 }
362

```

```

    // 取 euid。
363 int sys_geteuid(void)
364 {
365     return current->euid;
366 }
367
    // 取组号 gid。
368 int sys_getgid(void)
369 {
370     return current->gid;
371 }
372
    // 取 egid。
373 int sys_getegid(void)
374 {
375     return current->egid;
376 }
377
    // 系统调用功能 -- 降低对 CPU 的使用优先权（有人会用吗？☺）。
    // 应该限制 increment 大于 0，否则的话，可使优先权增大！！
378 int sys_nice(long increment)
379 {
380     if (current->priority-increment>0)
381         current->priority -= increment;
382     return 0;
383 }
384
    // 调度程序的初始化子程序。
385 void sched_init(void)
386 {
387     int i;
388     struct desc_struct * p; // 描述符表结构指针。
389
390     if (sizeof(struct sigaction) != 16)
391         panic("Struct sigaction MUST be 16 bytes");
    // 设置初始任务（任务 0）的任务状态段描述符和局部数据表描述符(include/asm/system.h, 65)。
392     set_tss_desc(gdt+FIRST_TSS_ENTRY, &(init_task.task.tss));
393     set_ldt_desc(gdt+FIRST_LDT_ENTRY, &(init_task.task.ldt));
    // 清任务数组和描述符表项（注意 i=1 开始，所以初始任务的描述符还在）。
394     p = gdt+2+FIRST_TSS_ENTRY;
395     for(i=1; i<NR_TASKS; i++) {
396         task[i] = NULL;
397         p->a=p->b=0;
398         p++;
399         p->a=p->b=0;
400         p++;
401     }
402     /* Clear NT, so that we won't have troubles with that later on */
    /* 清除标志寄存器中的位 NT，这样以后就不会有麻烦 */
    // NT 标志用于控制程序的递归调用(Nested Task)。当 NT 置位时，那么当前中断任务执行
    // iret 指令时就会引起任务切换。NT 指出 TSS 中的 back_link 字段是否有效。
403     __asm__ ("pushfl ; andl $0xffffbfff, (%esp) ; popfl"); // 复位 NT 标志。
404     ltr(0); // 将任务 0 的 TSS 加载到任务寄存器 tr。

```



```

405      lldt(0); // 将局部描述符表加载到局部描述符表寄存器。
// 注意！！是将 GDT 中相应 LDT 描述符的选择符加载到 ldtr。只明确加载这一次，以后新任务
// LDT 的加载，是 TSS 根据其中的 LDT 项自动加载。

// 下面代码用于初始化 8253 定时器。
406      outb\_p(0x36, 0x43); /* binary, mode 3, LSB/MSB, ch 0 */
407      outb\_p(LATCH & 0xff, 0x40); /* LSB */ // 定时值低字节。
408      outb(LATCH >> 8, 0x40); /* MSB */ // 定时值高字节。
// 设置时钟中断处理程序句柄（设置时钟中断门）。
409      set\_intr\_gate(0x20, &timer\_interrupt);
// 修改中断控制器屏蔽码，允许时钟中断。
410      outb(inb\_p(0x21) & ~0x01, 0x21);
// 设置系统调用中断门。
411      set\_system\_gate(0x80, &system\_call);
412 }
413

```

## 4.7.3 其它信息

### 4.7.3.1 软盘控制器编程

在编程时需要访问 4 个端口，分别对应一个或多个寄存器。对于 1.2M 的软盘控制器有以下一些端口。

表 软盘控制器端口

| I/O 端口 | 读写性 | 寄存器名称            |
|--------|-----|------------------|
| 0x3f2  | 只写  | 数字输出寄存器(数字控制寄存器) |
| 0x3f4  | 只读  | FDC 主状态寄存器       |
| 0x3f5  | 读/写 | FDC 数据寄存器        |
| 0x3f7  | 只读  | 数字输入寄存器          |
| 0x3f7  | 只写  | 磁盘控制寄存器(传输率控制)   |

数字输出端口（数字控制端口）是一个 8 位寄存器，它控制驱动器马达开启、驱动器选择、启动/复位 FDC 以及允许/禁止 DMA 及中断请求。

FDC 的主状态寄存器也是一个 8 位寄存器，用于反映软盘控制器 FDC 和软盘驱动器 FDD 的基本状态。通常，在 CPU 向 FDC 发送命令之前或从 FDC 获取操作结果之前，都要读取主状态寄存器的状态位，以判别当前 FDC 数据寄存器是否就绪，以及确定数据传送的方向。

FDC 的数据端口对应多个寄存器（只写型命令寄存器和参数寄存器、只读型结果寄存器），但任一时刻只能有一个寄存器出现在数据端口 0x3f5。在访问只写型寄存器时，主状态控制的 DIO 方向位必须为 0（CPU → FDC），访问只读型寄存器时则反之。在读取结果时只有在 FDC 不忙之后才算读完结果，通常结果数据最多有 7 个字节。

软盘控制器共可以接受 15 条命令。每个命令均经历三个阶段：命令阶段、执行阶段和结果阶段。

命令阶段是 CPU 向 FDC 发送命令字节和参数字节。每条命令的第一个字节总是命令字节（命令码）。其后跟着 0—8 字节的参数。

执行阶段是 FDC 执行命令规定的操作。在执行阶段 CPU 是不加干预的，一般是通过 FDC 发出中断请求获知命令执行的结束。如果 CPU 发出的 FDC 命令是传送数据，则 FDC 可以以中断方式或 DMA 方式进行。中断方式每次传送 1 字节。DMA 方式是在 DMA 控制器管理下，FDC 与内存进行数据的传输直至全部数据传送完。此时 DMA 控制器会将传输字节计数终止信号通知 FDC，最后由 FDC 发出中断请求信号告知 CPU 执行阶段结束。

结果阶段是由 CPU 读取 FDC 数据寄存器返回值，从而获得 FDC 命令执行的结果。返回结果数据的长度

为 0—7 字节。对于没有返回结果数据的命令，则应向 FDC 发送检测中断状态命令获得操作的状态。

## 4.8 signal.c 文件

### 4.8.1 功能描述

本程序给出了设置和获取进程信号阻塞码（屏蔽码）函数，`signal()` 和 `sigaction()` 系统调用以及在系统调用中断处理程序中处理信号的函数 `do_signal()`。

`signal()` 和 `sigaction()` 的功能比较类似，都是改变信号原处理句柄。但 `signal()` 会返回原信号处理句柄，并且在新句柄被调用一次后句柄就会恢复到默认值。而 `sigaction()` 则可以进行更自由的设置。

`do_signal()` 的处理原理是系统调用中断处理程序中真正的信号处理程序。该函数的主要作用是将信号的处理句柄插入到用户程序堆栈中。这样，在本系统调用结束返回后就会立刻执行信号句柄程序，然后继续执行用户的程序。

### 4.8.2 代码注释

列表 4.8 linux/kernel/signal.c 程序

```

1 /*
2  * linux/kernel/signal.c
3  *
4  * (C) 1991 Linus Torvalds
5  */
6
7 #include <linux/sched.h> // 调度程序头文件，定义了任务结构 task_struct、初始任务 0 的数据，
                        // 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
8 #include <linux/kernel.h> // 内核头文件。含有一些内核常用函数的原形定义。
9 #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。
10
11 #include <signal.h> // 信号头文件。定义信号符号常量，信号结构以及信号操作函数原型。
12
13 volatile void do_exit(int error_code);
14
15 // 获取当前任务信号屏蔽位图（屏蔽码）。
16 int sys_sgetmask()
17 {
18     return current->blocked;
19 }
20
21 // 设置新的信号屏蔽位图。SIGKILL 不能被屏蔽。返回值是原信号屏蔽位图。
22 int sys_ssetmask(int newmask)
23 {
24     int old=current->blocked;
25
26     current->blocked = newmask & ~(1<<(SIGKILL-1));
27     return old;
28 }
29
30 // 复制 sigaction 数据到 fs 数据段 to 处。。

```

```

28 static inline void save\_old(char * from, char * to)
29 {
30     int i;
31
32     verify\_area(to, sizeof(struct sigaction)); // 验证 to 处的内存是否足够。
33     for (i=0 ; i< sizeof(struct sigaction) ; i++) {
34         put\_fs\_byte(*from, to); // 复制到 fs 段。一般是用户数据段。
35         from++; // put\_fs\_byte() 在 include/asm/segment.h 中。
36         to++;
37     }
38 }
39
// 把 sigaction 数据从 fs 数据段复制到 to 处。
40 static inline void get\_new(char * from, char * to)
41 {
42     int i;
43
44     for (i=0 ; i< sizeof(struct sigaction) ; i++)
45         *(to++) = get\_fs\_byte(from++);
46 }
47
// signal() 系统调用。类似于 sigaction()。为指定的信号安装新的信号句柄(信号处理程序)。
// 信号句柄可以是用户指定的函数，也可以是 SIG_DFL (默认句柄) 或 SIG_IGN (忽略)。
// 参数 signal -- 指定的信号; handler -- 指定的句柄; restorer -- [??]
// 函数返回原信号句柄。
48 int sys\_signal(int signal, long handler, long restorer)
49 {
50     struct sigaction tmp;
51
52     if (signal<1 || signal>32 || signal==SIGKILL) // 信号值要在 (1-32) 范围内，
53         return -1; // 并且不得是 SIGKILL。
54     tmp.sa_handler = (void (*)(int)) handler; // 指定的信号处理句柄。
55     tmp.sa_mask = 0; // 执行时的信号屏蔽码。
56     tmp.sa_flags = SA\_ONESHOT | SA\_NOMASK; // 该句柄只使用 1 次后就恢复到默认值，
// 并允许信号在自己的处理句柄中收到。
57     tmp.sa_restorer = (void (*)(void)) restorer; // [??]
58     handler = (long) current->sigaction[signal-1].sa_handler;
59     current->sigaction[signal-1] = tmp;
60     return handler;
61 }
62
// sigaction() 系统调用。改变进程在收到一个信号时的操作。signal 是除了 SIGKILL 以外的任何
// 信号。[如果新操作(action)不为空]则新操作被安装。如果 oldaction 指针不为空，则原操作
// 被保留到 oldaction。成功则返回 0，否则为-1。
63 int sys\_sigaction(int signal, const struct sigaction * action,
64 struct sigaction * oldaction)
65 {
66     struct sigaction tmp;
67
68     // 信号值要在 (1-32) 范围内，并且信号 SIGKILL 的处理句柄不能被改变。
69     if (signal<1 || signal>32 || signal==SIGKILL)
70         return -1;
71
// 在信号的 sigaction 结构中设置新的操作(动作)。

```

```

70     tmp = current->sigaction[signum-1];
71     get\_new((char *) action,
72            (char *) (signum-1+current->sigaction));
// 如果 oldaction 指针不为空的话, 则将原操作指针保存到 oldaction 所指的位置。
73     if (oldaction)
74         save\_old((char *) &tmp, (char *) oldaction);
// 如果允许信号在自己的信号句柄中收到, 则令屏蔽码为 0, 否则设置屏蔽本信号。
75     if (current->sigaction[signum-1].sa_flags & SA\_NOMASK)
76         current->sigaction[signum-1].sa_mask = 0;
77     else
78         current->sigaction[signum-1].sa_mask |= (1<<(signum-1));
79     return 0;
80 }
81
// 系统调用中断处理程序中真正的信号处理程序 (在 kernel/system_call.s, 119 行)。
// 该段代码的主要作用是将信号的处理句柄插入到用户程序堆栈中, 并在本系统调用结束
// 返回后立刻执行信号句柄程序, 然后继续执行用户的程序。
82 void do\_signal(long signr, long eax, long ebx, long ecx, long edx,
83                long fs, long es, long ds,
84                long eip, long cs, long eflags,
85                unsigned long * esp, long ss)
86 {
87     unsigned long sa_handler;
88     long old_eip=eip;
89     struct sigaction * sa = current->sigaction + signr - 1; //current->sigaction[signu-1]。
90     int longs;
91     unsigned long * tmp_esp;
92
93     sa_handler = (unsigned long) sa->sa_handler;
// 如果信号句柄为 SIG_IGN(忽略), 则返回; 如果句柄为 SIG_DFL(默认处理), 则如果信号是
// SIGCHLD 则返回, 否则终止进程的执行
94     if (sa_handler==1)
95         return;
96     if (!sa_handler) {
97         if (signr==SIGCHLD)
98             return;
99         else
100             do\_exit(1<<(signr-1)); [?? 为什么以信号位图为参数? 不为什么!?! 😊]
101     }
// 如果该信号句柄只需使用一次, 则将该句柄置空(该信号句柄已经保存在 sa_handler 指针中)。
102     if (sa->sa_flags & SA\_ONESHOT)
103         sa->sa_handler = NULL;
// 下面这段代码将信号处理句柄插入到用户堆栈中, 同时也将 sa_restorer, signr, 进程屏蔽码(如果
// SA_NOMASK 设置位), eax, ecx, edx 作为参数以及原调用系统调用的程序返回指针及标志寄存器值
// 压入堆栈。因此在本次系统调用中断(0x80)返回用户程序时会首先执行用户的信号句柄程序, 然后
// 再继续执行用户程序。

// 将用户调用系统调用的代码指针 eip 指向该信号处理句柄。
104     *(&eip) = sa_handler;
// 如果允许信号自己的处理句柄收到信号自己, 则也需要将进程的阻塞码压入堆栈。
105     longs = (sa->sa_flags & SA\_NOMASK)?7:8;
// 将原调用程序的用户的堆栈指针向下扩展 7 (或 8) 个长字 (用来存放调用信号句柄的参数等),
// 并检查内存使用情况 (例如如果内存超界则分配新页等)。

```

```

106      *(&esp) -= longs;
107      verify\_area(esp, longs*4);
// 在用户堆栈中从下到上存放 sa_restorer, 信号 signr, 屏蔽码 blocked(如果 SA_NOMASK 置位),
// eax, ecx, edx, eflags 和用户程序原代码指针。
108      tmp_esp=esp;
109      put\_fs\_long((long) sa->sa_restorer, tmp_esp++);
110      put\_fs\_long(signr, tmp_esp++);
111      if (!(sa->sa_flags & SA\_NOMASK))
112          put\_fs\_long(current->blocked, tmp_esp++);
113      put\_fs\_long(eax, tmp_esp++);
114      put\_fs\_long(ecx, tmp_esp++);
115      put\_fs\_long(edx, tmp_esp++);
116      put\_fs\_long(eflags, tmp_esp++);
117      put\_fs\_long(old_eip, tmp_esp++);
118      current->blocked |= sa->sa_mask; // 进程阻塞码(屏蔽码)添上 sa_mask 中的码位。
119 }
120

```

## 4.8.3 其它信息

### 4.8.3.1 进程信号说明

| 标号 | 名称        | 说明  | 默认操作                          |
|----|-----------|---|-------------------------------|
| 1  | SIGHUP    | (Hangup) 当你不再控制终端时内核会产生该信号, 或者当你关闭 Xterm 或断开 modem。由于后台程序没有控制的终端, 因而它们常用 SIGUP 来发出需要重新读取其配置文件的信号。 | (Abort) 挂断控制中断或进程。            |
| 2  | SIGINT    | (Interrupt) 来自键盘的终端。通常终端驱动程序会将其与 ^C 绑定。   | (Abort) 终止程序。                 |
| 3  | SIGQUIT   | (Quit) 来自键盘的终端。通常终端驱动程序会将其与 ^\ 绑定。  | (Dump) 程序被终止并产生 dump core 文件。 |
| 4  | SIGILL    | (Illegal Instruction) 程序出错或者执行了一个非法的操作指令。   | (Dump) 程序被终止并产生 dump core 文件。 |
| 5  | SIGTRAP   | (Breakpoint/Trace Trap) 调试用, 跟踪断点。  |                               |
| 6  | SIGABRT   | (Abort) 放弃执行, 异常结束。   | (Dump) 程序被终止并产生 dump core 文件。 |
| 6  | SIGIOT    | (IO Trap) 同 SIGABRT   | (Dump) 程序被终止并产生 dump core 文件。 |
| 7  | SIGUNUSED | (Unused) 没有使用。  |                               |
| 8  | SIGFPE    | (Floating Point Exception) 浮点异常。  | (Dump) 程序被终止并产生 dump core 文件。 |
| 9  | SIGKILL   | (Kill) 程序被终止。该信号不能被捕获或者被忽略。想立刻终止一个进程, 就发送信号 9。注意程序将没有任何机会做清理工作。                                   | (Abort) 程序被终止。                |
| 10 | SIGUSR1   | (User defined Signal 1) 用户定义的信号。  | (Abort) 进程被终止。                |
| 11 | SIGSEGV   | (Segmentation Violation) 当程序引用无效的内存时会产生此信号。比如: 寻址没有映射的内存; 寻址未许可的内存。                               | (Dump) 程序被终止并产生 dump core 文件。 |
| 12 | SIGUSR2   | (User defined Signal 2) 保留给用户程序用于 IPC 或   | (Abort) 进程被终止。                |

|    |           |  |                         |
|----|-----------|--|-------------------------|
|    |           | 其它目的。  |                         |
| 13 | SIGPIPE   | (Pipe) 当程序向一个套接字或管道写时由于没有读者而产生该信号。   | (Abort) 进程被终止。          |
| 14 | SIGALRM   | (Alarm) 该信号会在用户调用 alarm 系统调用所设置的延迟秒数到后产生。该信号常用判别于系统调用超时。                                       | (Abort) 进程被终止。          |
| 15 | SIGTERM   | (Terminate) 用于和善地要求一个程序终止。它是 kill 的默认信号。与 SIGKILL 不同, 该信号能被捕获, 这样就能在退出运行前做清理工作。                | (Abort) 进程被终止。          |
| 16 | SIGSTKFLT | (Stack fault on coprocessor) 协处理器堆栈错误。   | (Abort) 进程被终止。          |
| 17 | SIGCHLD   | (Child) 父进程发出。停止或终止子进程。可改变其含义挪作它用。   | (Ignore) 子进程停止或结束。      |
| 18 | SIGCONT   | (Continue) 该信号致使被 SIGSTOP 停止的进程恢复运行。可以被捕获。   | (Continue) 恢复进程的<br>执行。 |
| 19 | SIGSTOP   | (Stop) 停止进程的运行。该信号不可被捕获或忽略。  | (Stop) 停止进程运行。          |
| 20 | SIGTSTP   | (Terminal Stop) 向终端发送停止键序列。该信号可以被捕获或忽略。  | (Stop) 停止进程运行。          |
| 21 | SIGTTIN   | (Terminal Input on Background) 后台进程试图从一个不再被控制的终端上读取数据, 此时该进程将被停止, 直到收到 SIGCONT 信号。该信号可以被捕获或忽略。 | (Stop) 停止进程运行。          |
| 22 | SIGTTOU   | (TTY Output on Background) 后台进程试图向一个不再被控制的终端上输出数据, 此时该进程将被停止, 直到收到 SIGCONT 信号。该信号可被捕获或忽略。      | (Stop) 停止进程运行。          |

## 4.9 exit.c 文件

### 4.9.1 功能描述

### 4.9.2 代码注释

列表 linux/kernel/exit.c 程序

```

1 /*
2  * linux/kernel/exit.c
3  *
4  * (C) 1991 Linus Torvalds
5  */
6
7 #include <errno.h>          // 错误号头文件。包含系统中各种出错号。(Linus 从 minix 中引进的)
8 #include <signal.h>        // 信号头文件。定义信号符号常量, 信号结构以及信号操作函数原型。
9 #include <sys/wait.h>       // 等待调用头文件。定义系统调用 wait() 和 waitpid() 及相关常数符号。
10
11 #include <linux/sched.h>    // 调度程序头文件, 定义了任务结构 task_struct、初始任务 0 的数据,
12                             // 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
13 #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 <linux/tty.h>    // tty 头文件, 定义了有关 tty_io, 串行通信方面的参数、常数。
14 #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。
15
16 int sys_pause(void);
17 int sys_close(int fd);
18
19 // 释放指定进程(任务)。
20 void release(struct task_struct * p)
21 {
22     int i;
23
24     if (!p)
25         return;
26     for (i=1 ; i<NR_TASKS ; i++)    // 扫描任务数组, 寻找指定任务。
27         if (task[i]==p) {
28             task[i]=NULL;    // 置空该任务项并释放相关内存页。
29             free_page((long)p);
30             schedule();    // 重新调度。
31             return;
32         }
33     panic("trying to release non-existent task"); // 指定任务若不存在则死机。
34 }
35
36 // 向指定任务(*p)发送信号(sig), 权限为 priv。
37 static inline int send_sig(long sig, struct task_struct * p, int priv)
38 {
39     // 若信号不正确或任务指针为空则出错退出。
40     if (!p || sig<1 || sig>32)
41         return -EINVAL;
42     // 若有权或进程有效用户标识符(euid)就是指定进程的 euid 或者是超级用户, 则在进程位图中添加
43     // 该信号, 否则出错退出。其中 suser() 定义为 (current->euid==0), 用于判断是否超级用户。
44     if (priv || (current->euid==p->euid) || suser())
45         p->signal |= (1<<(sig-1));
46     else
47         return -EPERM;
48     return 0;
49 }
50
51 // 终止会话(session)。
52 static void kill_session(void)
53 {
54     struct task_struct **p = NR_TASKS + task; // 指针的指针首先指向任务数组最末端。
55
56     // 对于所有的任务 (除任务 0 以外), 如果其会话等于当前进程的会话就其发送挂断进程信号。
57     while (--p > &FIRST_TASK) {
58         if (*p && (*p)->session == current->session)
59             (*p)->signal |= 1<<(SIGHUP-1); // 发送挂断进程信号。
60     }
61 }
62
63 /*
64  * XXX need to check permissions needed to send signals to process
65  * groups, etc. etc. kill() permissions semantics are tricky!

```



```

59  */
   /*
   * 为了向进程组等发送信号，XXX 需要检查许可。kill() 的许可机制非常巧妙！
   * /
   // kill() 系统调用可用于向任何进程组或进程发送任何信号。
   // 如果 pid 值>0，则信号被发送给 pid。
   // 如果 pid=0，那么信号就会被发送给当前进程的进程组中的所有进程。
   // 如果 pid=-1，则信号 sig 就会发送给除第一个进程外的所有进程。
   // 如果 pid < -1，则信号 sig 将发送给进程组-pid 的所有进程。
   // 如果信号 sig 为 0，则不发送信号，但仍会进行错误检查。如果成功则返回 0。
60 int sys_kill(int pid,int sig)
61 {
62     struct task_struct **p = NR_TASKS + task;
63     int err, retval = 0;
64
65     if (!pid) while (--p > &FIRST_TASK) {
66         if (*p && (*p)->pggrp == current->pid)
67             if (err=send_sig(sig,*p,1))
68                 retval = err;
69     } else if (pid>0) while (--p > &FIRST_TASK) {
70         if (*p && (*p)->pid == pid)
71             if (err=send_sig(sig,*p,0))
72                 retval = err;
73     } else if (pid == -1) while (--p > &FIRST_TASK)
74         if (err = send_sig(sig,*p,0))
75             retval = err;
76     else while (--p > &FIRST_TASK)
77         if (*p && (*p)->pggrp == -pid)
78             if (err = send_sig(sig,*p,0))
79                 retval = err;
80     return retval;
81 }
82
   // 通知父进程 -- 向进程 pid 发送信号 SIGCHLD：子进程将停止或终止。
   // 如果没有找到父进程，则自己释放。
83 static void tell_father(int pid)
84 {
85     int i;
86
87     if (pid)
88         for (i=0;i<NR_TASKS;i++) {
89             if (!task[i])
90                 continue;
91             if (task[i]->pid != pid)
92                 continue;
93             task[i]->signal |= (1<<(SIGCHLD-1));
94             return;
95         }
96     /* if we don't find any fathers, we just release ourselves */
97     /* This is not really OK. Must change it to make father 1 */
98     printk("BAD BAD - no father found\n\r");
99     release(current); // 如果没有找到父进程，则自己释放。
100 }

```



```

101
102 // 释放当前进程代码段和数据段所占的内存页 (free_page_tables() 在 mm/memory.c, 105 行)。
103 int do_exit(long code)    // code 是错误码。
104 {
105     int i;
106     free_page_tables(get_base(current->ldt[1]), get_limit(0x0f));
107     free_page_tables(get_base(current->ldt[2]), get_limit(0x17));
108     // 如果当前进程有子进程，就将子进程的 father 置为 1 (其父进程改为进程 1)。如果该子进程已经
109     // 处于僵死 (ZOMBIE) 状态，则向进程 1 发送子进程终止信号 SIGCHLD。
110     for (i=0 ; i<NR_TASKS ; i++)
111         if (task[i] && task[i]->father == current->pid) {
112             task[i]->father = 1;
113             if (task[i]->state == TASK_ZOMBIE)
114                 /* assumption task[1] is always init */
115                 (void) send_sig(SIGCHLD, task[1], 1);
116         }
117     // 关闭当前进程打开着的所有文件。
118     for (i=0 ; i<NR_OPEN ; i++)
119         if (current->filp[i])
120             sys_close(i);
121     // 同步当前进程工作目录 pwd、根目录 root 以及运行程序的 i 节点并分别置空。
122     iput(current->pwd);
123     current->pwd=NULL;
124     iput(current->root);
125     current->root=NULL;
126     iput(current->executable);
127     current->executable=NULL;
128     // 如果当前进程是领头 (leader) 进程并且其有控制的终端，则释放该终端。
129     if (current->leader && current->tty >= 0)
130         tty_table[current->tty].pgrp = 0;
131     // 如果当前进程上次使用过协处理器，则将 last_task_used_math 置空。
132     if (last_task_used_math == current)
133         last_task_used_math = NULL;
134     // 如果当前进程是 leader 进程，则终止所有相关进程。
135     if (current->leader)
136         kill_session();
137     // 当前进程置为僵死状态，并设置退出码。
138     current->state = TASK_ZOMBIE;
139     current->exit_code = code;
140     // 通知父进程，也即向父进程发送信号 SIGCHLD -- 子进程将停止或终止。
141     tell_father(current->father);
142     schedule();    // 重新调度进程的运行。
143     return (-1);    /* just to suppress warnings */
144 }
145
146 // 系统调用 exit()。终止进程。
147 int sys_exit(int error_code)
148 {
149     return do_exit((error_code&0xff)<<8);
150 }
151

```

```

//// 系统调用 waitpid()。挂起当前进程，直到 pid 指定的子进程退出（终止）或者收到要求终止
// 该进程的信号，或者是需要调用一个信号句柄（信号处理程序）。如果 pid 所指的子进程早已
// 退出（已成所谓的僵死进程），则本调用将立刻返回。子进程使用的所有资源将释放。
// 如果 pid > 0，表示等待进程号等于 pid 的子进程。
// 如果 pid = 0，表示等待进程组号等于当前进程的任何子进程。
// 如果 pid < -1，表示等待进程组号等于 pid 绝对值的任何子进程。
// [ 如果 pid = -1，表示等待任何子进程。]
// 若 options = WUNTRACED，表示如果子进程是停止的，也马上返回。
// 若 options = WNOHANG，表示如果没有子进程退出或终止就马上返回。
// 如果 stat_addr 不为空，则就将状态信息保存到那里。
142 int sys_waitpid(pid_t pid,unsigned long * stat_addr, int options)
143 {
144     int flag, code;
145     struct task_struct ** p;
146
147     verify_area(stat_addr, 4);
148 repeat:
149     flag=0;
150     for(p = &LAST_TASK ; p > &FIRST_TASK ; --p) { // 从任务数组末端开始扫描所有任务。
151         if (!*p || *p == current) // 跳过空项和本进程项。
152             continue;
153         if ((*p)->father != current->pid) // 如果不是当前进程的子进程则跳过。
154             continue;
155         if (pid>0) { // 如果指定的 pid>0，但扫描的进程 pid
156             if ((*p)->pid != pid) // 与之不等，则跳过。
157                 continue;
158         } else if (!pid) { // 如果指定的 pid=0，但扫描的进程组号
159             if ((*p)->pgrp != current->pgrp) // 与当前进程的组号不等，则跳过。
160                 continue;
161         } else if (pid != -1) { // 如果指定的 pid<-1，但扫描的进程组号
162             if ((*p)->pgrp != -pid) // 与其绝对值不等，则跳过。
163                 continue;
164         }
165         switch ((*p)->state) {
166             case TASK_STOPPED:
167                 if (!(options & WUNTRACED))
168                     continue;
169                 put_fs_long(0x7f, stat_addr); // 置状态信息为 0x7f。
170                 return (*p)->pid; // 退出，返回子进程的进程号。
171             case TASK_ZOMBIE:
172                 current->cutime += (*p)->utime; // 更新当前进程的子进程用户
173                 current->cstime += (*p)->stime; // 态和核心态运行时间。
174                 flag = (*p)->pid;
175                 code = (*p)->exit_code; // 取子进程的退出码。
176                 release(*p); // 释放该子进程。
177                 put_fs_long(code, stat_addr); // 置状态信息为退出码值。
178                 return flag; // 退出，返回子进程的 pid.
179             default:
180                 flag=1; // 如果子进程不在停止或僵死状态，则 flag=1。
181                 continue;
182         }
183     }
184     if (flag) { // 如果子进程没有处于退出或僵死状态，

```

```

185         if (options & WNOHANG)                // 并且 options = WNOHANG, 则立刻返回。
186             return 0;
187         current->state=TASK_INTERRUPTIBLE; // 置当前进程为可中断等待状态。
188         schedule();                          // 重新调度。
189         if (!(current->signal &= ~(1<<(SIGCHLD-1)))) // 又开始执行本进程时,
190             goto repeat; // 如果进程没有收到除 SIGCHLD 的信号, 则还是重复处理。
191     else
192         return -EINTR; // 退出, 返回出错码。
193     }
194     return -ECHILD;
195 }
196
197
198

```

### 4.9.3 其它信息

## 4.10 fork.c 程序

### 4.10.1 功能描述

fork() 系统调用用于创建子进程。Linux 中所有进程都是进程 0 (任务 0) 的子进程。

### 4.10.2 代码注释

列表 linux/kernel/fork.c 程序

```

1  /*
2   * linux/kernel/fork.c
3   *
4   * (C) 1991 Linus Torvalds
5   */
6
7  /*
8   * 'fork.c' contains the help-routines for the 'fork' system call
9   * (see also system_call.s), and some misc functions ('verify_area').
10  * Fork is rather simple, once you get the hang of it, but the memory
11  * management can be a bitch. See 'mm/mm.c': 'copy_page_tables()'
12  */
13 /*
14  * 'fork.c' 中含有系统调用 'fork' 的辅助子程序 (参见 system_call.s), 以及一些其它函数
15  * ('verify_area')。一旦你了解了 fork, 就会发现它是非常简单的, 但内存管理却有些难度。
16  * 参见 'mm/mm.c' 中的 'copy_page_tables()'。
17  */
18 #include <errno.h> // 错误号头文件。包含系统中各种出错号。(Linus 从 minix 中引进的)。
19 #include <linux/sched.h> // 调度程序头文件, 定义了任务结构 task_struct、初始任务 0 的数据,
20 // 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
21 #include <linux/kernel.h> // 内核头文件。含有一些内核常用函数的原形定义。
22 #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。

```

```

18 #include <asm/system.h> // 系统头文件。定义了设置或修改描述符/中断门等的嵌入式汇编宏。
19
20 extern void write_verify(unsigned long address);
21
22 long last_pid=0;
23
24 // 验证内存区域。对于给定起始地址和大小的内存块进行检查，如果发现超出分配给进程的页，就
25 // 再另行分配新内存页。
26 void verify_area(void * addr, int size)
27 {
28     unsigned long start;
29
30     start = (unsigned long) addr;
31     size += start & 0xfff;
32     start &= 0xfffff000;
33     start += get_base(current->ldt[2]);
34     while (size>0) {
35         size -= 4096;
36         write_verify(start);
37         start += 4096;
38     }
39
40 // 设置新任务的代码和数据段基址、限长并复制页表。
41 // nr 为新任务号；p 是新任务数据结构的指针。
42 int copy_mem(int nr, struct task_struct * p)
43 {
44     unsigned long old_data_base, new_data_base, data_limit;
45     unsigned long old_code_base, new_code_base, code_limit;
46
47     code_limit=get_limit(0x0f); // 取局部描述符表中代码段描述符项中段限长。
48     data_limit=get_limit(0x17); // 取局部描述符表中数据段描述符项中段限长。
49     old_code_base = get_base(current->ldt[1]); // 取原代码段基址。
50     old_data_base = get_base(current->ldt[2]); // 取原数据段基址。
51     if (old_data_base != old_code_base) // 0.11 版不支持代码和数据段分立的情况。
52         panic("We don't support separate I&D");
53     if (data_limit < code_limit) // 如果数据段长度 < 代码段长度也不对。
54         panic("Bad data_limit");
55     new_data_base = new_code_base = nr * 0x4000000; // 新基址=任务号*64Mb(任务大小)。
56     p->start_code = new_code_base;
57     set_base(p->ldt[1], new_code_base); // 设置代码段描述符中基址域。
58     set_base(p->ldt[2], new_data_base); // 设置数据段描述符中基址域。
59     if (copy_page_tables(old_data_base, new_data_base, data_limit)) { // 复制代码和数据段。
60         free_page_tables(new_data_base, data_limit); // 如果出错则释放申请的内存。
61         return -ENOMEM;
62     }
63     return 0;
64 }
65
66 /*
67  * Ok, this is the main fork-routine. It copies the system process
68  * information (task[nr]) and sets up the necessary registers. It
69  * also copies the data segment in it's entirety.

```

```

67  */
68  /*
69  * OK, 下面是主要的 fork 子程序。它复制系统进程信息(task[n])并且设置必要的寄存器。
70  * 它还整个地复制数据段。
71  */
72  // 复制进程。
73  int copy_process(int nr, long ebp, long edi, long esi, long gs, long none,
74                  long ebx, long ecx, long edx,
75                  long fs, long es, long ds,
76                  long eip, long cs, long eflags, long esp, long ss)
77  {
78      struct task_struct *p;
79      int i;
80      struct file *f;
81
82      p = (struct task_struct *) get_free_page(); // 为新任务数据结构分配内存。
83      if (!p) // 如果内存分配出错, 则返回出错码并退出。
84          return -EAGAIN;
85      task[nr] = p; // 将新任务结构指针放入任务数组中。
86                  // 其中 nr 为任务号, 由前面 find_empty_process() 返回。
87      *p = *current; /* NOTE! this doesn't copy the supervisor stack */
88                  /* 注意! 这样做不会复制超级用户的堆栈 */ (只复制当前进程内容)。
89      p->state = TASK_UNINTERRUPTIBLE; // 将新进程的状态先置为不可中断等待状态。
90      p->pid = last_pid; // 新进程号。由前面调用 find_empty_process() 得到。
91      p->father = current->pid; // 设置父进程号。
92      p->counter = p->priority;
93      p->signal = 0; // 信号位图置 0。
94      p->alarm = 0;
95      p->leader = 0; /* process leadership doesn't inherit */
96                  /* 进程的领导力是不能继承的 */
97      p->utime = p->stime = 0; // 初始化用户态时间和核心态时间。
98      p->cutime = p->cstime = 0; // 初始化子进程用户态和核心态时间。
99      p->start_time = jiffies; // 当前滴答数时间。
100  // 以下设置任务状态段 TSS 所需的数据 (参见列表后说明)。
101  p->tss.back_link = 0;
102  p->tss.esp0 = PAGE_SIZE + (long) p; // 堆栈指针 (由于是给任务结构 p 分配了 1 页
103                                      // 新内存, 所以此时 esp0 正好指向该页顶端)。
104  p->tss.ss0 = 0x10; // 堆栈段选择符 (内核数据段) [??]。
105  p->tss.eip = eip; // 指令代码指针。
106  p->tss.eflags = eflags; // 标志寄存器。
107  p->tss.eax = 0;
108  p->tss.ecx = ecx;
109  p->tss.edx = edx;
110  p->tss.ebx = ebx;
111  p->tss.esp = esp;
112  p->tss.ebp = ebp;
113  p->tss.esi = esi;
114  p->tss.edi = edi;
115  p->tss.es = es & 0xffff; // 段寄存器仅 16 位有效。
116  p->tss.cs = cs & 0xffff;
117  p->tss.ss = ss & 0xffff;
118  p->tss.ds = ds & 0xffff;
119  p->tss.fs = fs & 0xffff;

```

```

110     p->tss.gs = gs & 0xffff;
111     p->tss.ldt = LDT(nr); // 该新任务 nr 的局部描述符表选择符 (LDT 的描述符在 GDT 中)。
112     p->tss.trace_bitmap = 0x80000000; (高 16 位有效)。
// 如果当前任务使用了协处理器, 就保存其上下文。
113     if (last_task_used_math == current)
114         __asm__ ("clds ; fnsave %0"::"m" (p->tss.i387));
// 设置新任务的代码和数据段基址、限长并复制页表。如果出错 (返回值不是 0), 则复位任务数组中
// 相应项并释放为该新任务分配的内存页。
115     if (copy_mem(nr,p)) { // 返回不为 0 表示出错。
116         task[nr] = NULL;
117         free_page((long) p);
118         return -EAGAIN;
119     }
// 如果父进程中有文件是打开的, 则将对对应文件的打开次数增 1。
120     for (i=0; i<NR_OPEN;i++)
121         if (f=p->filp[i])
122             f->f_count++;
// 将当前进程 (父进程) 的 pwd, root 和 executable 引用次数均增 1。
123     if (current->pwd)
124         current->pwd->i_count++;
125     if (current->root)
126         current->root->i_count++;
127     if (current->executable)
128         current->executable->i_count++;
// 设置新任务的 TSS 和 LDT 描述符项 (在 GDT 中), 数据从 task 结构中取。
// 在任务切换时, 任务寄存器 tr 由 CPU 自动加载。局部描述符表寄存器 ldtr 已在 task0 时加载。
129     set_tss_desc(gdt+(nr<<1)+FIRST_TSS_ENTRY, &(p->tss));
130     set_ldt_desc(gdt+(nr<<1)+FIRST_LDT_ENTRY, &(p->ldt));
131     p->state = TASK_RUNNING; // do this last, just in case
// 最后再将新任务设置成可运行状态, 以防万一 */
132     return last_pid; // 返回新进程号 (与任务号是不同的)。
133 }
134
// 为新进程取得不重复的进程号 last_pid, 并返回在任务数组中的任务号 (数组 index)。
135 int find_empty_process(void)
136 {
137     int i;
138
139     repeat:
140         if ((++last_pid)<0) last_pid=1;
141         for(i=0 ; i<NR_TASKS ; i++)
142             if (task[i] && task[i]->pid == last_pid) goto repeat;
143         for(i=1 ; i<NR_TASKS ; i++) // 任务 0 排除在外。
144             if (!task[i])
145                 return i;
146         return -EAGAIN;
147 }
148

```

## 4.10.3 其它信息

### 4.10.3.1 任务状态段信息

Figure 7-1. 80386 32-Bit Task State Segment

| 31              | 23 | 15                        | 7 | 0               |
|-----------------|----|---------------------------|---|-----------------|
| I/O MAP BASE    |    | 0 0 0 0 0 0 0 0           |   | 0 0 0 0 0 0 0 0 |
| 0 0 0 0 0 0 0 0 |    | 0 0 0 0 0 0 0 0           |   | T               |
| 0 0 0 0 0 0 0 0 |    | 0 0 0 0 0 0 0 0           |   |                 |
| 0 0 0 0 0 0 0 0 |    | 0 0 0 0 0 0 0 0           |   |                 |
| 0 0 0 0 0 0 0 0 |    | 0 0 0 0 0 0 0 0           |   |                 |
| 0 0 0 0 0 0 0 0 |    | 0 0 0 0 0 0 0 0           |   |                 |
| 0 0 0 0 0 0 0 0 |    | 0 0 0 0 0 0 0 0           |   |                 |
| 0 0 0 0 0 0 0 0 |    | 0 0 0 0 0 0 0 0           |   |                 |
|                 |    | EDI                       |   |                 |
|                 |    | ESI                       |   |                 |
|                 |    | EBP                       |   |                 |
|                 |    | ESP                       |   |                 |
|                 |    | EBX                       |   |                 |
|                 |    | EDX                       |   |                 |
|                 |    | ECX                       |   |                 |
|                 |    | EAX                       |   |                 |
|                 |    | EFLAGS                    |   |                 |
|                 |    | INSTRUCTION POINTER <EIP> |   |                 |
|                 |    | CR3 <PDPR>                |   |                 |
| 0 0 0 0 0 0 0 0 |    | 0 0 0 0 0 0 0 0           |   |                 |
|                 |    | ESP2                      |   |                 |
| 0 0 0 0 0 0 0 0 |    | 0 0 0 0 0 0 0 0           |   |                 |
|                 |    | ESP1                      |   |                 |
| 0 0 0 0 0 0 0 0 |    | 0 0 0 0 0 0 0 0           |   |                 |
|                 |    | ESP0                      |   |                 |
| 0 0 0 0 0 0 0 0 |    | 0 0 0 0 0 0 0 0           |   |                 |
|                 |    | BACK LINK TO PREVIOUS TSS |   |                 |

## NOTE

0 MEANS INTEL RESERVED. DO NOT DEFINE.

CPU 管理任务需要的所有信息被存储于一个特殊类型的段中，任务状态段 (task state segment - TSS)。图中显示出执行 80386 任务的 TSS 格式。

TSS 中的字段可以分为两类：

1. CPU 在进行任务切换时更新的动态信息集。这些字段有：
  - o 通用寄存器 (EAX, ECX, EDX, EBX, ESP, EBP, ESI, EDI);
  - o 段寄存器 (ES, CS, SS, DS, FS, GS);
  - o 标志寄存器 (EIP);

- o 指令指针 (EIP);
- 前一个执行任务的 TSS 的选择符 (仅当返回时才更新)。
- 2. CPU 读取但不会更改的静态信息集。这些字段有:
  - o 任务的 LDT 的选择符;
  - o 含有任务页目录基地址的寄存器 (PDBR);
  - o 特权级 0-2 的堆栈指针;
  - o 当任务进行切换时导致 CPU 产生一个调试(debug)异常的 T-比特位 (调试跟踪位);
  - o I/O 比特位图基地址 (其长度上限就是 TSS 的长度上限, 在 TSS 描述符中说明)。

任务状态段可以存放在线性空间的任何地方。与其它各类段相似, 任务状态段也是由描述符来定义的。当前正在执行任务的 TSS 是由任务寄存器 (TR) 来指示的。指令 LTR 和 STR 用来修改和读取任务寄存器中的选择符 (任务寄存器的可见部分)。

I/O 比特位图中的每 1 比特对应 1 个 I/O 端口。比如端口 41 的比特位就是 I/O 位图基地址+5, 位偏移 1 处。在保护模式中, 当遇到 1 个 I/O 指令时 (IN, INS, OUT, OUTS), CPU 首先就会检查当前特权级是否小于标志寄存器的 IOPL, 如果这个条件满足, 就执行该 I/O 操作。如果不满足, 那么 CPU 就会检查 TSS 中的 I/O 比特位图。如果相应比特位是置位的, 就会产生一般保护性异常, 否则就会执行该 I/O 操作。

## 4.11 sys.c 程序

### 4.11.1 功能描述

sys.c 程序主要包含有各个系统调用功能的实现函数。其中, 若返回值为 -ENOSYS, 则表示本版的 linux 还没有实现该功能, 可以参考目前的代码来了解它们的实现方法。

### 4.11.2 代码注释

列表 linux/kernel/sys.c 程序

---

```

1  /*
2   *  linux/kernel/sys.c
3   *
4   *  (C) 1991  Linus Torvalds
5   */
6
7  #include <errno.h>          // 错误号头文件。包含系统中各种出错号。(Linus 从 minix 中引进的)。
8
9  #include <linux/sched.h>    // 调度程序头文件, 定义了任务结构 task_struct、初始任务 0 的数据,
                             // 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
10 #include <linux/tty.h>      // tty 头文件, 定义了有关 tty_io, 串行通信方面的参数、常数。
11 #include <linux/kernel.h>    // 内核头文件。含有一些内核常用函数的原形定义。
12 #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。
13 #include <sys/times.h>       // 定义了进程中运行时间的结构 tms 以及 times() 函数原型。
14 #include <sys/utsname.h>     // 系统名称结构头文件。
15
16 // 返回日期和时间。
17 int sys_fptime()
18 {
```



```
18         return -ENOSYS;
19     }
20
21     //
22     int sys_break()
23     {
24         return -ENOSYS;
25     }
26
27     // 用于当前进程对子进程进行调试(debugging)。
28     int sys_ptrace()
29     {
30         return -ENOSYS;
31     }
32
33     // 改变并打印终端行设置。
34     int sys_stty()
35     {
36         return -ENOSYS;
37     }
38
39     // 取终端行设置信息。
40     int sys_gtty()
41     {
42         return -ENOSYS;
43     }
44
45     // 修改文件名。
46     int sys_rename()
47     {
48         return -ENOSYS;
49     }
50
51     //
52     int sys_prof()
53     {
54         return -ENOSYS;
55     }
56
57     // 设置当前任务的实际以及/或者有效组 ID (gid) 。如果任务没有超级用户特权，
58     // 那么只能互换其实际组 ID 和有效组 ID。如果任务具有超级用户特权，就能任意设置有效的和实际
59     // 的组 ID。保留的 gid (saved gid) 被设置成与有效 gid 同值。
60     int sys_setregid(int rgid, int egid)
61     {
62         if (rgid>0) {
63             if ((current->gid == rgid) ||
64                 suser())
65                 current->gid = rgid;
66             else
67                 return(-EPERM);
68         }
69         if (egid>0) {
70             if ((current->gid == egid) ||
```

```

62         (current->egid == egid) ||
63         (current->sgid == egid) ||
64         suser())
65         current->egid = egid;
66     else
67         return(-EPERM);
68     }
69     return 0;
70 }
71
// 设置进程组号(gid)。如果任务没有超级用户特权，它可以使用 setgid()将其有效 gid
// (effective gid) 设置为成其保留 gid(saved gid)或其实际 gid(real gid)。如果任务有
// 超级用户特权，则实际 gid、有效 gid 和保留 gid 都被设置成参数指定的 gid。
72 int sys_setgid(int gid)
73 {
74     return(sys_setregid(gid, gid));
75 }
76
// 打开或关闭进程计帐功能。
77 int sys_acct()
78 {
79     return -ENOSYS;
80 }
81
// 映射任意物理内存到进程的虚拟地址空间。
82 int sys_phys()
83 {
84     return -ENOSYS;
85 }
86
87 int sys_lock()
88 {
89     return -ENOSYS;
90 }
91
92 int sys_mpx()
93 {
94     return -ENOSYS;
95 }
96
97 int sys_ulimit()
98 {
99     return -ENOSYS;
100 }
101
// 返回从 1970 年 1 月 1 日 00:00:00 GMT 开始计时的时间值（秒）。如果 tloc 不为 null，则时间值
// 也存储在那里。
102 int sys_time(long * tloc)
103 {
104     int i;
105
106     i = CURRENT_TIME;
107     if (tloc) {

```

```

108         verify\_area(tloc, 4);    // 验证内存容量是否够（这里是 4 字节）。
109         put\_fs\_long(i, (unsigned long *)tloc);    // 也放入用户数据段 tloc 处。
110     }
111     return i;
112 }
113
114 /*
115  * Unprivileged users may change the real user id to the effective uid
116  * or vice versa.
117  */
118 /*
119  * 无特权的用户可以见实际用户标识符(real uid)改成有效用户标识符(effective uid)，反之亦然。
120  */
121 // 设置任务的实际以及/或者有效用户 ID (uid)。如果任务没有超级用户特权，那么只能互换其
122 // 实际用户 ID 和有效用户 ID。如果任务具有超级用户特权，就能任意设置有效的和实际的用户 ID。
123 // 保留的 uid (saved uid) 被设置成与有效 uid 同值。
124 int sys\_setreuid(int ruid, int euid)
125 {
126     int old_ruid = current->uid;
127
128     if (ruid>0) {
129         if ((current->euid==ruid) ||
130             (old_ruid == ruid) ||
131             suser())
132             current->uid = ruid;
133         else
134             return(-EPERM);
135     }
136     if (euid>0) {
137         if ((old_ruid == euid) ||
138             (current->euid == euid) ||
139             suser())
140             current->euid = euid;
141         else {
142             current->uid = old_ruid;
143             return(-EPERM);
144         }
145     }
146     return 0;
147 }
148
149 // 设置任务用户号(uid)。如果任务没有超级用户特权，它可以使用 setuid() 将其有效 uid
150 // (effective uid) 设置成其保留 uid(saved uid) 或其实际 uid(real uid)。如果任务有
151 // 超级用户特权，则实际 uid、有效 uid 和保留 uid 都被设置成参数指定的 uid。
152 int sys\_setuid(int uid)
153 {
154     return(sys\_setreuid(uid, uid));
155 }
156
157 // 设置系统时间和日期。参数 tptr 是从 1970 年 1 月 1 日 00:00:00 GMT 开始计时的时间值（秒）。
158 // 调用进程必须具有超级用户权限。
159 int sys\_stime(long * tptr)
160 {

```

```

150     if (!suser())                // 如果不是超级用户则出错返回（许可）。
151         return -EPERM;
152     startup_time = get_fs_long((unsigned long *)tptr) - jiffies/HZ;
153     return 0;
154 }
155
// 获取当前任务时间。tms 结构中包括用户时间、系统时间、子进程用户时间、子进程系统时间。
156 int sys_times(struct tms * tbuf)
157 {
158     if (tbuf) {
159         verify_area(tbuf, sizeof *tbuf);
160         put_fs_long(current->utime, (unsigned long *)&tbuf->tms_utime);
161         put_fs_long(current->stime, (unsigned long *)&tbuf->tms_stime);
162         put_fs_long(current->cutime, (unsigned long *)&tbuf->tms_cutime);
163         put_fs_long(current->cstime, (unsigned long *)&tbuf->tms_cstime);
164     }
165     return jiffies;
166 }
167
// 当参数 end_data_seg 数值合理，并且系统确实有足够的内存，而且进程没有超越其最大数据段大小
// 时，该函数设置数据段末尾为 end_data_seg 指定的值。该值必须大于代码结尾并且要小于堆栈
// 结尾 16KB。返回值是数据段的新结尾值（如果返回值与要求值不同，则表明有错发生）。
// 该函数并不被用户直接调用，而由 libc 库函数进行包装，并且返回值也不一样。
168 int sys_brk(unsigned long end_data_seg)
169 {
170     if (end_data_seg >= current->end_code &&                // 如果参数>代码结尾，并且
171         end_data_seg < current->start_stack - 16384) // 小于堆栈-16KB，
172         current->brk = end_data_seg;                        // 则设置新数据段结尾值。
173     return current->brk;                                     // 返回进程当前的数据段结尾值。
174 }
175
176 /*
177  * This needs some heave checking ...
178  * I just haven't get the stomach for it. I also don't fully
179  * understand sessions/prgrp etc. Let somebody who does explain it.
180  */
181 /*
182  * 下面代码需要某些严格的检查...
183  * 我只是没有胃口来做这些。我也不完全明白 sessions/prgrp 等。还是让了解它们的人来做吧。
184  */
185 // 将参数 pid 指定进程的进程组 ID 设置成 pgid。如果参数 pid=0，则使用当前进程号。如果
186 // pgid 为 0，则使用参数 pid 指定的进程的组 ID 作为 pgid。如果该函数用于将进程从一个
187 // 进程组移到另一个进程组，则这两个进程组必须属于同一个会话(session)。在这种情况下，
188 // 参数 pgid 指定了要加入的现有进程组 ID，此时该组的会话 ID 必须与将要加入进程的相同(193 行)。
189 int sys_setpgid(int pid, int pgid)
190 {
191     int i;
192
193     if (!pid)                // 如果参数 pid=0，则使用当前进程号。
194         pid = current->pid;
195     if (!pgid)                // 如果 pgid 为 0，则使用当前进程 pid 作为 pgid。
196         pgid = current->pid;  // [??这里与 POSIX 的描述有出入]
197     for (i=0 ; i<NR_TASKS ; i++) // 扫描任务数组，查找指定进程号的任务。

```

```

190         if (task[i] && task[i]->pid==pid) {
191             if (task[i]->leader)           // 如果该任务已经是首领，则出错返回。
192                 return -EPERM;
193             if (task[i]->session != current->session) // 如果该任务的会话 ID
194                 return -EPERM;           // 与当前进程的不同，则出错返回。
195             task[i]->pgrp = pgid;           // 设置该任务的 pgrp。
196             return 0;
197         }
198     return -ESRCH;
199 }
200
// 返回当前进程的组号。与 getpgid(0) 等同。
201 int sys_getpgrp(void)
202 {
203     return current->pgrp;
204 }
205
// 创建一个会话(session) (即设置其 leader=1)，并且设置其会话=其组号=其进程号。
206 int sys_setsid(void)
207 {
208     if (current->leader && !suser()) // 如果当前进程已是会话首领并且不是超级用户
209         return -EPERM;           // 则出错返回。
210     current->leader = 1;           // 设置当前进程为新会话首领。
211     current->session = current->pgrp = current->pid; // 设置本进程 session = pid。
212     current->tty = -1;           // 表示当前进程没有控制终端。
213     return current->pgrp;           // 返回会话 ID。
214 }
215
// 获取系统信息。其中 utsname 结构包含 5 个字段，分别是：本版本操作系统的名称、网络节点名称、
// 当前发行级别、版本级别和硬件类型名称。
216 int sys_uname(struct utsname * name)
217 {
218     static struct utsname thisname = { // 这里给出了结构中的信息，这种编码肯定会改变。
219         "linux.0", "nodename", "release ", "version ", "machine "
220     };
221     int i;
222
223     if (!name) return -ERROR;           // 如果存放信息的缓冲区指针为空则出错返回。
224     verify_area(name, sizeof *name); // 验证缓冲区大小是否超限 (超出已分配的内存等)。
225     for(i=0; i<sizeof *name; i++) // 将 utsname 中的信息逐字节复制到用户缓冲区中。
226         put_fs_byte((char *) &thisname[i], i+(char *) name);
227     return 0;
228 }
229
// 设置当前进程创建文件属性屏蔽码为 mask & 0777。并返回原屏蔽码。
230 int sys_umask(int mask)
231 {
232     int old = current->umask;
233
234     current->umask = mask & 0777;
235     return (old);
236 }
237

```

### 4.11.3 其它信息

## 4.12 vsprintf.c 程序

### 4.12.1 功能描述

### 4.12.2 代码注释

列表 linux/kernel/vsprintf.c 程序

```

1  /*
2  *  linux/kernel/vsprintf.c
3  *
4  *  (C) 1991  Linus Torvalds
5  */
6
7  /* vsprintf.c -- Lars Wirzenius & Linus Torvalds. */
8  /*
9  * Wirzenius wrote this portably, Torvalds fucked it up :-)
10 */
11
12 #include <stdarg.h>          // 标准参数头文件。以宏的形式定义变量参数列表。主要说明了一个
                               // 类型(va_list)和三个宏(va_start, va_arg 和 va_end)，用于
                               // vsprintf、vprintf、vfprintf 函数。
13 #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。
14
15 /* we use this so that we can do without the ctype library */
16 /* 我们使用下面的定义，这样我们就可以不使用 ctype 库了 */
17 #define is_digit(c)          ((c) >= '0' && (c) <= '9') // 判断字符是否数字字符。
18
19 // 该函数将字符数字串转换成整数。输入是数字串指针的指针，返回是结果数值。另外指针将前移。
20 static int skip_atoi(const char **s)
21 {
22     int i=0;
23     while (is_digit(**s))
24         i = i*10 + *((**s)++) - '0';
25     return i;
26 }
27
28 // 这里定义转换类型的各种符号常数。
29 #define ZEROPAD 1             /* pad with zero */           /* 填充零 */
30 #define SIGN 2                /* unsigned/signed long */         /* 无符号/符号长整数 */
31 #define PLUS 4                /* show plus */                   /* 显示加 */
32 #define SPACE 8               /* space if plus */               /* 如是加，则置空格 */
33 #define LEFT 16               /* left justified */               /* 左调整 */

```

```

32 #define SPECIAL 32          /* 0x */          /* 0x */
33 #define SMALL 64           /* use 'abcdef' instead of 'ABCDEF' */ /* 使用小写字母 */
34
// 除操作。输入：n 为被除数，base 为除数；结果：n 为商，函数返回值为余数。
// 参见 4.5.3 节有关嵌入汇编的信息。
35 #define do_div(n, base) ({ \
36 int __res; \
37 __asm__("divl %4": "=a" (n), "=d" (__res): "" (n), "I" (0), "r" (base)); \
38 __res; })
39
// 将整数转换为指定进制的字符串。
// 输入：num-整数；base-进制；size-字符串长度；precision-数字长度(精度)；type-类型选项。
// 输出：str 字符串指针。
40 static char * number(char * str, int num, int base, int size, int precision
41 , int type)
42 {
43     char c, sign, tmp[36];
44     const char * digits = "0123456789ABCDEFGHIJKLMNOPQRSTUVWXYZ";
45     int i;
46
// 如果类型 type 指出用小写字母，则定义小写字母集。
// 如果类型指出要左调整（靠左边界），则屏蔽类型中的填零标志。
// 如果进制基数小于 2 或大于 36，则退出处理，也即本程序只能处理基数在 2-32 之间的数。
47     if (type & SMALL) digits = "0123456789abcdefghijklmnopqrstuvwxyz";
48     if (type & LEFT) type &= ~ZEROPAD;
49     if (base < 2 || base > 36)
50         return 0;
// 如果类型指出要填零，则置字符变量 c='0'（也即''），否则 c 等于空格字符。
// 如果类型指出是带符号数并且数值 num 小于 0，则置符号变量 sign=负号，并使 num 取绝对值。
// 否则如果类型指出是加号，则置 sign=加号，否则若类型带空格标志则 sign=空格，否则置 0。
51     c = (type & ZEROPAD) ? ' ' : '0';
52     if (type & SIGN && num < 0) {
53         sign = '-';
54         num = -num;
55     } else
56         sign = (type & PLUS) ? '+' : ((type & SPACE) ? ' ' : 0);
// 若带符号，则宽度值减 1。若类型指出是特殊转换，则对于十六进制宽度再减少 2 位(用于 0x)，
// 对于八进制宽度减 1（用于八进制转换结果前放一个零）。
57     if (sign) size--;
58     if (type & SPECIAL)
59         if (base == 16) size -= 2;
60         else if (base == 8) size--;
// 如果数值 num 为 0，则临时字符串='0'；否则根据给定的基数将数值 num 转换成字符形式。
61     i = 0;
62     if (num == 0)
63         tmp[i++] = '0';
64     else while (num != 0)
65         tmp[i++] = digits[do_div(num, base)];
// 若数值字符个数大于精度值，则精度值扩展为数字个数。
// 宽度值 size 减去用于存放数值字符的个数。
66     if (i > precision) precision = i;
67     size -= precision;

```

```

// 从这里真正开始形成所需要的转换结果，并暂时放在字符串 str 中。
// 若类型中没有填零(ZEROPAD)和左靠齐（左调整）标志，则在 str 中首先
// 填充剩余宽度值指出的空格数。若需带符号位，则存入符号。
68     if (!(type&(ZEROPAD+LEFT)))
69         while(size-->0)
70             *str++ = ' ';
71     if (sign)
72         *str++ = sign;
// 若类型指出是特殊转换，则对于八进制转换结果头一位放置一个'0'；而对于十六进制则存放'0x'。
73     if (type&SPECIAL)
74         if (base==8)
75             *str++ = '0';
76         else if (base==16) {
77             *str++ = '0';
78             *str++ = digits[33]; // 'X' 或 'x'
79         }
// 若类型中没有左调整（左靠齐）标志，则在剩余宽度中存放 c 字符（'0' 或空格），见 51 行。
80     if (!(type&LEFT))
81         while(size-->0)
82             *str++ = c;
// 此时 i 存有数值 num 的数字个数。若数字个数小于精度值，则 str 中放入（精度值-i）个'0'。
83     while(i<precision--)
84         *str++ = '0';
// 将转换数值换好的数字字符填入 str 中。共 i 个。
85     while(i-->0)
86         *str++ = tmp[i];
// 若宽度值仍大于零，则表示类型标志中有左靠齐标志标志。则在剩余宽度中放入空格。
87     while(size-->0)
88         *str++ = ' ';
89     return str; // 返回转换好的字符串。
90 }
91
// 下面函数是送格式化输出到字符串中。
// 为了能在内核中使用格式化的输出，Linux 在内核实现了该 C 标准函数。
// 其中参数 fmt 是格式字符串；args 是个数变化的值；buf 是输出字符串缓冲区。
// 请参见本代码列表后的有关格式转换字符的介绍。
92 int vsprintf(char *buf, const char *fmt, va_list args)
93 {
94     int len;
95     int i;
96     char *str; // 用于存放转换过程中的字符串。
97     char *s;
98     int *ip;
99
100     int flags; /* flags to number() */
101               /* number() 函数使用的标志 */
102     int field_width; /* width of output field */
103               /* 输出字段宽度 */
104     int precision; /* min. # of digits for integers; max
105                  number of chars for from string */
106               /* min. 整数数字个数; max. 字符串中字符个数 */
107     int qualifier; /* 'h', 'l', or 'L' for integer fields */

```



```

106                                     /* 'h', 'l', 或 'L' 用于整数字段 */
// 首先将字符指针指向 buf, 然后扫描格式字符串, 对各个格式转换指示进行相应的处理。
107     for (str=buf ; *fmt ; ++fmt) {
// 格式转换指示字符串均以 '%' 开始, 这里从 fmt 格式字符串中扫描 '%', 寻找格式转换字符串的开始。
// 不是格式指示的一般字符均被依次存入 str。
108         if (*fmt != '%') {
109             *str++ = *fmt;
110             continue;
111         }
112
// 下面取得格式指示字符串中的标志域, 并将标志常量放入 flags 变量中。
113         /* process flags */
114         flags = 0;
115         repeat:
116             ++fmt;          /* this also skips first '%' */
117             switch (*fmt) {
118                 case '-': flags |= LEFT; goto repeat;    // 左靠齐调整。
119                 case '+': flags |= PLUS; goto repeat;    // 放加号。
120                 case ' ': flags |= SPACE; goto repeat;    // 放空格。
121                 case '#': flags |= SPECIAL; goto repeat; // 是特殊转换。
122                 case '0': flags |= ZEROPAD; goto repeat; // 要填零(即'0')。
123             }
124
// 取当前参数字段宽度域值, 放入 field_width 变量中。如果宽度域中是数值则直接取其为宽度值。
// 如果宽度域中是字符 '*', 表示下一个参数指定宽度。因此调用 va_arg 取宽度值。若此时宽度值
// 小于 0, 则该负数表示其带有标志域 '-' 标志 (左靠齐), 因此还需在标志变量中添入该标志, 并
// 将字段宽度值取为其绝对值。
125         /* get field width */
126         field_width = -1;
127         if (is_digit(*fmt))
128             field_width = skip_atoi(&fmt);
129         else if (*fmt == '*') {
130             /* it's the next argument */
131             field_width = va_arg(args, int);
132             if (field_width < 0) {
133                 field_width = -field_width;
134                 flags |= LEFT;
135             }
136         }
137
// 下面这段代码, 取格式转换串的精度域, 并放入 precision 变量中。精度域开始的标志是 '.'。
// 其处理过程与上面宽度域的类似。如果精度域中是数值则直接取其为精度值。如果精度域中是
// 字符 '*', 表示下一个参数指定精度。因此调用 va_arg 取精度值。若此时宽度值小于 0, 则
// 将字段精度值取为其绝对值。
138         /* get the precision */
139         precision = -1;
140         if (*fmt == '.') {
141             ++fmt;
142             if (is_digit(*fmt))
143                 precision = skip_atoi(&fmt);
144             else if (*fmt == '*') {
145                 /* it's the next argument */
146                 precision = va_arg(args, int);

```

```

147         }
148         if (precision < 0)
149             precision = 0;
150     }
151
152     // 下面这段代码分析长度修饰符，并将其存入 qualifier 变量。（h, l, L 的含义参见列表后的说明）。
153     /* get the conversion qualifier */
154     qualifier = -1;
155     if (*fmt == 'h' || *fmt == 'l' || *fmt == 'L') {
156         qualifier = *fmt;
157         ++fmt;
158     }
159
160     // 下面分析转换指示符。
161     switch (*fmt) {
162     // 如果转换指示符是 'c'，则表示对应参数应是字符。此时如果标志域表明不是左靠齐，则该字段前面
163     // 放入宽度域值-1 个空格字符，然后再放入参数字符。如果宽度域还大于 0，则表示为左靠齐，则在
164     // 参数字符后面添加宽度值-1 个空格字符。
165     case 'c':
166         if (!(flags & LEFT))
167             while (--field_width > 0)
168                 *str++ = ' ';
169         *str++ = (unsigned char) va_arg(args, int);
170         while (--field_width > 0)
171             *str++ = ' ';
172         break;
173
174     // 如果转换指示符是 's'，则表示对应参数是字符串。首先取参数字符串的长度，若其超过了精度域值，
175     // 则扩展精度域=字符串长度。此时如果标志域表明不是左靠齐，则该字段前放入(宽度值-字符串长度)
176     // 个空格字符。然后再放入参数字符串。如果宽度域还大于 0，则表示为左靠齐，则在参数字符串后面
177     // 添加(宽度值-字符串长度)个空格字符。
178     case 's':
179         s = va_arg(args, char *);
180         len = strlen(s);
181         if (precision < 0)
182             precision = len;
183         else if (len > precision)
184             len = precision;
185
186         if (!(flags & LEFT))
187             while (len < field_width--)
188                 *str++ = ' ';
189         for (i = 0; i < len; ++i)
190             *str++ = *s++;
191         while (len < field_width--)
192             *str++ = ' ';
193         break;
194
195     // 如果格式转换符是 'o'，表示需将对应的参数转换成八进制数的字符串。调用 number() 函数处理。
196     case 'o':
197         str = number(str, va_arg(args, unsigned long), 8,
198                     field_width, precision, flags);
199         break;

```

```

190 // 如果格式转换符是'p'，表示对应参数的一个指针类型。此时若该参数没有设置宽度域，则默认宽度
191 // 为8，并且需要添零。然后调用 number() 函数进行处理。
192 case 'p':
193     if (field_width == -1) {
194         field_width = 8;
195         flags |= ZEROPAD;
196     }
197     str = number(str,
198                 (unsigned long) va_arg(args, void *), 16,
199                 field_width, precision, flags);
200     break;
201 // 若格式转换指示是'x'或'X'，则表示对应参数需要打印成十六进制数输出。'x'表示用小写字母表示。
202 case 'x':
203     flags |= SMALL;
204 case 'X':
205     str = number(str, va_arg(args, unsigned long), 16,
206                 field_width, precision, flags);
207     break;
208 // 如果格式转换字符是'd','i'或'u'，则表示对应参数是整数，'d','i'代表符号整数，因此需要加上
209 // 带符号标志。'u'代表无符号整数。
210 case 'd':
211 case 'i':
212     flags |= SIGN;
213 case 'u':
214     str = number(str, va_arg(args, unsigned long), 10,
215                 field_width, precision, flags);
216     break;
217 // 若格式转换指示符是'n'，则表示要把到目前为止转换输出的字符数保存到对应参数指针指定的位置中。
218 // 首先利用 va_arg() 取得该参数指针，然后将已经转换好的字符数存入该指针所指的位置。
219 case 'n':
220     ip = va_arg(args, int *);
221     *ip = (str - buf);
222     break;
223 // 若格式转换符不是'%', 则表示格式字符串有错，直接将一个 '%' 写入输出串中。
224 // 如果格式转换符的位置处还有字符，则也直接将该字符写入输出串中，并返回到 107 行继续处理
225 // 格式字符串。否则表示已经处理到格式字符串的结尾处，则退出循环。
226 default:
227     if (*fmt != '%')
228         *str++ = '%';
229     if (*fmt)
230         *str++ = *fmt;
231     else
232         --fmt;
233     break;
234 }
235 }
236 *str = '\0'; // 最后在转换好的字符串结尾处添上 null。
237 return str - buf; // 返回转换好的字符串长度值。

```

[233](#) }  
[234](#)

## 4.12.3 其它信息

### 4.12.3.1 vsprintf()的格式字符串

```
int vsprintf(char *buf, const char *fmt, va_list args)
```

vsprintf()函数是printf()系列函数之一。这些函数都产生格式化的输出：接受确定输出格式的格式字符串fmt，用格式字符串对个数变化的参数进行格式化，产生格式化的输出。

printf 直接把输出送到标准输出句柄 stdout。cprintf 把输出送到控制台。fprintf 把输出送到文件句柄。printf 前带'v'字符的(例如vfprintf)表示参数是从va\_arg数组的va\_list args中接受。printf 前面带's'字符则表示把输出送到以null结尾的字符串buf中(此时用户应确保buf有足够的空间存放字符串)。下面详细说明格式字符串的使用方法。

#### 1. 格式字符串

printf 系列函数中的格式字符串用于控制函数转换方式、格式化和输出其参数。对于每个格式，必须有对应的参数，参数过多将被忽略。格式字符串中含有两类成份，一种是将被直接复制到输出中的简单字符；另一种是用于对对应参数进行格式化的转换指示字符串。

#### 2. 格式指示字符串

格式指示串的形式如下：

```
%[flags][width][.prec][h|l|L][type]
```

每一个转换指示串均需要以百分号(%)开始。其中

|         |                       |
|---------|-----------------------|
| [flags] | 是可选的标志字符序列；           |
| [width] | 是可选的宽度指示符；            |
| [.prec] | 是可选的精度(precision)指示符； |
| [h l L] | 是可选的输入长度修饰符；          |
| [type]  | 是转换类型字符(或称为转换指示符)。    |

flags 控制输出对齐方式、数值符号、小数点、尾零、二进制、八进制或十六进制等，参见上面列表27-33行的注释。标志字符及其含义如下：

# 表示需要将相应参数转换为“特殊形式”。对于八进制(o)，则转换后的字符串的首位必须是一个零。对于十六进制(x或X)，则转换后的字符串需以'0x'或'0X'开头。对于e, E, f, F, g以及G，则即使没有小数位，转换结果也将总是有一个小数点。对于g或G，后拖的零也不会删除。

0 转换结果应该是附零的。对于d, i, o, u, x, X, e, E, f, g和G，转换结果的左边将用零填空而不是用空格。如果同时出现0和-标志，则0标志将被忽略。对于数值转换，如果给出了精度域，0标志也被忽略。

- 转换后的结果在相应字段边界内将作左调整(靠左)。(默认是作右调整—靠右)。n转换例外，转换结果将在右面填空格。

' ' 表示带符号转换产生的一个正数结果前应该留一个空格。

+ 表示在一个符号转换结果之前总需要放置一个符号(+或-)。对于默认情况，只有负数使用负号。

width 指定了输出字符串宽度，即指定了字段的最小宽度值。如果被转换的结果要比指定的宽度小，

则在其左边（或者右边，如果给出了左调整标志）需要填充空格或零（由 flags 标志确定）的个数等。除了使用数值来指定宽度域以外，也可以使用 '\*' 来指出字段的宽度由下一个整型参数给出。当转换值宽度大于 width 指定的宽度时，在任何情况下小宽度值都不会截断结果。字段宽度会扩充以包含完整结果。

precision 是说明输出数字起码的个数。对于 d, I, o, u, x 和 X 转换，精度值指出了起码出现数字的个数。对于 e, E, f 和 F，该值指出在小数点之后出现的数字的个数。对于 g 或 G，指出最大有效数字个数。对于 s 或 S 转换，精度值说明输出字符串的最大字符数。

长度修饰指示符说明了整型数转换后的输出类型形式。下面叙述中‘整型数转换’代表 d, i, o, u, x 或 X 转换。

- hh 说明后面的整型数转换对应于一个带符号字符或无符号字符参数。
- h 说明后面的整型数转换对应于一个带符号整数或无符号短整数参数。
- l 说明后面的整型数转换对应于一个长整数或无符号长整数参数。
- ll 说明后面的整型数转换对应于一个长长整数或无符号长长整数参数。
- L 说明 e, E, f, F, g 或 G 转换结果对应于一个长双精度参数。

type 是说明接受的输入参数类型和输出的格式。各个转换指示符的含义如下：

d, I 整型参数将被转换为带符号整数。如果有精度 (precision) 的话，则给出了需要输出的最少数字个数。如果被转换的值数字个数较少，就会在其左边添零。默认的精度值是 1。

o, u, x, X 会将无符号的整数转换为无符号八进制(o)、无符号十进制(u)或者是无符号十六进制(x 或 X)表示方式输出。x 表示要使用小写字母 (abcdef) 来表示十六进制数，X 表示用大写字母 (ABCDEF) 表示十六进制数。如果存在精度域的话，说明需要输出的最少数字个数。如果被转换的值数字个数较少，就会在其左边添零。默认的精度值是 1。

e, E 这两个转换字符用于经四舍五入将参数转换成[-]d.ddde±dd 的形式。小数点之后的数字个数等于精度。如果没有精度域，就取默认值 6。如果精度是 0，则没有小数出现。E 表示用大写字母 E 来表示指数。指数部分总是用 2 位数字表示。如果数值为 0，那么指数就是 00。

f, F 这两个转换字符用于经四舍五入将参数转换成[-]ddd.ddd 的形式。小数点之后的数字个数等于精度。如果没有精度域，就取默认值 6。如果精度是 0，则没有小数出现。如果有小数点，那么后面起码会有 1 位数字。

g, G 这两个转换字符将参数转换为 f 或 e 的格式（如果是 G，则是 F 或 E 格式）。精度值指定了整数的个数。如果没有精度域，则其默认值为 6。如果精度为 0，则作为 1 来对待。如果转换时指数小于-4 或大于等于精度，则采用 e 格式。小数部分后拖的零将被删除。仅当起码有一位小数时才会出现小数点。

c 参数将被转换成无符号字符并输出转换结果。

s 要求输入为指向字符串的指针，并且该字符串要以 null 结尾。如果有精度域，则只输出精度所要求的字符个数，并且字符串无须以 null 结尾。

p 以指针形式输出十六进制数。

n 用于把到目前为止转换输出的字符个数保存到由对应输入指针指定的位置中。不对参数进行转换。

% 输出一个百分号%，不进行转换。也即此时整个转换指示为%%。

#### 4.12.4 与当前版本的区别

由于该文件也属于库函数，所以从 1.2 版内核开始就直接使用库中的函数了。也即删除了该文件。

## 4.13 printk.c 程序

### 4.13.1 功能描述

### 4.13.2 代码注释

列表 linux/kernel/printk.c 程序

```

1  /*
2   * linux/kernel/printk.c
3   *
4   * (C) 1991 Linus Torvalds
5   */
6
7  /*
8   * When in kernel-mode, we cannot use printf, as fs is liable to
9   * point to 'interesting' things. Make a printf with fs-saving, and
10  * all is well.
11  */
12  /*
13   * 当处于内核模式时，我们不能使用 printf，因为寄存器 fs 指向其它不感兴趣的地方。
14   * 自己编制一个 printf 并在使用前保存 fs，一切就解决了。
15  */
16  #include <stdarg.h>          // 标准参数头文件。以宏的形式定义变量参数列表。主要说明了一个
17                               // 类型(va_list)和三个宏(va_start, va_arg 和 va_end)，用于
18                               // vsprintf、vprintf、vfprintf 函数。
19  #include <stddef.h>          // 标准定义头文件。定义了 NULL, offsetof(TYPE, MEMBER)。
20
21  #include <linux/kernel.h> // 内核头文件。含有一些内核常用函数的原形定义。
22
23  static char buf[1024];
24
25  // 下面该函数 vsprintf() 在 linux/kernel/vsprintf.c 中 92 行开始。
26  extern int vsprintf(char * buf, const char * fmt, va_list args);
27
28  int printk(const char *fmt, ...)
29  {
30      va_list args;                // va_list 实际上是一个字符指针类型。
31      int i;
32
33      va_start(args, fmt);         // 参数处理开始函数。在 (include/stdarg.h, 13)
34      i=vsprintf(buf, fmt, args);  // 使用格式串 fmt 将参数列表 args 输出到 buf 中。
35                                   // 返回值 i 等于输出字符串的长度。
36      va_end(args);               // 参数处理结束函数。
37      __asm__ ("push %%fs\n\t"    // 保存 fs。
38              "push %%ds\n\t"
39              "pop %%fs\n\t"      // 令 fs = ds。
40              "pushl %0\n\t"      // 将字符串长度压入堆栈(这三个入栈是调用参数)。

```

```

33         "pushl $_buf\n\t"      // 将 buf 的地址压入堆栈。
34         "pushl $0\n\t"        // 将数值 0 压入堆栈。是通道号 channel。
35         "call _tty_write\n\t"  // 调用 tty_write 函数。(kernel/chr_drv/tty_io.c, 290)。
36         "addl $8, %%esp\n\t"   // 跳过（丢弃）两个入栈参数(buf, channel)。
37         "popl %0\n\t"         // 弹出字符串长度值，作为返回值。
38         "pop %%fs"            // 恢复原 fs 寄存器。
39         "::"r"(i): "ax", "cx", "dx"); // 通知编译器，寄存器 ax, cx, dx 值可能已经改变。
40     return i;                // 返回字符串长度。
41 }
42

```

## 4.14 panic.c 程序

### 4.14.1 功能描述

当内核程序出错时，则调用函数 panic()，显示错误信息并使系统进入死循环。

在内核程序的许多地方，若出现严重出错时就要调用到该函数。在很多情况下，调用 panic() 函数是一种简明的处理方法。这样做很好地遵循了 UNIX “尽量简明” 的原则。

panic 是“惊慌，恐慌”的意思。在 Douglas Adams 的小说《Hitch hikers Guide to the Galaxy》（《银河徒步旅行者指南》），该书最有名的一句话就是“Don't Panic!”。该系列小说是 linux 骇客最常阅读的书之一。

### 4.14.2 代码注释

列表 linux/kernel/panic.c 程序

```

1  /*
2   * linux/kernel/panic.c
3   *
4   * (C) 1991 Linus Torvalds
5   */
6
7  /*
8   * This function is used through-out the kernel (includeinh mm and fs)
9   * to indicate a major problem.
10  */
11  /*
12   * 该函数在整个内核中使用（包括在 头文件*.h，内存管理程序 mm 和文件系统 fs 中），
13   * 用以指出主要的出错问题。
14   */
15  #include <linux/kernel.h> // 内核头文件。含有一些内核常用函数的原形定义。
16  #include <linux/sched.h> // 调度程序头文件，定义了任务结构 task_struct、初始任务 0 的数据，
17                          // 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
18
19  void sys_sync(void); /* it's really int */ /* 实际上是整型 int (fs/buffer.c, 44) */
20
21  // 该函数用来显示内核中出现的重大错误信息，并运行文件系统同步函数，然后进入死循环 -- 死机。
22  // 如果当前进程是任务 0 的话，还说明是交换任务出错，并且还没有运行文件系统同步函数。
23  volatile void panic(const char * s)

```

```
17 {  
18     printk("Kernel panic: %s\n\r", s);  
19     if (current == task[0])  
20         printk("In swapper task - not syncing\n\r");  
21     else  
22         sys\_sync();  
23     for(;;);  
24 }  
25
```







---



## 第 5 章 块设备驱动程序(block driver)

### 5.1 概述

列表 linux/kernel/blk\_drv 目录

| Name  | Size        | Last modified (GMT) | Description |
|---|-------------|---------------------|-------------|
|  <a href="#">Makefile</a>      | 1951 bytes  | 1991-12-05 19:59:42 | m           |
|  <a href="#">blk.h</a>         | 3464 bytes  | 1991-12-05 19:58:01 | m           |
|  <a href="#">floppy.c</a>      | 11429 bytes | 1991-12-07 00:00:38 | m           |
|  <a href="#">hd.c</a>          | 7807 bytes  | 1991-12-05 19:58:17 | m           |
|  <a href="#">ll_rw_blk.c</a> | 3539 bytes  | 1991-12-04 13:41:42 | m           |
|  <a href="#">ramdisk.c</a>   | 2740 bytes  | 1991-12-06 03:08:06 | m           |

### 5.2 Makefile 文件

#### 5.2.1 功能描述

#### 5.2.2 代码注释

列表 linux/kernel/blk\_drv/Makefile 文件

```

1 #
2 # Makefile for the FREAX-kernel block device drivers.
3 #
4 # Note! Dependencies are done automagically by 'make dep', which also
5 # removes any old dependencies. DON'T put your own dependencies here
6 # unless it's something special (ie not a .c file).
7 #
  # FREAX 内核块设备驱动程序的 Makefile 文件
  # 注意！依赖关系是由'make dep'自动进行的，它也会自动去除原来的依赖信息。不要把你自己的
  # 依赖关系信息放在这里，除非是特别文件的（也即不是一个.c文件的信息）。

```

# (Linux 最初的名字叫 FREAX, 后来被 ftp.funet.fi 的管理员改成 Linux 这个名字)。

```

8
9 AR      =gar      # GNU 的二进制文件处理程序, 用于创建、修改以及从归档文件中抽取文件。
10 AS      =gas      # GNU 的汇编程序。
11 LD      =gld      # GNU 的连接程序。
12 LDFLAGS =-s -x    # 连接程序所有的参数, -s 输出文件中省略所有符号信息。-x 删除所有局部符号。
13 CC      =gcc      # GNU C 语言编译器。
    # 下一行是 C 编译程序选项。-Wall 显示所有的警告信息; -O 优化选项, 优化代码长度和执行时间;
    # -fstrength-reduce 优化循环执行代码, 排除重复变量; -fomit-frame-pointer 省略保存不必要
    # 的框架指针; -fcombine-regs 合并寄存器, 减少寄存器类的使用; -finline-functions 将所有简
    # 单短小的函数代码嵌入调用程序中; -mstring-insns Linux 自己添加的优化选项, 以后不再使用;
    # -nostdinc -I../include 不使用默认路径中的包含文件, 而使用指定目录中的(.././include)。
14 CFLAGS  =-Wall -O -fstrength-reduce -fomit-frame-pointer -fcombine-regs \
15          -finline-functions -mstring-insns -nostdinc -I.././include
    # C 前处理选项。-E 只运行 C 前处理, 对所有指定的 C 程序进行预处理并将处理结果输出到标准输
    # 出设备或指定的输出文件中; -nostdinc -I.././include 同前。
16 CPP     =gcc -E -nostdinc -I.././include
17
    # 下面的规则指示 make 利用下面的命令将所有的.c 文件编译生成.s 汇编程序。该规则的命令
    # 指使 gcc 采用 CFLAGS 所指定的选项对 C 代码编译后不进行汇编就停止 (-S), 从而产生与
    # 输入的各个 C 文件对应的汇编代码文件。默认情况下所产生的汇编程序文件名是原 C 文件名
    # 去掉.c 而加上.s 后缀。-o 表示其后是输出文件的名称。其中$.s (或$@) 是自动目标变量,
    # $<代表第一个先决条件, 这里即是符合条件*.c 的文件。
18 .c.s:
19     $(CC) $(CFLAGS) \
20     -S -o $.s $<
    # 下面规则表示将所有.s 汇编程序文件编译成.o 目标文件。22 行是实现该操作的具体命令。
21 .s.o:
22     $(AS) -c -o $.o $<
23 .c.o:                                # 类似上面, *.c 文件→*.o 目标文件。不进行连接。
24     $(CC) $(CFLAGS) \
25     -c -o $.o $<
26
27 OBJS   = ll_rw_blk.o floppy.o hd.o ramdisk.o          # 定义目标文件变量 OBJS。
28
    # 在有了先决条件 OBJS 后使用下面的命令连接成目标 blk_drv.a 库文件。
29 blk_drv.a: $(OBJS)
30     $(AR) rcs blk_drv.a $(OBJS)
31     sync
32
    # 下面的规则用于清理工作。当执行'make clean'时, 就会执行 34--35 行上的命令, 去除所有编译
    # 连接生成的文件。'rm' 是文件删除命令, 选项-f 含义是忽略不存在的文件, 并且不显示删除信息。
33 clean:
34     rm -f core *.o *.a tmp_make
35     for i in *.c;do rm -f `basename $$i .c`.s;done
36
    # 下面得目标或规则用于检查各文件之间的依赖关系。方法如下:
    # 使用字符串编辑程序 sed 对 Makefile 文件 (即是本文件) 进行处理, 输出为删除 Makefile
    # 文件中'### Dependencies' 行后面的所有行 (下面从 44 开始的行), 并生成 tmp_make
    # 临时文件 (38 行的作用)。然后对 kernel/blk_drv/目录下的每个 C 文件执行 gcc 预处理操作。
    # -M 标志告诉预处理程序输出描述每个目标文件相关性的规则, 并且这些规则符合 make 语法。
    # 对于每一个源文件, 预处理程序输出一个 make 规则, 其结果形式是相应源程序文件的目标
    # 文件名加上其依赖关系--该源文件中包含的所有头文件列表。把预处理结果都添加到临时

```

# 文件 tmp\_make 中, 然后将该临时文件复制成新的 Makefile 文件。

```

37 dep:
38     sed '/\#\#\# Dependencies/q' < Makefile > tmp_make
39     (for i in *.c;do echo -n `echo $$i | sed 's,\.c,\.s,``" "; \
40         $(CPP) -M $$i;done) >> tmp_make
41     cp tmp_make Makefile
42
43 ### Dependencies:
44 floppy.s floppy.o : floppy.c ../../include/linux/sched.h ../../include/linux/head.h \
45     ../../include/linux/fs.h ../../include/sys/types.h ../../include/linux/mm.h \
46     ../../include/signal.h ../../include/linux/kernel.h \
47     ../../include/linux/fdreg.h ../../include/asm/system.h \
48     ../../include/asm/io.h ../../include/asm/segment.h blk.h
49 hd.s hd.o : hd.c ../../include/linux/config.h ../../include/linux/sched.h \
50     ../../include/linux/head.h ../../include/linux/fs.h \
51     ../../include/sys/types.h ../../include/linux/mm.h ../../include/signal.h \
52     ../../include/linux/kernel.h ../../include/linux/hdreg.h \
53     ../../include/asm/system.h ../../include/asm/io.h \
54     ../../include/asm/segment.h blk.h
55 ll_rw_blk.s ll_rw_blk.o : ll_rw_blk.c ../../include/errno.h ../../include/linux/sched.h \
56     ../../include/linux/head.h ../../include/linux/fs.h \
57     ../../include/sys/types.h ../../include/linux/mm.h ../../include/signal.h \
58     ../../include/linux/kernel.h ../../include/asm/system.h blk.h

```

## 5.3 blk.h 文件

### 5.3.1 功能描述

这是有关硬盘块设备参数的头文件, 因为只用于块设备, 所以与块设备代码放在同一个地方。

### 5.3.2 代码注释

列表 linux/kernel/blk\_drv/blk.h 文件

```

1 #ifndef BLK\_H
2 #define BLK\_H
3
4 #define NR\_BLK\_DEV      7          // 块设备的数量。
5 /*
6  * NR_REQUEST is the number of entries in the request-queue.
7  * NOTE that writes may use only the low 2/3 of these: reads
8  * take precedence.
9  *
10 * 32 seems to be a reasonable number: enough to get some benefit
11 * from the elevator-mechanism, but not so much as to lock a lot of
12 * buffers when they are in the queue. 64 seems to be too many (easily
13 * long pauses in reading when heavy writing/syncing is going on)
14 */
15 /*
16  * 下面定义的 NR_REQUEST 是请求队列中所包含的项数。

```

```

* 注意，读操作仅使用这些项低端的 2/3；读操作优先处理。
*
* 32 项好像是一个合理的数字：已经足够从电梯算法中获得好处，
* 但当缓冲区在队列中而锁住时又不显得是很大的数。64 就看上
* 去太大了（当大量的写/同步操作运行时很容易引起长时间的暂停）。
*/
15 #define NR_REQUEST 32
16
17 /*
18  * Ok, this is an expanded form so that we can use the same
19  * request for paging requests when that is implemented. In
20  * paging, 'bh' is NULL, and 'waiting' is used to wait for
21  * read/write completion.
22  */
/*
* OK，下面是 request 结构的一个扩展形式，因而当实现以后我们就可以在分页请求中
* 使用同样的请求。在分页处理中，'bh' 是 NULL，而'waiting' 则用于等待读/写的完成。
*/
// 下面是请求队列中项的结构。其中如果 dev=-1，则表示没有使用该项。
23 struct request {
24     int dev;                /* -1 if no request */    // 使用的设备号。
25     int cmd;                /* READ or WRITE */      // 命令(READ 或 WRITE)。
26     int errors;             // 操作时产生的错误次数。
27     unsigned long sector;   // 起始扇区。(1 块=2 扇区)
28     unsigned long nr_sectors; // 读/写扇区数。
29     char * buffer;         // 数据缓冲区。
30     struct task_struct * waiting; // 任务等待操作执行完成的地方。
31     struct buffer_head * bh;    // 缓冲区头指针(include/linux/fs.h, 68)。
32     struct request * next;     // 指向下一请求项。
33 };
34
35 /*
36  * This is used in the elevator algorithm: Note that
37  * reads always go before writes. This is natural: reads
38  * are much more time-critical than writes.
39  */
/*
* 下面的定义用于电梯算法：注意读操作总是在写操作之前进行。
* 这是很自然的：读操作对时间的要求要比写严格得多。
*/
40 #define IN_ORDER(s1,s2) \
41 ((s1)->cmd<(s2)->cmd || (s1)->cmd==(s2)->cmd && \
42 ((s1)->dev < (s2)->dev || ((s1)->dev == (s2)->dev && \
43 (s1)->sector < (s2)->sector)))
44
// 块设备结构。
45 struct blk_dev_struct {
46     void (*request_fn)(void);    // 请求操作的函数指针。
47     struct request * current_request; // 请求信息结构。
48 };
49
50 extern struct blk_dev_struct blk_dev[NR_BLK_DEV]; // 块设备数组，每种块设备占用一项。
51 extern struct request request[NR_REQUEST]; // 请求队列数组。

```

```

52 extern struct task\_struct * wait\_for\_request;           // 等待请求的任务结构。
53
54 #ifdef MAJOR\_NR           // 主设备号。
55
56 /*
57  * Add entries as needed. Currently the only block devices
58  * supported are hard-disks and floppies.
59  */
60 /*
61  * 需要时加入条目。目前块设备仅支持硬盘和软盘。
62  */
63
64 #if (MAJOR\_NR == 1) // RAM 盘的主设备号是 1。根据这里的定义可以推理内存块主设备号也为 1。
65 /* ram disk */      /* RAM 盘（内存虚拟盘）*/
66 #define DEVICE\_NAME "ramdisk"           // 设备名称 ramdisk。
67 #define DEVICE\_REQUEST do\_rd\_request    // 设备请求函数 do_rd_request()。
68 #define DEVICE\_NR(device) ((device) & 7) // 设备号 (0--7)。
69 #define DEVICE\_ON(device)               // 开启设备。虚拟盘无须开启和关闭。
70 #define DEVICE\_OFF(device)             // 关闭设备。
71
72 #elif (MAJOR\_NR == 2) // 软驱的主设备号是 2。
73 /* floppy */
74 #define DEVICE\_NAME "floppy"           // 设备名称 floppy。
75 #define DEVICE\_INTR do\_floppy         // 设备中断处理程序 do_floppy()。
76 #define DEVICE\_REQUEST do\_fd\_request    // 设备请求函数 do_fd_request()。
77 #define DEVICE\_NR(device) ((device) & 3) // 设备号 (0--3)。
78 #define DEVICE\_ON(device) floppy\_on(DEVICE\_NR(device)) // 开启设备函数 floppyon()。
79 #define DEVICE\_OFF(device) floppy\_off(DEVICE\_NR(device)) // 关闭设备函数 floppyoff()。
80
81 #elif (MAJOR\_NR == 3) // 硬盘主设备号是 3。
82 /* harddisk */
83 #define DEVICE\_NAME "harddisk"         // 硬盘名称 harddisk。
84 #define DEVICE\_INTR do\_hd              // 设备中断处理程序 do_hd()。
85 #define DEVICE\_REQUEST do\_hd\_request    // 设备请求函数 do_hd_request()。
86 #define DEVICE\_NR(device) (MINOR(device)/5) // 设备号 (0--1)。每个硬盘可以有 4 个分区。
87 #define DEVICE\_ON(device)              // 硬盘一直在工作，无须开启和关闭。
88 #define DEVICE\_OFF(device)
89
90 #elif
91 /* unknown blk device */ /* 未知块设备 */
92 #error "unknown blk device"
93
94 #endif
95
96 #define CURRENT (blk\_dev[MAJOR\_NR].current_request) // CURRENT 为指定主设备号的当前请求结构。
97 #define CURRENT\_DEV DEVICE\_NR(CURRENT->dev) // CURRENT_DEV 为 CURRENT 的设备号。
98
99 #ifdef DEVICE\_INTR
100 void (*DEVICE\_INTR)(void) = NULL;
101 #endif
102 static void (DEVICE\_REQUEST)(void);
103
104 // 释放锁定的缓冲区。

```

```

101 extern inline void unlock_buffer(struct buffer_head * bh)
102 {
103     if (!bh->b_lock)                // 如果指定的缓冲区 bh 没有被上锁，则显示警告信息。
104         printk(DEVICE_NAME ": free buffer being unlocked\n");
105     bh->b_lock=0;                    // 否则将该缓冲区解锁。
106     wake_up(&bh->b_wait);            // 唤醒等待该缓冲区的进程。
107 }
108
109 // 结束请求。
110 extern inline void end_request(int uptodate)
111 {
112     DEVICE_OFF(CURRENT->dev);        // 关闭设备。
113     if (CURRENT->bh) {
114         CURRENT->bh->b_uptodate = uptodate; // 置更新标志。
115         unlock_buffer(CURRENT->bh);        // 解锁缓冲区。
116     }
117     if (!uptodate) {                // 如果更新标志为 0 则显示设备错误信息。
118         printk(DEVICE_NAME " I/O error\n|r");
119         printk("dev %04x, block %d\n|r", CURRENT->dev,
120             CURRENT->bh->b_blocknr);
121     }
122     wake_up(&CURRENT->waiting);        // 唤醒等待该请求项的进程。
123     wake_up(&wait_for_request);        // 唤醒等待请求的进程。
124     CURRENT->dev = -1;                // 释放该请求项。
125     CURRENT = CURRENT->next;           // 从请求链表中删除该请求项。
126 }
127 // 定义初始化请求宏。
128 #define INIT_REQUEST \
129 repeat: \
130     if (!CURRENT) \                                // 如果当前请求结构指针为 null 则返回。
131         return; \
132     if (MAJOR(CURRENT->dev) != MAJOR_NR) \           // 如果当前设备的主设备号不对则死机。
133         panic(DEVICE_NAME ": request list destroyed"); \
134     if (CURRENT->bh) { \
135         if (!CURRENT->bh->b_lock) \                // 如果在进行请求操作时缓冲区没锁定则死机。
136             panic(DEVICE_NAME ": block not locked"); \
137     }
138 #endif
139
140 #endif
141

```

### 5.3.3 其它信息

## 5.4 hd.c 文件

### 5.4.1 功能描述

### 5.4.2 代码注释

列表 linux/kernel/blk\_drv/hd.c 程序

```

1  /*
2  *  linux/kernel/hd.c
3  *
4  *  (C) 1991  Linus Torvalds
5  */
6
7  /*
8  *  This is the low-level hd interrupt support. It traverses the
9  *  request-list, using interrupts to jump between functions. As
10 *  all the functions are called within interrupts, we may not
11 *  sleep. Special care is recommended.
12 *
13 *  modified by Drew Eckhardt to check nr of hd's from the CMOS.
14 */
15 /*
16 *  本程序是底层硬盘中断辅助程序。主要用于扫描请求列表，使用中断在函数之间跳转。
17 *  由于所有的函数都是在中断里调用的，所以这些函数不可以睡眠。请特别注意。
18 *  由 Drew Eckhardt 修改，利用 CMOS 信息检测硬盘数。
19 */
20
21 #include <linux/config.h> // 内核配置头文件。定义键盘语言和硬盘类型 (HD_TYPE) 可选项。
22 #include <linux/sched.h> // 调度程序头文件，定义了任务结构 task_struct、初始任务 0 的数据，
23 // 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
24 #include <linux/fs.h> // 文件系统头文件。定义文件表结构 (file, buffer_head, m_inode 等)。
25 #include <linux/kernel.h> // 内核头文件。含有一些内核常用函数的原形定义。
26 #include <linux/hdreg.h> // 硬盘参数头文件。定义访问硬盘寄存器端口，状态码，分区表等信息。
27 #include <asm/system.h> // 系统头文件。定义了设置或修改描述符/中断门等的嵌入式汇编宏。
28 #include <asm/io.h> // io 头文件。定义硬件端口输入/输出宏汇编语句。
29 #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。
30
31 #define MAJOR_NR 3 // 硬盘主设备号是 3。
32 #include "blk.h" // 块设备头文件。定义请求数据结构、块设备数据结构和宏函数等信息。
33
34 #define CMOS_READ(addr) ({ \ // 读 CMOS 参数宏函数。
35 outb_p(0x80|addr, 0x70); \
36 inb_p(0x71); \
37 })
38
39 /* Max read/write errors/sector */
40 #define MAX_ERRORS 7 // 读/写一个扇区时最大出错次数。
41 #define MAX_HD 2 // 最多硬盘数。
42
43 static void recal_intr(void); // 硬盘中断程序在复位操作时会调用的重新校正函数 (287 行)。

```



```

38
39 static int recalibrate = 1;    // 重新校正。
40 static int reset = 1;        // 复位。
41
42 /*
43  * This struct defines the HD's and their types.
44  */
45 /* 下面结构定义了硬盘参数及类型 */
46 // 各字段分别是磁头数、每磁道扇区数、柱面数、写前预补偿柱面号、磁头着陆区柱面号、控制字节。
47 struct hd_i_struct {
48     int head, sect, cyl, wpcom, lzone, ctl;
49 };
50 #ifdef HD_TYPE    // 如果已经在 include/linux/config.h 中定义了 HD_TYPE...
51 struct hd_i_struct hd_info[] = { HD_TYPE };    // 取定义好的参数作为 hd_info[] 的数据。
52 #define NR_HD ((sizeof (hd_info))/(sizeof (struct hd_i_struct)))    // 计算硬盘数。
53 #else
54 // 否则，都设为 0 值。
55 struct hd_i_struct hd_info[] = { {0, 0, 0, 0, 0, 0}, {0, 0, 0, 0, 0, 0} };
56 static int NR_HD = 0;
57 #endif
58
59 // 定义硬盘分区结构。给出每个分区的物理起始扇区号、分区扇区总数。
60 // 其中 5 的倍数处的项（例如 hd[0] 和 hd[5] 等）代表整个硬盘中的参数。
61 static struct hd_struct {
62     long start_sect;
63     long nr_sects;
64 } hd[5*MAX_HD] = {{0, 0}, };
65
66 // 读端口 port，共读 nr 字，保存在 buf 中。
67 #define port_read(port, buf, nr) \
68 __asm__ ("cld;rep;insw"::"d" (port), "D" (buf), "c" (nr): "cx", "di")
69
70 // 写端口 port，共写 nr 字，从 buf 中取数据。
71 #define port_write(port, buf, nr) \
72 __asm__ ("cld;rep;outsw"::"d" (port), "S" (buf), "c" (nr): "cx", "si")
73
74 extern void hd_interrupt(void);
75 extern void rd_load(void);
76
77 /* This may be used only once, enforced by 'static int callable' */
78 /* 下面该函数只在初始化时被调用一次。用静态变量 callable 作为可调用标志。*/
79 // 该函数的参数由初始化程序 init/main.c 的 init 子程序设置为指向 0x90080 处，此处存放着 setup.s
80 // 程序从 BIOS 取得的 2 个硬盘的基本参数表(32 字节)。硬盘参数表信息参见下面列表后的说明。
81 int sys_setup(void * BIOS)
82 {
83     static int callable = 1;
84     int i, drive;
85     unsigned char cmos_disks;
86     struct partition *p;
87     struct buffer_head *bh;
88
89     // 初始化时 callable=1，当运行该函数时将其设置为 0，使本函数只能执行一次。
90     if (!callable)
91         return -1;

```



```

81     callable = 0;
    // 如果没有在 config.h 中定义硬盘参数, 就从 0x90080 处读入。
82 #ifndef HD_TYPE
83     for (drive=0 ; drive<2 ; drive++) {
84         hd_info[drive].cyl = *(unsigned short *) BIOS;
85         hd_info[drive].head = *(unsigned char *) (2+BIOS);
86         hd_info[drive].wpcom = *(unsigned short *) (5+BIOS);
87         hd_info[drive].ctl = *(unsigned char *) (8+BIOS);
88         hd_info[drive].lzone = *(unsigned short *) (12+BIOS);
89         hd_info[drive].sect = *(unsigned char *) (14+BIOS);
90         BIOS += 16;          // 每个硬盘的参数表长 16 字节。
91     }
    // setup.s 程序在取 BIOS 中的硬盘参数表信息时, 如果只有 1 个硬盘, 就会将对应第 2 个硬盘的
    // 16 字节全部清零。因此这里只要判断第 2 个硬盘柱面数是否为 0 就可以知道有没有第 2 个硬盘了。
92     if (hd_info[1].cyl)
93         NR_HD=2;          // 硬盘数置为 2。
94     else
95         NR_HD=1;
96 #endif
    // 设置每个硬盘的起始扇区号和扇区总数。其中编号 i*5 含义参见本程序后的有关说明。
97     for (i=0 ; i<NR_HD ; i++) {
98         hd[i*5].start_sect = 0;          // 硬盘起始扇区号。
99         hd[i*5].nr_sects = hd_info[i].head*
100             hd_info[i].sect*hd_info[i].cyl; // 硬盘总扇区数。
101     }
102
103     /*
104         We query CMOS about hard disks : it could be that
105         we have a SCSI/ESDI/etc controller that is BIOS
106         compatable with ST-506, and thus showing up in our
107         BIOS table, but not register compatable, and therefore
108         not present in CMOS.
109
110         Furthurmore, we will assume that our ST-506 drives
111         <if any> are the primary drives in the system, and
112         the ones reflected as drive 1 or 2.
113
114         The first drive is stored in the high nibble of CMOS
115         byte 0x12, the second in the low nibble. This will be
116         either a 4 bit drive type or 0xf indicating use byte 0x19
117         for an 8 bit type, drive 1, 0x1a for drive 2 in CMOS.
118
119         Needless to say, a non-zero value means we have
120         an AT controller hard disk for that drive.
121
122         */
123     /*
    * 我们对 CMOS 有关硬盘的信息有些怀疑: 可能会出现这样的情况, 我们有一块 SCSI/ESDI/等的
    * 控制器, 它是以 ST-506 方式与 BIOS 兼容的, 因而会出现在我们的 BIOS 参数表中, 但却又不
    * 是寄存器兼容的, 因此这些参数在 CMOS 中又不存在。
    * 另外, 我们假设 ST-506 驱动器 (如果有的话) 是系统中的基本驱动器, 也即以驱动器 1 或 2
    * 出现的驱动器。

```

\* 第1个驱动器参数存放在CMOS字节0x12的高半字节中，第2个存放在低半字节中。该4位字节  
 \* 信息可以是驱动器类型，也可能仅是0xf。0xf表示使用CMOS中0x19字节作为驱动器1的8位  
 \* 类型字节，使用CMOS中0x1A字节作为驱动器2的类型字节。  
 \* 总之，一个非零值意味着我们有一个AT控制器硬盘兼容的驱动器。  
 \*/

```

124 // 这里根据上述原理来检测硬盘到底是否是AT控制器兼容的。
125 if ((cmos_disks = CMOS_READ(0x12)) & 0xf0)
126     if (cmos_disks & 0x0f)
127         NR_HD = 2;
128     else
129         NR_HD = 1;
130 else
131     NR_HD = 0;
132 // 若NR_HD=0，则两个硬盘都不是AT控制器兼容的，硬盘数据结构清零。
133 // 若NR_HD=1，则将第2个硬盘的参数清零。
134 for (i = NR_HD ; i < 2 ; i++) {
135     hd[i*5].start_sect = 0;
136     hd[i*5].nr_sects = 0;
137 }
138 // 读取每一个硬盘上第1块数据（第1个扇区有用），获取其中的分区表信息。
139 // 首先利用函数bread()读硬盘第1块数据(fs/buffer.c, 267)，参数中的0x300是硬盘的主设备号
140 // (参见列表后的说明)。然后根据硬盘头1个扇区位置0x1fe处的两个字节是否为'55AA'来判断
141 // 该扇区中位于0x1BE开始的分区表是否有效。最后将分区表信息放入硬盘分区数据结构hd中。
142 for (drive=0 ; drive<NR_HD ; drive++) {
143     if (!(bh = bread(0x300 + drive*5, 0))) { // 0x300, 0x305 逻辑设备号。
144         printk("Unable to read partition table of drive %d\n\r",
145             drive);
146         panic("");
147     }
148     if (bh->b_data[510] != 0x55 || (unsigned char)
149         bh->b_data[511] != 0xAA) { // 判断硬盘信息有效标志'55AA'。
150         printk("Bad partition table on drive %d\n\r", drive);
151         panic("");
152     }
153     p = 0x1BE + (void *)bh->b_data; // 分区表位于硬盘第1扇区的0x1BE处。
154     for (i=1; i<5; i++, p++) {
155         hd[i+5*drive].start_sect = p->start_sect;
156         hd[i+5*drive].nr_sects = p->nr_sects;
157     }
158     brelse(bh); // 释放为存放硬盘块而申请的内存缓冲区页。
159 }
160 if (NR_HD) // 如果有硬盘存在并且已读入分区表，则打印分区表正常信息。
161     printk("Partition table%s ok. \n\r", (NR_HD>1)? "s": "");
162 rd_load(); // 加载（创建）RAMDISK(kernel/blk_drv/ramdisk.c, 71)。
163 mount_root(); // 安装根文件系统(fs/super.c, 242)。
164 return (0);
165 }
166 // 判断并循环等待驱动器就绪。
167 // 读硬盘控制器状态寄存器端口HD_STATUS(0x1f7)，并循环检测驱动器就绪比特位和控制器忙位。
168 static int controller_ready(void)
169 {

```

```

163     int retries=10000;
164
165     while (--retries && (inb\_p(HD\_STATUS)&0xc0)!=0x40);
166     return (retries);        // 返回等待循环的次数。
167 }
168
169     // 检测硬盘执行命令后的状态。(win_表示温切斯特硬盘的缩写)
170     // 读取状态寄存器中的命令执行结果状态。返回 0 表示正常, 1 出错。如果执行命令错,
171     // 则再读错误寄存器 HD_ERROR(0x1f1)。
172 static int win\_result(void)
173 {
174     int i=inb\_p(HD\_STATUS); // 取状态信息。
175
176     if ((i & (BUSY\_STAT | READY\_STAT | WRERR\_STAT | SEEK\_STAT | ERR\_STAT))
177         == (READY\_STAT | SEEK\_STAT))
178         return(0); /* ok */
179     if (i&1) i=inb(HD\_ERROR); // 若 ERR_STAT 置位, 则读取错误寄存器。
180     return (1);
181 }
182
183     // 向硬盘控制器发送命令块 (参见列表后的说明)。
184     // 调用参数: drive - 硬盘号 (0-1);          nsect - 读写扇区数;
185     //          sect - 起始扇区;                head - 磁头号;
186     //          cyl - 柱面号;                  cmd - 命令码;
187     //          *intr_addr() - 硬盘中断处理程序中将调用的处理函数。
188 static void hd\_out(unsigned int drive,unsigned int nsect,unsigned int sect,
189                 unsigned int head,unsigned int cyl,unsigned int cmd,
190                 void (*intr_addr)(void))
191 {
192     register int port asm("dx"); // port 变量对应寄存器 dx。
193
194     if (drive>1 || head>15) // 如果驱动器号 (0,1)>1 或磁头号>15, 则程序不支持。
195         panic("Trying to write bad sector");
196     if (!controller\_ready()) // 如果等待一段时间后仍未就绪则出错, 死机。
197         panic("HD controller not ready");
198     do_hd = intr_addr; // do_hd 函数指针将在硬盘中断程序中被调用。
199     outb\_p(hd\_info[drive].ctl,HD\_CMD); // 向控制寄存器(0x3f6)输出控制字节。
200     port=HD\_DATA; // 置 dx 为数据寄存器端口(0x1f0)。
201     outb\_p(hd\_info[drive].wpcom>>2, ++port); // 参数: 写预补偿柱面号(需除 4)。
202     outb\_p(nsect, ++port); // 参数: 读/写扇区总数。
203     outb\_p(sect, ++port); // 参数: 起始扇区。
204     outb\_p(cyl, ++port); // 参数: 柱面号低 8 位。
205     outb\_p(cyl>>8, ++port); // 参数: 柱面号高 8 位。
206     outb\_p(0xA0 | (drive<<4) | head, ++port); // 参数: 驱动器号+磁头号。
207     outb(cmd, ++port); // 命令: 硬盘控制命令。
208 }
209
210     // 等待硬盘就绪。也即循环等待主状态控制器忙标志位复位。若仅有就绪或寻道结束标志
211     // 置位, 则成功, 返回 0。若经过一段时间仍为忙, 则返回 1。
212 static int drive\_busy(void)
213 {
214     unsigned int i;
215
216     while (i < 1000000)
217     {
218         if (win\_result() == 0)
219             return 0;
220         if (i % 100000 == 0)
221             outb\_p(0, HD\_DATA); // 向数据寄存器端口(0x1f0)写 0。
222         i++;
223     }
224     return 1;
225 }

```

```

206         for (i = 0; i < 10000; i++)                // 循环等待就绪标志位置位。
207             if (READY_STAT == (inb_p(HD_STATUS) & (BUSY_STAT | READY_STAT)))
208                 break;
209         i = inb(HD_STATUS);                        // 再取主控制器状态字节。
210         i &= BUSY_STAT | READY_STAT | SEEK_STAT;    // 检测忙位、就绪位和寻道结束位。
211         if (i == READY_STAT | SEEK_STAT)            // 若仅有就绪或寻道结束标志，则返回 0。
212             return(0);
213         printk("HD controller times out\n\r");      // 否则等待超时，显示信息。并返回 1。
214         return(1);
215     }
216
217     ///// 诊断复位（重新校正）硬盘控制器。
218     static void reset_controller(void)
219     {
220         int i;
221         outb(4, HD_CMD);                            // 向控制寄存器端口发送控制字节(4-复位)。
222         for(i = 0; i < 100; i++) nop();              // 等待一段时间（循环空操作）。
223         outb(hd_info[0].ctl & 0x0f, HD_CMD);        // 再发送正常的控制字节(不禁止重试、重读)。
224         if (drive_busy())                            // 若等待硬盘就绪超时，则显示出错信息。
225             printk("HD-controller still busy\n\r");
226         if ((i = inb(HD_ERROR)) != 1)              // 取错误寄存器，若不等于 1（无错误）则出错。
227             printk("HD-controller reset failed: %02x\n\r", i);
228     }
229
230     ///// 复位硬盘 nr。首先复位（重新校正）硬盘控制器。然后发送硬盘控制器命令“建立驱动器参数”，
231     // 其中 recal_intr() 是在硬盘中断处理程序中调用的重新校正处理函数。
232     static void reset_hd(int nr)
233     {
234         reset_controller();
235         hd_out(nr, hd_info[nr].sect, hd_info[nr].sect, hd_info[nr].head-1,
236             hd_info[nr].cyl, WIN_SPECIFY, &recal_intr);
237     }
238
239     ///// 意外硬盘中断调用函数。
240     // 发生意外硬盘中断时，硬盘中断处理程序中调用的默认 C 处理函数。在被调用函数指针为空时
241     // 调用该函数。参见(kernel/system_call.s, 241 行)。
242     void unexpected_hd_interrupt(void)
243     {
244         printk("Unexpected HD interrupt\n\r");
245     }
246
247     ///// 读写硬盘失败处理调用函数。
248     static void bad_rw_intr(void)
249     {
250         if (++CURRENT->errors >= MAX_ERRORS)        // 如果读扇区时的出错次数大于或等于 7 次时，
251             end_request(0);                          // 则结束请求并唤醒等待该请求的进程，而且
252                                                     // 对应缓冲区更新标志复位（没有更新）。
253         if (CURRENT->errors > MAX_ERRORS/2)        // 如果读一扇区时的出错次数已经大于 3 次，
254             reset = 1;                                // 则要求执行复位硬盘控制器操作。
255     }
256
257     ///// 读操作中中断调用函数。将在执行硬盘中断处理程序中被调用。

```

```

250 static void read_intr(void)
251 {
252     if (win_result()) {                // 若控制器忙、读写错或命令执行错，
253         bad_rw_intr();                // 则进行读写硬盘失败处理
254         do_hd_request();                // 然后再次请求硬盘作相应(复位)处理。
255         return;
256     }
257     port_read(HD_DATA, CURRENT->buffer, 256); // 将数据从数据寄存器端口读到请求缓冲区。
258     CURRENT->errors = 0;                // 清出错次数。
259     CURRENT->buffer += 512;             // 调整缓冲区指针，指向新的空区。
260     CURRENT->sector++;                  // 起始扇区号加1，
261     if (--CURRENT->nr_sectors) {        // 如果所需读出的扇区数还没有读完，则
262         do_hd = &read_intr;           // 再次置硬盘调用C函数指针为 read_intr()
263         return;                        // 因为硬盘中断处理程序每次调用 do_hd 时
264     }                                  // 都会将该函数指针置空。参见 system_call.s
265     end_request(1);                    // 若全部扇区数据已经读完，则处理请求结束事宜，
266     do_hd_request();                  // 执行其它硬盘请求操作。
267 }
268
269 // 写扇区中断调用函数。在硬盘中断处理程序中调用。
270 // 在写命令执行后，会产生硬盘中断信号，执行硬盘中断处理程序，此时在硬盘中断处理程序中调用的
271 // C函数指针 do_hd()已经指向 write_intr()，因此会在写操作完成（或出错）后，执行该函数。
272 static void write_intr(void)
273 {
274     if (win_result()) {                // 如果硬盘控制器返回错误信息，
275         bad_rw_intr();                // 则首先进行硬盘读写失败处理，
276         do_hd_request();                // 然后再次请求硬盘作相应(复位)处理，
277         return;                        // 然后返回（也退出了此次硬盘中断）。
278     }
279     if (--CURRENT->nr_sectors) {        // 否则将欲写扇区数减1，若还有扇区要写，则
280         CURRENT->sector++;              // 当前请求起始扇区号+1，
281         CURRENT->buffer += 512;         // 调整请求缓冲区指针，
282         do_hd = &write_intr;           // 置硬盘中断程序调用函数指针为 write_intr()，
283         port_write(HD_DATA, CURRENT->buffer, 256); // 再向数据寄存器端口写 256 字节。
284         return;                        // 返回等待硬盘再次完成写操作后的中断处理。
285     }
286     end_request(1);                    // 若全部扇区数据已经写完，则处理请求结束事宜，
287     do_hd_request();                  // 执行其它硬盘请求操作。
288 }
289
290 // 硬盘重新校正（复位）中断调用函数。在硬盘中断处理程序中调用。
291 // 如果硬盘控制器返回错误信息，则首先进行硬盘读写失败处理，然后请求硬盘作相应(复位)处理。
292 static void recal_intr(void)
293 {
294     if (win_result())
295         bad_rw_intr();
296     do_hd_request();
297 }
298
299 // 执行硬盘读写请求操作。
300 void do_hd_request(void)
301 {
302     int i, r;

```

```

297     unsigned int block, dev;
298     unsigned int sec, head, cyl;
299     unsigned int nsect;
300
301     INIT\_REQUEST; // 检测请求项的合法性(参见 kernel/blk_drv/blk.h, 127)。
302     dev = MINOR(CURRENT->dev); // 取设备号中的子设备号(见列表后对硬盘设备号的说明)。
                                    // CURRENT 定义为(blk_dev[MAJOR_NR].current_request)。
303     block = CURRENT->sector; // 请求的起始扇区。
// 如果子设备号不存在或者起始扇区大于该分区扇区数-2, 则结束该请求, 并跳转到标号 repeat 处
// (定义在 INIT\_REQUEST 开始处)。
304     if (dev >= 5*NR\_HD || block+2 > hd[dev].nr_sects) {
305         end\_request(0);
306         goto repeat;
307     }
308     block += hd[dev].start_sect; // 将所需读的块对应到硬盘绝对扇区号。
309     dev /= 5; // 此时 dev 代表硬盘号(0 或 1)。
// 下面嵌入汇编代码用来从硬盘信息结构中根据起始扇区号和每磁道扇区数计算在磁道中的
// 扇区号(sec)、所在柱面号(cyl)和磁头号(head)
310     __asm__ ("divl %4": "=a" (block), "=d" (sec): "" (block), "1" (0),
311             "r" (hd\_info[dev].sect));
312     __asm__ ("divl %4": "=a" (cyl), "=d" (head): "" (block), "1" (0),
313             "r" (hd\_info[dev].head));
314     sec++;
315     nsect = CURRENT->nr_sectors; // 欲读/写的扇区数。
// 如果 reset 置 1, 则执行复位操作。复位硬盘和控制器, 并置需要重新校正标志, 返回。
316     if (reset) {
317         reset = 0;
318         recalibrate = 1;
319         reset\_hd(CURRENT\_DEV);
320         return;
321     }
// 如果重新校正标志(recalibrate)置位, 则首先复位该标志, 然后向硬盘控制器发送重新校正命令。
322     if (recalibrate) {
323         recalibrate = 0;
324         hd\_out(dev, hd\_info[CURRENT\_DEV].sect, 0, 0, 0,
325             WIN\_RESTORE, &recal\_intr);
326         return;
327     }
// 如果当前请求是写扇区操作, 则发送写命令, 循环读取状态寄存器信息并判断请求服务标志
// DRQ_STAT 是否置位。
328     if (CURRENT->cmd == WRITE) {
329         hd\_out(dev, nsect, sec, head, cyl, WIN\_WRITE, &write\_intr);
330         for(i=0 ; i<3000 && !(r=inb\_p(HD\_STATUS)&DRQ\_STAT) ; i++)
331             /* nothing */;
// 如置位则退出循环。若等到循环结束 DRQ_STAT 也没有置位, 则此次写硬盘操作失败, 去处理下一个
// 硬盘请求。否则向硬盘控制器数据寄存器端口 HD_DATA 写入 1 个扇区的数据。
332         if (!r) {
333             bad\_rw\_intr();
334             goto repeat;
335         }
336         port\_write(HD\_DATA, CURRENT->buffer, 256);
// 如果当前请求是读硬盘扇区, 则向硬盘控制器发送读扇区命令。
337     } else if (CURRENT->cmd == READ) {

```



```

338         hd\_out(dev, nsect, sec, head, cyl, WIN_READ, &read\_intr);
339     } else
340         panic("unknown hd-command");
341 }
342 // 硬盘系统初始化。
343 void hd\_init(void)
344 {
345     blk\_dev[MAJOR_NR].request_fn = DEVICE\_REQUEST; // do_hd_request()。
346     set\_intr\_gate(0x2E, &hd\_interrupt); // 设置硬盘中断门向量 int 0x2E(46)。
                                     // hd\_interrupt 在(kernel/system_call.s, 221)。
347     outb\_p(inb\_p(0x21)&0xfb, 0x21); // 复位接联的主 8259A int2 的屏蔽位, 允许从片
                                     // 发出中断请求信号。
348     outb(inb\_p(0xA1)&0xbf, 0xA1); // 复位硬盘的中断请求屏蔽位(在从片上), 允许
                                     // 硬盘控制器发送中断请求信号。
349 }
350

```

## 5.4.3 其它信息

### 5.4.3.1 AT 硬盘接口寄存器

AT 硬盘控制器的编程寄存器端口见下表所示。

表 AT 硬盘控制器寄存器端口及作用

| I/O 端口 | 读操作  | 写操作      |
|--------|--|----------|
| 0x1f0  | 数据寄存器 -- 扇区数据(读、写、格式化)                         |          |
| 0x1f1  | 错误寄存器(错误状态)                                    | 写前预补偿寄存器 |
| 0x1f2  | 扇区数寄存器 -- 扇区数(读、写、检验、格式化)                      |          |
| 0x1f3  | 扇区号寄存器 -- 起始扇区(读、写、检验)                         |          |
| 0x1f4  | 柱面号寄存器 -- 柱面号低字节(读、写、检验、格式化)                   |          |
| 0x1f5  | 柱面号寄存器 -- 柱面号高字节(读、写、检验、格式化)                   |          |
| 0x1f6  | 驱动器/磁头寄存器 -- 驱动器号/磁头号(101dhhhh, d=驱动器号, h=磁头号) |          |
| 0x1f7  | 主状态寄存器   | 命令寄存器    |
| 0x3f6  | ---  | 硬盘控制寄存器  |
| 0x3f7  | 数字输入寄存器(与 1.2M 软盘合用)                           | ---      |

下面对各端口寄存器进行详细说明。

#### 1. 数据寄存器(HD\_DATA, 0x1f0)

这是一对 16 位高速 PIO 数据传输器, 用于扇区读、写和磁道格式化操作。CPU 通过该数据寄存器向硬盘写入或从硬盘读出 1 个扇区的数据, 也即要使用命令'rep outsw'或'rep insw'重复读/写 cx=256 字。

#### 2. 错误寄存器(读)/写前预补偿寄存器(写)(HD\_ERROR, 0x1f1)

在读时, 该寄存器存放有 8 位的错误状态。但只有当主状态寄存器(HD\_STATUS, 0x1f7)的位 0=1 时该寄存器中的数据才有效。执行控制器诊断命令时的含义与其它命令时的不同。下面分别列出:

表 硬盘控制器错误寄存器

| 值    | 诊断命令时 | 其它命令时  |
|------|-------|--------|
| 0x01 | 无错误   | 数据标志丢失 |

|      |         |        |
|------|---------|--------|
| 0x02 | 控制器出错   | 磁道 0 错 |
| 0x03 | 扇区缓冲区错  |        |
| 0x04 | ECC 部件错 | 命令放弃   |
| 0x05 | 控制处理器错  |        |
| 0x10 |         | ID 未找到 |
| 0x40 |         | ECC 错误 |
| 0x80 |         | 坏扇区    |

在写操作时，该寄存器即作为写前预补偿寄存器。它记录写预补偿起始柱面号。对应与硬盘基本参数表位移 0x05 处的一个字，需除 4 后输出。

### 3. 扇区数寄存器 (HD\_NSECTOR, 0x1f2)

该寄存器存放读、写、检验和格式化命令指定的扇区数。当用于多扇区操作时，每完成 1 扇区的操作该寄存器就自动减 1，直到为 0。若初值为 0，则表示传输最大扇区数 256。

### 4. 扇区号寄存器 (HD\_SECTOR, 0x1f3)

该寄存器存放读、写、检验操作命令指定的扇区号。在多扇区操作时，保存的是起始扇区号，而每完成 1 扇区的操作就自动增 1。

### 5. 柱面号寄存器 (HD\_LCYL, HD\_HCYL, 0x1f4, 0x1f5)

该两个柱面号寄存器分别存放有柱面号的低 8 位和高 2 位。

### 6. 驱动器/磁头寄存器 (HD\_CURRENT, 0x1f6)

该寄存器存放有读、写、检验、寻道和格式化命令指定的驱动器和磁头号。其位格式为 101dhhhh。其中 101 表示采用 ECC 校验码和每扇区为 512 字节；d 表示选择的驱动器 (0 或 1)；hhhh 表示选择的磁头。

### 7. 主状态寄存器 (读) / 命令寄存器 (写) (HD\_STATUS/HD\_COMMAND, 0x1f7)

在读时，对应一个 8 位主状态寄存器。反映硬盘控制器在执行命令前后的操作状态。各位的含义如下：

|            |      |                 |
|------------|------|-----------------|
| ERR_STAT   | 0x01 | // 命令执行错误。      |
| INDEX_STAT | 0x02 | // 收到索引。        |
| ECC_STAT   | 0x04 | // ECC 校验错。     |
| DRQ_STAT   | 0x08 | // 请求服务。        |
| SEEK_STAT  | 0x10 | // 寻道结束。        |
| WRERR_STAT | 0x20 | // 驱动器故障。       |
| READY_STAT | 0x40 | // 驱动器准备好 (就绪)。 |
| BUSY_STAT  | 0x80 | // 控制器忙碌。       |

当写时，该端口对应命令寄存器，接受 CPU 发出的硬盘控制命令，共有 8 种命令，见下表所示。

表 AT 硬盘控制器命令列表

| 命令名称         | 命令码字节 |    |    |       | 命令执行结束形式 |
|--------------|-------|----|----|-------|----------|
|              | 高 4 位 | D3 | D2 | D1 D0 |          |
| 驱动器重新校正 (复位) | 0x1   | R  | R  | R R   | 中断       |
| 读扇区          | 0x2   | 0  | 0  | L T   | 中断       |
| 写扇区          | 0x3   | 0  | 0  | L T   | 中断       |
| 扇区检验         | 0x4   | 0  | 0  | 0 T   | 中断       |



|         |     |         |       |
|---------|-----|---------|-------|
| 格式化磁道   | 0x5 | 0 0 0 0 | 中断    |
| 控制器初始化  | 0x6 | 0 0 0 0 | 中断    |
| 寻道      | 0x7 | R R R R | 中断    |
| 控制器诊断   | 0x9 | 0 0 0 0 | 控制器空闲 |
| 建立驱动器参数 | 0x9 | 0 0 0 1 | 控制器空闲 |

表中命令码字节的低 4 位是附加参数，其含义为：

R 是步进速率。R=0，则步进速率为 35us；R=1 为 0.5ms，以此量递增。

L 是数据模式。L=0 表示读/写扇区为 512 字节；L=1 表示读/写扇区为 512 加 4 字节的 ECC 码。

T 是重试模式。T=0 表示允许重试；T=1 则禁止重试。

#### 8. 硬盘控制寄存器（写）（HD\_CMD, 0x3f6）

该寄存器是只写的。用于存放硬盘控制字节并控制复位操作。其定义参见下面硬盘基本参数表的位移 0x08 处的字节说明。

### 5.4.3.2 AT 硬盘控制器编程

在对硬盘控制器进行操作控制时，需要同时发送参数和命令。其命令格式如下表所示。首先发送 6 字节的参数，最后发出 1 字节的命令码。不管什么命令均需要完整输出这 7 字节的命令块，依次写入端口 0x1f1 -- 0x1f7。。

|       |           |
|-------|-----------|
| 0x1f1 | 写预补偿起始柱面号 |
| 0x1f2 | 扇区数       |
| 0x1f3 | 起始扇区号     |
| 0x1f4 | 柱面号低字节    |
| 0x1f5 | 柱面号高字节    |
| 0x1f6 | 驱动器号/磁头号  |
| 0x1f7 | 命令码       |

首先 CPU 向控制寄存器端口 (HD\_CMD)0x3f6 输出控制字节，建立相应的硬盘控制方式。方式建立后即可按上面顺序发送参数和命令。步骤为：

1. 检测控制器空闲状态：CPU 通过读主状态寄存器，若位 7 为 0，表示控制器空闲。若在规定时间内控制器一直处于忙状态，则判为超时出错。
2. 检测驱动器就绪：CPU 判断主状态寄存器位 6 是否为 1 来看驱动器是否就绪。为 1 则可输出参数和命令。
3. 输出命令块：按顺序输出分别向对应端口输出参数和命令。
4. CPU 等待中断产生：命令执行后，由硬盘控制器产生中断请求信号（IRQ14 -- 对应中断 int46）或置控制器状态为空闲，表明操作结束或表示请求扇区传输（多扇区读/写）。
5. 检测操作结果：CPU 再次读主状态寄存器，若位 0 等于 0 则表示命令执行成功，否则失败。若失败则可进一步查询错误寄存器 (HD\_ERROR) 取错误码。

### 5.4.3.3 硬盘基本参数表

中断向量表中，int 0x41 的中断向量位置（4 \* 0x41 = 0x0000:0x0104）存放的并不是中断程序的地址而是第一个硬盘的基本参数表。对于 100%兼容的 BIOS 来说，这里存放着硬盘参数表阵列的首地址 F000h:E401h。第二个硬盘的基本参数表入口地址存于 int 0x46 中断向量中。

表 硬盘基本参数信息表

| 位移   | 大小 | 说明  |
|------|----|---|
| 0x00 | 字  | 柱面数   |
| 0x02 | 字节 | 磁头数   |
| 0x03 | 字  | 开始减小写电流的柱面(仅 PC XT 使用, 其它为 0)   |
| 0x05 | 字  | 开始写前预补偿柱面号(乘 4)   |
| 0x07 | 字节 | 最大 ECC 猝发长度(仅 XT 使用, 其它为 0)   |
| 0x08 | 字节 | 控制字节(驱动器步进选择)<br>位 0 未用<br>位 1 保留(0)(关闭 IRQ)<br>位 2 允许复位<br>位 3 若磁头数大于 8 则置 1<br>位 4 未用(0)<br>位 5 若在柱面数+1 处有生产商的坏区图, 则置 1<br>位 6 禁止 ECC 重试<br>位 7 禁止访问重试。 |
| 0x09 | 字节 | 标准超时值(仅 XT 使用, 其它为 0)   |
| 0x0A | 字节 | 格式化超时值(仅 XT 使用, 其它为 0)  |
| 0x0B | 字节 | 检测驱动器超时值(仅 XT 使用, 其它为 0)  |
| 0x0C | 字  | 磁头着陆(停止)柱面号   |
| 0x0E | 字节 | 每磁道扇区数  |
| 0x0F | 字节 | 保留。   |

#### 5.4.3.4 硬盘设备号命名方式

硬盘的主设备号是 3。其它设备的主设备号分别为:

1-内存, 2-磁盘, 3-硬盘, 4-ttyx, 5-tty, 6-并行口, 7-非命名管道

由于 1 个硬盘中可以存在 1—4 个分区, 因此硬盘还依据分区不同用次设备号进行指定分区。因此硬盘的逻辑设备号由以下方式构成:

设备号=主设备号\*256 + 次设备号

也即  $\text{dev\_no} = (\text{major} \ll 8) + \text{minor}$

两个硬盘的所有逻辑设备号见下表所示。

表 硬盘逻辑设备号

| 逻辑设备号 | 对应设备文件   | 说明                |
|-------|----------|-------------------|
| 0x300 | /dev/hd0 | 代表整个第 1 个硬盘       |
| 0x301 | /dev/hd1 | 表示第 1 个硬盘的第 1 个分区 |
| 0x302 | /dev/hd2 | 表示第 1 个硬盘的第 2 个分区 |
| 0x303 | /dev/hd3 | 表示第 1 个硬盘的第 3 个分区 |
| 0x304 | /dev/hd4 | 表示第 1 个硬盘的第 4 个分区 |
| 0x305 | /dev/hd5 | 代表整个第 2 个硬盘       |
| 0x306 | /dev/hd6 | 表示第 2 个硬盘的第 1 个分区 |
| 0x307 | /dev/hd7 | 表示第 2 个硬盘的第 2 个分区 |
| 0x308 | /dev/hd8 | 表示第 2 个硬盘的第 3 个分区 |
| 0x309 | /dev/hd9 | 表示第 2 个硬盘的第 4 个分区 |

其中 0x300 和 0x305 并不与哪个分区对应，而是代表整个硬盘。

从 linux 内核 0.95 版后已经不使用这种烦琐的命名方式，而是使用与现在相同的命名方法了。

### 5.4.3.5 硬盘分区表

为了实现多个操作系统共享硬盘资源，硬盘可以在逻辑上分为 1—4 个分区。每个分区之间的扇区号是邻接的。分区表由 4 个表项组成，每个表项由 16 字节组成，对应一个分区的信息，存放有分区的大小和起止的柱面号、磁道号和扇区号，见下表所示。分区表存放在硬盘的 0 柱面 0 头第 1 个扇区的 0x1BE—0x1FD 处。

表 硬盘分区表结构

| 位置        | 名称         | 大小 | 说明  |
|-----------|------------|----|---|
| 0x00      | boot_ind   | 字节 | 引导标志。4 个分区中同时只能有一个分区是可引导的。0x00—不从该分区引导操作系统；0x80—从该分区引导操作系统。 |
| 0x01      | head       | 字节 | 分区起始磁头号。  |
| 0x02      | sector     | 字节 | 分区起始扇区号(位 0–5)和起始柱面号高 2 位(位 6–7)。                           |
| 0x03      | cyl        | 字节 | 分区起始柱面号低 8 位。   |
| 0x04      | sys_ind    | 字节 | 分区类型字节。0x0b—DOS；0x80—Old Minix；0x83—Linux …                 |
| 0x05      | end_head   | 字节 | 分区的结束磁头号。   |
| 0x06      | end_sector | 字节 | 结束扇区号(位 0–5)和结束柱面号高 2 位(位 6–7)。                             |
| 0x07      | end_cyl    | 字节 | 结束柱面号低 8 位。   |
| 0x08—0x0b | start_sect | 长字 | 分区起始物理扇区号。  |
| 0x0c—0x0f | nr_sects   | 长字 | 分区占用的扇区数。   |

## 5.5 ll\_rw\_blk.c 文件

### 5.5.1 功能描述

### 5.5.2 代码注释

列表 linux/kernel/blk\_drv/ll\_rw\_blk.c 程序

```

1 /*
2  * linux/kernel/blk_dev/ll_rw.c
3  *
4  * (C) 1991 Linus Torvalds
5  */
6
7 /*
8  * This handles all read/write requests to block devices
9  */
10
11 /*
12  * 该程序处理块设备的所有读/写操作。
13  */

```

```

10 #include <errno.h>           // 错误号头文件。包含系统中各种出错号。(Linux 从 minix 中引进的)
11 #include <linux/sched.h>     // 调度程序头文件，定义了任务结构 task_struct、初始任务 0 的数据，
                               // 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
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/system.h>      // 系统头文件。定义了设置或修改描述符/中断门等的嵌入式汇编宏。
14
15 #include "blk.h"             // 块设备头文件。定义请求数据结构、块设备数据结构和宏函数等信息。
16
17 /*
18  * The request-struct contains all necessary data
19  * to load a nr of sectors into memory
20  */
21 /*
22  * 请求结构中含有加载 nr 扇区数据到内存的所有必须的信息。
23  */
24 struct request request[NR_REQUEST];
25
26 /*
27  * used to wait on when there are no free requests
28  */
29 /* 是用于请求数组没有空闲项时的临时等待处 */
30 struct task_struct * wait_for_request = NULL;
31
32 /* blk_dev_struct is:
33  *      do_request-address
34  *      next-request
35  */
36 /* blk_dev_struct 块设备结构是：(kernel/blk_drv/blk.h, 23)
37  *      do_request-address    // 对应主设备号的请求处理程序指针。
38  *      current-request      // 该设备的下一个请求。
39  */
40 // 该数组使用主设备号作为索引（下标）。
41 struct blk_dev_struct blk_dev[NR_BLK_DEV] = {
42     { NULL, NULL },          /* no_dev */    // 0 - 无设备。
43     { NULL, NULL },          /* dev mem */   // 1 - 内存。
44     { NULL, NULL },          /* dev fd */    // 2 - 软驱设备。
45     { NULL, NULL },          /* dev hd */    // 3 - 硬盘设备。
46     { NULL, NULL },          /* dev ttyx */  // 4 - ttyx 设备。
47     { NULL, NULL },          /* dev tty */   // 5 - tty 设备。
48     { NULL, NULL }           /* dev lp */    // 6 - lp 打印机设备。
49 };
50
51 // 锁定指定的缓冲区 bh。如果指定的缓冲区已经被其它任务锁定，则使自己睡眠（不可中断地等待），
52 // 直到被执行解锁缓冲区的任务明确地唤醒。
53 static inline void lock_buffer(struct buffer_head * bh)
54 {
55     cli();                  // 清中断许可。
56     while (bh->b_lock)      // 如果缓冲区已被锁定，则睡眠，直到缓冲区解锁。
57         sleep_on(&bh->b_wait);
58     bh->b_lock=1;           // 立刻锁定该缓冲区。
59     sti();                  // 开中断。
60 }

```

```

// 释放（解锁）锁定的缓冲区。
51 static inline void unlock_buffer(struct buffer_head * bh)
52 {
53     if (!bh->b_lock)                // 如果该缓冲区并没有被锁定，则打印出错信息。
54         printk("ll_rw_block.c: buffer not locked\n\r");
55     bh->b_lock = 0;                  // 清锁定标志。
56     wake_up(&bh->b_wait);           // 唤醒等待该缓冲区的任务。
57 }
58
59 /*
60  * add-request adds a request to the linked list.
61  * It disables interrupts so that it can muck with the
62  * request-lists in peace.
63  */
64 /*
65  * add-request() 向连表中加入一项请求。它关闭中断，
66  * 这样就能安全地处理请求连表了 */
67 /*
68  * 向链表中加入请求项。参数 dev 指定块设备，req 是请求的结构信息。
69 static void add_request(struct blk_dev_struct * dev, struct request * req)
70 {
71     struct request * tmp;
72
73     req->next = NULL;
74     cli();                // 关中断。
75     if (req->bh)
76         req->bh->b_dirt = 0;    // 清缓冲区“脏”标志。
77     // 如果 dev 的当前请求(current_request)子段为空，则表示目前该设备没有请求项，本次是第 1 个
78     // 请求项，因此可将块设备当前请求指针直接指向请求项。
79     if (!(tmp = dev->current_request)) {
80         dev->current_request = req;
81         sti();                // 开中断。
82         (dev->request_fn)();    // 对于硬盘设备(设备号 3)，是 do_hd_request()。
83         return;
84     }
85     // 如果目前该设备已经有请求项在等待，则首先利用电梯算法搜索最佳位置，然后将当前请求插入
86     // 请求链表中。
87     for ( ; tmp->next ; tmp=tmp->next)
88         if ((IN_ORDER(tmp, req) ||
89             !IN_ORDER(tmp, tmp->next)) &&
90             IN_ORDER(req, tmp->next))
91             break;
92     req->next=tmp->next;
93     tmp->next=req;
94     sti();
95 }
96
97 // 创建请求项并插入请求队列。参数是：主设备号 major，命令 rw，存放数据的缓冲区头指针 bh。
98 static void make_request(int major, int rw, struct buffer_head * bh)
99 {
100     struct request * req;
101     int rw_ahead;
102

```

```

93  /* WRITEA/READA is special case - it is not really needed, so if the */
94  /* buffer is locked, we just forget about it, else it's a normal read */
    /* WRITEA/READA 是特殊的情况 - 它们并不是必要的，所以如果缓冲区已经上锁，*/
    /* 我们就不管它而退出，否则的话就执行一般的读/写操作。 */
95      if (rw_ahead = (rw == READA || rw == WRITEA)) {
96          if (bh->b_lock)
97              return;
98          if (rw == READA)
99              rw = READ;
100         else
101             rw = WRITE;
102     }
    // 如果命令不是 READ 或 WRITE 则表示内核程序有错，显示出错信息并死机。
103     if (rw!=READ && rw!=WRITE)
104         panic("Bad block dev command, must be R/W/RA/WA");
    // 锁定缓冲区，如果缓冲区已经上锁，则当前任务（进程）就会睡眠，直到被明确地唤醒。
105     lock\_buffer(bh);
    // 如果命令是写并且缓冲区数据不脏，或者命令是读并且缓冲区数据是更新过的，则不用添加
    // 这个请求。将缓冲区解锁并退出。
106     if ((rw == WRITE && !bh->b_dirt) || (rw == READ && bh->b_uptodate)) {
107         unlock\_buffer(bh);
108         return;
109     }
110 repeat:
111 /* we don't allow the write-requests to fill up the queue completely:
112 * we want some room for reads: they take precedence. The last third
113 * of the requests are only for reads.
114 */
    /* 我们不能让队列中全都是写请求项：我们需要为读请求保留一些空间：读操作
    /* 是优先的。请求队列的后三分之一空间是为读准备的。
    */
    // 请求项是从请求数组末尾开始搜索空项填入的。根据上述要求，对于读命令请求，可以直接
    // 从队列末尾开始操作，而写请求则只能从队列的 2/3 处向头上搜索空项填入。
115     if (rw == READ)
116         req = request+NR\_REQUEST;           // 对于读请求，将队列指针指向队列尾部。
117     else
118         req = request+((NR\_REQUEST*2)/3); // 对于写请求，队列指针指向队列 2/3 处。
119 /* find an empty request */
    /* 搜索一个空请求项 */
    // 从后向前搜索，当请求结构 request 的 dev 字段值=-1 时，表示该项未被占用。
120     while (--req >= request)
121         if (req->dev<0)
122             break;
123 /* if none found, sleep on new requests: check for rw_ahead */
    /* 如果没有找到空闲项，则让该次新请求睡眠：需检查是否提前读/写 */
    // 如果没有一项是空闲的（此时 request 数组指针已经搜索越过头部），则查看此次请求是否是
    // 提前读/写（READA 或 WRITEA），如果是则放弃此次请求。否则让本次请求睡眠（等待请求队列
    // 腾出空项），过一会再来搜索请求队列。
124     if (req < request) {
125         if (rw_ahead) {
126             unlock\_buffer(bh);
127             return;
128         }

```

```

129         sleep\_on(&wait\_for\_request); // 否则让本次请求睡眠，过会再查看请求队列。
130         goto repeat;
131     }
132     /* fill up the request-info, and add it to the queue */
133     /* 向空闲请求项中填写请求信息，并将其加入队列中 */
134     // 请求结构参见 (kernel/blk_drv/blk.h, 23)。
135     req->dev = bh->b_dev; // 设备号。
136     req->cmd = rw; // 命令 (READ/WRITE)。
137     req->errors=0; // 操作时产生的错误次数。
138     req->sector = bh->b_blocknr<<1; // 起始扇区。(1 块=2 扇区)
139     req->nr_sectors = 2; // 读写扇区数。
140     req->buffer = bh->b_data; // 数据缓冲区。
141     req->waiting = NULL; // 任务等待操作执行完成的地方。
142     req->bh = bh; // 缓冲区头指针。
143     req->next = NULL; // 指向下一请求项。
144     add\_request(major+blk\_dev, req); // 将请求项加入队列中 (blk_dev[major], req)。
145 }
146
147 // 低层读写数据块函数。
148 // 该函数主要是在 fs/buffer.c 中被调用。实际的读写操作是由设备的 request_fn() 函数完成。
149 // 对于硬盘操作，该函数是 do_hd_request()。(kernel/blk_drv/hd.c, 294)
150 void ll\_rw\_block(int rw, struct buffer\_head * bh)
151 {
152     unsigned int major; // 主设备号 (对于硬盘是 3)。
153
154     // 如果设备的主设备号不存在或者该设备的读写操作函数不存在，则显示出错信息，并返回。
155     if ((major=MAJOR(bh->b_dev)) >= NR\_BLK\_DEV ||
156         !(blk\_dev[major].request_fn)) {
157         printk("Trying to read nonexistent block-device\n\r");
158         return;
159     }
160     make\_request(major, rw, bh); // 创建请求项并插入请求队列。
161 }
162
163 // 块设备初始化函数，由初始化程序 main.c 调用 (init/main.c, 128)。
164 // 初始化请求数组，将所有请求项置为空闲项 (dev = -1)。有 32 项 (NR_REQUEST = 32)。
165 void blk\_dev\_init(void)
166 {
167     int i;
168
169     for (i=0 ; i<NR\_REQUEST ; i++) {
170         request[i].dev = -1;
171         request[i].next = NULL;
172     }
173 }

```

## 5.6 ramdisk.c 文件

### 5.6.1 功能描述

### 5.6.2 代码注释

列表 linux/kernel/blk\_drv/ramdisk.c 程序

```

1  /*
2  *  linux/kernel/blk_drv/ramdisk.c
3  *
4  *  Written by Theodore Ts'o, 12/2/91
5  */
/* 由 Theodore Ts'o 编制, 12/2/91
*/
// Theodore Ts'o (Ted Ts'o) 是 linux 社区中的著名人物。Linux 在世界范围内的流行也有他很大的
// 功劳, 早在 Linux 操作系统刚问世时, 他就怀着极大的热情为 linux 的发展提供了 maillist, 并
// 在北美地区最早设立了 linux 的 ftp 站点 (tsx-ll.mit.edu), 而且至今仍然为广大 linux 用户
// 提供服务。他对 linux 作出的最大贡献之一是提出并实现了 ext2 文件系统。该文件系统已成为
// linux 世界中事实上的文件系统标准。最近他又推出了 ext3 文件系统, 大大提高了文件系统的
// 稳定性和访问效率。作为对他的推崇, 第 97 期 (2002 年 5 月) 的 linuxjournal 期刊将他作为
// 了封面人物, 并对他进行了采访。目前, 他为 IBM linux 技术中心工作, 并从事着有关 LSB
// (Linux Standard Base) 等方面的工作。(他的主页: http://thunk.org/tytso/)

6
7 #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。
8
9 #include <linux/config.h>    // 内核配置头文件。定义键盘语言和硬盘类型 (HD_TYPE) 可选项。
10 #include <linux/sched.h>    // 调度程序头文件, 定义了任务结构 task_struct、初始任务 0 的数据,
// 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
11 #include <linux/fs.h>       // 文件系统头文件。定义文件表结构 (file, buffer_head, m_inode 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/system.h>     // 系统头文件。定义了设置或修改描述符/中断门等的嵌入式汇编宏。
14 #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。
15 #include <asm/memory.h>    // 内存拷贝头文件。含有 memcpy() 嵌入式汇编宏函数。
16
17 #define MAJOR_NR 1          // 内存主设备号是 1。
18 #include "blk.h"
19
20 char    *rd_start;          // 虚拟盘在内存中的起始位置。在 52 行初始化函数 rd_init() 中
// 确定。参见 (init/main.c, 124) (缩写 rd_代表 ramdisk_)。
21 int     rd_length = 0;      // 虚拟盘所占内存大小 (字节)。
22
// 执行虚拟盘 (ramdisk) 读写操作。程序结构与 do_hd_request() 类似 (kernel/blk_drv/hd.c, 294)。
23 void do_rd_request(void)
24 {
25     int     len;
26     char    *addr;
27
28     INIT_REQUEST;          // 检测请求的合法性 (参见 kernel/blk_drv/blk.h, 127)。
// 下面语句取得 ramdisk 的起始扇区对应的内存起始位置和内存长度。
// 其中 sector << 9 表示 sector * 512, CURRENT 定义为 (blk_dev[MAJOR_NR].current_request)。

```



```

29     addr = rd_start + (CURRENT->sector << 9);
30     len = CURRENT->nr_sectors << 9;
// 如果子设备号不为 1 或者对应内存起始位置>虚拟盘末尾, 则结束该请求, 并跳转到 repeat 处
// (定义在 28 行的 INIT_REQUEST 内开始处)。
31     if ((MINOR(CURRENT->dev) != 1) || (addr+len > rd_start+rd_length)) {
32         end_request(0);
33         goto repeat;
34     }
// 如果是写命令(WRITE), 则将请求项中缓冲区的内容复制到 addr 处, 长度为 len 字节。
35     if (CURRENT->cmd == WRITE) {
36         (void) memcpy(addr,
37                       CURRENT->buffer,
38                       len);
// 如果是读命令(READ), 则将 addr 开始的内容复制到请求项中缓冲区中, 长度为 len 字节。
39     } else if (CURRENT->cmd == READ) {
40         (void) memcpy(CURRENT->buffer,
41                       addr,
42                       len);
// 否则显示命令不存在, 死机。
43     } else
44         panic("unknown ramdisk-command");
// 请求项成功后处理, 置更新标志。并继续处理本设备的下一请求项。
45     end_request(1);
46     goto repeat;
47 }
48
49 /*
50  * Returns amount of memory which needs to be reserved.
51  */
// 返回内存虚拟盘 ramdisk 所需的内存量 */
// 虚拟盘初始化函数。确定虚拟盘在内存中的起始地址, 长度。并对整个虚拟盘区清零。
52 long rd_init(long mem_start, int length)
53 {
54     int i;
55     char *cp;
56
57     blk_dev[MAJOR_NR].request_fn = DEVICE_REQUEST; // do_rd_request()。
58     rd_start = (char *) mem_start;
59     rd_length = length;
60     cp = rd_start;
61     for (i=0; i < length; i++)
62         *cp++ = '\0';
63     return(length);
64 }
65
66 /*
67  * If the root device is the ram disk, try to load it.
68  * In order to do this, the root device is originally set to the
69  * floppy, and we later change it to be ram disk.
70  */
//
// * 如果根文件系统设备(root device)是 ramdisk 的话, 则尝试加载它。root device 原先是指向
// * 软盘的, 我们将它改成指向 ramdisk。

```

```

    */
    //// 加载根文件系统到 ramdisk。
71 void rd_load(void)
72 {
73     struct buffer_head *bh;
74     struct super_block s;
75     int block = 256;    /* Start at block 256 */
76     int i = 1;
77     int nblocks;
78     char *cp;           /* Move pointer */
79
80     if (!rd_length)      // 如果 ramdisk 的长度为零，则退出。
81         return;
82     printk("Ram disk: %d bytes, starting at 0x%x\n", rd_length,
83            (int) rd_start); // 显示 ramdisk 的大小以及内存起始位置。
84     if (MAJOR(ROOT_DEV) != 2) // 如果此时根文件设备不是软盘，则退出。
85         return;
86     // 读软盘块 256+1, 256, 256+2。breada()用于读取指定的数据块，并标出还需要读的块，然后返回
87     // 含有数据块的缓冲区指针。如果返回 NULL，则表示数据块不可读(fs/buffer.c, 322)。
88     // 这里 block+1 是指磁盘上的超级块。
89     bh = breada(ROOT_DEV, block+1, block, block+2, -1);
90     if (!bh) {
91         printk("Disk error while looking for ramdisk!\n");
92         return;
93     }
94     // 将 s 指向缓冲区中的磁盘超级块。(d_super_block 磁盘中超级块结构)。
95     *((struct d_super_block *) &s) = *((struct d_super_block *) bh->b_data);
96     brelse(bh); // [?? 为什么数据没有复制就立刻释放呢?]
97     if (s.s_magic != SUPER_MAGIC) // 如果超级块中魔数不对，则说明不是 minix 文件系统。
98         /* No ram disk image present, assume normal floppy boot */
99         /* 磁盘中没有 ramdisk 映像文件，退出执行通常的软盘引导 */
100         return;
101     // 块数 = 逻辑块数(区段数) * 2^(每区段块数的次方)。
102     // 如果数据块数大于内存中虚拟盘所能容纳的块数，则不能加载，显示出错信息并返回。否则显示
103     // 加载数据块信息。
104     nblocks = s.s_nzones << s.s_log_zone_size;
105     if (nblocks > (rd_length >> BLOCK_SIZE_BITS)) {
106         printk("Ram disk image too big! (%d blocks, %d avail)\n",
107            nblocks, rd_length >> BLOCK_SIZE_BITS);
108         return;
109     }
110     printk("Loading %d bytes into ram disk... 0000k",
111        nblocks << BLOCK_SIZE_BITS);
112     // cp 指向虚拟盘起始处，然后将磁盘上的根文件系统映像文件复制到虚拟盘上。
113     cp = rd_start;
114     while (nblocks) {
115         if (nblocks > 2) // 如果需读取的块数多于 3 块则采用超前预读方式读数据块。
116             bh = breada(ROOT_DEV, block, block+1, block+2, -1);
117         else // 否则就单块读取。
118             bh = bread(ROOT_DEV, block);
119         if (!bh) {
120             printk("I/O error on block %d, aborting load\n",
121                block);

```

```

113         return;
114     }
115     (void) memcpy(cp, bh->b_data, BLOCK_SIZE); // 将缓冲区中的数据复制到 cp 处。
116     brelse(bh); // 释放缓冲区。
117     printk("010|010|010|010|010%4dk", i); // 打印加载块计数值。
118     cp += BLOCK_SIZE; // 虚拟盘指针前移。
119     block++;
120     nblocks--;
121     i++;
122 }
123 printk("010|010|010|010|010done\n");
124 ROOT_DEV=0x0101; // 修改 ROOT_DEV 使其指向虚拟盘 ramdisk。
125 }
126

```

## 5.7 floppy.c 文件

### 5.7.1 功能描述

### 5.7.2 代码注释

列表 linux/kernel/blk\_drv/floppy.c 程序

```

1  /*
2  *  linux/kernel/floppy.c
3  *
4  *  (C) 1991  Linus Torvalds
5  */
6
7  /*
8  *  02.12.91 - Changed to static variables to indicate need for reset
9  *  and recalibrate. This makes some things easier (output_byte reset
10 *  checking etc), and means less interrupt jumping in case of errors,
11 *  so the code is hopefully easier to understand.
12 */
13 /*
14 *  02.12.91 - 修改成静态变量，以适应复位和重新校正操作。这使得某些事情
15 *  做起来较为方便（output_byte 复位检查等），并且意味着在出错时中断跳转
16 *  要少一些，所以希望代码能更容易被理解。
17 */
18 /*
19 *  This file is certainly a mess. I've tried my best to get it working,
20 *  but I don't like programming floppies, and I have only one anyway.
21 *  Urgel. I should check for more errors, and do more graceful error
22 *  recovery. Seems there are problems with several drives. I've tried to
23 *  correct them. No promises.
24 */
25

```

```

* 这个文件当然比较混乱。我已经尽我所能使其能够工作，但我不喜欢软驱编程，
* 而且我也只有一个软驱。另外，我应该做更多的查错工作，以及改正更多的错误。
* 对于某些软盘驱动器好象还存在一些问题。我已经尝试着进行纠正了，但不能保证
* 问题已消失。
*/

```

```

21
22 /*
23  * As with hd.c, all routines within this file can (and will) be called
24  * by interrupts, so extreme caution is needed. A hardware interrupt
25  * handler may not sleep, or a kernel panic will happen. Thus I cannot
26  * call "floppy-on" directly, but have to set a special timer interrupt
27  * etc.
28  *
29  * Also, I'm not certain this works on more than 1 floppy. Bugs may
30  * abund.
31  */
/*
* 如同 hd.c 文件一样，该文件中的所有子程序都能够被中断调用，所以需要特别
* 地小心。硬件中断处理程序是不能睡眠的，否则内核就会傻掉(死机)☹。因此不能
* 直接调用“floppy-on”，而只能设置一个特殊的时间中断等。
*
* 另外，我不能保证该程序能在多于 1 个软驱的系统上工作，有可能存在错误。
*/

32
33 #include <linux/sched.h> // 调度程序头文件，定义了任务结构 task_struct、初始任务 0 的数据，
// 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
34 #include <linux/fs.h> // 文件系统头文件。定义文件表结构 (file, buffer_head, m_inode 等)。
35 #include <linux/kernel.h> // 内核头文件。含有一些内核常用函数的原形定义。
36 #include <linux/fdreg.h> // 软驱头文件。含有软盘控制器参数的一些定义。
37 #include <asm/system.h> // 系统头文件。定义了设置或修改描述符/中断门等的嵌入式汇编宏。
38 #include <asm/io.h> // io 头文件。定义硬件端口输入/输出宏汇编语句。
39 #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。
40
41 #define MAJOR_NR 2 // 软驱的主设备号是 2。
42 #include "blk.h" // 块设备头文件。定义请求数据结构、块设备数据结构和宏函数等信息。
43
44 static int recalibrate = 0; // 标志：需要重新校正。
45 static int reset = 0; // 标志：需要进行复位操作。
46 static int seek = 0; // 寻道。
47
48 extern unsigned char current_DOR; // 当前数字输出寄存器(Digital Output Register)。
49
50 #define immoutb_p(val,port) \ // 字节直接输出（嵌入汇编语言宏）。
51 __asm__( "outb %0,%1\n\tjmp 1f\n1:\tjmp 1f\n1:": "a" ((char) (val)), "i" (port))
52
// 这两个定义用于计算软驱的设备号。次设备号 = TYPE*4 + DRIVE。计算方法参见列表后。
53 #define TYPE(x) ((x)>>2) // 软驱类型 (2--1.2Mb, 7--1.44Mb)。
54 #define DRIVE(x) ((x)&0x03) // 软驱序号 (0--3 对应 A--D)。
55 /*
56  * Note that MAX_ERRORS=8 doesn't imply that we retry every bad read
57  * max 8 times - some types of errors increase the errorcount by 2,
58  * so we might actually retry only 5-6 times before giving up.
59  */

```

```

/*
 * 注意，下面定义 MAX_ERRORS=8 并不表示对每次读错误尝试最多 8 次 - 有些类型
 * 的错误将把出错计数值乘 2，所以我们实际上在放弃操作之前只需尝试 5-6 遍即可。
 */
60 #define MAX_ERRORS 8
61
62 /*
63  * globals used by 'result()'
64  */
/* 下面是函数 'result()' 使用的全局变量 */
// 这些状态字节中各比特位的含义请参见 include/linux/fdreg.h 头文件。
65 #define MAX_REPLIES 7 // FDC 最多返回 7 字节的结果信息。
66 static unsigned char reply_buffer[MAX_REPLIES]; // 存放 FDC 返回的结果信息。
67 #define ST0 (reply_buffer[0]) // 返回结果状态字节 0。
68 #define ST1 (reply_buffer[1]) // 返回结果状态字节 1。
69 #define ST2 (reply_buffer[2]) // 返回结果状态字节 2。
70 #define ST3 (reply_buffer[3]) // 返回结果状态字节 3。
71
72 /*
73  * This struct defines the different floppy types. Unlike minix
74  * linux doesn't have a "search for right type"-type, as the code
75  * for that is convoluted and weird. I've got enough problems with
76  * this driver as it is.
77  *
78  * The 'stretch' tells if the tracks need to be boubled for some
79  * types (ie 360kB diskette in 1.2MB drive etc). Others should
80  * be self-explanatory.
81  */
/*
 * 下面的软盘结构定义了不同的软盘类型。与 minix 不同的是，linux 没有
 * "搜索正确的类型"-类型，因为对其处理的代码令人费解且怪怪的。本程序
 * 已经让我遇到了许多的问题了。
 *
 * 对某些类型的软盘（例如在 1.2MB 驱动器中的 360kB 软盘等），'stretch' 用于
 * 检测磁道是否需要特殊处理。其它参数应该是自明的。
 */
// 软盘参数有：
// size      大小(扇区数)；
// sect      每磁道扇区数；
// head      磁头数；
// track     磁道数；
// stretch  对磁道是否要特殊处理（标志）；
// gap      扇区间隙长度(字节数)；
// rate      数据传输速率；
// spec1     参数（高 4 位步进速率，低四位磁头卸载时间）。
82 static struct floppy_struct {
83     unsigned int size, sect, head, track, stretch;
84     unsigned char gap, rate, spec1;
85 } floppy_type[] = {
86     { 0, 0, 0, 0, 0, 0x00, 0x00, 0x00 }, // no testing */
87     { 720, 9, 2, 40, 0, 0x2A, 0x02, 0xDF }, // 360kB PC diskettes */
88     { 2400, 15, 2, 80, 0, 0x1B, 0x00, 0xDF }, // 1.2 MB AT-diskettes */
89     { 720, 9, 2, 40, 1, 0x2A, 0x02, 0xDF }, // 360kB in 720kB drive */

```

```

90     { 1440, 9, 2, 80, 0, 0x2A, 0x02, 0xDF },      /* 3.5" 720kB diskette */
91     { 720, 9, 2, 40, 1, 0x23, 0x01, 0xDF },      /* 360kB in 1.2MB drive */
92     { 1440, 9, 2, 80, 0, 0x23, 0x01, 0xDF },      /* 720kB in 1.2MB drive */
93     { 2880, 18, 2, 80, 0, 0x1B, 0x00, 0xCF },      /* 1.44MB diskette */
94 };
95 /*
96  * Rate is 0 for 500kb/s, 2 for 300kbps, 1 for 250kbps
97  * Spec1 is 0xSH, where S is stepping rate (F=1ms, E=2ms, D=3ms etc),
98  * H is head unload time (1=16ms, 2=32ms, etc)
99  *
100  * Spec2 is (HLD<<1 | ND), where HLD is head load time (1=2ms, 2=4 ms etc)
101  * and ND is set means no DMA. Hardcoded to 6 (HLD=6ms, use DMA).
102  */
103 /*
104  * 上面速率 rate: 0 表示 500kb/s, 1 表示 300kbps, 2 表示 250kbps。
105  * 参数 spec1 是 0xSH, 其中 S 是步进速率 (F=1 毫秒, E=2ms, D=3ms 等),
106  * H 是磁头卸载时间 (1=16ms, 2=32ms 等)
107  *
108  * spec2 是 (HLD<<1 | ND), 其中 HLD 是磁头加载时间 (1=2ms, 2=4ms 等)
109  * ND 置位表示不使用 DMA (No DMA), 在程序中硬编码成 6 (HLD=6ms, 使用 DMA)。
110  */
111
112 extern void floppy_interrupt(void);
113 extern char tmp_floppy_area[1024];
114
115 /*
116  * These are global variables, as that's the easiest way to give
117  * information to interrupts. They are the data used for the current
118  * request.
119  */
120 /*
121  * 下面是一些全局变量, 因为这是将信息传给中断程序最简单的方式。它们是
122  * 用于当前请求的数据。
123  */
124
125 static int cur_spec1 = -1;
126 static int cur_rate = -1;
127 static struct floppy_struct * floppy = floppy_type;
128 static unsigned char current_drive = 0;
129 static unsigned char sector = 0;
130 static unsigned char head = 0;
131 static unsigned char track = 0;
132 static unsigned char seek_track = 0;
133 static unsigned char current_track = 255;
134 static unsigned char command = 0;
135 unsigned char selected = 0;
136 struct task_struct * wait_on_floppy_select = NULL;
137
138 // 释放 (取消选定的) 软盘 (软驱)。
139 // 数字输出寄存器 (DOR) 的低 2 位用于指定选择的软驱 (0-3 对应 A-D)。
140 void floppy_deselect(unsigned int nr)
141 {
142     if (nr != (current_DOR & 3))
143         printk("floppy_deselect: drive not selected\n");
144 }

```

```

129     selected = 0;
130     wake\_up(&wait\_on\_floppy\_select);
131 }
132
133 /*
134  * floppy-change is never called from an interrupt, so we can relax a bit
135  * here, sleep etc. Note that floppy-on tries to set current_DOR to point
136  * to the desired drive, but it will probably not survive the sleep if
137  * several floppies are used at the same time: thus the loop.
138  */
139
140 /*
141  * floppy-change() 不是从中断程序中调用的，所以这里我们可以轻松一下，睡觉等。
142  * 注意 floppy-on() 会尝试设置 current_DOR 指向所需的驱动器，但当同时使用几个
143  * 软盘时不能睡眠：因此此时只能使用循环方式。
144  */
145
146 // 检测指定软驱中软盘更换情况。如果软盘更换了则返回 1，否则返回 0。
139 int floppy\_change(unsigned int nr)
140 {
141     repeat:
142         floppy\_on(nr); // 开启指定软驱 nr (kernel/sched.c, 251)。
143         // 如果当前选择的软驱不是指定的软驱 nr，并且已经选择其它了软驱，则让当前任务进入可中断
144         // 等待状态。
145         while ((current\_DOR & 3) != nr && selected)
146             interruptible\_sleep\_on(&wait\_on\_floppy\_select);
147         // 如果当前没有选择其它软驱或者当前任务被唤醒时，当前软驱仍然不是指定的软驱 nr，则循环等待。
148         if ((current\_DOR & 3) != nr)
149             goto repeat;
150         // 取数字输入寄存器值，如果最高位（位 7）置位，则表示软盘已更换，此时关闭马达并退出返回 1。
151         // 否则关闭马达退出返回 0。
152         if (inb(FD\_DIR) & 0x80) {
153             floppy\_off(nr);
154             return 1;
155         }
156         floppy\_off(nr);
157         return 0;
158     }
159 }
160
161 // 复制内存块。
155 #define copy\_buffer(from,to) \
156 __asm__ ("cld ; rep ; movsl" \
157         ::: "c" (BLOCK\_SIZE/4), "S" ((long)(from)), "D" ((long)(to)) \
158         : "cx", "di", "si")
159
160 // 设置（初始化）软盘 DMA 通道。
160 static void setup\_DMA(void)
161 {
162     long addr = (long) CURRENT->buffer; // 当前请求项缓冲区所处内存中位置（地址）。
163
164     cli();
165     // 如果缓冲区处于内存 1M 以上的地方，则将 DMA 缓冲区设在临时缓冲区域(tmp_floppy_area 数组)
166     // (因为 8237A 芯片只能在 1M 地址范围内寻址)。如果是写盘命令，则还需将数据复制到该临时区域。
167     if (addr >= 0x100000) {
168         addr = (long) tmp\_floppy\_area;

```



```

167         if (command == FD\_WRITE)
168             copy\_buffer(CURRENT->buffer, tmp\_floppy\_area);
169     }
170     /* mask DMA 2 */ /* 屏蔽 DMA 通道 2 */
171     // 单通道屏蔽寄存器端口为 0x10。位 0-1 指定 DMA 通道(0-3)，位 2: 1 表示屏蔽，0 表示允许请求。
172     immoutb\_p(4|2, 10);
173     /* output command byte. I don't know why, but everyone (minix, */
174     /* sanches & canton) output this twice, first to 12 then to 11 */
175     /* 输出命令字节。我是不知道为什么，但是每个人 (minix, */
176     /* sanches 和 canton) 都输出两次，首先是 12 口，然后是 11 口 */
177     // 下面嵌入汇编代码向 DMA 控制器端口 12 和 11 写方式字 (读盘 0x46，写盘 0x4A)。
178     __asm__ ("outb %%al, $12\n\tjmp 1f\n1:\tjmp 1f\n1:\t"
179             "outb %%al, $11\n\tjmp 1f\n1:\tjmp 1f\n1:::"
180             "a" ((char) ((command == FD\_READ) ? DMA\_READ : DMA\_WRITE)));
181     /* 8 low bits of addr */ /* 地址低 0-7 位 */
182     // 向 DMA 通道 2 写入基/当前地址寄存器 (端口 4)。
183     immoutb\_p(addr, 4);
184     addr >>= 8;
185     /* bits 8-15 of addr */ /* 地址高 8-15 位 */
186     immoutb\_p(addr, 4);
187     addr >>= 8;
188     /* bits 16-19 of addr */ /* 地址 16-19 位 */
189     // DMA 只可以在 1M 内存空间内寻址，其高 16-19 位地址需放入页面寄存器 (端口 0x81)。
190     immoutb\_p(addr, 0x81);
191     /* low 8 bits of count-1 (1024-1=0x3ff) */ /* 计数器低 8 位 (1024-1=0x3ff) */
192     // 向 DMA 通道 2 写入基/当前字节计数器值 (端口 5)。
193     immoutb\_p(0xff, 5);
194     /* high 8 bits of count-1 */ /* 计数器高 8 位 */
195     // 一次共传输 1024 字节 (两个扇区)。
196     immoutb\_p(3, 5);
197     /* activate DMA 2 */ /* 开启 DMA 通道 2 的请求 */
198     // 复位对 DMA 通道 2 的屏蔽，开放 DMA2 请求 DREQ 信号。
199     immoutb\_p(0|2, 10);
200     sti();
201 }
202
203     //// 向软盘控制器输出一个字节数据 (命令或参数)。
204     static void output\_byte(char byte)
205     {
206         int counter;
207         unsigned char status;
208
209         if (reset)
210             return;
211         // 循环读取主状态控制器 FD_STATUS(0x3f4) 的状态。如果状态是 STATUS_READY 并且 STATUS_DIR=0
212         // (CPU→FDC)，则向数据端口输出指定字节。
213         for(counter = 0 ; counter < 10000 ; counter++) {
214             status = inb\_p(FD\_STATUS) & (STATUS\_READY | STATUS\_DIR);
215             if (status == STATUS\_READY) {
216                 outb(byte, FD\_DATA);
217                 return;
218             }
219         }
220     }
221 }

```



```

// 如果到循环 1 万次结束还不能发送，则置复位标志，并打印出错信息。
208     reset = 1;
209     printk("Unable to send byte to FDC\n\r");
210 }
211
///// 读取 FDC 执行的结果信息。
// 结果信息最多 7 个字节，存放在 reply_buffer[] 中。返回读入的结果字节数，若返回值=-1
// 表示出错。
212 static int result(void)
213 {
214     int i = 0, counter, status;
215
216     if (reset)
217         return -1;
218     for (counter = 0 ; counter < 10000 ; counter++) {
219         status = inb_p(FD_STATUS) & (STATUS_DIR | STATUS_READY | STATUS_BUSY);
220         if (status == STATUS_READY)
221             return i;
222         if (status == (STATUS_DIR | STATUS_READY | STATUS_BUSY)) {
223             if (i >= MAX_REPLIES)
224                 break;
225             reply_buffer[i++] = inb_p(FD_DATA);
226         }
227     }
228     reset = 1;
229     printk("Getstatus times out\n\r");
230     return -1;
231 }
232
///// 软盘操作出错中断调用函数。由软驱中断处理程序调用。
233 static void bad flp intr(void)
234 {
235     CURRENT->errors++; // 当前请求项出错次数增 1。
// 如果当前请求项出错次数大于最大允许出错次数，则取消选定当前软驱，并结束该请求项（不更新）。
236     if (CURRENT->errors > MAX_ERRORS) {
237         floppy_deselect(current_drive);
238         end_request(0);
239     }
// 如果当前请求项出错次数大于最大允许出错次数的一半，则置复位标志，需对软驱进行复位操作，
// 然后再试。否则软驱需重新校正一下，再试。
240     if (CURRENT->errors > MAX_ERRORS/2)
241         reset = 1;
242     else
243         recalibrate = 1;
244 }
245
246 /*
247  * Ok, this interrupt is called after a DMA read/write has succeeded,
248  * so we check the results, and copy any buffers.
249  */
250 /*
251  * OK, 下面该中断处理函数是在 DMA 读/写成功后调用的，这样我们就可以检查执行结果，
252  * 并复制缓冲区中的数据。

```

```

    */
    //// 软盘读写操作成功中断调用函数。。
250 static void rw_interrupt(void)
251 {
    // 如果返回结果字节数不等于 7，或者状态字节 0、1 或 2 中存在出错标志，则若是写保护
    // 就显示出错信息，释放当前驱动器，并结束当前请求项。否则就执行出错计数处理。
    // 然后继续执行软盘请求操作。
    // ( 0xf8 = ST0_INTR | ST0_SE | ST0_ECE | ST0_NR )
    // ( 0xbf = ST1_EOC | ST1_CRC | ST1_OR | ST1_ND | ST1_WP | ST1_MAM, 应该是 0xb7)
    // ( 0x73 = ST2_CM | ST2_CRC | ST2_WC | ST2_BC | ST2_MAM )
252     if (result() != 7 || (ST0 & 0xf8) || (ST1 & 0xbf) || (ST2 & 0x73)) {
253         if (ST1 & 0x02) { // 0x02 = ST1_WP - Write Protected.
254             printk("Drive %d is write protected\n\r", current_drive);
255             floppy_deselect(current_drive);
256             end_request(0);
257         } else
258             bad_flp_intr();
259         do_fd_request();
260         return;
261     }
    // 如果当前请求项的缓冲区位于 1M 地址以上，则说明此次软盘读操作的内容还放在临时缓冲区内，
    // 需要复制到请求项的缓冲区中（因为 DMA 只能在 1M 地址范围寻址）。
262     if (command == FD_READ && (unsigned long)(CURRENT->buffer) >= 0x100000)
263         copy_buffer(tmp_floppy_area, CURRENT->buffer);
    // 释放当前软盘，结束当前请求项（置更新标志），再继续执行其它软盘请求项。
264     floppy_deselect(current_drive);
265     end_request(1);
266     do_fd_request();
267 }
268
    //// 设置 DMA 并输出软盘操作命令和参数（输出 1 字节命令+ 0~7 字节参数）。
269 inline void setup_rw_floppy(void)
270 {
271     setup_DMA(); // 初始化软盘 DMA 通道。
272     do_floppy = rw_interrupt; // 置软盘中断调用函数指针。
273     output_byte(command); // 发送命令字节。
274     output_byte(head<<2 | current_drive); // 发送参数（磁头号+驱动器号）。
275     output_byte(track); // 发送参数（磁道号）。
276     output_byte(head); // 发送参数（磁头号）。
277     output_byte(sector); // 发送参数（起始扇区号）。
278     output_byte(2); /* sector size = 512 */ // 发送参数（字节数(N=2) 512 字节）。
279     output_byte(floppy->sect); // 发送参数（每磁道扇区数）。
280     output_byte(floppy->gap); // 发送参数（扇区间隔长度）。
281     output_byte(0xFF); /* sector size (0xff when n!=0 ?) */
    // 发送参数（当 N=0 时，扇区定义的字节长度），这里无用。
    // 若在发送命令和参数时发生错误，则继续执行下一软盘操作请求。
282     if (reset)
283         do_fd_request();
284 }
285
286 /*
287  * This is the routine called after every seek (or recalibrate) interrupt
288  * from the floppy controller. Note that the "unexpected interrupt" routine

```

```

289  * also does a recalibrate, but doesn't come here.
290  */
/*
 * 该子程序是在每次软盘控制器寻道（或重新校正）中断后被调用的。注意
 * "unexpected interrupt"(意外中断)子程序也会执行重新校正操作，但并不在此地。
 */
///// 寻道处理中断调用函数。
// 首先发送检测中断状态命令，获得状态信息 ST0 和磁头所在磁道信息。若出错则执行错误计数
// 检测处理或取消本次软盘操作请求项。否则根据状态信息设置当前磁道变量，然后调用函数
// setup_rw_floppy() 设置 DMA 并输出软盘读写命令和参数。
291 static void seek_interrupt(void)
292 {
293  /* sense drive status */ /* 检测中断状态 */
// 发送检测中断状态命令，该命令不带参数。返回结果信息两个字节：ST0 和磁头当前磁道号。
294  output_byte(FD_SENSEI);
// 如果返回结果字节数不等于 2，或者 ST0 不为寻道结束，或者磁头所在磁道(ST1)不等于设定磁道，
// 则说明发生了错误，于是执行检测错误计数处理，然后继续执行软盘请求项，并退出。
295  if (result() != 2 || (ST0 & 0xF8) != 0x20 || ST1 != seek_track) {
296      bad_flp_intr();
297      do_fd_request();
298      return;
299  }
300  current_track = ST1; // 设置当前磁道。
301  setup_rw_floppy(); // 设置 DMA 并输出软盘操作命令和参数。
302 }
303
304 /*
305  * This routine is called when everything should be correctly set up
306  * for the transfer (ie floppy motor is on and the correct floppy is
307  * selected).
308  */
/*
 * 该函数是在传输操作的所有信息都正确设置好后被调用的（也即软驱马达已开启
 * 并且已选择了正确的软盘（软驱））。
 */
///// 读写数据传输函数。

309 static void transfer(void)
310 {
// 首先看当前驱动器参数是否就是指定驱动器的参数，若不是就发送设置驱动器参数命令及相应
// 参数（参数 1：高 4 位步进速率，低四位磁头卸载时间；参数 2：磁头加载时间）。
311  if (cur_spec1 != floppy->spec1) {
312      cur_spec1 = floppy->spec1;
313      output_byte(FD_SPECIFY); // 发送设置磁盘参数命令。
314      output_byte(cur_spec1); // hut etc */ // 发送参数。
315      output_byte(6); // Head load time =6ms, DMA */
316  }
// 判断当前数据传输速率是否与指定驱动器的一致，若不是就发送指定软驱的速率值到数据传输
// 速率控制寄存器(FD_DCR)。
317  if (cur_rate != floppy->rate)
318      outb_p(cur_rate = floppy->rate, FD_DCR);
// 若返回结果信息表明出错，则再调用软盘请求函数，并返回。
319  if (reset) {

```

```

320         do_fd_request();
321         return;
322     }
    // 若寻道标志为零（不需要寻道），则设置 DMA 并发送相应读写操作命令和参数，然后返回。
323     if (!seek) {
324         setup_rw_floppy();
325         return;
326     }
    // 否则执行寻道处理。置软盘中断处理调用函数为寻道中断函数。
327     do_floppy = seek_interrupt;
    // 如果器始磁道号不等于零则发送磁头寻道命令和参数
328     if (seek_track) {
329         output_byte(FD_SEEK);    // 发送磁头寻道命令。
330         output_byte(head<<2 | current_drive); //发送参数：磁头号+当前软驱号。
331         output_byte(seek_track);    // 发送参数：磁道号。
332     } else {
333         output_byte(FD_RECALIBRATE);    // 发送重新校正命令。
334         output_byte(head<<2 | current_drive); //发送参数：磁头号+当前软驱号。
335     }
    // 如果复位标志已置位，则继续执行软盘请求项。
336     if (reset)
337         do_fd_request();
338 }
339
340 /*
341  * Special case - used after a unexpected interrupt (or reset)
342  */
343 /*
344  * 特殊情况 - 用于意外中断（或复位）处理后。
345  */
346 // 软驱重新校正中断调用函数。
347 // 首先发送检测中断状态命令（无参数），如果返回结果表明出错，则置复位标志，否则复位重新
348 // 校正标志。然后再次执行软盘请求。
349 static void recal_interrupt(void)
350 {
351     output_byte(FD_SENSEI);    // 发送检测中断状态命令。
352     if (result()!=2 || (ST0 & 0xE0) == 0x60) // 如果返回结果字节数不等于 2 或命令
353         reset = 1;    // 异常结束，则置复位标志。
354     else    // 否则复位重新校正标志。
355         recalibrate = 0;
356     do_fd_request();    // 执行软盘请求项。
357 }
358
359 // 意外软盘中断请求中断调用函数。
360 // 首先发送检测中断状态命令（无参数），如果返回结果表明出错，则置复位标志，否则置重新
361 // 校正标志。
362 void unexpected_floppy_interrupt(void)
363 {
364     output_byte(FD_SENSEI);    // 发送检测中断状态命令。
365     if (result()!=2 || (ST0 & 0xE0) == 0x60) // 如果返回结果字节数不等于 2 或命令
366         reset = 1;    // 异常结束，则置复位标志。
367     else    // 否则置重新校正标志。
368         recalibrate = 1;

```

```

360 }
361
362 // 软盘重新校正处理函数。
363 // 向软盘控制器 FDC 发送重新校正命令和参数，并复位重新校正标志。
364 static void recalibrate_floppy(void)
365 {
366     recalibrate = 0; // 复位重新校正标志。
367     current_track = 0; // 当前磁道号归零。
368     do_floppy = recal_interrupt; // 置软盘中断调用函数指针指向重新校正调用函数。
369     output_byte(FD_RECALIBRATE); // 发送命令：重新校正。
370     output_byte(head<<2 | current_drive); // 发送参数：（磁头号加）当前驱动器号。
371     if (reset) // 如果出错（复位标志被置位）则继续执行软盘请求。
372         do_fd_request();
373 }
374
375 // 软盘控制器 FDC 复位中断调用函数。在软盘中断处理程序中调用。
376 // 首先发送检测中断状态命令（无参数），然后读出返回的结果字节。接着发送设定软驱参数命令
377 // 和相关参数，最后再次调用执行软盘请求。
378 static void reset_interrupt(void)
379 {
380     output_byte(FD_SENSEI); // 发送检测中断状态命令。
381     (void) result(); // 读取命令执行结果字节。
382     output_byte(FD_SPECIFY); // 发送设定软驱参数命令。
383     output_byte(cur_spec1); /* hut etc */ // 发送参数。
384     output_byte(6); /* Head load time =6ms, DMA */
385     do_fd_request(); // 调用执行软盘请求。
386 }
387
388 /*
389  * reset is done by pulling bit 2 of DOR low for a while.
390  */
391 /* FDC 复位是通过将数字输出寄存器 (DOR) 位 2 置 0 一会儿实现的 */
392 // 复位软盘控制器。
393 static void reset_floppy(void)
394 {
395     int i;
396
397     reset = 0; // 复位标志置 0。
398     cur_spec1 = -1;
399     cur_rate = -1;
400     recalibrate = 1; // 重新校正标志置位。
401     printk("Reset-floppy called\n\r"); // 显示执行软盘复位操作信息。
402     cli(); // 关中断。
403     do_floppy = reset_interrupt; // 设置在软盘中断处理程序中调用的函数。
404     outb_p(current_DOR & ~0x04, FD_DOR); // 对软盘控制器 FDC 执行复位操作。
405     for (i=0 ; i<100 ; i++) // 空操作，延迟。
406         __asm__("nop");
407     outb(current_DOR, FD_DOR); // 再启动软盘控制器。
408     sti(); // 开中断。
409 }
410
411 // 软驱启动定时中断调用函数。
412 // 首先检查数字输出寄存器 (DOR)，使其选择当前指定的驱动器。然后调用执行软盘读写传输

```

```

// 函数 transfer()。
404 static void floppy_on_interrupt(void)
405 {
406 /* We cannot do a floppy-select, as that might sleep. We just force it */
/* 我们不能任意设置选择的软驱，因为这样做可能会引起进程睡眠。我们只是迫使它自己选择 */
407     selected = 1; // 置已选择当前驱动器标志。
// 如果当前驱动器号与数字输出寄存器 DOR 中的不同，则重新设置 DOR 为当前驱动器 current_drive。
// 定时延迟 2 个滴答时间，然后调用软盘读写传输函数 transfer()。否则直接调用软盘读写传输函数。
408     if (current_drive != (current_DOR & 3)) {
409         current_DOR &= 0xFC;
410         current_DOR |= current_drive;
411         outb(current_DOR, FD_DOR); // 向数字输出寄存器输出当前 DOR。
412         add_timer(2, &transfer); // 添加定时器并执行传输函数。
413     } else
414         transfer(); // 执行软盘读写传输函数。
415 }
416
///// 软盘读写请求项处理函数。
//
417 void do_fd_request(void)
418 {
419     unsigned int block;
420
421     seek = 0;
// 如果复位标志已置位，则执行软盘复位操作，并返回。
422     if (reset) {
423         reset_floppy();
424         return;
425     }
// 如果重新校正标志已置位，则执行软盘重新校正操作，并返回。
426     if (recalibrate) {
427         recalibrate_floppy();
428         return;
429     }
// 检测请求项的合法性(参见 kernel/blk_drv/blk.h, 127)。
430     INIT_REQUEST;
// 将请求项结构中软盘设备号中的软盘类型(MINOR(CURRENT->dev)>>2)作为索引取得软盘参数块。
431     floppy = (MINOR(CURRENT->dev)>>2) + floppy_type;
// 如果当前驱动器不是请求项中指定的驱动器，则置标志 seek，表示需要进行寻道操作。
// 然后置请求项设备为当前驱动器。
432     if (current_drive != CURRENT_DEV)
433         seek = 1;
434     current_drive = CURRENT_DEV;
// 设置读写起始扇区。因为每次读写是以块为单位(1 块 2 个扇区)，所以起始扇区需要起码比
// 磁盘总扇区数小 2 个扇区。否则结束该次软盘请求项，执行下一个请求项。
435     block = CURRENT->sector; // 取当前软盘请求项中起始扇区号→block。
436     if (block+2 > floppy->size) { // 如果 block+2 大于磁盘扇区总数，则
437         end_request(0); // 结束本次软盘请求项。
438         goto repeat;
439     }
// 求对应在磁道上的扇区号，磁头号，磁道号，搜寻磁道号(对于软驱读不同格式的盘)。
440     sector = block % floppy->sect; // 起始扇区对每磁道扇区数取模，得磁道上扇区号。
441     block /= floppy->sect; // 起始扇区对每磁道扇区数取整，得起始磁道数。

```



```

442     head = block % floppy->head;    // 起始磁道数对磁头数取模，得操作的磁头号。
443     track = block / floppy->head;    // 起始磁道数对磁头数取整，得操作的磁道号。
444     seek_track = track << floppy->stretch; // 相应于驱动器中盘类型进行调整，得寻道号。
// 如果寻道号与当前磁头所在磁道不同，则置需要寻道标志 seek。
445     if (seek_track != current_track)
446         seek = 1;
447     sector++;                        // 磁盘上实际扇区计数是从1算起。
448     if (CURRENT->cmd == READ)        // 如果请求项中是读操作，则置软盘读命令码。
449         command = FD_READ;
450     else if (CURRENT->cmd == WRITE)   // 如果请求项中是写操作，则置软盘写命令码。
451         command = FD_WRITE;
452     else
453         panic("do_fd_request: unknown command");
// 添加定时器，用于指定驱动器到能正常运行所需延迟的时间（滴答数），当定时时间到时就调用
// 函数 floppy_on_interrupt(),
454     add_timer(ticks to floppy_on(current drive), &floppy_on_interrupt);
455 }
456
///// 软盘系统初始化。
// 设置软盘块设备的请求处理函数(do_fd_request())，并设置软盘中断门(int 0x26，对应硬件
// 中断请求信号 IRQ6)，然后取消对该中断信号的屏蔽，允许软盘控制器 FDC 发送中断请求信号。
457 void floppy_init(void)
458 {
459     blk_dev[MAJOR_NR].request_fn = DEVICE_REQUEST; // = do_fd_request()。
460     set_trap_gate(0x26, &floppy_interrupt); // 设置软盘中断门 int 0x26(38)。
461     outb(inb_p(0x21)&~0x40, 0x21);           // 复位软盘的中断请求屏蔽位，允许
                                                // 软盘控制器发送中断请求信号。
462 }
463

```

## 5.7.3 其它信息

### 5.7.3.1 软盘驱动器的设备号

在 Linux 中，软驱的主设备号是 2，次设备号 = TYPE\*4 + DRIVE，其中 DRIVE 为 0-3，分别对应软驱 A、B、C 或 D；TYPE 是软驱的类型，2 表示 1.2M 软驱，7 表示 1.44M 软驱，也即 floppy.c 中 85 行定义的软盘类型 (floppy\_type[]) 数组的索引值：

- |   |                       |
|---|-----------------------|
| 0 | 不用；                   |
| 1 | 360kB PC 软驱；          |
| 2 | 1.2MB AT 软驱；          |
| 3 | 360kB 在 720kB 驱动器中使用； |
| 4 | 3.5" 720kB 软盘；        |
| 5 | 360kB 在 1.2MB 驱动器中使用； |
| 6 | 720kB 在 1.2MB 驱动器中使用； |
| 7 | 1.44MB 软驱。            |

例如，因为  $7*4 + 0 = 28$ ，所以 /dev/PS0 (2, 28) 指的是 1.44M A 驱动器，其设备号是 0x021c。

同理 /dev/at0 (2, 8) 指的是 1.2M A 驱动器，其设备号是 0x0208。

### 5.7.3.2 软盘控制器编程方法

对软盘控制器的编程比较烦琐。在编程时需要访问 4 个端口，分别对应一个或多个寄存器。对于 1.2M 的软盘控制器有以下一些端口。

表 软盘控制器端口

| I/O 端口 | 读写性 | 寄存器名称                   |
|--------|-----|-------------------------|
| 0x3f2  | 只写  | 数字输出寄存器 (DOR) (数字控制寄存器) |
| 0x3f4  | 只读  | FDC 主状态寄存器 (STATUS)     |
| 0x3f5  | 读/写 | FDC 数据寄存器 (DATA)        |
| 0x3f7  | 只读  | 数字输入寄存器 (DIR)           |
| 0x3f7  | 只写  | 磁盘控制寄存器 (DCR) (传输率控制)   |

数字输出端口 DOR (数字控制端口) 是一个 8 位寄存器, 它控制驱动器马达开启、驱动器选择、启动/复位 FDC 以及允许/禁止 DMA 及中断请求。

表 数字输出寄存器定义

| D7         | D6         | D5         | D4         | D3       | D2        | D1   | D0 |
|------------|------------|------------|------------|----------|-----------|------|----|
| 启动<br>马达 D | 启动<br>马达 C | 启动<br>马达 B | 启动<br>马达 A | 允许<br>请求 | 启动<br>FDC | 软驱选择 |    |

D7D6D5D4 - 分别控制驱动器 D-A 的马达, 1 启动马达; 0 关闭马达。

D3 - 1 允许 DMA 和中断请求; 0 禁止 DMA 和中断请求;

D2 - 1 启动软驱; 0 复位软驱;

D1D0 - 00-11 用于选择软盘驱动器 A-D;

FDC 的主状态寄存器也是一个 8 位寄存器, 用于反映软盘控制器 FDC 和软盘驱动器 FDD 的基本状态。通常, 在 CPU 向 FDC 发送命令之前或从 FDC 获取操作结果之前, 都要读取主状态寄存器的状态位, 以判别当前 FDC 数据寄存器是否就绪, 以及确定数据传送的方向。

表 FDC 主状态控制器定义

| D7        | D6       | D5          | D4       | D3        | D2        | D1        | D0        |
|-----------|----------|-------------|----------|-----------|-----------|-----------|-----------|
| 数据口<br>就绪 | 传送<br>方向 | 非 DMA<br>方式 | FDC<br>忙 | 软驱 D<br>忙 | 软驱 C<br>忙 | 软驱 B<br>忙 | 软驱 A<br>忙 |
| RQM       | DIO      | NDM         | CB       | DDB       | DCB       | DBB       | DAB       |

D7 - 1 表示 FDC 数据寄存器已准备就绪;

D6 - 1 表示数据从 FDC 到 CPU; 0 表示数据从 CPU 到 FDC;

D5 - 1 表示 FDC 工作在非 DMA 方式;

D4 - 1 表示 FDC 正处于命令执行忙碌状态;

D3D2D1D0 - 分别代表驱动器 D-A 忙碌状态。

FDC 的数据端口对应多个寄存器 (只写型命令寄存器和参数寄存器、只读型结果寄存器), 但任一时刻只能有一个寄存器出现在数据端口 0x3f5。在访问只写型寄存器时, 主状态控制的 DIO 方向位必须为 0 (CPU → FDC), 访问只读型寄存器时则反之。在读取结果时只有在 FDC 不忙之后才算读完结果, 通常结果数据最多有 7 个字节。

数据输入寄存器 (DIR) 只有位 7 (D7) 对软盘有效, 用来表示盘片更换状态。其余七位用于硬盘控制器接口。



磁盘控制寄存器(DCR)用于选择盘片在不同类型驱动器上使用的数据传输率。仅使用低 2 位 (D1D0), 00 - 500kbps, 01 - 300kbps, 10 - 250kbps。

软盘控制器共可以接受 15 条命令。每个命令均经历三个阶段：命令阶段、执行阶段和结果阶段。

命令阶段是 CPU 向 FDC 发送命令字节和参数字节。每条命令的第一个字节总是命令字节（命令码）。其后跟着 0—8 字节的参数。执行阶段是 FDC 执行命令规定的操作。在执行阶段 CPU 是不加干预的，一般是通过 FDC 发出中断请求获知命令执行的结束。如果 CPU 发出的 FDC 命令是传送数据，则 FDC 可以以中断方式或 DMA 方式进行。中断方式每次传送 1 字节。DMA 方式是在 DMA 控制器管理下，FDC 与内存进行数据的传输直至全部数据传送完。此时 DMA 控制器会将传输字节计数终止信号通知 FDC，最后由 FDC 发出中断请求信号告知 CPU 执行阶段结束。结果阶段是由 CPU 读取 FDC 数据寄存器返回值，从而获得 FDC 命令执行的结果。返回结果数据的长度为 0—7 字节。对于没有返回结果数据的命令，则应向 FDC 发送检测中断状态命令获得操作的状态。








### 5.7.3.3 DMA 控制器编程



## 第 6 章 字符设备驱动程序(char driver)

### 6.1 概述

列表 linux/kernel/chr\_drv 目录

| Name  | Size        | Last modified (GMT) | Description |
|---|-------------|---------------------|-------------|
|  <a href="#">Makefile</a>      | 2443 bytes  | 1991-12-02 03:21:41 | m           |
|  <a href="#">console.c</a>     | 14568 bytes | 1991-11-23 18:41:21 | m           |
|  <a href="#">keyboard.S</a>    | 12780 bytes | 1991-12-04 15:07:58 | m           |
|  <a href="#">rs_io.s</a>       | 2718 bytes  | 1991-10-02 14:16:30 | m           |
|  <a href="#">serial.c</a>    | 1406 bytes  | 1991-11-17 21:49:05 | m           |
|  <a href="#">tty_io.c</a>    | 7634 bytes  | 1991-12-08 18:09:15 | m           |
|  <a href="#">tty_ioctl.c</a> | 4979 bytes  | 1991-11-25 19:59:38 | m           |

### 6.2 Makefile 文件

#### 6.2.1 功能描述

#### 6.2.2 代码注释

列表 linux/kernel/chr\_drv/Makefile 文件

```

1 #
2 # Makefile for the FREAX-kernel character device drivers.
3 #
4 # Note! Dependencies are done automagically by 'make dep', which also
5 # removes any old dependencies. DON'T put your own dependencies here
6 # unless it's something special (ie not a .c file).
7 #

```

# FREAX(Linux)内核字符设备驱动程序的 Makefile 文件。  
 # 注意！依赖关系是由'make dep'自动进行的，它也会自动去除原来的依赖信息。不要把你自己的  
 # 依赖关系信息放在这里，除非是特别文件的（也即不是一个.c文件的信息）。

```

8
9 AR      =gar      # GNU 的二进制文件处理程序，用于创建、修改以及从归档文件中抽取文件。
10 AS      =gas      # GNU 的汇编程序。
11 LD      =gld      # GNU 的连接程序。
12 LDFLAGS =-s -x    # 连接程序所有的参数，-s 输出文件中省略所有符号信息。-x 删除所有局部符号。
13 CC      =gcc      # GNU C 语言编译器。
    # 下一行是 C 编译程序选项。-Wall 显示所有的警告信息；-O 优化选项，优化代码长度和执行时间；
    # -fstrength-reduce 优化循环执行代码，排除重复变量；-fomit-frame-pointer 省略保存不必要
    # 的框架指针；-fcombine-regs 合并寄存器，减少寄存器类的使用；-finline-functions 将所有简
    # 单短小的函数代码嵌入调用程序中；-mstring-insns Linus 自己添加的优化选项，以后不再使用；
    # -nostdinc -I../include 不使用默认路径中的包含文件，而使用指定目录中的(../../include)。
14 CFLAGS  =-Wall -O -fstrength-reduce -fomit-frame-pointer -fcombine-regs \
15         -finline-functions -mstring-insns -nostdinc -I../include
    # C 前处理选项。-E 只运行 C 前处理，对所有指定的 C 程序进行预处理并将处理结果输出到标准输
    # 出设备或指定的输出文件中；-nostdinc -I../include 同前。
16 CPP     =gcc -E -nostdinc -I../include
17
    # 下面的规则指示 make 利用下面的命令将所有的.c 文件编译生成.s 汇编程序。该规则的命令
    # 指使 gcc 采用 CFLAGS 所指定的选项对 C 代码编译后不进行汇编就停止（-S），从而产生与
    # 输入的各个 C 文件对应的汇编代码文件。默认情况下所产生的汇编程序文件名是原 C 文件名
    # 去掉.c 而加上.s 后缀。-o 表示其后是输出文件的名称。其中$.s（或$@）是自动目标变量，
    # $<代表第一个先决条件，这里即是符合条件*.c 的文件。
18 .c.s:
19     $(CC) $(CFLAGS) \
20     -S -o $.s $<
    # 下面规则表示将所有.s 汇编程序文件编译成.o 目标文件。22 行是实现该操作的具体命令。
21 .s.o:
22     $(AS) -c -o $.o $<
23 .c.o:                                # 类似上面，*.c 文件→*.o 目标文件。不进行连接。
24     $(CC) $(CFLAGS) \
25     -c -o $.o $<
26
27 OBJS = tty_io.o console.o keyboard.o serial.o rs_io.o \    # 定义目标文件变量 OBJS。
28        tty_ioctl.o
29
30 chr_drv.a: $(OBJS)          # 在有了先决条件 OBJS 后使用下面的命令连接成目标 chr_drv.a 库文件。
31     $(AR) rcs chr_drv.a $(OBJS)
32     sync
33
    # 对 keyboard.S 汇编程序进行预处理。-traditional 选项用来对程序作修改使其支持传统的 C 编译器。
    # 处理后的程序改名为 kernboard.s。
34 keyboard.s: keyboard.S ../include/linux/config.h
35     $(CPP) -traditional keyboard.S -o keyboard.s
36
    # 下面的规则用于清理工作。当执行'make clean'时，就会执行下面的命令，去除所有编译
    # 连接生成的文件。'rm'是文件删除命令，选项-f 含义是忽略不存在的文件，并且不显示删除信息。
37 clean:
38     rm -f core *.o *.a tmp_make keyboard.s
39     for i in *.c;do rm -f `basename $$i .c`.s;done
40

```

# 下面得目标或规则用于检查各文件之间的依赖关系。方法如下：  
 # 使用字符串编辑程序 sed 对 Makefile 文件（即是本文件）进行处理，输出为删除 Makefile  
 # 文件中'### Dependencies'行后面的所有行（下面从 48 开始的行），并生成 tmp\_make  
 # 临时文件（44 行的作用）。然后对 kernel/chr\_drv/目录下的每个 C 文件执行 gcc 预处理操作。  
 # -M 标志告诉预处理程序输出描述每个目标文件相关性的规则，并且这些规则符合 make 语法。  
 # 对于每一个源文件，预处理程序输出一个 make 规则，其结果形式是相应源程序文件的目标  
 # 文件名加上其依赖关系—该源文件中包含的所有头文件列表。把预处理结果都添加到临时  
 # 文件 tmp\_make 中，然后将该临时文件复制成新的 Makefile 文件。

```

41 dep:
42     sed '/\#\#\# Dependencies/q' < Makefile > tmp_make
43     (for i in *.c;do echo -n `echo $$i | sed 's,\.c,\.s,`" "; \
44         $(CPP) -M $$i;done) >> tmp_make
45     cp tmp_make Makefile
46
47 ### Dependencies:
48 console.s console.o : console.c ../../include/linux/sched.h \
49     ../../include/linux/head.h ../../include/linux/fs.h \
50     ../../include/sys/types.h ../../include/linux/mm.h ../../include/signal.h \
51     ../../include/linux/tty.h ../../include/termios.h ../../include/asm/io.h \
52     ../../include/asm/system.h
53 serial.s serial.o : serial.c ../../include/linux/tty.h ../../include/termios.h \
54     ../../include/linux/sched.h ../../include/linux/head.h \
55     ../../include/linux/fs.h ../../include/sys/types.h ../../include/linux/mm.h \
56     ../../include/signal.h ../../include/asm/system.h ../../include/asm/io.h
57 tty_io.s tty_io.o : tty_io.c ../../include/ctype.h ../../include/errno.h \
58     ../../include/signal.h ../../include/sys/types.h \
59     ../../include/linux/sched.h ../../include/linux/head.h \
60     ../../include/linux/fs.h ../../include/linux/mm.h ../../include/linux/tty.h \
61     ../../include/termios.h ../../include/asm/segment.h \
62     ../../include/asm/system.h
63 tty_ioctl.s tty_ioctl.o : tty_ioctl.c ../../include/errno.h ../../include/termios.h \
64     ../../include/linux/sched.h ../../include/linux/head.h \
65     ../../include/linux/fs.h ../../include/sys/types.h ../../include/linux/mm.h \
66     ../../include/signal.h ../../include/linux/kernel.h \
67     ../../include/linux/tty.h ../../include/asm/io.h \
68     ../../include/asm/segment.h ../../include/asm/system.h

```

## 6.3 serial.c 文件

### 6.3.1 功能描述

### 6.3.2 代码注释

列表 linux/kernel/chr\_drv/serial.c 程序

```

1 /*
2  * linux/kernel/serial.c
3  *
4  * (C) 1991 Linus Torvalds

```

```

5  */
6
7  /*
8  *      serial.c
9  *
10 * This module implements the rs232 io functions
11 *      void rs_write(struct tty_struct * queue);
12 *      void rs_init(void);
13 * and all interrupts pertaining to serial IO.
14 */
15
16 /*
17 *      serial.c
18 * 该程序用于实现 rs232 的输入输出功能
19 *      void rs_write(struct tty_struct *queue);
20 *      void rs_init(void);
21 * 以及与传输 IO 有关系的所有中断处理程序。
22 */
23
24 #include <linux/tty.h>    // tty 头文件, 定义了有关 tty_io, 串行通信方面的参数、常数。
25 #include <linux/sched.h> // 调度程序头文件, 定义了任务结构 task_struct、初始任务 0 的数据,
26 // 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
27 #include <asm/system.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。
28 #include <asm/io.h>       // io 头文件。定义硬件端口输入/输出宏汇编语句。
29
30 #define WAKEUP_CHARS (TTY_BUF_SIZE/4) // 当写队列中含有 WAKEUP_CHARS 个字符时, 就开始发送。
31
32 extern void rs1_interrupt(void);        // 串行口 1 的中断处理程序(rs_io.s, 34)。
33 extern void rs2_interrupt(void);        // 串行口 2 的中断处理程序(rs_io.s, 38)。
34
35 // 初始化串行端口
36 // port: 串口 1 - 0x3F8, 串口 2 - 0x2F8。
37 static void init(int port)
38 {
39     outb_p(0x80, port+3);    /* set DLAB of line control reg */
40                               /* 设置线路控制寄存器的 DLAB 位(位 7) */
41     outb_p(0x30, port);      /* LS of divisor (48 -> 2400 bps) */
42                               /* 发送波特率因子低字节, 0x30->2400bps */
43     outb_p(0x00, port+1);     /* MS of divisor */
44                               /* 发送波特率因子高字节, 0x00 */
45     outb_p(0x03, port+3);     /* reset DLAB */
46                               /* 复位 DLAB 位, 数据位为 8 位 */
47     outb_p(0x0b, port+4);     /* set DTR, RTS, OUT_2 */
48                               /* 设置 DTR, RTS, 辅助用户输出 2 */
49     outb_p(0x0d, port+1);     /* enable all intrs but writes */
50                               /* 除了写(写保持空)以外, 允许所有中断源中断 */
51     (void)inb(port);          /* read data port to reset things (?) */
52                               /* 读数据口, 以进行复位操作(?) */
53 }
54
55 // 初始化串行中断程序和串行接口。
56 void rs_init(void)
57 {
58     set_intr_gate(0x24, rs1_interrupt); // 设置串行口 1 的中断门向量(硬件 IRQ4 信号)。

```

```

40     set_intr_gate(0x23, rs2_interrupt); // 设置串行口 2 的中断门向量(硬件 IRQ3 信号)。
41     init(tty_table[1].read_q.data);    // 初始化串行口 1(.data 是端口号)。
42     init(tty_table[2].read_q.data);    // 初始化串行口 2。
43     outb(inb_p(0x21)&0xE7, 0x21);      // 允许主 8259A 芯片的 IRQ3, IRQ4 中断信号请求。
44 }
45
46 /*
47  * This routine gets called when tty_write has put something into
48  * the write_queue. It must check wheter the queue is empty, and
49  * set the interrupt register accordingly
50  *
51  * void _rs_write(struct tty_struct * tty);
52  */
53
54 /*
55  * 在 tty_write() 已将数据放入输出(写)队列时会调用下面的子程序。必须首先
56  * 检查写队列是否为空, 并相应设置中断寄存器。
57  */
58
59 // 串行数据发送输出。
60 // 实际上只是开启串行发送保持寄存器已空中断标志, 在 UART 将数据发送出去后允许发中断信号。
61 void rs_write(struct tty_struct * tty)
62 {
63     cli(); // 关中断。
64     // 如果写队列不空, 则从 0x3f9(或 0x2f9) 首先读取中断允许寄存器内容, 添上发送保持寄存器
65     // 中断允许标志(位 1)后, 再写回该寄存器。
66     if (!EMPTY(tty->write_q))
67         outb(inb_p(tty->write_q.data+1) | 0x02, tty->write_q.data+1);
68     sti(); // 开中断。
69 }
70

```

## 6.3.3 其它信息

### 6.3.3.1 异步串行通信芯片 UART

PC 微机的串行通信使用的异步串行通信芯片是 INS 8250 或 NS16450 兼容芯片, 统称为 UART(通用异步接收发送器)。对 UART 的编程实际上是对其内部寄存器执行读写操作。因此可将 UART 看作是一组寄存器集合, 包含发送、接收和控制三部分。UART 内部有 10 个寄存器, 供 CPU 通过 IN/OUT 指令对其进行访问。这些寄存器的端口和用途见下表所示。其中端口 0x3f8-0x3fe 用于微机上 COM1 串行口, 0x2f8-0x2fe 对应 COM2 端口。条件 DLAB(Divisor Latch Access Bit)是除数锁存访问位, 是指线路控制寄存器的位 7。

表 UART 内部寄存器对应端口及用途

| 端口            | 读/写 | 条件     | 用途  |
|---------------|-----|--------|---|
| 0x3f8 (0x2f8) | 写   | DLAB=0 | 写发送保持寄存器。含有将发送的字符。  |
|               | 读   | DLAB=0 | 读接收缓存寄存器。含有收到的字符。   |
|               | 读/写 | DLAB=1 | 读/写波特率因子低字节 (LSB)。  |
| 0x3f9 (0x2f9) | 读/写 | DLAB=1 | 读/写波特率因子高字节 (MSB)。  |
|               | 读/写 | DLAB=0 | 读/写中断允许寄存器。<br>位 7-4 全 0 保留不用;<br>位 3=1 modem 状态中断允许;<br>位 2=1 接收器线路状态中断允许; |

|               |   |  |  |
|---------------|---|--|--|
|               |   |  | 位 1=1 发送保持寄存器空中断允许;<br>位 0=1 已接收到数据中断允许。   |
| 0x3fa (0x2fa) | 读 |  | 读中断标识寄存器。中断处理程序用以判断此次中断是 4 种中的那一种。<br>位 7-3 全 0 (不用);<br>位 2-1 确定中断的优先级;<br>= 11 接收状态有错中断, 优先级最高;<br>= 10 已接收到数据中断, 优先级第 2;<br>= 01 发送保持寄存器空中断, 优先级第 3;<br>= 00 modem 状态改变中断, 优先级第 4。<br>位 0=0 有待处理中断; =1 无中断。   |
| 0x3fb (0x2fb) | 写 |  | 写线路控制寄存器。<br>位 7=1 除数锁存访问位 (DLAB)。<br>0 接收器, 发送保持或中断允许寄存器访问;<br>位 6=1 允许间断;<br>位 5=1 保持奇偶位;<br>位 4=1 偶校验; =0 奇校验;<br>位 3=1 允许奇偶校验; =0 无奇偶校验;<br>位 2=1 1 位停止位; =0 无停止位;<br>位 1-0 数据位长度:<br>= 00 5 位数据位;<br>= 01 6 位数据位;<br>= 10 7 位数据位;<br>= 11 8 位数据位。 |
| 0x3fc (0x2fc) | 写 |  | 写 modem 控制寄存器。<br>位 7-5 全 0 保留;<br>位 4=1 芯片处于循环反馈诊断操作模式;<br>位 3=1 辅助用户指定输出 2, 允许 INTRPT 到系统;<br>位 2=1 辅助用户指定输出 1, PC 机未用;<br>位 1=1 使请求发送 RTS 有效;<br>位 0=1 使数据终端就绪 DTR 有效。  |
| 0x3fd (0x2fd) | 读 |  | 读线路状态寄存器。<br>位 7=0 保留;<br>位 6=1 发送移位寄存器为空;<br>位 5=1 发送保持寄存器为空, 可以取字符发送;<br>位 4=1 接收到满足间断条件的位序列;<br>位 3=1 帧格式错误;<br>位 2=1 奇偶校验错误;<br>位 1=1 超越覆盖错误;<br>位 0=1 接收器数据准备好, 系统可读取。  |
| 0x3fe (0x2fe) | 读 |  | 读 modem 状态寄存器。δ 表示信号发生变化。<br>位 7=1 载波检测 (CD) 有效;<br>位 6=1 响铃指示 (RI) 有效;<br>位 5=1 数据设备就绪 (DSR) 有效;  |



|  |  |  |  |
|--|--|--|--|
|  |  |  | 位 4=1 清除发送 (CTS) 有效;<br>位 3=1 检测到 $\delta$ 载波;<br>位 2=1 检测到响铃信号边沿;<br>位 1=1 $\delta$ 数据设备就绪 (DSR);<br>位 0=1 $\delta$ 清除发送 (CTS)。 |
|--|--|--|--|

## 6.4 rs\_io.s 文件

### 6.4.1 功能描述

该汇编程序实现 rs232 串行通信中断处理程序。

### 6.4.2 代码注释

列表 linux/kernel/chr\_drv/rs\_io.s 程序

```
1 /*
2  * linux/kernel/rs_io.s
3  *
4  * (C) 1991 Linus Torvalds
5  */
6
7 /*
8  *      rs_io.s
9  *
10 * This module implements the rs232 io interrupts.
11 */
12 /*
13  * 该程序模块实现 rs232 输入输出中断处理程序。
14 */
15
16 // size 是读写队列缓冲区的字节长度。
17 size = 1024                                /* must be power of two !    必须是 2 的次方并且需
18                                         and must match the value 与 tty_io.c 中的值匹配!
19                                         in tty_io.c!!! */
20
21 /* these are the offsets into the read/write buffer structures */
22 /* 以下这些是读写缓冲结构中的偏移量 */
23 rs_addr = 0                                // 串行端口号字段偏移。
24 head = 4                                  // 缓冲区中头指针字段偏移。
25 tail = 8                                  // 缓冲区中尾指针字段偏移。
26 proc_list = 12                             // 等待该缓冲的进程字段偏移。
27 buf = 16                                  // 缓冲区字段偏移。
28
29 startup = 256                             /* chars left in write queue when we restart it */
30                                         /* 当写队列里还剩 256 个字符空间 (WAKEUP_CHARS) 时，我们就可以写 */
```

```

28
29 /*
30 * These are the actual interrupt routines. They look where
31 * the interrupt is coming from, and take appropriate action.
32 */
33 /*
34 * 这些是实际的中断程序。程序首先检查中断的来源，然后执行相应
35 * 的处理。
36 */
37 .align 2
38 // 串行端口 1 中断处理程序入口点。
39 _rs1_interrupt:
40     pushl $_table_list+8 // tty 表中对应串口 1 的读写缓冲指针的地址入栈(tty_io.c, 99)。
41     jmp rs_int
42 .align 2
43 // 串行端口 2 中断处理程序入口点。
44 _rs2_interrupt:
45     pushl $_table_list+16 // tty 表中对应串口 2 的读写缓冲队列指针的地址入栈。
46 rs_int:
47     pushl %edx
48     pushl %ecx
49     pushl %ebx
50     pushl %eax
51     push %es
52     push %ds // as this is an interrupt, we cannot */
53     pushl $0x10 // know that bs is ok. Load it */
54     pop %ds // 由于这是一个中断程序，我们不知道 ds 是否正确，*/
55     pushl $0x10 // 所以加载它们(让 ds、es 指向内核数据段 */
56     pop %es
57     movl 24(%esp), %edx // 将缓冲队列指针地址存入 edx 寄存器，
58                        // 也即 35 或 39 行上最先压入堆栈的地址。
59     movl (%edx), %edx // 取读队列指针(地址)→edx。
60     movl rs_addr(%edx), %edx // 取串口 1 的端口号→edx。
61     addl $2, %edx // interrupt ident. reg */ /* edx 指向中断标识寄存器 */
62 rep_int:
63     xorl %eax, %eax // eax 清零。
64     inb %dx, %al // 取中断标识字节，用以判断中断来源(有 4 种中断情况)。
65     testb $1, %al // 首先判断有无待处理的中断(位 0=1 无中断; =0 有中断)。
66     jne end // 若无待处理中断，则跳转至退出处理处 end。
67     cmpb $6, %al // this shouldn't happen, but ... */ /* 这不会发生，但是...*/
68     ja end // al 值>6? 是则跳转至 end (没有这种状态)。
69     movl 24(%esp), %ecx // 再取缓冲队列指针地址→ecx。
70     pushl %edx // 将端口号 0x3fa(0x2fa)入栈。
71     subl $2, %edx // 0x3f8(0x2f8)。
72     call jmp_table(, %eax, 2) /* NOTE! not *4, bit0 is 0 already */ /* 不乘 4，位 0 已是 0*/
73 // 上面语句是指，当有待处理中断时，al 中位 0=0，位 2-1 是中断类型，因此相当于已经将中断类型
74 // 乘了 2，这里再乘 2，得到跳转表对应各中断类型地址，并跳转到那里去作相应处理。
75     popl %edx // 弹出中断标识寄存器端口号 0x3fa (或 0x2fa)。
76     jmp rep_int // 跳转，继续判断有无待处理中断并继续处理。
77 end:     movb $0x20, %al // 向中断控制器发送结束中断指令 EOI。
78     outb %al, $0x20 // EOI */
79     pop %ds
80     pop %es

```

```

72     popl %eax
73     popl %ebx
74     popl %ecx
75     popl %edx
76     addl $4,%esp          # jump over _table_list entry # 丢弃缓冲队列指针地址。
77     iret
78
// 各中断类型处理程序地址跳转表，共有 4 种中断来源：
// modem 状态变化中断，写字符中断，读字符中断，线路状态有问题中断。
79 jmp_table:
80     .long modem_status,write_char,read_char,line_status
81
82 .align 2
83 modem_status:
84     addl $6,%edx          /* clear intr by reading modem status reg */
85     inb %dx,%al          /* 通过读 modem 状态寄存器进行复位(0x3fe) */
86     ret
87
88 .align 2
89 line_status:
90     addl $5,%edx          /* clear intr by reading line status reg. */
91     inb %dx,%al          /* 通过读线路状态寄存器进行复位(0x3fd) */
92     ret
93
94 .align 2
95 read_char:
96     inb %dx,%al          /* 读取字符→al。
97     movl %ecx,%edx        /* 当前串口缓冲队列指针地址→edx。
98     subl $_table_list,%edx // 缓冲队列指针表首址 - 当前串口队列指针地址→edx，
99     shr1 $3,%edx         // 差值/8。对于串口 1 是 1，对于串口 2 是 2。
100    movl (%ecx),%ecx      # read-queue # 取读缓冲队列结构地址→ecx。
101    movl head(%ecx),%ebx  // 取读队列中缓冲头指针→ebx。
102    movb %al,buf(%ecx,%ebx) // 将字符放在缓冲区中头指针所指的位置。
103    incl %ebx            // 将头指针前移一字节。
104    andl $size-1,%ebx    // 用缓冲区大小对头指针进行模操作。指针不能超过缓冲区大小。
105    cmpl tail(%ecx),%ebx // 缓冲区头指针与尾指针比较。
106    je 1f               // 若相等，表示缓冲区满，跳转到标号 1 处。
107    movl %ebx,head(%ecx) // 保存修改过的头指针。
108 1:    pushl %edx         // 将串口号压入堆栈(1- 串口 1, 2 - 串口 2)，作为参数，
109    call _do_tty_interrupt // 调用 tty 中断处理 C 函数 (。
110    addl $4,%esp         // 丢弃入栈参数，并返回。
111    ret
112
113 .align 2
114 write_char:
115     movl 4(%ecx),%ecx    # write-queue # 取写缓冲队列结构地址→ecx。
116     movl head(%ecx),%ebx // 取写队列头指针→ebx。
117     subl tail(%ecx),%ebx // 头指针 - 尾指针 = 队列中字符数。
118     andl $size-1,%ebx   # nr chars in queue # 对指针取模运算。
119     je write_buffer_empty // 如果头指针 = 尾指针，说明写队列无字符，跳转处理。
120     cmpl $startup,%ebx  // 队列中字符数超过 256 个？
121     ja 1f               // 超过，则跳转处理。
122     movl proc_list(%ecx),%ebx # wake up sleeping process # 唤醒等待的进程。

```

```

123      testl %ebx,%ebx                // 取等待该队列的进程的指针，并判断是否为空。
124      je 1f                          # is there any? # 有等待的进程吗？
125      movl $0, (%ebx)                // 是空的，则向前跳转到标号 1 处。
126 1:      movl tail(%ecx), %ebx        // 否则将进程置为可运行状态(唤醒进程)。。
127      movb buf(%ecx, %ebx), %al      // 取尾指针。
128      outb %al, %dx                  // 从缓冲中尾指针处取一字符→al。
129      incl %ebx                      // 向端口 0x3f8(0x2f8)送出到保持寄存器中。
130      andl $size-1, %ebx             // 尾指针前移。
131      movl %ebx, tail(%ecx)          // 尾指针若到缓冲区末端，则折回。
132      cmpl head(%ecx), %ebx          // 保存已修改过的尾指针。
133      je write_buffer_empty          // 尾指针与头指针比较，
134      ret                            // 若相等，表示队列已空，则跳转。
135      .align 2
136 write_buffer_empty:
137      movl proc_list(%ecx), %ebx     # wake up sleeping process # 唤醒等待的进程。
138      testl %ebx,%ebx                // 取等待该队列的进程的指针，并判断是否为空。
139      je 1f                          # is there any? # 有等待的进程吗？
140      movl $0, (%ebx)                # 无，则向前跳转到标号 1 处。
141 1:      incl %edx                    # 否则将进程置为可运行状态(唤醒进程)。
142      inb %dx, %al                   # 指向端口 0x3f9(0x2f9)。
143      jmp 1f                          # 读取中断允许寄存器。
144 1:      jmp 1f                      # 稍作延迟。
145 1:      andb $0xd, %al              /* disable transmit interrupt */
146      outb %al, %dx                  /* 屏蔽发送保持寄存器空中断（位 1） */
147      ret                            // 写入 0x3f9(0x2f9)。

```

## 6.5 tty\_io.c 文件

### 6.5.1 功能描述

### 6.5.2 代码注释

列表 linux/kernel/chr\_drv/tty\_io.c 程序

```

1  /*
2  *  linux/kernel/tty_io.c
3  *
4  *  (C) 1991  Linus Torvalds
5  */
6
7  /*
8  *  'tty_io.c' gives an orthogonal feeling to tty's, be they consoles
9  *  or rs-channels. It also implements echoing, cooked mode etc.
10 *
11 *  Kill-line thanks to John T Kohl.
12 */

```

```

/*
 * 'tty_io.c' 给 tty 一种非相关的感觉，是控制台还是串行通道。该程序同样
 * 实现了回显、规范(熟)模式等。
 *
 * Kill-line, 谢谢 John T Kahl.
 */
13 #include <ctype.h>           // 字符类型头文件。定义了一些有关字符类型判断和转换的宏。
14 #include <errno.h>          // 错误号头文件。包含系统中各种出错号。(Linux 从 minix 中引进的)。
15 #include <signal.h>         // 信号头文件。定义信号符号常量，信号结构以及信号操作函数原型。
16
// 下面给出相应信号在信号位图中的对应比特位。
17 #define ALRMASK (1<<(SIGALRM-1)) // 警告(alarm)信号屏蔽位。
18 #define KILLMASK (1<<(SIGKILL-1)) // 终止(kill)信号屏蔽位。
19 #define INTMASK (1<<(SIGINT-1)) // 键盘中断(int)信号屏蔽位。
20 #define QUITMASK (1<<(SIGQUIT-1)) // 键盘退出(quit)信号屏蔽位。
21 #define TSTPMASK (1<<(SIGTSTP-1)) // tty 发出的停止进程(tty stop)信号屏蔽位。
22
23 #include <linux/sched.h> // 调度程序头文件，定义了任务结构 task_struct、初始任务 0 的数据，
// 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
24 #include <linux/tty.h> // tty 头文件，定义了有关 tty_io，串行通信方面的参数、常数。
25 #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。
26 #include <asm/system.h> // 系统头文件。定义了设置或修改描述符/中断门等的嵌入式汇编宏。
27
28 #define L_FLAG(tty, f) ((tty)->termios.c_lflag & f) // 取 termios 结构中的本地模式标志。
29 #define I_FLAG(tty, f) ((tty)->termios.c_iflag & f) // 取 termios 结构中的输入模式标志。
30 #define O_FLAG(tty, f) ((tty)->termios.c_oflag & f) // 取 termios 结构中的输出模式标志。
31
// 取 termios 结构中本地模式标志集中的一个标志位。
32 #define L_CANON(tty) L_FLAG((tty), ICANON) // 取本地模式标志集中规范(熟)模式标志位。
33 #define L_ISIG(tty) L_FLAG((tty), ISIG) // 取信号标志位。
34 #define L_ECHO(tty) L_FLAG((tty), ECHO) // 取回显字符标志位。
35 #define L_ECHOE(tty) L_FLAG((tty), ECHOE) // 规范模式时，取回显擦出标志位。
36 #define L_ECHOK(tty) L_FLAG((tty), ECHOK) // 规范模式时，取 KILL 擦除当前行标志位。
37 #define L_ECHOCTL(tty) L_FLAG((tty), ECHOCTL) // 取回显控制字符标志位。
38 #define L_ECHOKL(tty) L_FLAG((tty), ECHOKL) // 规范模式时，取 KILL 擦除行并回显标志位。
39
// 取 termios 结构中输入模式标志集中的一个标志位。
40 #define I_UCLC(tty) I_FLAG((tty), IUCLC) // 取输入模式标志集中大写到低转换标志位。
41 #define I_NLCR(tty) I_FLAG((tty), INLCR) // 取换行符 NL 转回车符 CR 标志位。
42 #define I_CRNL(tty) I_FLAG((tty), ICRNL) // 取回车符 CR 转换行符 NL 标志位。
43 #define I_NOCR(tty) I_FLAG((tty), IGNCR) // 取忽略回车符 CR 标志位。
44
// 取 termios 结构中输出模式标志集中的一个标志位。
45 #define O_POST(tty) O_FLAG((tty), OPOST) // 取输出模式标志集中执行输出处理标志。
46 #define O_NLCR(tty) O_FLAG((tty), ONLCR) // 取换行符 NL 转回车换行符 CR-NL 标志。
47 #define O_CRNL(tty) O_FLAG((tty), OCRNL) // 取回车符 CR 转换行符 NL 标志。
48 #define O_NLRET(tty) O_FLAG((tty), ONLRET) // 取换行符 NL 执行回车功能的标志。
49 #define O_LCUC(tty) O_FLAG((tty), OLCUC) // 取小写转大写字母标志。
50
// tty 数据结构的 tty_table 数组。其中包含三个初始化项数据，分别对应控制台、串口终端 1 和
// 串口终端 2 的初始化数据。
51 struct tty_struct tty_table[] = {
52     {

```

```

53      {ICRNL,          /* change incoming CR to NL */ /* 将输入的 CR 转换为 NL */
54      OPOST|ONLCR,    /* change outgoing NL to CRNL */ /* 将输出的 NL 转 CRNL */
55      0,               // 控制模式标志初始化为 0。
56      ISIG | ICANON | ECHO | ECHOCTL | ECHOKE, // 本地模式标志。
57      0,               /* console termio */ // 控制台 termio, 不用线路规程。
58      INIT_C_CC},      // 控制字符数组。
59      0,               /* initial pgrp */ // 所属初始进程组。
60      0,               /* initial stopped */ // 初始停止标志。
61      con_write,       // tty 写函数指针。
62      {0, 0, 0, 0, ""}, /* console read-queue */ // tty 控制台读队列。
63      {0, 0, 0, 0, ""}, /* console write-queue */ // tty 控制台写队列。
64      {0, 0, 0, 0, ""} /* console secondary queue */ // tty 控制台辅助(第二)队列。
65  }, {
66      {0, /* no translation */ // 输入模式标志。0, 无须转换。
67      0, /* no translation */ // 输出模式标志。0, 无须转换。
68      B2400 | CS8,          // 控制模式标志。波特率 2400bps, 8 位数据位。
69      0,                    // 本地模式标志 0。
70      0,                    // 线路规程 0。
71      INIT_C_CC},          // 控制字符数组。
72      0,                    // 所属初始进程组。
73      0,                    // 初始停止标志。
74      rs_write,             // 串口 1 tty 写函数指针。
75      {0x3f8, 0, 0, 0, ""}, /* rs 1 */ // 串行终端 1 读缓冲队列。
76      {0x3f8, 0, 0, 0, ""}, // 串行终端 1 写缓冲队列。
77      {0, 0, 0, 0, ""}     // 串行终端 1 辅助缓冲队列。
78  }, {
79      {0, /* no translation */ // 输入模式标志。0, 无须转换。
80      0, /* no translation */ // 输出模式标志。0, 无须转换。
81      B2400 | CS8,          // 控制模式标志。波特率 2400bps, 8 位数据位。
82      0,                    // 本地模式标志 0。
83      0,                    // 线路规程 0。
84      INIT_C_CC},          // 控制字符数组。
85      0,                    // 所属初始进程组。
86      0,                    // 初始停止标志。
87      rs_write,             // 串口 2 tty 写函数指针。
88      {0x2f8, 0, 0, 0, ""}, /* rs 2 */ // 串行终端 2 读缓冲队列。
89      {0x2f8, 0, 0, 0, ""}, // 串行终端 2 写缓冲队列。
90      {0, 0, 0, 0, ""}     // 串行终端 2 辅助缓冲队列。
91  }
92 };
93
94 /*
95  * these are the tables used by the machine code handlers.
96  * you can implement pseudo-tty's or something by changing
97  * them. Currently not done.
98  */
99 /*
100  * 下面是汇编程序使用的缓冲队列地址表。通过修改你可以实现
101  * 伪 tty 终端或其它终端类型。目前还没有这样做。
102  */
103 // tty 缓冲队列地址表。rs_io.s 汇编程序使用, 用于取得读写缓冲队列地址。
104 struct tty_queue * table_list[]={
105     &tty_table[0].read_q, &tty_table[0].write_q, // 控制台终端读、写缓冲队列地址。

```

```

101     &tty_table[1].read_q, &tty_table[1].write_q, // 串行口 1 终端读、写缓冲队列地址。
102     &tty_table[2].read_q, &tty_table[2].write_q // 串行口 2 终端读、写缓冲队列地址。
103 };
104
105     ///// tty 终端初始化函数。
106     // 初始化串口终端和控制台终端。
107 void tty_init(void)
108 {
109     rs_init(); // 初始化串行中断程序和串行接口 1 和 2。(serial.c, 37)
110     con_init(); // 初始化控制台终端。(console.c, 617)
111 }
112
113     ///// tty 键盘终端字符处理函数。
114     // 参数: tty - 相应 tty 终端结构指针; mask - 信号屏蔽位。
115 void tty_intr(struct tty_struct * tty, int mask)
116 {
117     int i;
118
119     // 如果 tty 所属组号小于等于 0, 则退出。
120     if (tty->pgrp <= 0)
121         return;
122     // 扫描任务数组, 向 tty 相应组的所有任务发送指定的信号。
123     for (i=0; i<NR_TASKS; i++)
124     // 如果该项任务指针不为空, 并且其组号等于 tty 组号, 则设置该任务指定的信号 mask。
125     if (task[i] && task[i]->pgrp==tty->pgrp)
126         task[i]->signal |= mask;
127 }
128
129     ///// 如果队列缓冲区空则让进程进入可中断的睡眠状态。
130     // 参数: queue - 指定队列的指针。
131     // 进程在取队列缓冲区中字符时调用此函数。
132 static void sleep_if_empty(struct tty_queue * queue)
133 {
134     cli(); // 关中断。
135     // 若当前进程没有信号要处理并且指定的队列缓冲区空, 则让进程进入可中断睡眠状态, 并让
136     // 队列的进程等待指针指向该进程。
137     while (!current->signal && EMPTY(*queue))
138         interruptible_sleep_on(&queue->proc_list);
139     sti(); // 开中断。
140 }
141
142     ///// 若队列缓冲区满则让进程进入可中断的睡眠状态。
143     // 参数: queue - 指定队列的指针。
144     // 进程在往队列缓冲区中写入时调用此函数。
145 static void sleep_if_full(struct tty_queue * queue)
146 {
147     // 若队列缓冲区不满, 则返回退出。
148     if (!FULL(*queue))
149         return;
150     cli(); // 关中断。
151     // 如果进程没有信号需要处理并且队列缓冲区中空闲剩余区长度<128, 则让进程进入可中断睡眠状态,
152     // 并让该队列的进程等待指针指向该进程。
153     while (!current->signal && LEFT(*queue)<128)

```



```

136             interruptible sleep on(&queue->proc_list);
137         sti(); // 开中断。
138     }
139     // 等待按键。
140     // 如果控制台的读队列缓冲区空则让进程进入可中断的睡眠状态。
141     void wait_for_keypress(void)
142     {
143         sleep_if_empty(&tty_table[0].secondary);
144     }
145     // 复制成规范模式字符序列。
146     // 将指定 tty 终端队列缓冲区中的字符复制成规范(熟)模式字符并存放在辅助队列(规范模式队列)中。
147     // 参数: tty - 指定终端的 tty 结构。
148     void copy_to_cooked(struct tty_struct * tty)
149     {
150         signed char c;
151         // 如果 tty 的读队列缓冲区不空并且辅助队列缓冲区为空, 则循环执行下列代码。
152         while (!EMPTY(tty->read_q) && !FULL(tty->secondary)) {
153             // 从队列尾处取一字符到 c, 并前移尾指针。
154             GETCH(tty->read_q, c);
155             // 下面对输入字符, 利用输入模式标志集进行处理。
156             // 如果该字符是回车符 CR(13), 则: 若回车转换标志 CRNL 置位则将该字符转换为换行符 NL(10);
157             // 否则若忽略回车标志 NOCR 置位, 则忽略该字符, 继续处理其它字符。
158             if (c==13)
159                 if (I_CRNL(tty))
160                     c=10;
161                 else if (I_NOCR(tty))
162                     continue;
163                 else ;
164             // 如果该字符是换行符 NL(10) 并且换行转回车标志 NLCR 置位, 则将其转换为回车符 CR(13)。
165             else if (c==10 && I_NLCR(tty))
166                 c=13;
167             // 如果大写转小写标志 UCLC 置位, 则将该字符转换为小写字符。
168             if (I_UCLC(tty))
169                 c=tolower(c);
170             // 如果本地模式标志集中规范(熟)模式标志 CANON 置位, 则进行以下处理。
171             if (L_CANON(tty)) {
172                 // 如果该字符是键盘终止控制字符 KILL(^U), 则进行删除输入行处理。
173                 if (c==KILL_CHAR(tty)) {
174                     /* deal with killing the input line */ /* 删除输入行处理 */
175                     // 如果 tty 辅助队列不空, 或者辅助队列中最后一个字符是换行 NL(10), 或者该字符是文件结束字符
176                     // (^D), 则循环执行下列代码。
177                     while(!EMPTY(tty->secondary) ||
178                           (c=LAST(tty->secondary))==10 ||
179                           c==EOF_CHAR(tty))) {
180                         // 如果本地回显标志 ECHO 置位, 那么: 若字符是控制字符(值<32), 则往 tty 的写队列中放入擦除
181                         // 字符 ERASE。再放入一个擦除字符 ERASE, 并且调用该 tty 的写函数。
182                         if (L_ECHO(tty)) {
183                             if (c<32)
184                                 PUTCH(127, tty->write_q);
185                                 PUTCH(127, tty->write_q);

```



```

171         tty->write(tty);
172     }
    // 将 tty 辅助队列头指针后退 1 字节。
173     DEC(tty->secondary.head);
174 }
175     continue;    // 继续读取并处理其它字符。
176 }
    // 如果该字符是删除控制字符 ERASE(^H), 那么:
177     if (c==ERASE_CHAR(tty)) {
    // 若 tty 的辅助队列为空, 或者其最后一个字符是换行符 NL(10), 或者是文件结束符, 继续处理
    // 其它字符。
178         if (EMPTY(tty->secondary) ||
179             (c==LAST(tty->secondary))==10 ||
180             c==EOF_CHAR(tty))
181             continue;
    // 如果本地回显标志 ECHO 置位, 那么: 若字符是控制字符(值<32), 则往 tty 的写队列中放入擦除
    // 字符 ERASE。再放入一个擦除字符 ERASE, 并且调用该 tty 的写函数。
182         if (L_ECHO(tty)) {
183             if (c<32)
184                 PUTCH(127, tty->write_q);
185                 PUTCH(127, tty->write_q);
186                 tty->write(tty);
187         }
    // 将 tty 辅助队列头指针后退 1 字节, 继续处理其它字符。
188     DEC(tty->secondary.head);
189     continue;
190 }
    //如果该字符是停止字符(^S), 则置 tty 停止标志, 继续处理其它字符。
191     if (c==STOP_CHAR(tty)) {
192         tty->stopped=1;
193         continue;
194     }
    // 如果该字符是停止字符(^Q), 则复位 tty 停止标志, 继续处理其它字符。
195     if (c==START_CHAR(tty)) {
196         tty->stopped=0;
197         continue;
198     }
199 }
    // 若输入模式标志集中 ISIG 标志置位, 则在收到 INTR、QUIT、SUSP 或 DSUSP 字符时, 需要为进程
    // 产生相应的信号。
200     if (L_ISIG(tty)) {
    // 如果该字符是键盘中断符(^C), 则向当前进程发送键盘中断信号, 并继续处理下一字符。
201         if (c==INTR_CHAR(tty)) {
202             tty_intr(tty, INTRMASK);
203             continue;
204         }
    // 如果该字符是键盘中断符(^\\), 则向当前进程发送键盘退出信号, 并继续处理下一字符。
205         if (c==QUIT_CHAR(tty)) {
206             tty_intr(tty, QUITMASK);
207             continue;
208         }
209     }
    // 如果该字符是换行符 NL(10), 或者是文件结束符 EOF(^D), 辅助缓冲队列字符数加 1。[??]

```

```

210         if (c==10 || c==EOF_CHAR(tty))
211             tty->secondary.data++;
// 如果本地模式标志集中回显标志 ECHO 置位, 那么, 如果字符是换行符 NL(10), 则将换行符 NL(10)
// 和回车符 CR(13)放入 tty 写队列缓冲区中; 如果字符是控制字符(字符值<32)并且回显控制字符标志
// ECHOCTL 置位, 则将字符'~'和字符 c+64 放入 tty 写队列中(也即会显示^C、^H等); 否则将该字符
// 直接放入 tty 写缓冲队列中。最后调用该 tty 的写操作函数。
212         if (L_ECHO(tty)) {
213             if (c==10) {
214                 PUTCH(10, tty->write_q);
215                 PUTCH(13, tty->write_q);
216             } else if (c<32) {
217                 if (L_ECHOCTL(tty)) {
218                     PUTCH('~', tty->write_q);
219                     PUTCH(c+64, tty->write_q);
220                 }
221             } else
222                 PUTCH(c, tty->write_q);
223             tty->write(tty);
224         }
// 将该字符放入辅助队列中。
225         PUTCH(c, tty->secondary);
226     }
// 唤醒等待该辅助缓冲队列的进程(如果有的话)。
227     wake_up(&tty->secondary.proc_list);
228 }
229
//// tty 读函数。
// 参数: channel - 子设备号; buf - 缓冲区指针; nr - 欲读字节数。
// 返回已读字节数。
230 int tty_read(unsigned channel, char * buf, int nr)
231 {
232     struct tty_struct * tty;
233     char c, * b=buf;
234     int minimum, time, flag=0;
235     long oldalarm;
236
// 本版本 linux 内核的终端只有 3 个子设备, 分别是控制台(0)、串口终端 1(1)和串口终端 2(2)。
// 所以任何大于 2 的子设备号都是非法的。写的字节数当然也不能小于 0 的。
237     if (channel>2 || nr<0) return -1;
// tty 指针指向子设备号对应 ttb_table 表中的 tty 结构。
238     tty = &tty_table[channel];
// 下面首先保存进程原定时值, 然后根据控制字符 VTIME 和 VMIN 设置读字符操作的超时定时值。
// 在非规范模式下, 这两个值是超时定时值。MIN 表示为了满足读操作, 需要读取的最少字符数。
// TIME 是一个十分之一秒计数的计时值。
// 首先取进程中的(报警)定时值(滴答数)。
239     oldalarm = current->alarm;
// 并设置读操作超时定时值 time 和需要最少读取的字符个数 minimum。
240     time = 10L*tty->termios.c_cc[VTIME];
241     minimum = tty->termios.c_cc[VMIN];
// 如果设置了读超时定时值 time 但没有设置最少读取个数 minimum, 那么在读到至少一个字符或者
// 定时超时后读操作将立刻返回。所以这里置 minimum=1。
242     if (time && !minimum) {
243         minimum=1;

```

```

// 如果进程原定时值是 0 或者 time+当前系统时间值小于进程原定时值的话，则置重新设置进程定时
// 值为 time+当前系统时间，并置 flag 标志。
244         if (flag!=(oldalarm || time+jiffies<oldalarm))
245             current->alarm = time+jiffies;
246     }
// 如果设置的最少读取字符数>欲读的字符数，则令其等于此次欲读取的字符数。
247     if (minimum>nr)
248         minimum=nr;
// 当欲读的字节数>0，则循环执行以下操作。
249     while (nr>0) {
// 如果 flag 不为 0(即进程原定时值是 0 或者 time+当前系统时间值小于进程原定时值)并且进程有定
// 时信号 SIGALRM，则复位进程的定时信号并中断循环。
250         if (flag && (current->signal & ALRMMASK)) {
251             current->signal &= ~ALRMMASK;
252             break;
253         }
// 如果当前进程有信号要处理，则退出，返回 0。
254         if (current->signal)
255             break;
// 如果辅助缓冲队列(规范模式队列)为空，或者设置了规范模式标志并且辅助队列中字符数为 0 以及
// 辅助模式缓冲队列空闲空间>20，则进入可中断睡眠状态，返回后继续处理。
256         if (EMPTY(tty->secondary) || (L CANON(tty) &&
257             !tty->secondary.data && LEFT(tty->secondary)>20)) {
258             sleep_if_empty(&tty->secondary);
259             continue;
260         }
// 执行以下操作，直到 nr=0 或者辅助缓冲队列为空。
261         do {
// 取辅助缓冲队列字符 c。
262             GETCH(tty->secondary, c);
// 如果该字符是文件结束符(^D)或者是换行符 NL(10)，则辅助缓冲队列字符数减 1。
263             if (c==EOF_CHAR(tty) || c==10)
264                 tty->secondary.data--;
// 如果该字符是文件结束符(^D)并且规范模式标志置位，则返回已读字符数，并退出。
265             if (c==EOF_CHAR(tty) && L CANON(tty))
266                 return (b-buf);
// 否则将该字符放入用户数据段缓冲区 buf 中，欲读字符数减 1，如果欲读字符数已为 0，则中断循环。
267             else {
268                 put_fs_byte(c, b++);
269                 if (!--nr)
270                     break;
271             }
272         } while (nr>0 && !EMPTY(tty->secondary));
// 如果超时定时值 time 不为 0 并且规范模式标志没有置位(非规范模式)，那么：
273         if (time && !L CANON(tty))
// 如果进程原定时值是 0 或者 time+当前系统时间值小于进程原定时值的话，则置重新设置进程定时值
// 为 time+当前系统时间，并置 flag 标志。否则让进程的定时值等于进程原定时值。
274         if (flag!=(oldalarm || time+jiffies<oldalarm))
275             current->alarm = time+jiffies;
276         else
277             current->alarm = oldalarm;
// 如果规范模式标志置位，那么若没有读到 1 个字符则中断循环。否则若已读取数大于或等于最少要
// 求读取的字符数，则也中断循环。

```

```

278         if (L_CANON(tty)) {
279             if (b->buf)
280                 break;
281         } else if (b->buf >= minimum)
282             break;
283     }
    // 让进程的定时值等于进程原定值。
284     current->alarm = oldalarm;
    // 如果进程有信号并且没有读取任何字符，则返回出错号（超时）。
285     if (current->signal && !(b->buf))
286         return -EINTR;
287     return (b->buf);        // 返回已读取的字符数。
288 }
289
    //// tty 写函数。
    // 参数：channel - 子设备号；buf - 缓冲区指针；nr - 写字节数。
    // 返回已写字节数。
290 int tty_write(unsigned channel, char * buf, int nr)
291 {
292     static cr_flag=0;
293     struct tty_struct * tty;
294     char c, *b=buf;
295
    // 本版本 linux 内核的终端只有 3 个子设备，分别是控制台(0)、串口终端 1(1)和串口终端 2(2)。
    // 所以任何大于 2 的子设备号都是非法的。写的字节数当然也不能小于 0 的。
296     if (channel>2 || nr<0) return -1;
    // tty 指针指向子设备号对应 ttb_table 表中的 tty 结构。
297     tty = channel + tty_table;
    // 字符设备是一个一个字符进行处理的，所以这里对于 nr 大于 0 时对每个字符进行循环处理。
298     while (nr>0) {
    // 如果此时 tty 的写队列已满，则当前进程进入可中断的睡眠状态。
299         sleep_if_full(&tty->write_q);
    // 如果当前进程有信号要处理，则退出，返回 0。
300         if (current->signal)
301             break;
    // 当要写的字节数>0 并且 tty 的写队列不满时，循环执行以下操作。
302         while (nr>0 && !FULL(tty->write_q)) {
    // 从用户数据段内存中取一字节 c。
303             c=get_fs_byte(b);
    // 如果终端输出模式标志集中执行输出处理标志 OPOST 置位，则执行下列输出时处理过程。
304             if (O_POST(tty)) {
    // 如果该字符是回车符 '\r' (CR, 13) 并且回车符转换行符标志 OCRNL 置位，则将该字符换成换行符
    // '\n' (NL, 10)；否则如果该字符是换行符 '\n' (NL, 10) 并且换行转回车功能标志 ONLRET 置位的话，
    // 则将该字符换成回车符 '\r' (CR, 13)。
305             if (c=='\r' && O_CRNL(tty))
306                 c='\n';
307             else if (c=='\n' && O_NLRET(tty))
308                 c='\r';
    // 如果该字符是换行符 '\n' 并且回车标志 cr_flag 没有置位，换行转回车-换行标志 ONLCR 置位的话，
    // 则将 cr_flag 置位，并将一回车符放入写队列中。然后继续处理下一个字符。
309             if (c=='\n' && !cr_flag && O_NLCR(tty)) {
310                 cr_flag = 1;
311                 PUTCH(13, tty->write_q);

```

```

312                                     continue;
313                                     }
// 如果小写转大写标志 OLCUC 置位的话, 就将该字符转成大写字符。
314                                     if (0 LUCUC(tty))
315                                         c=toupper(c);
316                                     }
// 用户数据缓冲指针 b 前进 1 字节; 欲写字节数减 1 字节; 复位 cr_flag 标志, 并将该字节放入 tty
// 写队列中。
317                                     b++; nr--;
318                                     cr_flag = 0;
319                                     PUTCH(c, tty->write_q);
320                                     }
// 若字节全部写完, 或者写队列已满, 则程序执行到这里。调用对应 tty 的写函数, 若还有字节要写,
// 则等待写队列不满, 所以调用调度程序, 先去执行其它任务。
321                                     tty->write(tty);
322                                     if (nr>0)
323                                         schedule();
324                                     }
325                                     return (b-buf); // 返回写入的字节数。
326     }
327
328 /*
329  * Jeh, sometimes I really like the 386.
330  * This routine is called from an interrupt,
331  * and there should be absolutely no problem
332  * with sleeping even in an interrupt (I hope).
333  * Of course, if somebody proves me wrong, I'll
334  * hate intel for all time :-). We'll have to
335  * be careful and see to reinstating the interrupt
336  * chips before calling this, though.
337  *
338  * I don't think we sleep here under normal circumstances
339  * anyway, which is good, as the task sleeping might be
340  * totally innocent.
341  */
/*
 * 呵, 有时我是真得很喜欢 386。该子程序是从一个中断处理程序中调用的, 即使在
 * 中断处理程序中睡眠也应该绝对没有问题(我希望如此)。当然, 如果有人证明我是
 * 错的, 那么我将憎恨 intel 一辈子☺。但是我们必须小心, 在调用该子程序之前需
 * 要恢复中断。
 *
 * 我不认为在通常环境下会处在这里睡眠, 这样很好, 因为任务睡眠是完全任意的。
 */
//// tty 中断处理调用函数 - 执行 tty 中断处理。
// 参数: tty - 指定的 tty 终端号 (0, 1 或 2) 。
// 将指定 tty 终端队列缓冲区中的字符复制成规范(熟)模式字符并存放在辅助队列(规范模式队列)中。
// 在串口读字符中断(rs_io.s, 109)和键盘中断(kerboard.S, 69)中调用。
342 void do tty interrupt(int tty)
343 {
344     copy_to_cooked(tty_table+tty);
345 }
346
//// 字符设备初始化函数。空, 为以后扩展做准备。

```

```

347 void chr_dev_init(void)
348 {
349 }
350

```

## 6.5.3 其它信息

### 6.5.3.1 控制字符 VTIME、VMIN

在非规范模式下，这两个值是超时定时值。MIN 表示为了满足读操作，需要读取的最少字符数。TIME 是一个十分之一秒计数的计时值。当这两个都设置的话，读操作将等待，直到至少读到一个字符，然后在以读取 MIN 个字符或者时间 TIME 在读取最后一个字符后超时。如果仅设置了 MIN，那么在读取 MIN 个字符之前读操作将不返回。如果仅设置了 TIME，那么在读到至少一个字符或者定时超时后读操作将立刻返回。如果两个都没有设置，则读操作将立刻返回，仅给出目前已读的字节数。详细说明参见 `termios.h` 文件。

## 6.6 tty\_ioctl.c 文件

### 6.6.1 功能描述

### 6.6.2 代码注释

列表 linux/kernel/chr\_drv/tty\_ioctl.c 程序

```

1  /*
2  *  linux/kernel/chr_drv/tty_ioctl.c
3  *
4  *  (C) 1991  Linus Torvalds
5  */
6
7  #include <errno.h>          // 错误号头文件。包含系统中各种出错号。(Linus 从 minix 中引进的)。
8  #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。
9
10 #include <linux/sched.h>    // 调度程序头文件，定义了任务结构 task_struct、初始任务 0 的数据，
    // 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
11 #include <linux/kernel.h>   // 内核头文件。含有一些内核常用函数的原形定义。
12 #include <linux/tty.h>     // tty 头文件，定义了有关 tty_io，串行通信方面的参数、常数。
13
14 #include <asm/io.h>         // io 头文件。定义硬件端口输入/输出宏汇编语句。
15 #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。
16 #include <asm/system.h>    // 系统头文件。定义了设置或修改描述符/中断门等的嵌入式汇编宏。
17
    // 这是波特率因子数组（或称为除数数组）。波特率与波特率因子的对应关系参见列表后的说明。
18 static unsigned short quotient[] = {
19     0, 2304, 1536, 1047, 857,
20     768, 576, 384, 192, 96,
21     64, 48, 24, 12, 6, 3
22 };
23

```

```

24  static void change_speed(struct tty_struct * tty)
25  {
26      unsigned short port, quot;
27
28      // 对于串口终端，其 tty 结构的读缓冲队列 data 字段存放的是串行端口号 (0x3f8 或 0x2f8)。
29      if (!(port = tty->read_q.data))
30          return;
31      // 从 tty 的 termios 结构控制模式标志集中取得设置的波特率索引号，据此从波特率因子数组中取得
32      // 对应的波特率因子值。CBAUD 是控制模式标志集中波特率位屏蔽码。
33      quot = quotient[tty->termios.c_cflag & CBAUD];
34      cli(); // 关中断。
35      outb_p(0x80, port+3); // set DLAB // 首先设置除数锁定标志 DLAB。
36      outb_p(quot & 0xff, port); // LS of divisor // 输出因子低字节。
37      outb_p(quot >> 8, port+1); // MS of divisor // 输出因子高字节。
38      outb(0x03, port+3); // reset DLAB // 复位 DLAB。
39      sti(); // 开中断。
40  }
41
42  // 刷新 tty 缓冲队列。
43  // 参数：queue - 指定的缓冲队列指针。
44  // 令缓冲队列的头指针等于尾指针，从而达到清空缓冲区(零字符)的目的。
45  static void flush(struct tty_queue * queue)
46  {
47      cli();
48      queue->head = queue->tail;
49      sti();
50  }
51
52  // 等待字符发送出去。
53  static void wait_until_sent(struct tty_struct * tty)
54  {
55      // do nothing - not implemented // 什么都没做 - 还未实现
56  }
57
58  // 发送 BREAK 控制符。
59  static void send_break(struct tty_struct * tty)
60  {
61      // do nothing - not implemented // 什么都没做 - 还未实现
62  }
63
64  // 取终端 termios 结构信息。
65  // 参数：tty - 指定终端的 tty 结构指针；termios - 用户数据区 termios 结构缓冲区指针。
66  // 返回 0。
67  static int get_termios(struct tty_struct * tty, struct termios * termios)
68  {
69      int i;
70
71      // 首先验证一下用户的缓冲区指针所指内存区是否足够，如不够则分配内存。
72      verify_area(termios, sizeof (*termios));

```



```

// 复制指定 tty 结构中的 termios 结构信息到用户 termios 结构缓冲区。
61     for (i=0 ; i< (sizeof (*termios)) ; i++)
62         put_fs_byte( ((char *)&tty->termios)[i] , i+(char *)termios );
63     return 0;
64 }
65
//// 设置终端 termios 结构信息。
// 参数: tty - 指定终端的 tty 结构指针; termios - 用户数据区 termios 结构指针。
// 返回 0 。
66 static int set_termios(struct tty_struct * tty, struct termios * termios)
67 {
68     int i;
69
// 首先复制用户数据区中 termios 结构信息到指定 tty 结构中。
70     for (i=0 ; i< (sizeof (*termios)) ; i++)
71         ((char *)&tty->termios)[i]=get_fs_byte(i+(char *)termios);
// 用户有可能已修改了 tty 的串行口传输波特率, 所以根据 termios 结构中的控制模式标志 c_cflag
// 修改串行芯片 UART 的传输波特率。
72     change_speed(tty);
73     return 0;
74 }
75
//// 读取 termio 结构中的信息。
// 参数: tty - 指定终端的 tty 结构指针; termio - 用户数据区 termio 结构缓冲区指针。
// 返回 0。
76 static int get_termio(struct tty_struct * tty, struct termio * termio)
77 {
78     int i;
79     struct termio tmp_termio;
80
// 首先验证一下用户的缓冲区指针所指内存区是否足够, 如不够则分配内存。
81     verify_area(termio, sizeof (*termio));
// 将 termios 结构的信息复制到 termio 结构中。目的是为了其中模式标志集的类型进行转换, 也即
// 从 termios 的长整数类型转换为 termio 的短整数类型。
82     tmp_termio.c_iflag = tty->termios.c_iflag;
83     tmp_termio.c_oflag = tty->termios.c_oflag;
84     tmp_termio.c_cflag = tty->termios.c_cflag;
85     tmp_termio.c_lflag = tty->termios.c_lflag;
// 两种结构的 c_line 和 c_cc[] 字段是完全相同的。
86     tmp_termio.c_line = tty->termios.c_line;
87     for(i=0 ; i < NCC ; i++)
88         tmp_termio.c_cc[i] = tty->termios.c_cc[i];
// 最后复制指定 tty 结构中的 termio 结构信息到用户 termio 结构缓冲区。
89     for (i=0 ; i< (sizeof (*termio)) ; i++)
90         put_fs_byte( ((char *)&tmp_termio)[i] , i+(char *)termio );
91     return 0;
92 }
93
94 /*
95  * This only works as the 386 is low-byt-first
96  */
97 /*
98  * 下面的 termio 设置函数仅在 386 低字节在前的方式下可用。

```



```

    */
    //// 设置终端 termio 结构信息。
    // 参数: tty - 指定终端的 tty 结构指针; termio - 用户数据区 termio 结构指针。
    // 将用户缓冲区 termio 的信息复制到终端的 termios 结构中。返回 0。
97 static int set_termio(struct tty_struct * tty, struct termio * termio)
98 {
99     int i;
100     struct termio tmp_termio;
101
102     // 首先复制用户数据区中 termio 结构信息到临时 termio 结构中。
103     for (i=0 ; i< (sizeof (*termio)) ; i++)
104         ((char *)&tmp_termio)[i]=get_fs_byte(i+(char *)termio);
105     // 再将 termio 结构的信息复制到 tty 的 termios 结构中。目的是为了其中模式标志集的类型进行转换,
106     // 也即从 termio 的短整数类型转换成 termios 的长整数类型。
107     *(unsigned short *)&tty->termios.c_iflag = tmp_termio.c_iflag;
108     *(unsigned short *)&tty->termios.c_oflag = tmp_termio.c_oflag;
109     *(unsigned short *)&tty->termios.c_cflag = tmp_termio.c_cflag;
110     *(unsigned short *)&tty->termios.c_lflag = tmp_termio.c_lflag;
111     // 两种结构的 c_line 和 c_cc[] 字段是完全相同的。
112     tty->termios.c_line = tmp_termio.c_line;
113     for(i=0 ; i < NCC ; i++)
114         tty->termios.c_cc[i] = tmp_termio.c_cc[i];
115     // 用户可能已修改了 tty 的串行口传输波特率, 所以根据 termios 结构中的控制模式标志集 c_cflag
116     // 修改串行芯片 UART 的传输波特率。
117     change_speed(tty);
118     return 0;
119 }
120
121 //// tty 终端设备的 ioctl 函数。
122 // 参数: dev - 设备号; cmd - ioctl 命令; arg - 操作参数指针。
123 int tty_ioctl(int dev, int cmd, int arg)
124 {
125     struct tty_struct * tty;
126     // 首先取 tty 的子设备号。如果主设备号是 5(tty 终端), 则进程的 tty 字段即是子设备号; 如果进程
127     // 的 tty 子设备号是负数, 表明该进程没有控制终端, 也即不能发出该 ioctl 调用, 出错死机。
128     if (MAJOR(dev) == 5) {
129         dev=current->tty;
130         if (dev<0)
131             panic("tty_ioctl: dev<0");
132     }
133     // 否则直接从设备号中取出子设备号。
134     } else
135         dev=MINOR(dev);
136     // 子设备号可以是 0(控制台终端)、1(串口 1 终端)、2(串口 2 终端)。
137     // 让 tty 指向对应子设备号的 tty 结构。
138     tty = dev + tty_table;
139     // 根据 tty 的 ioctl 命令进行分别处理。
140     switch (cmd) {
141     case TCGETS:
142         //取相应终端 termios 结构中的信息。
143         return get_termios(tty, (struct termios *) arg);
144     case TCSETSF:
145         // 在设置 termios 的信息之前, 需要先等待输出队列中所有数据处理完, 并且刷新(清空)输入队列。
146         // 再设置。

```

```

129         flush(&tty->read_q); /* fallthrough */
130     case TCSETSW:
    // 在设置终端 termios 的信息之前, 需要先等待输出队列中所有数据处理完(耗尽)。对于修改参数
    // 会影响输出的情况, 就需要使用这种形式。
131         wait_until_sent(tty); /* fallthrough */
132     case TCSETS:
    // 设置相应终端 termios 结构中的信息。
133         return set_termios(tty, (struct termios *) arg);
134     case TCGETA:
    // 取相应终端 termio 结构中的信息。
135         return get_termio(tty, (struct termio *) arg);
136     case TCSETAF:
    // 在设置 termio 的信息之前, 需要先等待输出队列中所有数据处理完, 并且刷新(清空)输入队列。
    // 再设置。
137         flush(&tty->read_q); /* fallthrough */
138     case TCSETAW:
    // 在设置终端 termio 的信息之前, 需要先等待输出队列中所有数据处理完(耗尽)。对于修改参数
    // 会影响输出的情况, 就需要使用这种形式。
139         wait_until_sent(tty); /* fallthrough */ /* 继续执行 */
140     case TCSETA:
    // 设置相应终端 termio 结构中的信息。
141         return set_termio(tty, (struct termio *) arg);
142     case TCSBRK:
    // 等待输出队列处理完毕(空), 如果参数值是 0, 则发送一个 break。
143         if (!arg) {
144             wait_until_sent(tty);
145             send_break(tty);
146         }
147         return 0;
148     case TCXONC:
    // 开始/停止控制。如果参数值是 0, 则挂起输出; 如果是 1, 则重新开启挂起的输出; 如果是 2, 则挂起
    // 输入; 如果是 3, 则重新开启挂起的输入。
149         return -EINVAL; /* not implemented */ /* 未实现 */
150     case TCFLSH:
    // 刷新已写输出但还没发送或已收但还没有读数据。如果参数是 0, 则刷新(清空)输入队列; 如果是 1,
    // 则刷新输出队列; 如果是 2, 则刷新输入和输出队列。
151         if (arg==0)
152             flush(&tty->read_q);
153         else if (arg==1)
154             flush(&tty->write_q);
155         else if (arg==2) {
156             flush(&tty->read_q);
157             flush(&tty->write_q);
158         } else
159             return -EINVAL;
160         return 0;
161     case TIOCEXCL:
    // 设置终端串行线路专用模式。
162         return -EINVAL; /* not implemented */ /* 未实现 */
163     case TIOCNXCL:
    // 复位终端串行线路专用模式。
164         return -EINVAL; /* not implemented */ /* 未实现 */
165     case TIOCSTTY:

```

```

// 设置 tty 为控制终端。(TIOCNOTTY - 禁止 tty 为控制终端)。
166         return -EINVAL; /* set controlling term NI */ /* 设置控制终端 NI */
167     case TIOCGPGRP:                                     // NI - Not Implemented.
// 读取指定终端设备进程的组 id。首先验证用户缓冲区长度，然后复制 tty 的 pgrp 字段到用户缓冲区。
168         verify_area((void *) arg, 4);
169         put_fs_long(tty->pgrp, (unsigned long *) arg);
170         return 0;
171     case TIOCSPPGRP:
// 设置指定终端设备进程的组 id。
172         tty->pgrp=get_fs_long((unsigned long *) arg);
173         return 0;
174     case TIOCOUTQ:
// 返回输出队列中还未送出的字符数。首先验证用户缓冲区长度，然后复制队列中字符数给用户。
175         verify_area((void *) arg, 4);
176         put_fs_long(CHARS(tty->write_q), (unsigned long *) arg);
177         return 0;
178     case TIOCINQ:
// 返回输入队列中还未读取的字符数。首先验证用户缓冲区长度，然后复制队列中字符数给用户。
179         verify_area((void *) arg, 4);
180         put_fs_long(CHARS(tty->secondary),
181                     (unsigned long *) arg);
182         return 0;
183     case TIOCS TI:
// 模拟终端输入。该命令以一个指向字符的指针作为参数，并假装该字符是在终端上键入的。用户必须
// 在该控制终端上具有超级用户权限或具有读许可权限。
184         return -EINVAL; /* not implemented */ /* 未实现 */
185     case TIOCGWINSZ:
// 读取终端设备窗口大小信息（参见 termios.h 中的 winsize 结构）。
186         return -EINVAL; /* not implemented */ /* 未实现 */
187     case TIOCSWINSZ:
// 设置终端设备窗口大小信息（参见 winsize 结构）。
188         return -EINVAL; /* not implemented */ /* 未实现 */
189     case TIOCMGET:
// 返回 modem 状态控制引线的当前状态比特位标志集（参见 termios.h 中 185-196 行）。
190         return -EINVAL; /* not implemented */ /* 未实现 */
191     case TIOCMBS:
// 设置单个 modem 状态控制引线的状态(true 或 false)。
192         return -EINVAL; /* not implemented */ /* 未实现 */
193     case TIOCMBS:
// 复位单个 modem 状态控制引线的状态。
194         return -EINVAL; /* not implemented */ /* 未实现 */
195     case TIOCMSET:
// 设置 modem 状态引线的状态。如果某一比特位置位，则 modem 对应的状态引线将置为有效。
196         return -EINVAL; /* not implemented */ /* 未实现 */
197     case TIOCGSOFTCAR:
// 读取软件载波检测标志(1 - 开启; 0 - 关闭)。
198         return -EINVAL; /* not implemented */ /* 未实现 */
199     case TIOCSSOFTCAR:
// 设置软件载波检测标志(1 - 开启; 0 - 关闭)。
200         return -EINVAL; /* not implemented */ /* 未实现 */
201     default:
202         return -EINVAL;
203 }

```

[204](#) }  
[205](#)

### 6.6.3 其它信息

#### 6.6.3.1 波特率与波特率因子

波特率 = 1.8432MHz / (16 \* 波特率因子)。程序中波特率与波特率因子的对应关系见下表所示。

| 波特率   | 波特率因子      |      | 波特率   | 波特率因子      |     |
|-------|------------|------|-------|------------|-----|
|       | MSB, LSB   | 合并值  |       | MSB, LSB   | 合并值 |
| 50    | 0x09, 0x00 | 2304 | 1200  | 0x00, 0x60 | 96  |
| 75    | 0x06, 0x00 | 1536 | 1800  | 0x00, 0x40 | 64  |
| 110   | 0x04, 0x17 | 1047 | 2400  | 0x00, 0x30 | 48  |
| 134.5 | 0x03, 0x59 | 857  | 4800  | 0x00, 0x18 | 24  |
| 150   | 0x03, 0x00 | 768  | 9600  | 0x00, 0x1c | 12  |
| 200   | 0x02, 0x40 | 576  | 19200 | 0x00, 0x06 | 6   |
| 300   | 0x01, 0x80 | 384  | 38400 | 0x00, 0x03 | 3   |
| 600   | 0x00, 0xc0 | 192  |       |            |     |

## 6.7 keyboard.S 文件

### 6.7.1 功能描述

该键盘驱动汇编程序主要包括键盘中断处理程序。在英文惯用法中，make 表示键被按下；break 表示键被松开(放开)。

对于 AT 键盘的扫描码，当键按下时，则对应键的扫描码被送出，但当键松开时，将会发送两个字节，第一个是 0xf0，第 2 个还是按下时的扫描码。为了向下的兼容性，设计人员将 AT 键盘发出的扫描码转换成了老式 PC/XT 标准键盘的扫描码。因此这里仅对 PC/XT 的扫描码进行处理即可。

### 6.7.2 代码注释

列表 linux/kernel/chr\_drv/keyboard.S 文件

```
1 /*  
2  * linux/kernel/keyboard.S  
3  *  
4  * (C) 1991 Linus Torvalds  
5  */  
6  
7 /*  
8  * Thanks to Alfred Leung for US keyboard patches  
9  * Wolfgang Thiel for German keyboard patches  
10 * Marc Corsini for the French keyboard
```

```

11 */
12 /*
13  * 感谢 Alfred Leung 添加了 US 键盘补丁程序;
14  *      Wolfgang Thiel 添加了德语键盘补丁程序;
15  *      Marc Corsini 添加了法文键盘补丁程序。
16  */
17
18 #include <linux/config.h> // 内核配置头文件。定义键盘语言和硬盘类型 (HD_TYPE) 可选项。
19
20 .text
21 .globl _keyboard_interrupt
22
23 /*
24  * these are for the keyboard read functions
25  */
26 /*
27  * 以下这些是用于键盘读操作。
28  */
29 // size 是键盘缓冲区的长度 (字节数)。
30 size = 1024 /* must be a power of two ! And MUST be the same
31              as in tty_io.c !!!! */
32 /* 数值必须是 2 的次方! 并且与 tty_io.c 中的值匹配!!!! */
33 // 以下这些是缓冲队列结构中的偏移量 */
34 head = 4 // 缓冲区中头指针字段偏移。
35 tail = 8 // 缓冲区中尾指针字段偏移。
36 proc_list = 12 // 等待该缓冲队列的进程字段偏移。
37 buf = 16 // 缓冲区字段偏移。
38
39 // mode 是键盘特殊键的按下状态标志。
40 // 表示大小写转换键(caps)、交换键(alt)、控制键(ctrl)和换档键(shift)的状态。
41 // 位 7 caps 键按下;
42 // 位 6 caps 键的状态(应该与 leds 中的对应标志位一样);
43 // 位 5 右 alt 键按下;
44 // 位 4 左 alt 键按下;
45 // 位 3 右 ctrl 键按下;
46 // 位 2 左 ctrl 键按下;
47 // 位 1 右 shift 键按下;
48 // 位 0 左 shift 键按下。
49 mode: .byte 0 /* caps, alt, ctrl and shift mode */
50 // 数字锁定键(num-lock)、大小写转换键(caps-lock)和滚动锁定键(scroll-lock)的 LED 发光管状态。
51 // 位 7-3 全 0 不用;
52 // 位 2 caps-lock;
53 // 位 1 num-lock(初始置 1, 也即设置数字锁定键(num-lock)发光管为亮);
54 // 位 0 scroll-lock。
55 leds: .byte 2 /* num-lock, caps, scroll-lock mode (nom-lock on) */
56 // 当扫描码是 0xe0 或 0xe1 时, 置该标志。表示其后还跟随着 1 个或 2 个字符扫描码, 参见列表后说明。
57 // 位 1 =1 收到 0xe1 标志;
58 // 位 0 =1 收到 0xe0 标志。
59 e0: .byte 0
60
61 /*
62  * con_int is the real interrupt routine that reads the
63  * keyboard scan-code and converts it into the appropriate

```

```

35 *  ascii character(s).
36 */
/*
 * con_int 是实际的中断处理子程序，用于读键盘扫描码并将其转换
 * 成相应的 ascii 字符。
 */
///// 键盘中断处理程序入口点。
37 _keyboard_interrupt:
38     pushl %eax
39     pushl %ebx
40     pushl %ecx
41     pushl %edx
42     push %ds
43     push %es
44     movl $0x10,%eax        // 将 ds、es 段寄存器置为内核数据段。
45     mov %ax,%ds
46     mov %ax,%es
47     xorl %al,%al          /* %eax is scan code */ /* eax 中是扫描码 */
48     inb $0x60,%al         // 读取扫描码→al。
49     cmpb $0xe0,%al        // 该扫描码是 0xe0 吗？如果是则跳转到设置 e0 标志代码处。
50     je set_e0
51     cmpb $0xe1,%al        // 扫描码是 0xe1 吗？如果是则跳转到设置 e1 标志代码处。
52     je set_e1
53     call key_table(,%eax,4) // 调用键处理程序 ker_table + eax * 4（参见下面 502 行）。
54     movb $0,e0            // 复位 e0 标志。
// 下面这段代码(55-65 行)是针对使用 8255A 的 PC 标准键盘电路进行硬件复位处理。端口 0x61 是
// 8255A 输出口 B 的地址，该输出端口的第 7 位(PB7)用于禁止和允许对键盘数据的处理。
// 这段程序用于对收到的扫描码做出应答。方法是首先禁止键盘，然后立刻重新允许键盘工作。
55 e0_e1:  inb $0x61,%al      // 取 PPI 端口 B 状态，其位 7 用于允许/禁止(0/1)键盘。
56         jmp 1f            // 延迟一会。
57 1:      jmp 1f
58 1:      orb $0x80,%al      // al 位 7 置位(禁止键盘工作)。
59         jmp 1f            // 再延迟一会。
60 1:      jmp 1f
61 1:      outb %al,$0x61     // 使 PPI PB7 位置位。
62         jmp 1f            // 延迟一会。
63 1:      jmp 1f
64 1:      andb $0x7F,%al     // al 位 7 复位。
65         outb %al,$0x61     // 使 PPI PB7 位复位（允许键盘工作）。

66         movb $0x20,%al     // 向 8259 中断芯片发送 EOI(中断结束)信号。
67         outb %al,$0x20

68     pushl $0              // 控制台 tty 号=0，作为参数入栈。
69     call _do_tty_interrupt // 将收到的数据复制成规范模式数据并存放在规范字符缓冲队列中。
70     addl $4,%esp          // 丢弃入栈的参数，弹出保留的寄存器，并中断返回。
71     pop %es
72     pop %ds
73     popl %edx
74     popl %ecx
75     popl %ebx
76     popl %eax
77     iret

```

```

78 set_e0: movb $1,e0          // 收到扫描前导码 0xe0 时, 设置 e0 标志 (位 0)。
79         jmp e0_e1
80 set_e1: movb $2,e0          // 收到扫描前导码 0xe1 时, 设置 e1 标志 (位 1)。
81         jmp e0_e1
82
83 /*
84 * This routine fills the buffer with max 8 bytes, taken from
85 * %ebx:%eax. (%edx is high). The bytes are written in the
86 * order %al,%ah,%eal,%eah,%bl,%bh ... until %eax is zero.
87 */
/*
* 下面该子程序把 ebx:eax 中的最多 8 个字符添入缓冲队列中。(edx 是
* 所写入字符的顺序是 al, ah, eal, eah, bl, bh... 直到 eax 等于 0。
*/
88 put_queue:
89     pushl %ecx               // 保存 ecx, edx 内容。
90     pushl %edx               // 取控制台 tty 结构中读缓冲队列指针。
91     movl _table_list,%edx    # read-queue for console
92     movl head(%edx),%ecx     // 取缓冲队列中头指针→ecx。
93 1:     movb %al,buf(%edx,%ecx) // 将 al 中的字符放入缓冲队列头指针位置处。
94         incl %ecx            // 头指针前移 1 字节。
95         andl $size-1,%ecx    // 以缓冲区大小调整头指针(若超出则返回缓冲区开始)。
96         cmpl tail(%edx),%ecx # buffer full - discard everything
97         je 3f                // 头指针==尾指针吗(缓冲队列满)?
98         shrdl $8,%ebx,%eax    // 如果已满, 则后面未放入的字符全抛弃。
99         je 2f                // 将 ebx 中 8 位比特位右移 8 位到 eax 中, 但 ebx 不变。
100        shr $8,%ebx           // 还有字符吗? 若没有(等于 0)则跳转。
101        jmp 1b                // 将 ebx 中比特位右移 8 位, 并跳转到标号 1 继续操作。
102 2:     movl %ecx,head(%edx)    // 若已将所有字符都放入了队列, 则保存头指针。
103        movl proc_list(%edx),%ecx // 该队列的等待进程指针?
104        testl %ecx,%ecx        // 检测任务结构指针是否为空(有等待该队列的进程吗?)。
105        je 3f                // 无, 则跳转;
106        movl $0, (%ecx)        // 有, 则置该进程为可运行就绪状态(唤醒该进程)。
107 3:     popl %edx              // 弹出保留的寄存器并返回。
108        popl %ecx
109        ret
110
// 下面这段代码根据 ctrl 或 alt 的扫描码, 分别设置模式标志中相应位。如果该扫描码之前收到过
// 0xe0 扫描码(e0 标志置位), 则说明按下的是键盘右边的 ctrl 或 alt 键, 则对应设置 ctrl 或 alt
// 在模式标志 mode 中的比特位。
111 ctrl:  movb $0x04,%al        // 0x4 是模式标志 mode 中左 ctrl 键对应的比特位(位 2)。
112         jmp 1f
113 alt:    movb $0x10,%al        // 0x10 是模式标志 mode 中左 alt 键对应的比特位(位 4)。
114 1:      cmpb $0,e0            // e0 标志置位了吗(按下的是右边的 ctrl 或 alt 键吗)?
115         je 2f                // 不是则转。
116         addb %al,%al          // 是, 则改成置相应右键的标志位(位 3 或位 5)。
117 2:      orb %al,mode           // 设置模式标志 mode 中对应的比特位。
118         ret
// 这段代码处理 ctrl 或 alt 键松开的扫描码, 对应复位模式标志 mode 中的比特位。在处理时要根据
// e0 标志是否置位来判断是否是键盘右边的 ctrl 或 alt 键。
119 unctrl: movb $0x04,%al        // 模式标志 mode 中左 ctrl 键对应的比特位(位 2)。
120         jmp 1f

```



```

121 unalt:   movb $0x10,%al           // 0x10 是模式标志 mode 中左 alt 键对应的比特位(位 4)。
122 1:       cmpb $0,e0              // e0 标志置位了吗(释放的是右边的 ctrl 或 alt 键吗)?
123         je 2f                   // 不是, 则转。
124         addb %al,%al             // 是, 则该成复位相应右键的标志位(位 3 或位 5)。
125 2:       notb %al                // 复位模式标志 mode 中对应的比特位。
126         andb %al,mode
127         ret
128
129 lshift:
130         orb $0x01,mode           // 是左 shift 键按下, 设置 mode 中对应的标志位(位 0)。
131         ret
132 unlshift:
133         andb $0xfe,mode          // 是左 shift 键松开, 复位 mode 中对应的标志位(位 0)。
134         ret
135 rshift:
136         orb $0x02,mode           // 是右 shift 键按下, 设置 mode 中对应的标志位(位 1)。
137         ret
138 unrshift:
139         andb $0xfd,mode          // 是右 shift 键松开, 复位 mode 中对应的标志位(位 1)。
140         ret
141
142 caps:    testb $0x80,mode         // 测试模式标志 mode 中位 7 是否已经置位(按下状态)。
143         jne 1f                   // 如果已处于按下状态, 则返回(ret)。
144         xorb $4,leds              // 翻转 leds 标志中 caps-lock 比特位(位 2)。
145         xorb $0x40,mode           // 翻转 mode 标志中 caps 键按下的比特位(位 6)。
146         orb $0x80,mode           // 设置 mode 标志中 caps 键已按下标志位(位 7)。
// 这段代码根据 leds 标志, 开启或关闭 LED 指示器。
147 set_leds:
148         call kb_wait              // 等待键盘控制器输入缓冲空。
149         movb $0xed,%al            /* set leds command */ /* 设置 LED 的命令 */
150         outb %al,$0x60            // 发送键盘命令 0xed 到 0x60 端口。
151         call kb_wait              // 等待键盘控制器输入缓冲空。
152         movb leds,%al             // 取 leds 标志, 作为参数。
153         outb %al,$0x60            // 发送该参数。
154         ret
155 uncaps:   andb $0x7f,mode          // caps 键松开, 则复位模式标志 mode 中的对应位(位 7)。
156         ret
157 scroll:
158         xorb $1,leds              // scroll 键按下, 则翻转 leds 标志中的对应位(位 0)。
159         jmp set_leds              // 根据 leds 标志重新开启或关闭 LED 指示器。
160 num:     xorb $2,leds              // num 键按下, 则翻转 leds 标志中的对应位(位 1)。
161         jmp set_leds              // 根据 leds 标志重新开启或关闭 LED 指示器。
162
163 /*
164  * curosr-key/numeric keypad cursor keys are handled here.
165  * checking for numeric keypad etc.
166  */
/*
 * 这里处理方向键/数字小键盘方向键, 检测数字小键盘等。
 */
167 cursor:
168         subb $0x47,%al            // 扫描码是小数字键盘上的键(其扫描码>=0x47)发出的?
169         jb 1f                     // 如果小于则不处理, 返回。

```



```

170      cmpb $12,%al          // 如果扫描码 > 0x53(0x53 - 0x47= 12), 则
171      ja 1f                 // 扫描码值超过 83(0x53), 不处理, 返回。
172      jne cur2              /* check for ctrl-alt-del */ /* 检查是否 ctrl-alt-del */
                                // 如果等于 12, 则说明 del 键已被按下, 则继续判断 ctrl
                                // 和 alt 是否也同时按下。

173      testb $0x0c,mode     // 有 ctrl 键按下吗?
174      je cur2              // 无, 则跳转。
175      testb $0x30,mode     // 有 alt 键按下吗?
176      jne reboot          // 有, 则跳转到重启处理。
177 cur2:  cmpb $0x01,e0      /* e0 forces cursor movement */ /* e0 置位表示光标移动 */
                                // e0 标志置位了吗?

178      je cur               // 置位了, 则跳转光标移动处理处 cur。
179      testb $0x02,leds     /* not num-lock forces cursor */ /* num-lock 键则不许 */
                                // 测试 leds 中标志 num-lock 键标志是否置位。

180      je cur              // 如果没有置位(num 的 LED 不亮), 则也进行光标移动处理。
181      testb $0x03,mode     /* shift forces cursor */ /* shift 键也使光标移动 */
                                // 测试模式标志 mode 中 shift 按下标志。

182      jne cur              // 如果有 shift 键按下, 则也进行光标移动处理。
183      xorl %ebx,%ebx       // 否则查询扫描数字表(199 行), 取对应键的数字 ASCII 码。
184      movb num_table(%eax),%al // 以 eax 作为索引值, 取对应数字字符→al。
185      jmp put_queue        // 将该字符放入缓冲队列中。
186 1:      ret
187
// 这段代码处理光标的移动。
188 cur:    movb cur_table(%eax),%al // 取光标字符表中相应键的代表字符→al。
189      cmpb $'9',%al          // 若该字符<='9', 说明是上一页、下一页、插入或删除键,
190      ja ok_cur              // 则功能字符序列中要添入字符'~'。
191      movb $'~',%ah
192 ok_cur: shll $16,%eax       // 将 ax 中内容移到 eax 高字中。
193      movw $0x5b1b,%ax      // 在 ax 中放入'esc [' 字符, 与 eax 高字中字符组成移动序列。
194      xorl %ebx,%ebx
195      jmp put_queue         // 将该字符放入缓冲队列中。
196
197 #if defined(KBD_FR)
198 num_table:
199     .ascii "789 456 1230." // 数字小键盘上键对应的数字 ASCII 码表。
200 #else
201 num_table:
202     .ascii "789 456 1230,"
203 #endif
204 cur_table:
205     .ascii "HA5 DGC YB623" // 数字小键盘上方向键或插入删除键对应的移动表示字符表。
206
207 /*
208  * this routine handles function keys
209  */
// 下面子程序处理功能键。
210 func:
211     pushl %eax
212     pushl %ecx
213     pushl %edx
214     call _show_stat        // 调用显示各任务状态函数(kernel/sched.c, 37)。
215     popl %edx

```

```

216     popl %ecx
217     popl %eax
218     subb $0x3B,%al      // 功能键'F1'的扫描码是 0x3B, 因此此时 al 中是功能键索引号。
219     jb end_func         // 如果扫描码小于 0x3b, 则不处理, 返回。
220     cmpb $9,%al        // 功能键是 F1-F10?
221     jbe ok_func         // 是, 则跳转。
222     subb $18,%al       // 是功能键 F11, F12 吗?
223     cmpb $10,%al       // 是功能键 F11?
224     jb end_func         // 不是, 则不处理, 返回。
225     cmpb $11,%al       // 是功能键 F12?
226     ja end_func        // 不是, 则不处理, 返回。
227 ok_func:
228     cmpl $4,%ecx        /* check that there is enough room */ * 检查是否有足够空间*/
229     jl end_func         // 需要放入 4 个字符序列, 如果放不下, 则返回。
230     movl func_table(,%eax,4),%eax // 取功能键对应字符序列。
231     xorl %ebx,%ebx
232     jmp put_queue       // 放入缓冲队列中。
233 end_func:
234     ret
235
236 /*
237  * function keys send F1:'esc [ [ A' F2:'esc [ [ B' etc.
238  */
239 /*
240  * 功能键发送的扫描码, F1 键为: 'esc [ [ A', F2 键为: 'esc [ [ B' 等。
241  */
242 func_table:
243     .long 0x415b5b1b,0x425b5b1b,0x435b5b1b,0x445b5b1b
244     .long 0x455b5b1b,0x465b5b1b,0x475b5b1b,0x485b5b1b
245     .long 0x495b5b1b,0x4a5b5b1b,0x4b5b5b1b,0x4c5b5b1b
246
247 // 扫描码-ASCII 字符映射表。
248 // 根据在 config.h 中定义的键盘类型 (FINNISH, US, GERMEN, FRANCH), 将相应键的扫描码映射
249 // 到 ASCII 字符。
250 #if defined(KBD_FINNISH)
251 // 以下是芬兰语键盘的扫描码映射表。
252 key_map:
253     .byte 0,27          // 扫描码 0x00,0x01 对应的 ASCII 码;
254     .ascii "1234567890+' " // 扫描码 0x02,...0x0c,0x0d 对应的 ASCII 码, 以下类似。
255     .byte 127,9
256     .ascii "qwertyuiop}"
257     .byte 0,13,0
258     .ascii "asdfghjkl|{"
259     .byte 0,0
260     .ascii ""zxcvbnm,.-"
261     .byte 0,'*',0,32    /* 36-39 */ /* 扫描码 0x36-0x39 对应的 ASCII 码 */
262     .fill 16,1,0        /* 3A-49 */ /* 扫描码 0x3A-0x49 对应的 ASCII 码 */
263     .byte '-',0,0,0,0,'+' /* 4A-4E */ /* 扫描码 0x4A-0x4E 对应的 ASCII 码 */
264     .byte 0,0,0,0,0,0,0 /* 4F-55 */ /* 扫描码 0x4F-0x55 对应的 ASCII 码 */
265     .byte '<'
266     .fill 10,1,0
267
268 // shift 键同时按下时的映射表。

```

```

261 shift_map:
262     .byte 0, 27
263     .ascii "!\"#$%&/()=?`"
264     .byte 127, 9
265     .ascii "QWERTYUIOP]^"
266     .byte 13, 0
267     .ascii "ASDFGHJKL\\["
268     .byte 0, 0
269     .ascii "*ZXCVBNM;:_'"
270     .byte 0, '*', 0, 32      /* 36-39 */
271     .fill 16, 1, 0          /* 3A-49 */
272     .byte '-', 0, 0, 0, '+' /* 4A-4E */
273     .byte 0, 0, 0, 0, 0, 0 /* 4F-55 */
274     .byte '>'
275     .fill 10, 1, 0
276
    // alt 键同时按下时的映射表。
277 alt_map:
278     .byte 0, 0
279     .ascii "\0@\0$\0\0{[]}\0"
280     .byte 0, 0
281     .byte 0, 0, 0, 0, 0, 0, 0, 0, 0, 0
282     .byte '~', 13, 0
283     .byte 0, 0, 0, 0, 0, 0, 0, 0, 0, 0
284     .byte 0, 0
285     .byte 0, 0, 0, 0, 0, 0, 0, 0, 0, 0
286     .byte 0, 0, 0, 0      /* 36-39 */
287     .fill 16, 1, 0      /* 3A-49 */
288     .byte 0, 0, 0, 0, 0 /* 4A-4E */
289     .byte 0, 0, 0, 0, 0, 0 /* 4F-55 */
290     .byte '|'
291     .fill 10, 1, 0
292
293 #elif defined(KBD_US)
294
    // 以下是美式键盘的扫描码映射表。
295 key_map:
296     .byte 0, 27
297     .ascii "1234567890-="
298     .byte 127, 9
299     .ascii "qwertyuiop[]"
300     .byte 13, 0
301     .ascii "asdfghjkl;'"
302     .byte '~', 0
303     .ascii "\\zxcvbnm,./"
304     .byte 0, '*', 0, 32      /* 36-39 */
305     .fill 16, 1, 0          /* 3A-49 */
306     .byte '-', 0, 0, 0, '+' /* 4A-4E */
307     .byte 0, 0, 0, 0, 0, 0 /* 4F-55 */
308     .byte '<'
309     .fill 10, 1, 0
310
311

```

```

312 shift_map:
313     .byte 0, 27
314     .ascii "!@#$$%^&*()_+"
315     .byte 127, 9
316     .ascii "QWERTYUIOP{}"
317     .byte 13, 0
318     .ascii "ASDFGHJKL:\'"
319     .byte '~', 0
320     .ascii "|ZXCVBNM<>?"
321     .byte 0, '*', 0, 32      /* 36-39 */
322     .fill 16, 1, 0          /* 3A-49 */
323     .byte '-', 0, 0, 0, '+' /* 4A-4E */
324     .byte 0, 0, 0, 0, 0, 0 /* 4F-55 */
325     .byte '>'
326     .fill 10, 1, 0
327
328 alt_map:
329     .byte 0, 0
330     .ascii "\0@\0$\0\0{[]}\0\0"
331     .byte 0, 0
332     .byte 0, 0, 0, 0, 0, 0, 0, 0, 0, 0
333     .byte '~', 13, 0
334     .byte 0, 0, 0, 0, 0, 0, 0, 0, 0, 0
335     .byte 0, 0
336     .byte 0, 0, 0, 0, 0, 0, 0, 0, 0, 0
337     .byte 0, 0, 0, 0        /* 36-39 */
338     .fill 16, 1, 0          /* 3A-49 */
339     .byte 0, 0, 0, 0, 0     /* 4A-4E */
340     .byte 0, 0, 0, 0, 0, 0 /* 4F-55 */
341     .byte '| '
342     .fill 10, 1, 0
343
344 #elif defined(KBD_GR)
345
346 // 以下是德语键盘的扫描码映射表。
347 key_map:
348     .byte 0, 27
349     .ascii "1234567890\`\'"
350     .byte 127, 9
351     .ascii "qwertzuiop@+"
352     .byte 13, 0
353     .ascii "asdfghjkl[]^"
354     .byte 0, '#'
355     .ascii "yxcvbnm,.-"
356     .byte 0, '*', 0, 32      /* 36-39 */
357     .fill 16, 1, 0          /* 3A-49 */
358     .byte '-', 0, 0, 0, '+' /* 4A-4E */
359     .byte 0, 0, 0, 0, 0, 0 /* 4F-55 */
360     .byte '<'
361     .fill 10, 1, 0
362
363 shift_map:

```

```

364     .byte 0, 27
365     .ascii "!\"#$%&/'()=?`"
366     .byte 127, 9
367     .ascii "QWERTZUIOP\\*"
368     .byte 13, 0
369     .ascii "ASDFGHJKL {} ~"
370     .byte 0, ''
371     .ascii "YXCVBNM;:_ "
372     .byte 0, '*', 0, 32          /* 36-39 */
373     .fill 16, 1, 0              /* 3A-49 */
374     .byte '-', 0, 0, 0, '+      /* 4A-4E */
375     .byte 0, 0, 0, 0, 0, 0, 0   /* 4F-55 */
376     .byte '>'
377     .fill 10, 1, 0
378
379 alt_map:
380     .byte 0, 0
381     .ascii "\0@\0$\0\0{[]}\0\0"
382     .byte 0, 0
383     .byte '@, 0, 0, 0, 0, 0, 0, 0, 0, 0
384     .byte '~', 13, 0
385     .byte 0, 0, 0, 0, 0, 0, 0, 0, 0, 0
386     .byte 0, 0
387     .byte 0, 0, 0, 0, 0, 0, 0, 0, 0, 0
388     .byte 0, 0, 0, 0           /* 36-39 */
389     .fill 16, 1, 0            /* 3A-49 */
390     .byte 0, 0, 0, 0, 0       /* 4A-4E */
391     .byte 0, 0, 0, 0, 0, 0, 0 /* 4F-55 */
392     .byte '|'
393     .fill 10, 1, 0
394
395
396 #elif defined(KBD_FR)
397
398 // 以下是法语键盘的扫描码映射表。
399 key_map:
400     .byte 0, 27
401     .ascii "&{\'" (-)_/@="
402     .byte 127, 9
403     .ascii "azertyuiop`$"
404     .byte 13, 0
405     .ascii "qsd fghjklm|"
406     .byte '`', 0, 42          /* coin sup gauche, don't know, [*|mu] */
407     .ascii "wxcvbn,;:!"
408     .byte 0, '*', 0, 32          /* 36-39 */
409     .fill 16, 1, 0              /* 3A-49 */
410     .byte '-', 0, 0, 0, '+      /* 4A-4E */
411     .byte 0, 0, 0, 0, 0, 0, 0   /* 4F-55 */
412     .byte '<'
413     .fill 10, 1, 0
414
415 shift_map:
416     .byte 0, 27

```

```

416     .ascii "1234567890]+'"
417     .byte 127,9
418     .ascii "AZERTYUIOP<>"
419     .byte 13,0
420     .ascii "QSDFGHJKLM%"
421     .byte '~',0,'#'
422     .ascii "WXCVCBN?. /\\\"
423     .byte 0,'*',0,32      /* 36-39 */
424     .fill 16,1,0          /* 3A-49 */
425     .byte '-',0,0,0,'+'   /* 4A-4E */
426     .byte 0,0,0,0,0,0,0  /* 4F-55 */
427     .byte '>'
428     .fill 10,1,0
429
430 alt_map:
431     .byte 0,0
432     .ascii "\\0~#{[|`\\`@]}"
433     .byte 0,0
434     .byte '@',0,0,0,0,0,0,0,0,0,0
435     .byte '~',13,0
436     .byte 0,0,0,0,0,0,0,0,0,0,0
437     .byte 0,0
438     .byte 0,0,0,0,0,0,0,0,0,0,0
439     .byte 0,0,0,0        /* 36-39 */
440     .fill 16,1,0          /* 3A-49 */
441     .byte 0,0,0,0,0      /* 4A-4E */
442     .byte 0,0,0,0,0,0,0  /* 4F-55 */
443     .byte '|'
444     .fill 10,1,0
445
446 #else
447 #error "KBD-type not defined"
448 #endif
449 /*
450  * do_self handles "normal" keys, ie keys that don't change meaning
451  * and which have just one character returns.
452  */
453 do_self:
454     // 454-460 行用于根据模式标志 mode 选择 alt_map、shift_map 或 key_map 映射表之一。
455     lea alt_map,%ebx          // alt 键同时按下时的映射表基址 alt_map→ebx。
456     testb $0x20,mode         /* alt-gr */ /* 右 alt 键同时按下了? */
457     jne 1f                  // 是，则向前跳转到标号 1 处。
458     lea shift_map,%ebx       // shift 键同时按下时的映射表基址 shift_map→ebx。
459     testb $0x03,mode         // 有 shift 键同时按下了吗？
460     jne 1f                  // 有，则向前跳转到标号 1 处。
461     lea key_map,%ebx         // 否则使用普通映射表 key_map。
462     // 取映射表中对应扫描码的 ASCII 字符，若没有对应字符，则返回(转 none)。
463 1:     movb (%ebx,%eax),%al    // 将扫描码作为索引值，取对应的 ASCII 码→al。
464     orb %al,%al             // 检测看是否有对应的 ASCII 码。
465     je none                 // 若没有(对应的 ASCII 码=0)，则返回。

```

```

// 若 ctrl 键已按下或 caps 键锁定，并且字符在'a'-'}' (0x61-0x7D) 范围内，则将其转成大写字符
// (0x41-0x5D)。
464      testb $0x4c,mode          /* ctrl or caps */ /* 控制键已按下或 caps 亮? */
465      je 2f                    // 没有，则向前跳转标号 2 处。
466      cmpb $'a',%al            // 将 al 中的字符与'a'比较。
467      jb 2f                    // 若 al 值<'a'，则转标号 2 处。
468      cmpb $'}',%al           // 将 al 中的字符与'}'比较。
469      ja 2f                    // 若 al 值>'}'，则转标号 2 处。
470      subb $32,%al             // 将 al 转换为大写字符(减 0x20)。
// 若 ctrl 键已按下，并且字符在'`'-'_' (0x40-0x5F)之间(是大写字符)，则将其转换为控制字符
// (0x00-0x1F)。
471 2:      testb $0x0c,mode      /* ctrl */ /* ctrl 键同时按下了吗? */
472      je 3f                    // 若没有则转标号 3。
473      cmpb $64,%al            // 将 al 与'@' (64) 字符比较(即判断字符所属范围)。
474      jb 3f                    // 若值<'@'，则转标号 3。
475      cmpb $64+32,%al         // 将 al 与'`' (96) 字符比较(即判断字符所属范围)。
476      jae 3f                  // 若值>='`'，则转标号 3。
477      subb $64,%al            // 否则 al 值减 0x40，
// 即将字符转换为 0x00-0x1f 之间的控制字符。
// 若左 alt 键同时按下，则将字符的位 7 置位。
478 3:      testb $0x10,mode      /* left alt */ /* 左 alt 键同时按下? */
479      je 4f                    // 没有，则转标号 4。
480      orb $0x80,%al           // 字符的位 7 置位。
// 将 al 中的字符放入读缓冲队列中。
481 4:      andl $0xff,%eax        // 清 eax 的高字和 ah。
482      xorl %ebx,%ebx          // 清 ebx。
483      call put_queue          // 将字符放入缓冲队列中。
484 none:   ret
485
486 /*
487  * minus has a routine of it's own, as a 'E0h' before
488  * the scan code for minus means that the numeric keypad
489  * slash was pushed.
490  */
491 /*
492  * 减号有它自己的处理子程序，因为在减号扫描码之前的 0xe0
493  * 意味着按下了数字小键盘上的斜杠键。
494  */
491 minus:  cmpb $1,e0            // e0 标志置位了吗?
492          jne do_self           // 没有，则调用 do_self 对减号符进行普通处理。
493          movl $/,%eax          // 否则用 '/' 替换减号 '-' ➔ al。
494          xorl %ebx,%ebx
495          jmp put_queue         // 并将字符放入缓冲队列中。
496
497 /*
498  * This table decides which routine to call when a scan-code has been
499  * gotten. Most routines just call do_self, or none, depending if
500  * they are make or break.
501  */
// 下面是一张子程序地址跳转表。当取得扫描码后就根据此表调用相应的扫描码处理子程序。
// 大多数调用的子程序是 do_self，或者是 none，这取决于按键(make)还是释放键(break)。
//
502 key_table:

```

```

503 . long none, do_self, do_self, do_self /* 00-03 s0 esc 1 2 */
504 . long do_self, do_self, do_self, do_self /* 04-07 3 4 5 6 */
505 . long do_self, do_self, do_self, do_self /* 08-0B 7 8 9 0 */
506 . long do_self, do_self, do_self, do_self /* 0C-0F + ' bs tab */
507 . long do_self, do_self, do_self, do_self /* 10-13 q w e r */
508 . long do_self, do_self, do_self, do_self /* 14-17 t y u i */
509 . long do_self, do_self, do_self, do_self /* 18-1B o p } ^ */
510 . long do_self, ctrl, do_self, do_self /* 1C-1F enter ctrl a s */
511 . long do_self, do_self, do_self, do_self /* 20-23 d f g h */
512 . long do_self, do_self, do_self, do_self /* 24-27 j k l | */
513 . long do_self, do_self, lshift, do_self /* 28-2B { para lshift , */
514 . long do_self, do_self, do_self, do_self /* 2C-2F z x c v */
515 . long do_self, do_self, do_self, do_self /* 30-33 b n m , */
516 . long do_self, minus, rshift, do_self /* 34-37 . - rshift * */
517 . long alt, do_self, caps, func /* 38-3B alt sp caps f1 */
518 . long func, func, func, func /* 3C-3F f2 f3 f4 f5 */
519 . long func, func, func, func /* 40-43 f6 f7 f8 f9 */
520 . long func, num, scroll, cursor /* 44-47 f10 num scr home */
521 . long cursor, cursor, do_self, cursor /* 48-4B up pgup - left */
522 . long cursor, cursor, do_self, cursor /* 4C-4F n5 right + end */
523 . long cursor, cursor, cursor, cursor /* 50-53 dn pgdn ins del */
524 . long none, none, do_self, func /* 54-57 sysreq ? < f11 */
525 . long func, none, none, none /* 58-5B f12 ? ? ? */
526 . long none, none, none, none /* 5C-5F ? ? ? ? */
527 . long none, none, none, none /* 60-63 ? ? ? ? */
528 . long none, none, none, none /* 64-67 ? ? ? ? */
529 . long none, none, none, none /* 68-6B ? ? ? ? */
530 . long none, none, none, none /* 6C-6F ? ? ? ? */
531 . long none, none, none, none /* 70-73 ? ? ? ? */
532 . long none, none, none, none /* 74-77 ? ? ? ? */
533 . long none, none, none, none /* 78-7B ? ? ? ? */
534 . long none, none, none, none /* 7C-7F ? ? ? ? */
535 . long none, none, none, none /* 80-83 ? br br br */
536 . long none, none, none, none /* 84-87 br br br br */
537 . long none, none, none, none /* 88-8B br br br br */
538 . long none, none, none, none /* 8C-8F br br br br */
539 . long none, none, none, none /* 90-93 br br br br */
540 . long none, none, none, none /* 94-97 br br br br */
541 . long none, none, none, none /* 98-9B br br br br */
542 . long none, unctrl, none, none /* 9C-9F br unctrl br br */
543 . long none, none, none, none /* A0-A3 br br br br */
544 . long none, none, none, none /* A4-A7 br br br br */
545 . long none, none, unlshift, none /* A8-AB br br unlshift br */
546 . long none, none, none, none /* AC-AF br br br br */
547 . long none, none, none, none /* B0-B3 br br br br */
548 . long none, none, unrshift, none /* B4-B7 br br unrshift br */
549 . long unalt, none, uncaps, none /* B8-BB unalt br uncaps br */
550 . long none, none, none, none /* BC-BF br br br br */
551 . long none, none, none, none /* C0-C3 br br br br */
552 . long none, none, none, none /* C4-C7 br br br br */
553 . long none, none, none, none /* C8-CB br br br br */
554 . long none, none, none, none /* CC-CF br br br br */
555 . long none, none, none, none /* D0-D3 br br br br */

```



```

556      .long none, none, none, none          /* D4-D7 br br br br */
557      .long none, none, none, none          /* D8-DB br ? ? ? */
558      .long none, none, none, none          /* DC-DF ? ? ? ? */
559      .long none, none, none, none          /* E0-E3 e0 e1 ? ? */
560      .long none, none, none, none          /* E4-E7 ? ? ? ? */
561      .long none, none, none, none          /* E8-EB ? ? ? ? */
562      .long none, none, none, none          /* EC-EF ? ? ? ? */
563      .long none, none, none, none          /* F0-F3 ? ? ? ? */
564      .long none, none, none, none          /* F4-F7 ? ? ? ? */
565      .long none, none, none, none          /* F8-FB ? ? ? ? */
566      .long none, none, none, none          /* FC-FF ? ? ? ? */
567
568 /*
569  * kb_wait waits for the keyboard controller buffer to empty.
570  * there is no timeout - if the buffer doesn't empty, we hang.
571  */
572 /*
573  * 子程序 kb_wait 用于等待键盘控制器缓冲空。不存在超时处理 - 如果
574  * 缓冲永远不空的话，程序就会永远等待(死掉)。
575  */
576 kb_wait:
577     pushl %eax
578     1:    inb $0x64,%al                      // 读键盘控制器状态。
579     testb $0x02,%al                      // 测试输入缓冲器是否为空(等于 0)。
580     jne 1b                                // 若不空，则跳转循环等待。
581     popl %eax
582     ret
583 /*
584  * This routine reboots the machine by asking the keyboard
585  * controller to pulse the reset-line low.
586  */
587 /*
588  * 该子程序通过设置键盘控制器，向复位线输出负脉冲，使系统复位重启(reboot)。
589  */
590 reboot:
591     call kb_wait                          // 首先等待键盘控制器输入缓冲器空。
592     movw $0x1234,0x472                    /* don't do memory check */
593     movb $0xfc,%al                        /* pulse reset and A20 low */
594     outb %al,$0x64                        // 向系统复位和 A20 线输出负脉冲。
595 die:   jmp die                            // 死机。

```

## 6.7.3 其它信息

### 6.7.3.1 AT 键盘接口编程

主机系统板上所采用的键盘控制器是 intel 8042 芯片或其兼容芯片，其逻辑示意图，见下图所示。其中输出端口 P2 分别用于其它目的。位 0(P20 引脚)用于实现 CPU 的复位操作，位 1(P21 引脚)用于控制 A20 信号线的开启与否。当该输出端口位 0 为 1 时就开启(选通)了 A20 信号线，为 0 则禁止 A20 信号线。参见引导启动程序一章中对 A20 信号线的详细说明。

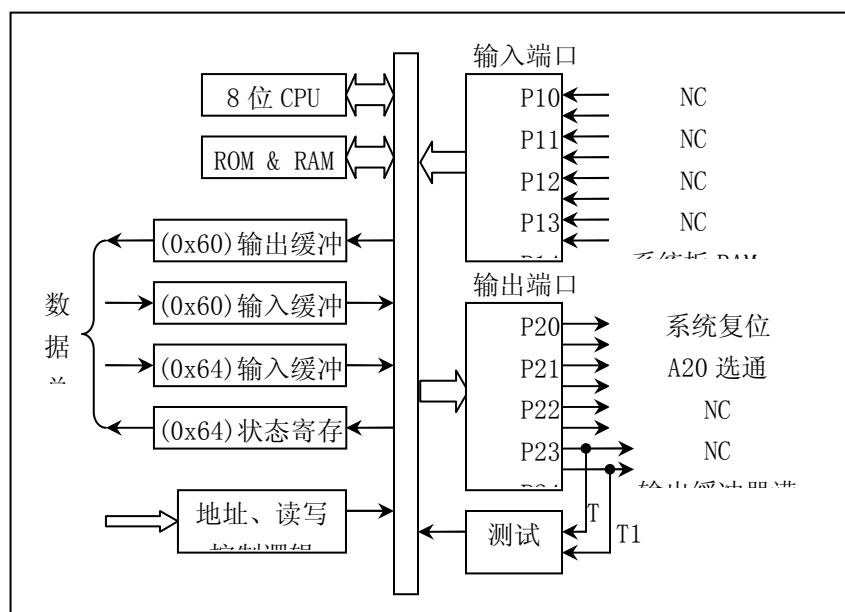


图. 键盘控制器 804X 逻辑示意图

分配给键盘控制器的 I/O 端口范围是 0x60-0x6f，但实际上 IBM CP/AT 使用的只有 0x60 和 0x64 两个口地址（0x61、0x62 和 0x63 用于与 XT 兼容目的）见下表所示，加上对端口的读和写操作含义不同，因此主要可有 4 种不同操作。对键盘控制器进行编程，将涉及芯片中的状态寄存器、输入缓冲器和输出缓冲器。

表 键盘控制器 804X 端口

| 端口   | 读/写 | 名称         | 用途   |
|------|-----|------------|--|
| 0x60 | 读   | 数据端口或输出缓冲器 | 是一个 8 位只读寄存器。当键盘控制器收到来自键盘的扫描码或命令响应时，一方面置状态寄存器位 0 = 1，另一方面产生中断 IRQ1。通常应该仅在状态端口位 0 = 1 时才读。  |
| 0x60 | 写   | 输入缓冲器      | 用于向键盘发送命令与/或随后的参数，或向键盘控制器写参数。键盘命令共有 10 多条，见表格后说明。通常都应该仅在状态端口位 1=0 时才写。   |
| 0x61 | 读/写 |            | 该端口 0x61 是 8255A 输出口 B 的地址，是针对使用/兼容 8255A 的 PC 标准键盘电路进行硬件复位处理。该端口用于对收到的扫描码做出应答。方法是首先禁止键盘，然后立刻重新允许键盘。所操作的数据为：<br>位 7=1 禁止键盘；=0 允许键盘；<br>位 6=0 迫使键盘时钟为低位，因此键盘不能发送任何数据。<br>位 5-0 这些位与键盘无关，是用于可编程并行接口 (PPI)。 |
| 0x64 | 读   | 状态寄存器      | 是一个 8 位只读寄存器，其位字段含义分别为：<br>位 7=1 来自键盘传输数据奇偶校验错；<br>位 6=1 接收超时（键盘传送未产生 IRQ1）；<br>位 5=1 发送超时（键盘无响应）；<br>位 4=1 键盘接口被键盘锁禁止；[??是=0 时]<br>位 3=1 写入输入缓冲器中的数据是命令（通过端口 0x64）；<br>=0 写入输入缓冲器中的数据是参数（通过端口 0x60）；    |

|      |   |       |   |
|------|---|-------|---|
|      |   |       | 位 2 系统标志状态: 0 = 上电启动或复位; 1 = 自检通过;<br>位 1=1 输入缓冲器满 (0x60/64 口有给 8042 的数据);<br>位 0=1 输出缓冲器满 (数据端口 0x60 有给系统的数据)。 |
| 0x64 | 写 | 输入缓冲器 | 向键盘控制器写命令。可带一参数, 参数从端口 0x60 写入。键盘控制器命令有 12 条, 见表格后说明。   |

### 6.7.3.2 键盘命令

系统在向端口 0x60 写入 1 字节, 便是发送键盘命令。键盘在接收到命令后 20ms 内应予以响应, 即回送一个命令响应。有的命令后还需要跟一参数 (也写到该端口)。命令列表见下表。注意, 如果没有另外指明, 所有命令均被回送一个 0xfa 响应码 (ACK)。

表 键盘命令一览表

| 命令码       | 参数 | 功能  |
|-----------|----|---|
| 0xed      | 有  | 设置/复位模式指示器。置 1 开启, 0 关闭。参数字节:<br>位 7-3 保留全为 0;<br>位 2 = caps-lock 键;<br>位 1 = num-lock 键;<br>位 0 = scroll-lock 键。   |
| 0xee      | 无  | 诊断回应。键盘应回送 0xee。  |
| 0xef      |    | 保留不用。   |
| 0xf0      | 有  | 读取/设置扫描码集。参数字节等于:<br>0x00 - 选择当前扫描码集;<br>0x01 - 选择扫描码集 1 (用于 PCs, PS/2 30 等);<br>0x02 - 选择扫描码集 2 (用于 AT, PS/2, 是缺省值);<br>0x03 - 选择扫描码集 3。   |
| 0xf1      |    | 保留不用。   |
| 0xf2      | 无  | 读取键盘标识号 (读取 2 个字节)。AT 键盘返回响应码 0xfa。   |
| 0xf3      | 有  | 设置扫描码连续发送时的速率和延迟时间。参数字节的含义为:<br>位 7 保留为 0;<br>位 6-5 延时值: 令 $C = \text{位 6-5}$ , 则有公式: 延时值 $= (1+C) * 250\text{ms}$ ;<br>位 4-0 扫描码连续发送的速率: 令 $B = \text{位 4-3}$ ; $A = \text{位 2-0}$ , 则有公式:<br>速率 $= 1 / ((8+A) * 2^B * 0.00417)$ 。<br>参数缺省值为 0x2c。 |
| 0xf4      | 无  | 开启键盘。   |
| 0xf5      | 无  | 禁止键盘。   |
| 0xf6      | 无  | 设置键盘默认参数。   |
| 0xf7-0xfd |    | 保留不用。   |
| 0xfe      | 无  | 重发扫描码。当系统检测到键盘传输数据有错, 则发此命令。  |
| 0xff      | 无  | 执行键盘上电复位操作, 称之为基本保证测试 (BAT)。操作过程为:<br>1. 键盘收到该命令后立刻响应发送 0xfa;<br>2. 键盘控制器使键盘时钟和数据线置为高电平;<br>3. 键盘开始执行 BAT 操作;<br>4. 若正常完成, 则键盘发送 0xaa; 否则发送 0xfd 并停止扫描。   |

### 6.7.3.3 键盘控制器命令

系统向输入缓冲(端口 0x64)写入 1 字节,即发送一键盘控制器命令。可带一参数。参数是通过写 0x60 端口发送的。

表 键盘控制器命令一览表

| 命令        | 参数 | 功能  |
|-----------|----|---|
| 0x20      | 无  | 读给键盘控制器的最后一个命令字节,放在端口 0x60 供系统读取。   |
| 0x21-0x3f | 无  | 读取由命令低 5 比特位指定的控制器内部 RAM 中的命令。  |
| 0x60-0x7f | 有  | 写键盘控制器命令字节。参数字节:(默认值为 0x5d)<br>位 7 保留为 0;<br>位 6 IBM PC 兼容模式(奇偶检验,转换为系统扫描码,单字节 PC 断开码);<br>位 5 PC 模式(对扫描码不进行奇偶校验;不转换成系统扫描码);<br>位 4 禁止键盘工作(使键盘时钟为低电平);<br>位 3 禁止超越(override),对键盘锁定转换不起作用;<br>位 2 系统标志;1 表示控制器工作正确;<br>位 1 保留为 0;<br>位 0 允许输出寄存器满中断。 |
| 0xaa      | 无  | 初始化键盘控制器自测试。成功返回 0x55;失败返回 0xfc。  |
| 0xab      | 无  | 初始化键盘接口测试。返回字节:<br>0x00 无错;<br>0x01 键盘时钟线为低(始终为低,低粘连);<br>0x02 键盘时钟线为高;<br>0x03 键盘数据线为低;<br>0x04 键盘数据线为高;   |
| 0xac      | 无  | 诊断转储。804x 的 16 字节 RAM、输出口、输入口状态依次输出给系统。   |
| 0xad      | 无  | 禁止键盘工作(设置命令字节位 4=1)。  |
| 0xae      | 无  | 允许键盘工作(复位命令字节位 4=0)。  |
| 0xc0      | 无  | 读 804x 的输入端口 P1,并放在 0x60 供读取;   |
| 0xd0      | 无  | 读 804x 的输出端口 P2,并放在 0x60 供读取;   |
| 0xd1      | 有  | 写 804x 的输出端口 P2,原 IBM PC 使用输出端口的位 2 控制 A20 门。注意,位 0(系统复位)应该总是置位的。   |
| 0xe0      | 无  | 读测试端 T0 和 T1 的输入送输出缓冲器供系统读取。<br>位 1 键盘数据;位 0 键盘时钟。  |
| 0xed      | 有  | 控制 LED 的状态。置 1 开启,0 关闭。参数字节:<br>位 7-3 保留全为 0;<br>位 2 = caps-lock 键;<br>位 1 = num-lock 键;<br>位 0 = scroll-lock 键。  |
| 0xf0-0xff | 无  | 送脉冲到输出端口。该命令序列控制输出端口 P20-23 线,参见键盘控制器逻辑示意图。欲让哪一位输出负脉冲(6 微秒),即置该位为 0。也即该命令的低 4 位分别控制负脉冲的输出。例如,若要复位系统,则需发出命令 0xfe(P20 低)即可。   |

### 6.7.3.4 键盘扫描码

PC 机采用的均是非编码键盘。键盘上每个键都有一个位置编号，是从左到右从上到下。并且 PC XT 机与 AT 机键盘的位置码差别很大。键盘内的微处理机向系统发送的是键对应的扫描码。当键按下时，键盘输出的扫描码称为接通(make)扫描码，而该键松开时发送的则称为断开(break)扫描码。

键盘上的每个键都有一个包含在字节低 7 位（位 6-0）中相应的扫描码。在高位（位 7）表示是按键还是松开按键。位 7=0 表示刚将键按下的扫描码，位 7=1 表示键松开的扫描码。例如，如果某人刚把 ESC 键按下，则传输给系统的扫描码将是 1（1 是 ESC 键的扫描码），当该键释放时将产生  $1+0x80=129$  扫描码。

对于 PC、PC/XT 的标准 83 键键盘，接通扫描码与键号（键的位置码）是一样的。并用 1 字节表示。例如“A”键，键位置号是 30，接通码是扫描码是 0x1e。而其断开码是是接通扫描码加上 0x80，即 0x9e。对于 AT 机使用的 84/101/102 扩展键盘，则与 PC/XT 标准键盘区别较大。

对于某些“扩展”的键，则情况有些不同。当一个扩展键被按下时，将产生一个中断并且键盘端口将输出一个“扩展的”的扫描码前缀 0xe0，而在下一个“中断”中将给出。比如，对于 PC/XT 标准键盘，左边的控制键 ctrl 的扫描码是 29，而右边的“扩展的”控制键 ctrl 则具有一个扩展的扫描码 29。这个规则同样适合于 alt、箭头键。

另外，还有两个键的处理是非常特殊的。PrtScn 键和 Pause/Break 键。按下 PrtScn 键将会向键盘中断程序发送\*2\*个扩展字符，42(0x2a)和 55(0x37)，所以实际的字节序列将是 0xe0, 0x2a, 0xe0, 0x37。但在键重复产生时将只发送扩展码 0x37。当键松开时，又重新发送两个扩展的加上 0x80 的码(0xe0, 0xaa, 0xe0, 0xb7)。当 prtscn 键按下时，如果 shift 或 ctrl 键也按下了，则仅发送 0xe0, 0x37，并且在松开时仅发送 0xe0, 0xb7)。

对于 Pause/Break 键。如果你在按下该键的同时也按下了控制键，则将行如扩展键 70，而在其它情况下它将发送字符序列 0xe1, 0x1d, 0x45, 0xe1, 0x9d, 0xc5。将键按下并不会产生重复的扫描码，而松开键也并不会产生任何扫描码。

因此，可以这样来看待和处理：扫描码 0xe0 意味着还有一个字符跟随其后，而扫描码 0xe1 则表示后面跟随着 2 个字符。

对于 AT 键盘的扫描码，与 PC/XT 的略有不同。当键按下时，则对应键的扫描码被送出，但当键松开时，将会发送两个字节，第一个是 0xf0，第 2 个还是相同的键扫描码。现在键盘设计者使用 8049 作为 AT 键盘的输入处理器，为了向下的兼容性将 AT 键盘发出的扫描码转换成了老式 PC/XT 标准键盘的扫描码。

AT 键盘有三种独立的扫描码集：一种是我们上面说明的(83 键映射，而增加的键有多余的 0xe0 码)，一种几乎是顺序的，还有一种却只有 1 个字节！最后一种所带来的问题是只有左 shift, caps, 左 ctrl 和左 alt 键的松开码被发送。键盘的默认扫描码集是扫描码集 2，可以利用命令更改。

对于扫描码集 1 和 2，有特殊码 0xe0 和 0xe1。它们用于具有相同功能的键。比如：左控制键 ctrl 位置是 0x1d(对于 PC/XT)，则右边的控制键就是 0xe0, 0x1d。这是为了与 PC/XT 程序兼容。请注意唯一使用 0xe1 的时候是当它表示临时控制键时，对此情况同时也有一个 0xe0 的版本。

The XT Scancodes



```

10  * This module implements the console io functions
11  *      'void con_init(void)'
12  *      'void con_write(struct tty_queue * queue)'
13  * Hopefully this will be a rather complete VT102 implementation.
14  *
15  * Beeping thanks to John T Kohl.
16  */
/*
 * 该模块实现控制台输入输出功能
 *      'void con_init(void)'
 *      'void con_write(struct tty_queue * queue)'
 * 希望这是一个非常完整的 VT102 实现。
 *
 * 感谢 John T Kohl 实现了蜂鸣指示。
 */
17
18 /*
19  * NOTE!!! We sometimes disable and enable interrupts for a short while
20  * (to put a word in video IO), but this will work even for keyboard
21  * interrupts. We know interrupts aren't enabled when getting a keyboard
22  * interrupt, as we use trap-gates. Hopefully all is well.
23  */
/*
 * 注意!!! 我们有时短暂地禁止和允许中断(在将一个字(word)放到视频 IO), 但即使
 * 对于键盘中断这也是可以工作的。因为我们使用陷阱门, 所以我们知道在获得一个
 * 键盘中断时中断是不允许的。希望一切均正常。
 */
24
25 /*
26  * Code to check for different video-cards mostly by Galen Hunt,
27  * <g-hunt@ee.utah.edu>
28  */
/*
 * 检测不同显示卡的代码大多数是 Galen Hunt 编写的,
 * <g-hunt@ee.utah.edu>
 */
29
30 #include <linux/sched.h> // 调度程序头文件, 定义了任务结构 task_struct、初始任务 0 的数据,
                        // 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
31 #include <linux/tty.h> // tty 头文件, 定义了有关 tty_io, 串行通信方面的参数、常数。
32 #include <asm/io.h> // io 头文件。定义硬件端口输入/输出宏汇编语句。
33 #include <asm/system.h> // 系统头文件。定义了设置或修改描述符/中断门等的嵌入式汇编宏。
34
35 /*
36  * These are set up by the setup-routine at boot-time:
37  */
/*
 * 这些是设置子程序 setup 在引导启动系统时设置的参数:
 */
38
// 参见对 boot/setup.s 的注释, 和 setup 程序读取并保留的参数表。
39 #define ORIG_X (* (unsigned char *) 0x90000) // 光标列号。
40 #define ORIG_Y (* (unsigned char *) 0x90001) // 光标行号。

```

```

41 #define ORIG_VIDEO_PAGE          (*(unsigned short *)0x90004)    // 显示页面。
42 #define ORIG_VIDEO_MODE          ((*(unsigned short *)0x90006) & 0xff) // 显示模式。
43 #define ORIG_VIDEO_COLS          (((*(unsigned short *)0x90006) & 0xff00) >> 8) // 字符列数。
44 #define ORIG_VIDEO_LINES         (25)                          // 显示行数。
45 #define ORIG_VIDEO_EGA_AX        (*(unsigned short *)0x90008)    // [??]
46 #define ORIG_VIDEO_EGA_BX        (*(unsigned short *)0x9000a)    // 显示内存大小和色彩模式。
47 #define ORIG_VIDEO_EGA_CX        (*(unsigned short *)0x9000c)    // 显示卡特性参数。
48
// 定义显示器单色/彩色显示模式类型符号常数。
49 #define VIDEO_TYPE_MDA            0x10    /* Monochrome Text Display */ /* 单色文本 */
50 #define VIDEO_TYPE_CGA            0x11    /* CGA Display */ /* CGA 显示器 */
51 #define VIDEO_TYPE_EGAM           0x20    /* EGA/VGA in Monochrome Mode */ /* EGA/VGA 单色 */
52 #define VIDEO_TYPE_EGAC           0x21    /* EGA/VGA in Color Mode */ /* EGA/VGA 彩色 */
53
54 #define NPAR 16
55
56 extern void keyboard_interrupt(void); // 键盘中断处理程序(keyboard.S)。
57
58 static unsigned char    video_type;      /* Type of display being used */
/* 使用的显示类型 */
59 static unsigned long    video_num_columns; /* Number of text columns */
/* 屏幕文本列数 */
60 static unsigned long    video_size_row;   /* Bytes per row */
/* 每行使用的字节数 */
61 static unsigned long    video_num_lines; /* Number of test lines */
/* 屏幕文本行数 */
62 static unsigned char    video_page;      /* Initial video page */
/* 初始显示页面 */
63 static unsigned long    video_mem_start; /* Start of video RAM */
/* 显示内存起始地址 */
64 static unsigned long    video_mem_end;   /* End of video RAM (sort of) */
/* 显示内存结束(末端)地址 */
65 static unsigned short   video_port_reg;  /* Video register select port */
/* 显示控制索引寄存器端口 */
66 static unsigned short   video_port_val;  /* Video register value port */
/* 显示控制数据寄存器端口 */
67 static unsigned short   video_erase_char; /* Char+Attrib to erase with */
/* 擦除字符属性与字符(0x0720) */
68
// 以下这些变量用于屏幕滚屏操作。
69 static unsigned long    origin;           /* Used for EGA/VGA fast scroll */ // scr_start。
/* 用于 EGA/VGA 快速滚屏 */ // 滚屏起始内存地址。
70 static unsigned long    scr_end;         /* Used for EGA/VGA fast scroll */
/* 用于 EGA/VGA 快速滚屏 */ // 滚屏末端内存地址。
71 static unsigned long    pos;             // 当前光标对应的显示内存位置。
72 static unsigned long    x,y;            // 当前光标位置。
73 static unsigned long    top,bottom;      // 滚动时顶行行号; 底行行号。
// state 用于标明处理 ESC 转义序列时的当前步骤。npar, par[] 用于存放 ESC 序列的中间处理参数。
74 static unsigned long    state=0;         // ANSI 转义字符序列处理状态。
75 static unsigned long    npar, par[NPAR]; // ANSI 转义字符序列参数个数和参数数组。
76 static unsigned long    ques=0;
77 static unsigned char    attr=0x07;      // 字符属性(黑底白字)。
78

```



```

79 static void sysbeep(void);           // 系统蜂鸣函数。
80
81 /*
82  * this is what the terminal answers to a ESC-Z or csi0c
83  * query (= vt100 response).
84  */
85 /*
86  * 下面是终端回应 ESC-Z 或 csi0c 请求的应答(=vt100 响应)。
87  */
88 // csi - 控制序列引导码(Control Sequence Introducer)。
89 #define RESPONSE "\033[?1;2c"
90
91 /* NOTE! gotoxy thinks x==video_num_columns is ok */
92 /* 注意! gotoxy 函数认为 x==video_num_columns, 这是正确的 */
93 // 跟踪光标当前位置。
94 // 参数: new_x - 光标所在列号; new_y - 光标所在行号。
95 // 更新当前光标位置变量 x, y, 并修正 pos 指向光标在显示内存中的对应位置。
96 static inline void gotoxy(unsigned int new_x, unsigned int new_y)
97 {
98     // 如果输入的光标行号超出显示器列数, 或者光标行号超出显示的最大行数, 则退出。
99     if (new_x > video\_num\_columns || new_y >= video\_num\_lines)
100         return;
101     // 更新当前光标变量; 更新光标位置对应的在显示内存中位置变量 pos。
102     x=new_x;
103     y=new_y;
104     pos=origin + y*video\_size\_row + (x<<1);
105 }
106
107 // 设置滚屏起始显示内存地址。
108 static inline void set\_origin(void)
109 {
110     cli();
111     // 首先选择显示控制数据寄存器 r12, 然后写入滚屏起始地址高字节。向右移动 9 位, 表示向右移动
112     // 8 位, 再除以 2(2 字节代表屏幕上 1 字符)。是相对于默认显示内存操作的。
113     outb\_p(12, video\_port\_reg);
114     outb\_p(0xff&((origin-video\_mem\_start)>>9), video\_port\_val);
115     // 再选择显示控制数据寄存器 r13, 然后写入滚屏起始地址底字节。向右移动 1 位表示除以 2。
116     outb\_p(13, video\_port\_reg);
117     outb\_p(0xff&((origin-video\_mem\_start)>>1), video\_port\_val);
118     sti();
119 }
120
121 // 向上卷动一行(屏幕窗口向下移动)。
122 // 将屏幕窗口向下移动一行。参见程序列表后说明。
123 static void scrup(void)
124 {
125     // 如果显示类型是 EGA, 则执行以下操作。
126     if (video\_type == VIDEO\_TYPE\_EGAC || video\_type == VIDEO\_TYPE\_EGAM)
127     {
128         // 如果移动起始行 top=0, 移动最底行 bottom=video_num_lines=25, 则表示整屏窗口向下移动。
129         if (!top && bottom == video\_num\_lines) {
130             // 调整屏幕显示对应内存的起始位置指针 origin 为向下移一行屏幕字符对应的内存位置, 同时也调整
131             // 当前光标对应的内存位置以及屏幕末行末端字符指针 scr_end 的位置。

```

```

112         origin += video_size_row;
113         pos += video_size_row;
114         scr_end += video_size_row;
// 如果屏幕末端最后一个显示字符所对应的显示内存指针 scr_end 超出了实际显示内存的末端, 则将
// 屏幕内容内存数据移动到显示内存的起始位置 video_mem_start 处, 并在出现的新行上填入空格字符。
115         if (scr_end > video_mem_end) {
// %0 - eax(擦除字符+属性); %1 - ecx((显示器字符行数-1)所对应的字符数/2, 是以长字移动);
// %2 - edi(显示内存起始位置 video_mem_start); %3 - esi(屏幕内容对应的内存起始位置 origin)。
// 移动方向: [edi]→[esi], 移动 ecx 个长字。
116             __asm__( "cld\n\t"           // 清方向位。
117                     "rep\n\t"           // 重复操作, 将当前屏幕内存数据
118                     "movsl\n\t"         // 移动到显示内存起始处。
119                     "movl _video_num_columns, %1\n\t" // ecx=1 行字符数。
120                     "rep\n\t"           // 在新行上填入空格字符。
121                     "stosw"
122                     :: "a" (video_erase_char),
123                     "c" ((video_num_lines-1)*video_num_columns>>1),
124                     "D" (video_mem_start),
125                     "S" (origin)
126                     : "cx", "di", "si");
// 根据屏幕内存数据移动后的情况, 重新调整当前屏幕对应内存的起始指针、光标位置指针和屏幕末端
// 对应内存指针 scr_end。
127             scr_end -= origin-video_mem_start;
128             pos -= origin-video_mem_start;
129             origin = video_mem_start;
130         } else {
// 如果调整后的屏幕末端对应的内存指针 scr_end 没有超出显示内存的末端 video_mem_end, 则只需在
// 新行上填入擦除字符(空格字符)。
// %0 - eax(擦除字符+属性); %1 - ecx(显示器字符行数); %2 - edi(屏幕对应内存最后一行开始处);
131         __asm__( "cld\n\t"           // 清方向位。
132                 "rep\n\t"           // 重复操作, 在新出现行上
133                 "stosw"             // 填入擦除字符(空格字符)。
134                 :: "a" (video_erase_char),
135                 "c" (video_num_columns),
136                 "D" (scr_end-video_size_row)
137                 : "cx", "di");
138     }
// 向显示控制器中写入新的屏幕内容对应的内存起始位置值。
139     set_origin();
// 则表示不是整屏移动。也即表示从指定行 top 开始的所有行向上移动 1 行(删除 1 行)。此时直接
// 将屏幕从指定行 top 到屏幕末端所有行对应的显示内存数据向上移动 1 行, 并在新出现的行上填入擦
// 除字符。
// %0-eax(擦除字符+属性); %1-ecx(top 行下 1 行开始到屏幕末行的行数所对应的内存长字数);
// %2-edi(top 行所处的内存位置); %3-esi(top+1 行所处的内存位置)。
140     } else {
141         __asm__( "cld\n\t"           // 清方向位。
142                 "rep\n\t"           // 循环操作, 将 top+1 到 bottom 行
143                 "movsl\n\t"         // 所对应的内存块移到 top 行开始处。
144                 "movl _video_num_columns, %%ecx\n\t" // ecx = 1 行字符数。
145                 "rep\n\t"           // 在新行上填入擦除字符。
146                 "stosw"
147                 :: "a" (video_erase_char),
148                 "c" ((bottom-top-1)*video_num_columns>>1),

```

```

149         "D" (origin+video_size row*top),
150         "S" (origin+video_size row*(top+1))
151         : "cx", "di", "si");
152     }
153 }
// 如果显示类型不是 EGA(是 MDA), 则执行下面移动操作。因为 MDA 显示控制卡会自动调整超出显示范围
// 的情况, 也即会自动翻卷指针, 所以这里不对屏幕内容对应内存超出显示内存的情况单独处理。处理
// 方法与 EGA 非整屏移动情况完全一样。
154     else                /* Not EGA/VGA */
155     {
156         __asm__( "cld\n\t"
157                 "rep\n\t"
158                 "movsl\n\t"
159                 "movl _video_num_columns, %%ecx\n\t"
160                 "rep\n\t"
161                 "stosw"
162                 :: "a" (video_erase_char),
163                 "c" ((bottom-top-1)*video_num_columns>>1),
164                 "D" (origin+video_size row*top),
165                 "S" (origin+video_size row*(top+1))
166                 : "cx", "di", "si");
167     }
168 }
169
//// 向下卷动一行(屏幕窗口向上移动)。
// 将屏幕窗口向上移动一行, 屏幕显示的内容向下移动 1 行, 在被移动开始行的上方出现一新行。参见
// 程序列表后说明。处理方法与 scrup() 相似, 只是为了在移动显示内存数据时不出现数据覆盖错误情
// 况, 复制是以反方向进行的, 也即从屏幕倒数第 2 行的最后一个字符开始复制
170 static void scrdown(void)
171 {
// 如果显示类型是 EGA, 则执行下列操作。
// [??好象 if 和 else 的操作完全一样啊! 为什么还要分别处理呢? 难道与任务切换有关?]
172     if (video_type == VIDEO_TYPE EGAC || video_type == VIDEO_TYPE EGAM)
173     {
// %0-eax(擦除字符+属性); %1-ecx(top 行开始到屏幕末行-1 行的行数所对应的内存长字数);
// %2-edi(屏幕右下角最后一个长字位置); %3-esi(屏幕倒数第 2 行最后一个长字位置)。
// 移动方向: [esi]→[edi], 移动 ecx 个长字。
174         __asm__( "std\n\t"                // 置方向位。
175                 "rep\n\t"                // 重复操作, 向下移动从 top 行到 bottom-1 行
176                 "movsl\n\t"              // 对应的内存数据。
177                 "addl $2, %%edi\n\t"      /* %edi has been decremented by 4 */
178                                           /* %edi 已经减 4, 因为也是方向填擦除字符 */
179                 "movl _video_num_columns, %%ecx\n\t" // 置 ecx=1 行字符数。
180                 "rep\n\t"                // 将擦除字符填入上方新行中。
181                 "stosw"
182                 :: "a" (video_erase_char),
183                 "c" ((bottom-top-1)*video_num_columns>>1),
184                 "D" (origin+video_size row*bottom-4),
185                 "S" (origin+video_size row*(bottom-1)-4)
186                 : "ax", "cx", "di", "si");
187     }
// 如果不是 EGA 显示类型, 则执行以下操作(目前与上面完全一样)。
188     else                /* Not EGA/VGA */

```

```

188     {
189         __asm__( "std\n\t"
190                "rep\n\t"
191                "movsl\n\t"
192                "addl $2,%edi\n\t" /* %edi has been decremented by 4 */
193                "movl _video_num_columns,%ecx\n\t"
194                "rep\n\t"
195                "stosw"
196                ":: \"a\" (video_erase_char),
197                \"c\" ((bottom-top-1)*video_num_columns>>1),
198                \"D\" (origin+video_size_row*bottom-4),
199                \"S\" (origin+video_size_row*(bottom-1)-4)
200                : \"ax\", \"cx\", \"di\", \"si\" );
201     }
202 }
203
204 // 光标位置下移一行(lf - line feed 换行)。
205 static void lf(void)
206 {
207     // 如果光标没有处在倒数第2行之后，则直接修改光标当前行变量 y++，并调整光标对应显示内存位置
208     // pos(加上屏幕一行字符所对应的内存长度)。
209     if (y+1<bottom) {
210         y++;
211         pos += video_size_row;
212         return;
213     }
214     // 否则需要将屏幕内容上移一行。
215     scrup();
216 }
217
218 // 光标上移一行(ri - reverse line feed 反向换行)。
219 static void ri(void)
220 {
221     // 如果光标不在第1行上，则直接修改光标当前行标量 y--，并调整光标对应显示内存位置 pos，减去
222     // 屏幕上一行字符所对应的内存长度字节数。
223     if (y>top) {
224         y--;
225         pos -= video_size_row;
226         return;
227     }
228     // 否则需要将屏幕内容下移一行。
229     scrdwn();
230 }
231
232 // 光标回到第1列(0列)左端(cr - carriage return 回车)。
233 static void cr(void)
234 {
235     // 光标所在的列号*2 即 0 列到光标所在列对应的内存字节长度。
236     pos -= x<<1;
237     x=0;
238 }
239
240 // 擦除光标前一字符(用空格替代)(del - delete 删除)。

```

```

230 static void del(void)
231 {
    // 如果光标没有处在 0 列，则将光标对应内存位置指针 pos 后退 2 字节(对应屏幕上一个字符)，然后
    // 将当前光标变量列值减 1，并将光标所在位置字符擦除。
232     if (x) {
233         pos -= 2;
234         x--;
235         *(unsigned short *)pos = video_erase_char;
236     }
237 }
238
    //// 删除屏幕上与光标位置相关的部分，以屏幕为单位。csi - 控制序列引导码(Control Sequence
    // Introducer)。
    // ANSI 转义序列: 'ESC [sJ' (s = 0 删除光标到屏幕底端; 1 删除屏幕开始到光标处; 2 整屏删除)。
    // 参数: par - 对应上面 s。
239 static void csi_J(int par)
240 {
241     long count __asm__("cx");    // 设为寄存器变量。
242     long start __asm__("di");
243
    // 首先根据三种情况分别设置需要删除的字符数和删除开始的显示内存位置。
244     switch (par) {
245         case 0: /* erase from cursor to end of display */ /* 擦除光标到屏幕底端 */
246             count = (scr_end-pos)>>1;
247             start = pos;
248             break;
249         case 1: /* erase from start to cursor */ /* 删除从屏幕开始到光标处的字符 */
250             count = (pos-origin)>>1;
251             start = origin;
252             break;
253         case 2: /* erase whole display */ /* 删除整个屏幕上的字符 */
254             count = video_num_columns * video_num_lines;
255             start = origin;
256             break;
257         default:
258             return;
259     }
    // 然后使用擦除字符填写删除字符的地方。
    // %0 - ecx(要删除的字符数 count); %1 - edi(删除操作开始地址); %2 - eax(填入的擦除字符)。
260     __asm__("cld\n\t"
261            "rep\n\t"
262            "stosw\n\t"
263            "::\"c\" (count),
264            \"D\" (start), \"a\" (video_erase_char)
265            :\"cx\", \"di\");
266 }
267
    //// 删除行内与光标位置相关的部分，以一行为单位。
    // ANSI 转义字符序列: 'ESC [sK' (s = 0 删除到行尾; 1 从开始删除; 2 整行都删除)。
268 static void csi_K(int par)
269 {
270     long count __asm__("cx");    // 设置寄存器变量。
271     long start __asm__("di");

```

```

272 // 首先根据三种情况分别设置需要删除的字符数和删除开始的显示内存位置。
273 switch (par) {
274     case 0: /* erase from cursor to end of line */ /* 删除光标到行尾字符 */
275         if (x >= video_num_columns)
276             return;
277         count = video_num_columns - x;
278         start = pos;
279         break;
280     case 1: /* erase from start of line to cursor */ /* 删除从行开始到光标处 */
281         start = pos - (x << 1);
282         count = (x < video_num_columns) ? x : video_num_columns;
283         break;
284     case 2: /* erase whole line */ /* 将整行字符全删除 */
285         start = pos - (x << 1);
286         count = video_num_columns;
287         break;
288     default:
289         return;
290 }
291 // 然后使用擦除字符填写删除字符的地方。
292 // %0 - ecx(要删除的字符数 count); %1 - edi(删除操作开始地址); %2 - eax(填入的擦除字符)。
293 __asm__ ("cld\n|t"
294         "rep\n|t"
295         "stosw\n|t"
296         ":: \"c\" (count),
297         "D\" (start), \"a\" (video_erase_char)
298         : \"cx\", \"di\");
299
300 //// 允许翻译(重显)(允许重新设置字符显示方式, 比如加粗、加下划线、闪烁、反显等)。
301 // ANSI 转义字符序列: 'ESC [nm'。n = 0 正常显示; 1 加粗; 4 加下划线; 7 反显; 27 正常显示。
302 void csi_m(void)
303 {
304     int i;
305     for (i=0; i<=npar; i++)
306         switch (par[i]) {
307             case 0: attr=0x07; break;
308             case 1: attr=0x0f; break;
309             case 4: attr=0x0f; break;
310             case 7: attr=0x70; break;
311             case 27: attr=0x07; break;
312         }
313
314     //// 根据设置显示光标。
315     // 根据显示内存光标对应位置 pos, 设置显示控制器光标的显示位置。
316 static inline void set_cursor(void)
317 {
318     cli();
319     // 首先使用索引寄存器端口选择显示控制数据寄存器 r14(光标当前显示位置高字节), 然后写入光标
320     // 当前位置高字节(向右移动 9 位表示高字节移到低字节再除以 2)。是相对于默认显示内存操作的。

```

```

316     outb_p(14, video_port_reg);
317     outb_p(0xff & ((pos-video_mem_start)>>9), video_port_val);
    // 再使用索引寄存器选择 r15, 并将光标当前位置低字节写入其中。
318     outb_p(15, video_port_reg);
319     outb_p(0xff & ((pos-video_mem_start)>>1), video_port_val);
320     sti();
321 }
322
    // 发送对终端 VT100 的响应序列。
    // 将响应序列放入读缓冲队列中。
323 static void respond(struct tty_struct * tty)
324 {
325     char * p = RESPONSE;
326
327     cli();                // 关中断。
328     while (*p) {          // 将字符序列放入写队列。
329         PUTCH(*p, tty->read_q);
330         p++;
331     }
332     sti();                // 开中断。
333     copy_to_cooked(tty);  // 转换成规范模式(放入辅助队列中)。
334 }
335
    // 在光标处插入一空格字符。
336 static void insert_char(void)
337 {
338     int i=x;
339     unsigned short tmp, old = video_erase_char;
340     unsigned short * p = (unsigned short *) pos;
341
    // 光标开始的所有字符右移一格, 并将擦除字符插入在光标所在处。
    // 若一行上都有字符的话, 则行最后一个字符将不会更动☺?
342     while (i++<video_num_columns) {
343         tmp=*p;
344         *p=old;
345         old=tmp;
346         p++;
347     }
348 }
349
    // 在光标处插入一行(则光标将处在新的空行上)。
    // 将屏幕从光标所在行到屏幕底向下卷动一行。
350 static void insert_line(void)
351 {
352     int oldtop, oldbottom;
353
354     oldtop=top;            // 保存原 top, bottom 值。
355     oldbottom=bottom;
356     top=y;                // 设置屏幕卷动开始行。
357     bottom = video_num_lines; // 设置屏幕卷动最后行。
358     scrdown();            // 从光标开始处, 屏幕内容向下滚动一行。
359     top=oldtop;            // 恢复原 top, bottom 值。
360     bottom=oldbottom;

```

```

361 }
362
363 // 删除光标处的一个字符。
364 static void delete_char(void)
365 {
366     int i;
367     unsigned short * p = (unsigned short *) pos;
368
369     // 如果光标超出屏幕最右列，则返回。
370     if (x >= video_num_columns)
371         return;
372     // 从光标右一个字符开始到行末所有字符左移一格。
373     i = x;
374     while (++i < video_num_columns) {
375         *p = *(p+1);
376         p++;
377     }
378     // 最后一个字符处填入擦除字符(空格字符)。
379     *p = video_erase_char;
380 }
381
382 // 删除光标所在行。
383 // 从光标所在行开始屏幕内容上卷一行。
384 static void delete_line(void)
385 {
386     int oldtop, oldbottom;
387
388     oldtop = top; // 保存原 top, bottom 值。
389     oldbottom = bottom;
390     top = y; // 设置屏幕卷动开始行。
391     bottom = video_num_lines; // 设置屏幕卷动最后行。
392     scrup(); // 从光标开始处，屏幕内容向上滚动一行。
393     top = oldtop; // 恢复原 top, bottom 值。
394     bottom = oldbottom;
395 }
396
397 // 在光标处插入 nr 个字符。
398 // ANSI 转义字符序列: 'ESC [n@'。
399 // 参数 nr = 上面 n。
400 static void csi_at(unsigned int nr)
401 {
402     // 如果插入的字符数大于一行字符数，则截为一行字符数；若插入字符数 nr 为 0，则插入 1 个字符。
403     if (nr > video_num_columns)
404         nr = video_num_columns;
405     else if (!nr)
406         nr = 1;
407     // 循环插入指定的字符数。
408     while (nr--)
409         insert_char();
410 }
411
412 // 在光标位置处插入 nr 行。
413 // ANSI 转义字符序列'ESC [nL'。

```



```

401 static void csi L(unsigned int nr)
402 {
403     // 如果插入的行数大于屏幕最多行数，则截为屏幕显示行数；若插入行数 nr 为 0，则插入 1 行。
404     if (nr > video\_num\_lines)
405         nr = video\_num\_lines;
406     else if (!nr)
407         nr = 1;
408     // 循环插入指定行数 nr。
409     while (nr--)
410         insert\_line();
411 }
412
413 // 删除光标处的 nr 个字符。
414 // ANSI 转义序列：'ESC [nP'。
415 static void csi P(unsigned int nr)
416 {
417     // 如果删除的字符数大于一行字符数，则截为一行字符数；若删除字符数 nr 为 0，则删除 1 个字符。
418     if (nr > video\_num\_columns)
419         nr = video\_num\_columns;
420     else if (!nr)
421         nr = 1;
422     // 循环删除指定字符数 nr。
423     while (nr--)
424         delete\_char();
425 }
426
427 // 删除光标处的 nr 行。
428 // ANSI 转义序列：'ESC [nM'。
429 static void csi M(unsigned int nr)
430 {
431     // 如果删除的行数大于屏幕最多行数，则截为屏幕显示行数；若删除的行数 nr 为 0，则删除 1 行。
432     if (nr > video\_num\_lines)
433         nr = video\_num\_lines;
434     else if (!nr)
435         nr=1;
436     // 循环删除指定行数 nr。
437     while (nr--)
438         delete\_line();
439 }
440
441 static int saved\_x=0;           // 保存的光标列号。
442 static int saved\_y=0;           // 保存的光标行号。
443
444 // 保存当前光标位置。
445 static void save\_cur(void)
446 {
447     saved\_x=x;
448     saved\_y=y;
449 }
450
451 // 恢复保存的光标位置。
452 static void restore\_cur(void)
453 {

```

```

442         gotoxy(saved_x, saved_y);
443     }
444
445     // 控制台写函数。
446     // 从终端对应的 tty 写缓冲队列中取字符，并显示在屏幕上。
447 void con_write(struct tty_struct * tty)
448 {
449     int nr;
450     char c;
451
452     // 首先取得写缓冲队列中现有字符数 nr，然后针对每个字符进行处理。
453     nr = CHARS(tty->write_q);
454     while (nr-->0) {
455         // 从写队列中取一字符 c，根据前面所处理字符的状态 state 分别处理。状态之间的转换关系为：
456         // state = 0: 初始状态；或者原是状态 4；或者原是状态 1，但字符不是 '['；
457         // 1: 原是状态 0，并且字符是转义字符 ESC(0x1b = 033 = 27)；
458         // 2: 原是状态 1，并且字符是 '['；
459         // 3: 原是状态 2；或者原是状态 3，并且字符是 ';' 或数字。
460         // 4: 原是状态 3，并且字符不是 ';' 或数字；
461         GETCH(tty->write_q, c);
462         switch(state) {
463             case 0:
464                 // 如果字符不是控制字符(c>31)，并且也不是扩展字符(c<127)，则
465                 if (c>31 && c<127) {
466                     // 若当前光标处在行末端或末端以外，则将光标移到下行头列。并调整光标位置对应的内存指针 pos。
467                     if (x>=video_num_columns) {
468                         x -= video_num_columns;
469                         pos -= video_size_row;
470                         lf();
471                     }
472                     // 将字符 c 写到显示内存中 pos 处，并将光标右移 1 列，同时也将 pos 对应地移动 2 个字节。
473                     __asm__( "movb _attr, %%ah\n\t"
474                             "movw %%ax, %I\n\t"
475                             :: "a" (c), "m" (*(short *)pos)
476                             : "ax");
477                     pos += 2;
478                     x++;
479                     // 如果字符 c 是转义字符 ESC，则转换状态 state 到 1。
480                     } else if (c==27)
481                         state=1;
482                     // 如果字符 c 是换行符(10)，或是垂直制表符 VT(11)，或者是换页符 FF(12)，则移动光标到下一行。
483                     else if (c==10 || c==11 || c==12)
484                         lf();
485                     // 如果字符 c 是回车符 CR(13)，则将光标移动到头列(0 列)。
486                     else if (c==13)
487                         cr();
488                     // 如果字符 c 是 DEL(127)，则将光标右边一字符擦除(用空格字符替代)，并将光标移到被擦除位置。
489                     else if (c==ERASE_CHAR(tty))
490                         del();
491                     // 如果字符 c 是 BS(backspace, 8)，则将光标右移 1 格，并相应调整光标对应内存位置指针 pos。
492                     else if (c==8) {
493                         if (x) {
494                             x--;

```

```

478                                     pos -= 2;
479                                     }
// 如果字符 c 是水平制表符 TAB(9)，则将光标移到 8 的倍数列上。若此时光标列数超出屏幕最大列数，
// 则将光标移到下一行上。
480                                     } else if (c==9) {
481                                         c=8-(x&7);
482                                         x += c;
483                                         pos += c<<1;
484                                         if (x>video_num_columns) {
485                                             x -= video_num_columns;
486                                             pos -= video_size_row;
487                                             lf();
488                                         }
489                                         c=9;
// 如果字符 c 是响铃符 BEL(7)，则调用蜂鸣函数，是扬声器发声。
490                                     } else if (c==7)
491                                         sysbeep();
492                                     break;
// 如果原状态是 0，并且字符是转义字符 ESC(0x1b = 033 = 27)，则转到状态 1 处理。
493                                     case 1:
494                                         state=0;
// 如果字符 c 是 '['，则将状态 state 转到 2。
495                                         if (c=='[')
496                                             state=2;
// 如果字符 c 是 'E'，则光标移到下一行开始处(0 列)。
497                                         else if (c=='E')
498                                             gotoxy(0, y+1);
// 如果字符 c 是 'M'，则光标上移一行。
499                                         else if (c=='M')
500                                             ri();
// 如果字符 c 是 'D'，则光标下移一行。
501                                         else if (c=='D')
502                                             lf();
// 如果字符 c 是 'Z'，则发送终端应答字符序列。
503                                         else if (c=='Z')
504                                             respond(tty);
// 如果字符 c 是 '7'，则保存当前光标位置。注意这里代码写错！应该是(c=='7')。
505                                         else if (x=='7')
506                                             save_cur();
// 如果字符 c 是 '8'，则恢复到原保存的光标位置。注意这里代码写错！应该是(c=='8')。
507                                         else if (x=='8')
508                                             restore_cur();
509                                     break;
// 如果原状态是 1，并且上一字符是 '['，则转到状态 2 来处理。
510                                     case 2:
// 首先对 ESC 转义字符序列参数使用的处理数组 par[]清零，索引变量 npar 指向首项，并且设置状态
// 为 3。若此时字符不是 '?'，则直接转到状态 3 去处理，否则去读一字符，再到状态 3 处理代码处。
511                                     for(npar=0;npar<NPAR;npar++)
512                                         par[npar]=0;
513                                     npar=0;
514                                     state=3;
515                                     if (ques==(c=='?'))
516                                         break;

```

```

// 如果原来是状态 2; 或者原来就是状态 3, 但原字符是 ';' 或数字, 则在下面处理。
517         case 3:
// 如果字符 c 是分号 ';', 并且数组 par 未滿, 则索引值加 1。
518             if (c==';' && npar<NPAR-1) {
519                 npar++;
520                 break;
// 如果字符 c 是数字字符 '0'-'9', 则将该字符转换成数值并与 npar 所索引的项组成 10 进制数。
521             } else if (c>='0' && c<='9') {
522                 par[npar]=10*par[npar]+c-'0';
523                 break;
// 否则转到状态 4。
524             } else state=4;
// 如果原状态是状态 3, 并且字符不是 ';' 或数字, 则转到状态 4 处理。首先复位状态 state=0。
525         case 4:
526             state=0;
527             switch(c) {
// 如果字符 c 是 'G' 或 '`', 则 par[] 中第一个参数代表列号。若列号不为零, 则将光标右移一格。
528                 case 'G': case '`':
529                     if (par[0]) par[0]--;
530                     gotoxy(par[0], y);
531                     break;
// 如果字符 c 是 'A', 则第一个参数代表光标上移的行数。若参数为 0 则上移一行。
532                 case 'A':
533                     if (!par[0]) par[0]++;
534                     gotoxy(x, y-par[0]);
535                     break;
// 如果字符 c 是 'B' 或 'e', 则第一个参数代表光标下移的行数。若参数为 0 则下移一行。
536                 case 'B': case 'e':
537                     if (!par[0]) par[0]++;
538                     gotoxy(x, y+par[0]);
539                     break;
// 如果字符 c 是 'C' 或 'a', 则第一个参数代表光标右移的格数。若参数为 0 则右移一格。
540                 case 'C': case 'a':
541                     if (!par[0]) par[0]++;
542                     gotoxy(x+par[0], y);
543                     break;
// 如果字符 c 是 'D', 则第一个参数代表光标左移的格数。若参数为 0 则左移一格。
544                 case 'D':
545                     if (!par[0]) par[0]++;
546                     gotoxy(x-par[0], y);
547                     break;
// 如果字符 c 是 'E', 则第一个参数代表光标向下移动的行数, 并回到 0 列。若参数为 0 则下移一行。
548                 case 'E':
549                     if (!par[0]) par[0]++;
550                     gotoxy(0, y+par[0]);
551                     break;
// 如果字符 c 是 'F', 则第一个参数代表光标向上移动的行数, 并回到 0 列。若参数为 0 则上移一行。
552                 case 'F':
553                     if (!par[0]) par[0]++;
554                     gotoxy(0, y-par[0]);
555                     break;
// 如果字符 c 是 'd', 则第一个参数代表光标所需的行号(从 0 计数)。
556                 case 'd':

```

```

557         if (par[0]) par[0]--;
558         gotoxy(x, par[0]);
559         break;
// 如果字符 c 是'H' 或'f'，则第一个参数代表光标移到的行号，第二个参数代表光标移到的列号。
560         case 'H': case 'f':
561             if (par[0]) par[0]--;
562             if (par[1]) par[1]--;
563             gotoxy(par[1], par[0]);
564             break;
// 如果字符 c 是'J'，则第一个参数代表以光标所处位置清屏的方式：
// ANSI 转义序列：'ESC [sJ' (s = 0 删除光标到屏幕底端；1 删除屏幕开始到光标处；2 整屏删除)。
565         case 'J':
566             csi_J(par[0]);
567             break;
// 如果字符 c 是'K'，则第一个参数代表以光标所在位置对行中字符进行删除处理的方式。
// ANSI 转义字符序列：'ESC [sK' (s = 0 删除到行尾；1 从开始删除；2 整行都删除)。
568         case 'K':
569             csi_K(par[0]);
570             break;
// 如果字符 c 是'L'，表示在光标位置处插入 n 行(ANSI 转义字符序列'ESC [nL')。
571         case 'L':
572             csi_L(par[0]);
573             break;
// 如果字符 c 是'M'，表示在光标位置处删除 n 行(ANSI 转义字符序列'ESC [nM')。
574         case 'M':
575             csi_M(par[0]);
576             break;
// 如果字符 c 是'P'，表示在光标位置处删除 n 个字符(ANSI 转义字符序列'ESC [nP')。
577         case 'P':
578             csi_P(par[0]);
579             break;
// 如果字符 c 是'@'，表示在光标位置处插入 n 个字符(ANSI 转义字符序列'ESC [n@')。
580         case '@':
581             csi_at(par[0]);
582             break;
// 如果字符 c 是'm'，表示改变光标处字符的显示属性，比如加粗、加下划线、闪烁、反显等。
// ANSI 转义字符序列：'ESC [nm'。n = 0 正常显示；1 加粗；4 加下划线；7 反显；27 正常显示。
583         case 'm':
584             csi_m();
585             break;
// 如果字符 c 是'r'，则表示用两个参数设置滚屏的起始行号和终止行号。
586         case 'r':
587             if (par[0]) par[0]--;
588             if (!par[1]) par[1] = video_num_lines;
589             if (par[0] < par[1] &&
590                 par[1] <= video_num_lines) {
591                 top=par[0];
592                 bottom=par[1];
593             }
594             break;
// 如果字符 c 是's'，则表示保存当前光标所在位置。
595         case 's':
596             save_cur();

```

```

597                                     break;
// 如果字符 c 是 'u'，则表示恢复光标到原保存的位置处。
598                                     case 'u':
599                                         restore_cur();
600                                     break;
601                                 }
602                            }
603                    }
// 最后根据上面设置的光标位置，向显示控制器发送光标显示位置。
604        set_cursor();
605    }
606
607    /*
608     * void con_init(void);
609     *
610     * This routine initializes console interrupts, and does nothing
611     * else. If you want the screen to clear, call tty_write with
612     * the appropriate escape-sequence.
613     *
614     * Reads the information preserved by setup.s to determine the current display
615     * type and sets everything accordingly.
616     */
617    void con_init(void)
618    {
619        register unsigned char a;
620        char *display_desc = "????";
621        char *display_ptr;
622
623        video_num_columns = ORIG_VIDEO_COLS;    // 显示器显示字符列数。
624        video_size_row = video_num_columns * 2; // 每行需使用字节数。
625        video_num_lines = ORIG_VIDEO_LINES;    // 显示器显示字符行数。
626        video_page = ORIG_VIDEO_PAGE;          // 当前显示页面。
627        video_erase_char = 0x0720;             // 擦除字符(0x20 显示字符， 0x07 是属性)。
628
629        // 如果原始显示模式等于 7，则表示是单色显示器。
630        if (ORIG_VIDEO_MODE == 7)                /* Is this a monochrome display? */
631        {
632            video_mem_start = 0xb0000;           // 设置单显映象内存起始地址。
633            video_port_reg = 0x3b4;              // 设置单显索引寄存器端口。
634            video_port_val = 0x3b5;              // 设置单显数据寄存器端口。
635
636            // 根据 BIOS 中断 int 0x10 功能 0x12 获得的显示模式信息，判断显示卡单色显示卡还是彩色显示卡。
637            // 如果使用上述中断功能所得到的 BX 寄存器返回值不等于 0x10，则说明是 EGA 卡。因此初始
638            // 显示类型为 EGA 单色；所使用映象内存末端地址为 0xb8000；并置显示器描述字符串为 'EGAm'。
639            // 在系统初始化期间显示器描述字符串将显示在屏幕的右上角。
640            if ((ORIG_VIDEO_EGA_BX & 0xff) != 0x10)
641            {

```

```

636         video_type = VIDEO_TYPE_EGAM; // 设置显示类型(EGA 单色)。
637         video_mem_end = 0xb8000;      // 设置显示内存末端地址。
638         display_desc = "EGAm";        // 设置显示描述字符串。
639     }
    // 如果 BX 寄存器的值等于 0x10, 则说明是单色显示卡 MDA。则设置相应参数。
640     else
641     {
642         video_type = VIDEO_TYPE_MDA; // 设置显示类型(MDA 单色)。
643         video_mem_end = 0xb2000;      // 设置显示内存末端地址。
644         display_desc = "MDA";         // 设置显示描述字符串。
645     }
646 }
    // 如果显示模式不为 7, 则为彩色模式。此时所用的显示内存起始地址为 0xb800; 显示控制索引寄存
    // 器端口地址为 0x3d4; 数据寄存器端口地址为 0x3d5。
647     else /* If not, it is color. */
648     {
649         video_mem_start = 0xb8000;    // 显示内存起始地址。
650         video_port_reg = 0x3d4;        // 设置彩色显示索引寄存器端口。
651         video_port_val = 0x3d5;        // 设置彩色显示数据寄存器端口。
    // 再判断显示卡类别。如果 BX 不等于 0x10, 则说明是 EGA 显示卡。
652         if ((ORIG_VIDEO_EGA_BX & 0xff) != 0x10)
653         {
654             video_type = VIDEO_TYPE_EGAC; // 设置显示类型(EGA 彩色)。
655             video_mem_end = 0xbc000;      // 设置显示内存末端地址。
656             display_desc = "EGAc";        // 设置显示描述字符串。
657         }
    // 如果 BX 寄存器的值等于 0x10, 则说明是 CGA 显示卡。则设置相应参数。
658         else
659         {
660             video_type = VIDEO_TYPE_CGA; // 设置显示类型(CGA)。
661             video_mem_end = 0xba000;      // 设置显示内存末端地址。
662             display_desc = "CGA";         // 设置显示描述字符串。
663         }
664     }
665
666     /* Let the user known what kind of display driver we are using */
    /* 让用户知道我们正在使用哪一类显示驱动程序 */
667
    // 在屏幕的右上角显示显示描述字符串。采用的方法是直接将字符串写到显示内存的相应位置处。
    // 首先将显示指针 display_ptr 指到屏幕第一行右端差 4 个字符处(每个字符需 2 个字节, 因此减 8)。
668     display_ptr = ((char *)video_mem_start) + video_size_row - 8;
    // 然后循环复制字符串中的字符, 并且每复制一个字符都空开一个属性字节。
669     while (*display_desc)
670     {
671         *display_ptr++ = *display_desc++; // 复制字符。
672         display_ptr++; // 空开属性字节位置。
673     }
674
675     /* Initialize the variables used for scrolling (mostly EGA/VGA) */
    /* 初始化用于滚屏的变量(主要用于 EGA/VGA) */
676
677     origin = video_mem_start; // 滚屏起始显示内存地址。
678     scr_end = video_mem_start + video_num_lines * video_size_row; // 滚屏结束内存地址。

```

```

679         top      = 0;                                // 最顶行号。
680         bottom   = video_num_lines;                  // 最底行号。
681
682         gotoxy(ORIG_X, ORIG_Y);                      // 初始化光标位置 x, y 和对应的内存位置 pos。
683         set_trap_gate(0x21, &keyboard_interrupt);    // 设置键盘中断陷阱门。
684         outb_p(inb_p(0x21)&0xfd, 0x21);              // 取消 8259A 中对键盘中断的屏蔽, 允许 IRQ1。
685         a=inb_p(0x61);                                // 延迟读取键盘端口 0x61 (8255A 端口 PB)。
686         outb_p(a|0x80, 0x61);                        // 设置禁止键盘工作(位 7 置位),
687         outb(a, 0x61);                                // 再允许键盘工作, 用以复位键盘操作。
688     }
689     /* from bsd-net-2: */
690
691     // 停止蜂鸣。
692     // 复位 8255A PB 端口的位 1 和位 0。
693     void sysbeepstop(void)
694     {
695         /* disable counter 2 */ /* 禁止定时器 2 */
696         outb(inb_p(0x61)&0xFC, 0x61);
697     }
698     int beepcount = 0;
699
700     // 开通蜂鸣。
701     // 8255A 芯片 PB 端口的位 1 用作扬声器的开门信号; 位 0 用作 8253 定时器 2 的门信号, 该定时器的
702     // 输出脉冲送往扬声器, 作为扬声器发声的频率。因此要使扬声器蜂鸣, 需要两步: 首先开启 PB 端口
703     // 位 1 和位 0 (置位), 然后设置定时器发送一定的定时频率即可。
704     static void sysbeep(void)
705     {
706         /* enable counter 2 */ /* 开启定时器 2 */
707         outb_p(inb_p(0x61)|3, 0x61);
708         /* set command for counter 2, 2 byte write */ /* 送设置定时器 2 命令 */
709         outb_p(0xB6, 0x43);
710         /* send 0x637 for 750 HZ */ /* 设置频率为 750HZ, 因此送定时值 0x637 */
711         outb_p(0x37, 0x42);
712         outb(0x06, 0x42);
713         /* 1/8 second */ /* 蜂鸣时间为 1/8 秒 */
714         beepcount = HZ/8;
715     }
716 }

```

## 6.8.3 其它信息

### 6.8.3.1 显示控制卡编程

这里仅给出和说明兼容显示卡端口的说明。描述了 MDA、CGA、EGA 和 VGA 显示控制卡的通用编程端口, 这些端口都是与 CGA 使用的 MC6845 芯片兼容, 其名称和用途见下面列表。其中以 CGA/EGA/VGA 的端口 (0x3d0-0x3df) 为例进行说明, MDA 的端口是 0x3b0 - 0x3bf。

对显示控制卡进行编程的基本步骤是: 首先写显示卡的索引寄存器, 选择要进行设置的显示控制内部寄存器之一(r0-r17), 然后将参数写到其数据寄存器端口。也即显示卡的数据寄存器端口每次只能对显示卡中的一个内部寄存器进行操作。



表 CGA 端口寄存器名称及作用

| 端口    | 读/写 | 名称和用途   |
|-------|-----|---|
| 0x3d4 | 写   | CRT (6845) 索引寄存器。用于选择通过端口 0x3b5 访问的各个数据寄存器 (r0-r17)。  |
| 0x3d5 | 写   | CRT (6845) 数据寄存器。其中数据寄存器 r12-r15 还可以读。<br>各个数据寄存器的功能说明见下表。  |
| 0x3d8 | 读/写 | 模式控制寄存器。<br>位 7-6 未用;<br>位 5=1 允许闪烁;<br>位 4=1 640*200 图形模式;<br>位 3=1 允许视频;<br>位 2=1 单色显示;<br>位 1=1 图形模式; =0 文本模式;<br>位 0=1 80*25 文本模式; =0 40*25 文本模式。   |
| 0x3d9 | 读/写 | CGA 调色板寄存器。选择所采用的色彩。<br>位 7-6 未用;<br>位 5=1 激活色彩集: 青(cyan)、紫(magenta)、白(white);<br>=0 激活色彩集: 红(red)、绿(green)、蓝(blue);<br>位 4=1 增强显示图形、文本背景色彩;<br>位 3=1 增强显示 40*25 的边框、320*200 的背景、640*200 的前景颜色;<br>位 2=1 显示红色: 40*25 的边框、320*200 的背景、640*200 的前景;<br>位 1=1 显示绿色: 40*25 的边框、320*200 的背景、640*200 的前景;<br>位 0=1 显示蓝色: 40*25 的边框、320*200 的背景、640*200 的前景; |
| 0x3da | 读   | CGA 显示状态寄存器。<br>位 7-4 未用;<br>位 3=1 在垂直回扫阶段;<br>位 2=1 光笔开关关闭; =0 光笔开关接通;<br>位 1=1 光笔选通有效;<br>位 0=1 可以不干扰显示访问显示内存; =0 此时不要使用显示内存。   |
| 0x3db | 写   | 清除光笔锁存 (复位光笔寄存器)。   |
| 0x3dc | 读/写 | 预设置光笔锁存 (强制光笔选通有效)。   |

MC6845 内部数据寄存器及初始值

| 编号 | 名称       | 单位  | 读/写 | 40*25 模式 | 80*25 模式 | 图形模式 |
|----|----------|-----|-----|----------|----------|------|
| r0 | 水平字符总数   | 字符  | 写   | 0x38     | 0x71     | 0x38 |
| r1 | 水平显示字符数  | 字符  | 写   | 0x28     | 0x50     | 0x28 |
| r2 | 水平同步位置   | 字符  | 写   | 0x2d     | 0x5a     | 0x2d |
| r3 | 水平同步脉冲宽度 | 字符  | 写   | 0x0a     | 0x0a     | 0x0a |
| r4 | 垂直字符总数   | 字符行 | 写   | 0x1f     | 0x1f     | 0x7f |
| r5 | 垂直同步脉冲宽度 | 扫描行 | 写   | 0x06     | 0x06     | 0x06 |
| r6 | 垂直显示字符数  | 字符行 | 写   | 0x19     | 0x19     | 0x64 |
| r7 | 垂直同步位置   | 字符行 | 写   | 0x1c     | 0x1c     | 0x70 |
| r8 | 隔行/逐行选择  |     | 写   | 0x02     | 0x02     | 0x02 |
| r9 | 最大扫描行数   | 扫描行 | 写   | 0x07     | 0x07     | 0x01 |

|     |             |     |     |      |      |      |
|-----|-------------|-----|-----|------|------|------|
| r10 | 光标开始位置      | 扫描行 | 写   | 0x06 | 0x06 | 0x06 |
| r11 | 光标结束位置      | 扫描行 | 写   | 0x07 | 0x07 | 0x07 |
| r12 | 显示内存起始位置(高) |     | 写   | 0x00 | 0x00 | 0x00 |
| r13 | 显示内存末端位置(低) |     | 写   | 0x00 | 0x00 | 0x00 |
| r14 | 光标当前位置(高)   |     | 读/写 | 可变   |      |      |
| r15 | 光标当前位置(低)   |     | 读/写 |      |      |      |
| r16 | 光笔当前位置(高)   |     | 读   | 可变   |      |      |
| r17 | 光笔当前位置(低)   |     | 读   |      |      |      |

### 6.8.3.2 滚屏操作原理

滚屏操作是指将指定开始行和结束行的一块文本内容向上移动(向上卷动 scroll up)或向下移动(向下卷动 scroll down), 如果将屏幕看作是显示内存上对应屏幕内容的一个窗口的话, 那么将屏幕内容向上移动即是窗口沿显示内存向下移动; 将屏幕内容向下移动即是窗口向下移动。在程序中就是重新设置显示控制器中显示内存的起始位置 origin 以及调整程序中相应的变量。对于这两种操作各自都有两种情况。

对于向上卷动, 当屏幕对应的显示内存窗口在向下移动后仍然在显示内存范围之内, 也即对应当前屏幕的内存块位置始终在显示内存起始位置(video\_mem\_start)和末端位置 video\_mem\_end 之间, 那么只需要调整显示控制器中起始显示内存位置即可。但是当对对应屏幕的内存块位置在向下移动时超出了实际显示内存的末端(video\_mem\_end)这种情况, 就需要移动对应显示内存中的数据, 以保证所有当前屏幕数据都落在显示内存范围内。在这第二种情况, 程序中是将屏幕对应的内存数据移动到实际显示内存的开始位置处(video\_mem\_start)。

程序中实际的处理过程分三步进行。首先调整屏幕显示起始位置 origin; 然后判断对应屏幕内存数据是否超出显示内存下界(video\_mem\_end), 如果超出就将屏幕对应的内存数据移动到实际显示内存的开始位置处(video\_mem\_start); 最后对移动后屏幕上出现的新行用空格字符填满。见下面图中所示。其中图(a)对应第一种简单情况, 图(b)对应需要移动内存数据时的情况。

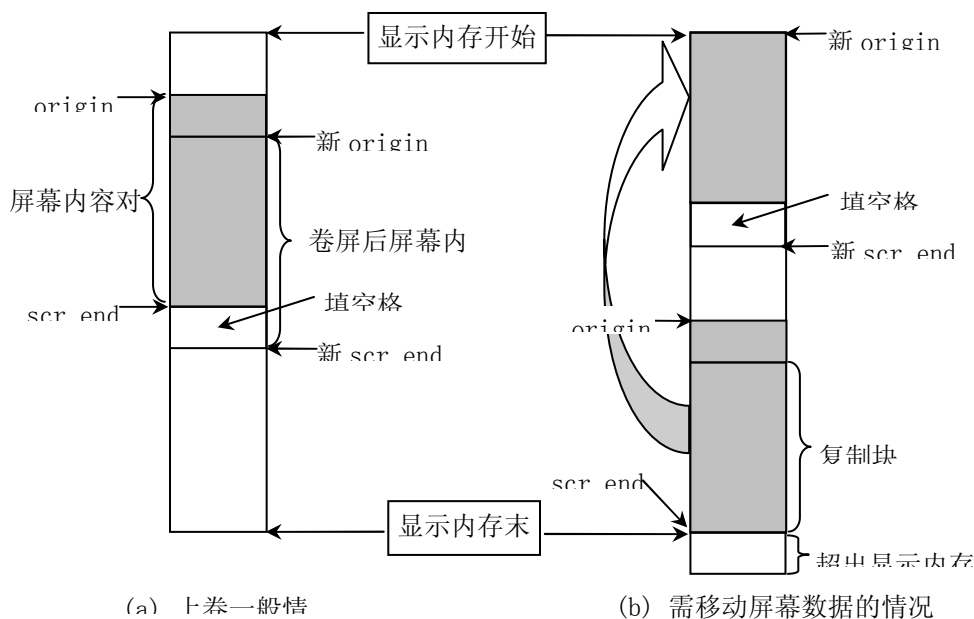


图 向上卷屏(scroll up)操作示意图

向下卷动屏幕的操作与向上卷屏相似，也会遇到这两种类似情况，只是由于屏幕窗口上移，因此会在屏幕上方出现一空行，并且在屏幕内容所对应的内存超出显示内存范围时需要将屏幕数据内存块往下移动到显示内存的末端位置。

### 6.8.3.3 ANSI 转义控制序列

终端通常有两部分功能，分别作为计算机信息的输入设备(键盘)和输出设置(显示器)。终端可有许多控制命令，使得终端执行一定的操作而不是仅仅在屏幕上显示一个字符。使用这种方式，计算机就可以命令终端执行移动光标、切换显示模式和响铃等操作。为了能理解程序的执行处理过程，下面对终端控制命令进行一些简单描述。首先说明控制字符和控制序列的含义。

控制字符是指 ASCII 码表开头的 32 个字符(0x00 - 0x1f 或 0-31)以及字符 DEL(0x7f 或 127)，参见附录中的 ASCII 码表。通常一个指定类型的终端都会采用其中的一个子集作为控制字符，而其它的控制字符将不起作用。例如，对于 VT100 终端所采用的控制字符见下表所示。

表 控制字符

| 控制字符 | 八进制 | 十六进制 | Action Taken  |
|------|-----|------|---|
| NUL  | 000 | 0x00 | Ignored on input (not stored in input buffer; see full duplex protocol).  |
| ENQ  | 005 | 0x05 | Transmit answerback message.  |
| BEL  | 007 | 0x07 | Sound bell tone from keyboard.  |
| BS   | 010 | 0x08 | Move the cursor to the left one character position, unless it is at the left margin, in which case no action occurs.                            |
| HT   | 011 | 0x09 | Move the cursor to the next tab stop, or to the right margin if no further tab stops are present on the line.                                   |
| LF   | 012 | 0x0a | This code causes a line feed or a new line operation. (See new line mode).  |
| VT   | 013 | 0x0b | Interpreted as LF.  |
| FF   | 014 | 0x0c | Interpreted as LF.  |
| CR   | 015 | 0x0d | Move cursor to the left margin on the current line.   |
| SO   | 016 | 0x0e | Invoke G1 character set, as designated by SCS control sequence.   |
| SI   | 017 | 0x0f | Select G0 character set, as selected by ESC ( sequence.   |
| XON  | 021 | 0x11 | Causes terminal to resume transmission.   |
| XOFF | 023 | 0x13 | Causes terminal to stop transmitted all codes except XOFF and XON.  |
| CAN  | 030 | 0x18 | If sent during a control sequence, the sequence is immediately terminated and not executed. It also causes the error character to be displayed. |
| SUB  | 032 | 0x1a | Interpreted as CAN.   |
| ESC  | 033 | 0x1b | Invokes a control sequence.   |
| DEL  | 177 | 0x7f | Ignored on input (not stored in input buffer).  |

控制序列已经由 ANSI(美国国家标准局 American National Standards Institute)制定为标准：X3.64-1977。控制序列是指由一些非控制字符构成的一个特殊字符序列，终端在收到这个序列时并不是将它们直接显示在屏幕上，而是采取一定的控制操作，比如，移动光标、删除字符、删除行、插入字符或插入行等操作。ANSI 控制序列由以下一些基本元素组成：

控制序列引入码(Control Sequence Introducer - CSI): 表示一个转移序列, 提供辅助的控制并且本身是影响随后一系列连续字符含义解释的前缀。通常, 一般 CSI 都使用 ESC [。

参数(Parameter): 零个或多个数字字符组成的一个数值。

数值参数(Numeric Parameter): 表示一个数的参数, 使用 n 表示。

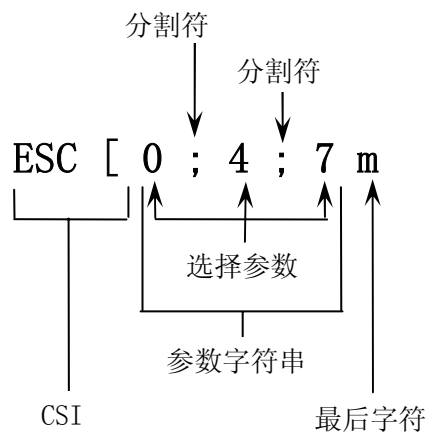
选择参数(Selective Parameter): 用于从一功能子集中选择一个子功能, 一般用 s 表示。通常, 具有多个选择参数的一个控制序列所产生的作用, 如同分立的几个控制序列。例如: CSI sa;sb;sc F 的作用是与 CSI sa F CSI sb F CSI sc F 完全一样的。

参数字符串(Parameter String): 用分号';'隔开的参数字符串。

默认值(Default): 当没有明确指定一个值或者值是 0 的话, 就会指定一个与功能相关的值。

最后字符(Final character): 用于结束一个转义或控制序列。

下面是一个控制序列的例子: 取消所有字符的属性, 然后开启下划线和反显属性。ESC [ 0;4;7m



下面是常用的一些控制序列列表。

Esc Seq Function

```
=====
E[nA    move cursor up n lines
E[nB    move cursor down n lines
E[nC    move cursor right n characters
E[nD    move cursor left n characters
E[n`    move cursor to character position n
E[na    move cursor right n characters
```

---



E[nd     move cursor to line n  
 E[ne     move cursor down n lines  
 E[nF     move cursor to start of line, n lines up  
 E[nE     move cursor to start of line, n lines down  
 E[y;xH   Move cursor to x, y   E[H homes cursor  
 E[y;xf   Move cursor to x, y  
 E[nZ     Move cursor back n tab stops  
 E[nL     Insert n blank lines  
 E[n@     Insert n blank characters  
 E[nM     Delete n lines  
 E[nP     Delete n characters  
 E[nJ     Erase part or all of display:  
           n = 0 from cursor to end of display,  
           n = 1 from begin of display to cursor,  
           n = 2 entire display.  
 E[nK     Erase part or all of line:  
           n = 0 from cursor to end of line,  
           n = 1 from begin of line to cursor,  
           n = 2 entire line.  
 E[nX     Erase n characters  
 E[nS     Scroll display n lines up (forward)  
 E[nT     Scroll display n lines down (reverse)  
 E[nm     Set character attributes:  
           n = 0 normal attribute (all off)  
           n = 1 bold  
           n = 4 underscore  
           n = 5 blink  
           n = 7 reverse  
           n = 3X set foreground color  
           n = 4X set background color  
           X = 0 black     X = 1 red  
           X = 2 green     X = 3 brown  
           X = 4 blue      X = 5 magenta  
           X = 6 cyan      X = 7 white  
           You can set more than one thing by separating them with a  
           semi-colon. eg. E[0;1;33;40m  
 E[s     Save cursor position  
 E[u     Restore saved cursor position  
 E means OX1B  
 if n is 0 then it can also be left off  
 E[0J == E[J



## 第 7 章 数学协处理器 (math)

### 7.1 概述

列表 linux/kernel/math 目录

| Name   | Size       | Last modified (GMT) | Description |
|--|------------|---------------------|-------------|
|  <a href="#">Makefile</a>       | 936 bytes  | 1991-11-18 00:21:45 |             |
|  <a href="#">math_emulate.c</a> | 1023 bytes | 1991-11-23 15:36:34 |             |

### 7.2 Makefile 文件

#### 7.2.1 功能描述

#### 7.2.2 代码注释

列表 linux/kernel/math/Makefile 文件

```
1 #
2 # Makefile for the FREAX-kernel character device drivers.
3 #
4 # Note! Dependencies are done automagically by 'make dep', which also
5 # removes any old dependencies. DON'T put your own dependencies here
6 # unless it's something special (ie not a .c file).
7 #
8 # FREAX(Linux)内核字符设备驱动程序的 Makefile 文件。
9 # 注意！依赖关系是由'make dep'自动进行的，它也会自动去除原来的依赖信息。不要把你自己的
10 # 依赖关系信息放在这里，除非是特别文件的（也即不是一个.c文件的信息）。
11
12 8
13 9 AR      =gar      # GNU 的二进制文件处理程序，用于创建、修改以及从归档文件中抽取文件。
14 10 AS      =gas      # GNU 的汇编程序。
15 11 LD      =gld      # GNU 的连接程序。
16 12 LDFLAGS =-s -x    # 连接程序所有的参数，-s 输出文件中省略所有符号信息。-x 删除所有局部符号。
17 13 CC      =gcc      # GNU C 语言编译器。
18 # 下一行是 C 编译程序选项。-Wall 显示所有的警告信息；-O 优化选项，优化代码长度和执行时间；
19 # -fstrength-reduce 优化循环执行代码，排除重复变量；-fomit-frame-pointer 省略保存不必要
20 # 的框架指针；-fcombine-regs 合并寄存器，减少寄存器类的使用；-finline-functions 将所有简
21 # 单短小的函数代码嵌入调用程序中；-mstring-insns Linus 自己添加的优化选项，以后不再使用；
```

```

# -nostdinc -I../include 不使用默认路径中的包含文件，而使用指定目录中的(../../include)。
14 CFLAGS =-Wall -O -fstrength-reduce -fomit-frame-pointer -fcombine-regs \
15      -finline-functions -mstring-insns -nostdinc -I../../include
# C 前处理选项。-E 只运行 C 前处理，对所有指定的 C 程序进行预处理并将处理结果输出到标准输
# 出设备或指定的输出文件中；-nostdinc -I../../include 同前。
16 CPP =gcc -E -nostdinc -I../../include
17
# 下面的规则指示 make 利用下面的命令将所有的.c 文件编译生成.s 汇编程序。该规则的命令
# 指使 gcc 采用 CFLAGS 所指定的选项对 C 代码编译后不进行汇编就停止 (-S)，从而产生与
# 输入的各个 C 文件对应的汇编代码文件。默认情况下所产生的汇编程序文件名是原 C 文件名
# 去掉.c 而加上.s 后缀。-o 表示其后是输出文件的名称。其中$.s (或$@) 是自动目标变量，
# $<代表第一个先决条件，这里即是符合条件*.c 的文件。
18 .c.s:
19     $(CC) $(CFLAGS) \
20     -S -o $.s $<
# 下面规则表示将所有.s 汇编程序文件编译成.o 目标文件。22 行是实现该操作的具体命令。
21 .s.o:
22     $(AS) -c -o $.o $<
23 .c.o:                                # 类似上面，*.c 文件→*.o 目标文件。不进行连接。
24     $(CC) $(CFLAGS) \
25     -c -o $.o $<
26
27 OBJS = math_emulate.o    # 定义目标文件变量 OBJS。
28
29 math.a: $(OBJS)           # 在有了先决条件 OBJS 后使用下面的命令连接成目标 math.a 库文件。
30     $(AR) rcs math.a $(OBJS)
31     sync
32
# 下面的规则用于清理工作。当执行'make clean'时，就会执行下面的命令，去除所有编译
# 连接生成的文件。'rm' 是文件删除命令，选项-f 含义是忽略不存在的文件，并且不显示删除信息。
33 clean:
34     rm -f core *.o *.a tmp_make
35     for i in *.c;do rm -f `basename $$i .c`.s;done
36
# 下面得目标或规则用于检查各文件之间的依赖关系。方法如下：
# 使用字符串编辑程序 sed 对 Makefile 文件（即是本文件）进行处理，输出为删除 Makefile
# 文件中'### Dependencies' 行后面的所有行，并生成 tmp_make 临时文件。然后对 kernel/math/
# 目录下的每个 C 文件执行 gcc 预处理操作。
# -M 标志告诉预处理程序输出描述每个目标文件相关性的规则，并且这些规则符合 make 语法。
# 对于每一个源文件，预处理程序输出一个 make 规则，其结果形式是相应源程序文件的目标
# 文件名加上其依赖关系——该源文件中包含的所有头文件列表。把预处理结果都添加到临时
# 文件 tmp_make 中，然后将该临时文件复制成新的 Makefile 文件。
37 dep:
38     sed '/\#\#\# Dependencies/q' < Makefile > tmp_make
39     (for i in *.c;do echo -n `echo $$i | sed 's,\.c,\.s,`" "; \
40         $(CPP) -M $$i;done) >> tmp_make
41     cp tmp_make Makefile
42
43 ### Dependencies:

```



## 7.3 math-emulation.c 程序

### 7.3.1 功能描述

### 7.3.2 代码注释

列表 linux/kernel/math/math\_emulate.c 程序

```

1  /*
2   * linux/kernel/math/math_emulate.c
3   *
4   * (C) 1991 Linus Torvalds
5   */
6
7  /*
8   * This directory should contain the math-emulation code.
9   * Currently only results in a signal.
10  */
11  /*
12   * 该目录里应该包含数学仿真代码。目前仅产生一个信号。
13   */
14
15 #include <signal.h>          // 信号头文件。定义信号符号常量，信号结构以及信号操作函数原型。
16
17 #include <linux/sched.h>     // 调度程序头文件，定义了任务结构 task_struct、初始任务 0 的数据，
18                               // 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
19
20 #include <linux/kernel.h>    // 内核头文件。含有一些内核常用函数的原形定义。
21
22 #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。
23
24
25 // 协处理器仿真函数。
26 // 中断处理程序调用的 C 函数，参见 (kernel/math/system_call.s, 169 行)。
27 void math_emulate(long edi, long esi, long ebp, long sys_call_ret,
28                   long eax, long ebx, long ecx, long edx,
29                   unsigned short fs, unsigned short es, unsigned short ds,
30                   unsigned long eip, unsigned short cs, unsigned long eflags,
31                   unsigned short ss, unsigned long esp)
32 {
33     unsigned char first, second;
34
35     /* 0x0007 means user code space */
36     /* 0x0007 表示用户代码空间 */
37     // 选择符 0x000F 表示在局部描述符表中描述符索引值=1，即代码空间。如果段寄存器 cs 不等于 0x000F
38     // 则表示 cs 一定是内核代码选择符，是在内核代码空间，则出错，显示此时的 cs:eip 值，并显示信息
39     // “内核中需要数学仿真”，然后进入死机状态。
40     if (cs != 0x000F) {
41         printk("math_emulate: %04x:%08x\n|r", cs, eip);
42         panic("Math emulation needed in kernel");
43     }
44     // 取用户数据区堆栈数据 first 和 second，显示这些数据，并给进程设置浮点异常信号 SIGFPE。
45     first = get_fs_byte((char *)(&eip++));
46     second = get_fs_byte((char *)(&eip++));

```








---



```
33     printk( "%04x:%08x %02x %02x\n|r", cs, eip-2, first, second);
34     current->signal |= 1<<(SIGFPE-1);
35 }
36
37     ///// 协处理器出错处理函数。
38     // 中断处理程序调用的 C 函数，参见 (kernel/math/system_call.s, 145 行)。
39 void math_error(void)
40 {
41     // 协处理器指令。(以非等待形式)清除所有异常标志、忙标志和状态字位 7。
42     __asm__( "fnclex");
43     // 如果上个任务使用过协处理器，则向上个任务发送协处理器异常信号。
44     if (last_task_used_math)
45         last_task_used_math->signal |= 1<<(SIGFPE-1);
46 }
```

---

## 第 8 章 文件系统(fs)

### 8.1 概述

| Name   | Size        | Last modified (GMT) | Description |
|--|-------------|---------------------|-------------|
|  <a href="#">Makefile</a>       | 5053 bytes  | 1991-12-02 03:21:31 | m           |
|  <a href="#">bitmap.c</a>       | 4042 bytes  | 1991-11-26 21:31:53 | m           |
|  <a href="#">block_dev.c</a>    | 1422 bytes  | 1991-10-31 17:19:55 | m           |
|  <a href="#">buffer.c</a>       | 9072 bytes  | 1991-12-06 20:21:00 | m           |
|  <a href="#">char_dev.c</a>   | 2103 bytes  | 1991-11-19 09:10:22 |             |
|  <a href="#">exec.c</a>       | 9134 bytes  | 1991-12-01 20:01:01 | m           |
|  <a href="#">fcntl.c</a>      | 1455 bytes  | 1991-10-02 14:16:29 | m           |
|  <a href="#">file_dev.c</a>   | 1852 bytes  | 1991-12-01 19:02:43 | m           |
|  <a href="#">file_table.c</a> | 122 bytes   | 1991-10-02 14:16:29 | m           |
|  <a href="#">inode.c</a>      | 6933 bytes  | 1991-12-06 20:16:35 | m           |
|  <a href="#">ioctl.c</a>      | 977 bytes   | 1991-11-19 09:13:05 |             |
|  <a href="#">namei.c</a>      | 16562 bytes | 1991-11-25 19:19:59 | m           |
|  <a href="#">open.c</a>       | 4340 bytes  | 1991-11-25 19:21:01 | m           |
|  <a href="#">pipe.c</a>       | 2385 bytes  | 1991-10-18 19:02:33 | m           |
|  <a href="#">read_write.c</a> | 2802 bytes  | 1991-11-25 15:47:20 | m           |
|  <a href="#">stat.c</a>       | 1175 bytes  | 1991-10-02 14:16:29 | m           |

|  |            |                     |   |
|--|------------|---------------------|---|
|  <a href="#">super.c</a>    | 5628 bytes | 1991-12-06 20:10:12 | m |
|  <a href="#">truncate.c</a> | 1148 bytes | 1991-10-02 14:16:29 | m |

## 8.2 Makefile 文件

### 8.2.1 功能描述

### 8.2.2 代码注释

列表 linux/fs/Makefile 文件

```

1 AR      =gar      # GNU 的二进制文件处理程序，用于创建、修改以及从归档文件中抽取文件。
2 AS      =gas      # GNU 的汇编程序。
3 CC      =gcc      # GNU C 语言编译器。
4 LD      =gld      # GNU 的连接程序。

# C 编译程序选项。-Wall 显示所有的警告信息；-O 优化选项，优化代码长度和执行时间；
# -fstrength-reduce 优化循环执行代码，排除重复变量；-fomit-frame-pointer 省略保存不必要
# 的框架指针；-fcombine-regs 合并寄存器，减少寄存器类的使用；-mstring-insns Linus 自己
# 填加的优化选项，以后不再使用；-nostdinc -I../include 不使用默认路径中的包含文件，而使
# 用这里指定目录中的(../include)。
5 CFLAGS  =-Wall -O -fstrength-reduce -fcombine-regs -fomit-frame-pointer \
6         -mstring-insns -nostdinc -I../include

# C 前处理选项。-E 只运行 C 前处理，对所有指定的 C 程序进行预处理并将处理结果输出到标准输
# 出设备或指定的输出文件中；-nostdinc -I../include 同前。
7 CPP     =gcc -E -nostdinc -I../include
8
# 下面的规则指示 make 利用下面的命令将所有的.c 文件编译生成.s 汇编程序。该规则的命令
# 指使 gcc 采用 CFLAGS 所指定的选项对 C 代码编译后不进行汇编就停止（-S），从而产生与
# 输入的各个 C 文件对应的汇编代码文件。默认情况下所产生的汇编程序文件名是原 C 文件名
# 去掉.c 而加上.s 后缀。-o 表示其后是输出文件的名称。其中$.s（或$o）是自动目标变量，
# $<代表第一个先决条件，这里即是符合条件*.c 的文件。
9 .c.s:
10      $(CC) $(CFLAGS) \
11      -S -o $.s $<
# 将所有*.c 文件编译成*.o 目标文件。不进行连接。
12 .c.o:
13      $(CC) $(CFLAGS) \
14      -c -o $.o $<
# 下面规则表示将所有.s 汇编程序文件编译成.o 目标文件。16 行是实现该操作的具体命令。
15 .s.o:
16      $(AS) -o $.o $<
17

```

```

# 定义目标文件变量 OBJS。
18 OBJS= open.o read_write.o inode.o file_table.o buffer.o super.o \
19       block_dev.o char_dev.o file_dev.o stat.o exec.o pipe.o namei.o \
20       bitmap.o fcntl.o ioctl.o truncate.o
21
# 在有了先决条件 OBJS 后使用下面的命令连接成目标 fs.o
22 fs.o: $(OBJS)
23       $(LD) -r -o fs.o $(OBJS)
24
# 下面的规则用于清理工作。当执行'make clean'时,就会执行 26--27 行上的命令,去除所有编译
# 连接生成的文件。'rm'是文件删除命令,选项-f 含义是忽略不存在的文件,并且不显示删除信息。
25 clean:
26       rm -f core *.o *.a tmp_make
27       for i in *.c;do rm -f `basename $$i .c`.s;done
28
# 下面得目标或规则用于检查各文件之间的依赖关系。方法如下:
# 使用字符串编辑程序 sed 对 Makefile 文件(这里即是自己)进行处理,输出为删除 Makefile
# 文件中'### Dependencies'行后面的所有行(下面从 35 开始的行),并生成 tmp_make
# 临时文件(30 行的作用)。然后对 fs/目录下的每一个 C 文件执行 gcc 预处理操作。
# -M 标志告诉预处理程序输出描述每个目标文件相关性的规则,并且这些规则符合 make 语法。
# 对于每一个源文件,预处理程序输出一个 make 规则,其结果形式是相应源程序文件的目标
# 文件名加上其依赖关系--该源文件中包含的所有头文件列表。把预处理结果都添加到临时
# 文件 tmp_make 中,然后将该临时文件复制成新的 Makefile 文件。
29 dep:
30       sed '/\#\#\# Dependencies/q' < Makefile > tmp_make
31       (for i in *.c;do $(CPP) -M $$i;done) >> tmp_make
32       cp tmp_make Makefile
33
34 ### Dependencies:
35 bitmap.o : bitmap.c ../include/string.h ../include/linux/sched.h \
36 ..../include/linux/head.h ../include/linux/fs.h ../include/sys/types.h \
37 ..../include/linux/mm.h ../include/signal.h ../include/linux/kernel.h
38 block_dev.o : block_dev.c ../include/errno.h ../include/linux/sched.h \
39 ..../include/linux/head.h ../include/linux/fs.h ../include/sys/types.h \
40 ..../include/linux/mm.h ../include/signal.h ../include/linux/kernel.h \
41 ..../include/asm/segment.h ../include/asm/system.h
42 buffer.o : buffer.c ../include/stdarg.h ../include/linux/config.h \
43 ..../include/linux/sched.h ../include/linux/head.h ../include/linux/fs.h \
44 ..../include/sys/types.h ../include/linux/mm.h ../include/signal.h \
45 ..../include/linux/kernel.h ../include/asm/system.h ../include/asm/io.h
46 char_dev.o : char_dev.c ../include/errno.h ../include/sys/types.h \
47 ..../include/linux/sched.h ../include/linux/head.h ../include/linux/fs.h \
48 ..../include/linux/mm.h ../include/signal.h ../include/linux/kernel.h \
49 ..../include/asm/segment.h ../include/asm/io.h
50 exec.o : exec.c ../include/errno.h ../include/string.h \
51 ..../include/sys/stat.h ../include/sys/types.h ../include/a.out.h \
52 ..../include/linux/fs.h ../include/linux/sched.h ../include/linux/head.h \
53 ..../include/linux/mm.h ../include/signal.h ../include/linux/kernel.h \
54 ..../include/asm/segment.h
55 fcntl.o : fcntl.c ../include/string.h ../include/errno.h \
56 ..../include/linux/sched.h ../include/linux/head.h ../include/linux/fs.h \
57 ..../include/sys/types.h ../include/linux/mm.h ../include/signal.h \
58 ..../include/linux/kernel.h ../include/asm/segment.h ../include/fcntl.h \

```

```

59  ../include/sys/stat.h
60 file_dev.o : file_dev.c ../include/errno.h ../include/fcntl.h \
61  ../include/sys/types.h ../include/linux/sched.h ../include/linux/head.h \
62  ../include/linux/fs.h ../include/linux/mm.h ../include/signal.h \
63  ../include/linux/kernel.h ../include/asm/segment.h
64 file_table.o : file_table.c ../include/linux/fs.h ../include/sys/types.h
65 inode.o : inode.c ../include/string.h ../include/sys/stat.h \
66  ../include/sys/types.h ../include/linux/sched.h ../include/linux/head.h \
67  ../include/linux/fs.h ../include/linux/mm.h ../include/signal.h \
68  ../include/linux/kernel.h ../include/asm/system.h
69 ioctl.o : ioctl.c ../include/string.h ../include/errno.h \
70  ../include/sys/stat.h ../include/sys/types.h ../include/linux/sched.h \
71  ../include/linux/head.h ../include/linux/fs.h ../include/linux/mm.h \
72  ../include/signal.h
73 namei.o : namei.c ../include/linux/sched.h ../include/linux/head.h \
74  ../include/linux/fs.h ../include/sys/types.h ../include/linux/mm.h \
75  ../include/signal.h ../include/linux/kernel.h ../include/asm/segment.h \
76  ../include/string.h ../include/fcntl.h ../include/errno.h \
77  ../include/const.h ../include/sys/stat.h
78 open.o : open.c ../include/string.h ../include/errno.h ../include/fcntl.h \
79  ../include/sys/types.h ../include/utime.h ../include/sys/stat.h \
80  ../include/linux/sched.h ../include/linux/head.h ../include/linux/fs.h \
81  ../include/linux/mm.h ../include/signal.h ../include/linux/tty.h \
82  ../include/termios.h ../include/linux/kernel.h ../include/asm/segment.h
83 pipe.o : pipe.c ../include/signal.h ../include/sys/types.h \
84  ../include/linux/sched.h ../include/linux/head.h ../include/linux/fs.h \
85  ../include/linux/mm.h ../include/asm/segment.h
86 read_write.o : read_write.c ../include/sys/stat.h ../include/sys/types.h \
87  ../include/errno.h ../include/linux/kernel.h ../include/linux/sched.h \
88  ../include/linux/head.h ../include/linux/fs.h ../include/linux/mm.h \
89  ../include/signal.h ../include/asm/segment.h
90 stat.o : stat.c ../include/errno.h ../include/sys/stat.h \
91  ../include/sys/types.h ../include/linux/fs.h ../include/linux/sched.h \
92  ../include/linux/head.h ../include/linux/mm.h ../include/signal.h \
93  ../include/linux/kernel.h ../include/asm/segment.h
94 super.o : super.c ../include/linux/config.h ../include/linux/sched.h \
95  ../include/linux/head.h ../include/linux/fs.h ../include/sys/types.h \
96  ../include/linux/mm.h ../include/signal.h ../include/linux/kernel.h \
97  ../include/asm/system.h ../include/errno.h ../include/sys/stat.h
98 truncate.o : truncate.c ../include/linux/sched.h ../include/linux/head.h \
99  ../include/linux/fs.h ../include/sys/types.h ../include/linux/mm.h \
100  ../include/signal.h ../include/sys/stat.h

```

### 8.2.3 其它信息

## 8.3 bitmap.c 程序

### 8.3.1 功能描述

该程序主要用于处理 i 节点和逻辑块（磁盘块或区段）的位图。

### 8.3.2 代码注释

列表 linux/fs/bitmap.c 程序

```

1  /*
2   *  linux/fs/bitmap.c
3   *
4   *  (C) 1991  Linus Torvalds
5   */
6
7  /* bitmap.c contains the code that handles the inode and block bitmaps */
   /* bitmap.c 程序含有处理 i 节点和磁盘块位图的代码 */
8  #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。
9                              // 主要使用了其中的 memset() 函数。
10 #include <linux/sched.h>     // 调度程序头文件，定义了任务结构 task_struct、初始任务 0 的数据，
                              // 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
11 #include <linux/kernel.h>    // 内核头文件。含有一些内核常用函数的原形定义。
12
   //// 将指定地址(addr)处的一块内存清零。嵌入汇编程序宏。
   // 输入: eax = 0, ecx = 数据块大小 BLOCK_SIZE/4, edi = addr。
13 #define clear_block(addr) \
14 __asm__( "cld\n\t" \           // 清方向位。
15         "rep\n\t" \           // 重复执行存储数据 (0)。
16         "stosl" \
17         :: "a" (0), "c" (BLOCK_SIZE/4), "D" ((long) (addr)): "cx", "di")
18
   //// 置位指定地址开始的第 nr 个位偏移处的比特位(nr 可以大于 32!)。返回原比特位 (0 或 1)。
   // 输入: %0 - eax (返回值), %1 - eax(0); %2 - nr, 位偏移值; %3 - (addr), addr 的内容。
19 #define set_bit(nr,addr) ({\
20 register int res __asm__("ax"); \
21 __asm__ __volatile__ ("btsl %2,%3\n\tsetb %%al": \
22 "=a" (res): "" (0), "r" (nr), "m" (*(addr))); \
23 res;})
24
   //// 复位指定地址开始的第 nr 位偏移处的比特位。返回原比特位的反码 (1 或 0)。
   // 输入: %0 - eax (返回值), %1 - eax(0); %2 - nr, 位偏移值; %3 - (addr), addr 的内容。
25 #define clear_bit(nr,addr) ({\
26 register int res __asm__("ax"); \
27 __asm__ __volatile__ ("btrl %2,%3\n\tsetnb %%al": \
28 "=a" (res): "" (0), "r" (nr), "m" (*(addr))); \
29 res;})
30
   //// 从 addr 开始寻找第 1 个 0 值比特位。
   // 输入: %0 - ecx(返回值); %1 - ecx(0); %2 - esi(addr)。
   // 在 addr 指定地址开始的位图中寻找第 1 个是 0 的比特位，并将其距离 addr 的比特位偏移值返回。
31 #define find_first_zero(addr) ({ \
32 int __res; \
33 __asm__( "cld\n\t" \           // 清方向位。
34         "l:|t lodsl\n\t" \     // 取[esi]→eax。
35         "notl %%eax\n\t" \     // eax 中每位取反。
36         "bsfl %%eax,%%edx\n\t" \ // 从位 0 扫描 eax 中是 1 的第 1 个位，其偏移值→edx。
37         "je 2f\n\t" \         // 如果 eax 中全是 0，则向前跳转到标号 2 处(40 行)。

```

```

38     "addl %%edx, %%ecx|n|t" \    // 偏移值加入 ecx(ecx 中是位图中首个是 0 的比特位的偏移值)
39     "jmp 3f|n" \                // 向前跳转到标号 3 处 (结束)。
40     "2:|taddl $32, %%ecx|n|t" \ // 没有找到 0 比特位, 则将 ecx 加上 1 个长字的位偏移量 32。
41     "cmpl $8192, %%ecx|n|t" \   // 已经扫描了 8192 位 (1024 字节) 了吗?
42     "j1 1b|n" \                // 若还没有扫描完 1 块数据, 则向前跳转到标号 1 处, 继续。
43     "3:" \                      // 结束。此时 ecx 中是位偏移量。
44     : "=c" (__res): "c" (0), "S" (addr): "ax", "dx", "si"; \
45     __res;})
46
47     /// 释放设备 dev 上数据区中的逻辑块 block。
48     /// 复位指定逻辑块 block 的逻辑块位图比特位。
49 void free_block(int dev, int block)
50 {
51     struct super_block * sb;
52     struct buffer_head * bh;
53
54     // 从设备 dev 取超级块, 如果指定设备不存在, 则出错死机。
55     if (!(sb = get_super(dev)))
56         panic("trying to free block on nonexistent device");
57     // 若逻辑块号小于首个逻辑块号或者大于设备上总逻辑块数, 则出错, 死机。
58     if (block < sb->s_firstdatazone || block >= sb->s_nzones)
59         panic("trying to free block not in datazone");
60     // 从 hash 表中寻找该块数据。若找到了则判断其有效性, 并清已修改和更新标志, 释放该数据块。
61     bh = get_hash_table(dev, block);
62     if (bh) {
63         if (bh->b_count != 1) {
64             printk("trying to free block (%04x:%d), count=%d\n",
65                 dev, block, bh->b_count);
66             return;
67         }
68         bh->b_dirt=0;           // 复位脏 (已修改) 标志位。
69         bh->b_uptodate=0;       // 复位更新标志。
70         brelse(bh);
71     }
72     // 计算 block 在数据区开始算起的逻辑块号 (从 1 开始计数)。然后对逻辑块 (区段) 位图进行操作,
73     // 复位对应的比特位。若对应比特位原来即是 0, 则出错, 死机。
74     block -= sb->s_firstdatazone - 1; // block = block - (-1);
75     if (clear_bit(block&8191, sb->s_zmap[block/8192]->b_data)) {
76         printk("block (%04x:%d)", dev, block+sb->s_firstdatazone-1);
77         panic("free_block: bit already cleared");
78     }
79     // 置相应逻辑块位图所在缓冲区已修改标志。
80     sb->s_zmap[block/8192]->b_dirt = 1;
81 }
82
83     /// 向设备 dev 申请一个逻辑块 (磁盘块, 区段)。返回逻辑块号。
84     /// 置位指定逻辑块 block 的逻辑块位图比特位。
85 int new_block(int dev)
86 {
87     struct buffer_head * bh;
88     struct super_block * sb;
89     int i, j;
90

```



```

// 从设备 dev 取超级块，如果指定设备不存在，则出错死机。
81     if (!(sb = get\_super(dev)))
82         panic("trying to get new block from nonexistant device");
// 扫描逻辑块位图，寻找首个 0 比特位，寻找空闲逻辑块，获取放置该逻辑块的块号。。
83     j = 8192;
84     for (i=0 ; i<8 ; i++)
85         if (bh=sb->s_zmap[i])
86             if ((j=find\_first\_zero(bh->b_data))<8192)
87                 break;
// 如果全部扫描完还没找到(i>=8 或 j>=8192)或者位图所在的缓冲块无效(bh=NULL)则 返回 0，
// 退出（没有空闲逻辑块）。
88     if (i>=8 || !bh || j>=8192)
89         return 0;
// 设置新逻辑块对应逻辑块位图中的比特位，若对应比特位已经置位，则出错，死机。
90     if (set\_bit(j, bh->b_data))
91         panic("new block: bit already set");
// 置对应缓冲区块的已修改标志。如果新逻辑块大于该设备上的总逻辑块数，则说明指定逻辑块在
// 对应设备上不存在。申请失败，返回 0，退出。
92     bh->b_dirt = 1;
93     j += i*8192 + sb->s_firstdatazone-1;
94     if (j >= sb->s_nzones)
95         return 0;
// 读取设备上的该新逻辑块数据（验证）。如果失败则死机。
96     if (!(bh=getblk(dev, j)))
97         panic("new block: cannot get block");
// 新块的引用计数应为 1。否则死机。
98     if (bh->b_count != 1)
99         panic("new block: count is != 1");
// 将该新逻辑块清零，并置位更新标志和已修改标志。然后释放对应缓冲区，返回逻辑块号。
100    clear\_block(bh->b_data);
101    bh->b_uptodate = 1;
102    bh->b_dirt = 1;
103    brelse(bh);
104    return j;
105 }
106
///// 释放指定的 i 节点。
// 复位对应 i 节点位图比特位。
107 void free\_inode(struct m\_inode * inode)
108 {
109     struct super\_block * sb;
110     struct buffer\_head * bh;
111
// 如果 i 节点指针=NULL，则退出。
112     if (!inode)
113         return;
// 如果 i 节点上的设备号字段为 0，说明该节点无用，则用 0 清空对应 i 节点所占内存区，并返回。
114     if (!inode->i_dev) {
115         memset(inode, 0, sizeof(*inode));
116         return;
117     }
// 如果此 i 节点还有其它程序引用，则不能释放，说明内核有问题，死机。
118     if (inode->i_count>1) {

```

```

119         printk("trying to free inode with count=%d\n", inode->i_count);
120         panic("free_inode");
121     }
    // 如果文件目录项连接数不为0, 则表示还有其它文件目录项在使用该节点, 不应释放, 而应该放回等。
122     if (inode->i_nlinks)
123         panic("trying to free inode with links");
    // 取 i 节点所在设备的超级块, 测试设备是否存在。
124     if (!(sb = get_super(inode->i_dev)))
125         panic("trying to free inode on nonexistent device");
    // 如果 i 节点号=0 或大于该设备上 i 节点总数, 则出错 (0 号 i 节点保留没有使用)。
126     if (inode->i_num < 1 || inode->i_num > sb->s_ninodes)
127         panic("trying to free inode 0 or nonexistant inode");
    // 如果该 i 节点对应的节点位图不存在, 则出错。
128     if (!(bh=sb->s_imap[inode->i_num>>13]))
129         panic("nonexistent imap in superblock");
    // 复位 i 节点对应的节点位图中的比特位, 如果该比特位已经等于 0, 则出错。
130     if (clear_bit(inode->i_num&8191, bh->b_data))
131         printk("free_inode: bit already cleared. \n|r");
    // 置 i 节点位图所在缓冲区已修改标志, 并清空该 i 节点结构所占内存区。
132     bh->b_dirt = 1;
133     memset(inode, 0, sizeof(*inode));
134 }
135
    //// 为设备 dev 建立一个新 i 节点。返回该新 i 节点的指针。
    // 在内存 i 节点表中获取一个空闲 i 节点表项, 并从 i 节点位图中找一个空闲 i 节点。
136 struct m_inode * new_inode(int dev)
137 {
138     struct m_inode * inode;
139     struct super_block * sb;
140     struct buffer_head * bh;
141     int i, j;
142
    // 从内存 i 节点表(inode_table)中获取一个空闲 i 节点项(inode)。
143     if (!(inode=get_empty_inode()))
144         return NULL;
    // 读取指定设备的超级块结构。
145     if (!(sb = get_super(dev)))
146         panic("new_inode with unknown device");
    // 扫描 i 节点位图, 寻找首个 0 比特位, 寻找空闲节点, 获取放置该 i 节点的节点号。
147     j = 8192;
148     for (i=0 ; i<8 ; i++)
149         if (bh=sb->s_imap[i])
150             if ((j=find_first_zero(bh->b_data))<8192)
151                 break;
    // 如果全部扫描完还没找到, 或者位图所在的缓冲块无效(bh=NULL)则 返回 0, 退出(没有空闲 i 节点)。
152     if (!bh || j >= 8192 || j+i*8192 > sb->s_ninodes) {
153         iput(inode);
154         return NULL;
155     }
    // 置位对应新 i 节点的 i 节点位图相应比特位, 如果已经置位, 则出错。
156     if (set_bit(j, bh->b_data))
157         panic("new_inode: bit already set");
    // 置 i 节点位图所在缓冲区已修改标志。

```

```

158     bh->b_dirt = 1;
    // 初始化该 i 节点结构。
159     inode->i_count=1;           // 引用计数。
160     inode->i_nlinks=1;         // 文件目录项链接数。
161     inode->i_dev=dev;          // i 节点所在的设备号。
162     inode->i_uid=current->euid; // i 节点所属用户 id。
163     inode->i_gid=current->egid; // 组 id。
164     inode->i_dirt=1;           // 已修改标志置位。
165     inode->i_num = j + i*8192; // 对应设备中的 i 节点号。
166     inode->i_mtime = inode->i_atime = inode->i_ctime = CURRENT_TIME; // 设置时间。
167     return inode;             // 返回该 i 节点指针。
168 }
169

```

### 8.3.3 其它信息

## 8.4 inode.c 程序

### 8.4.1 功能描述

该程序含有处理 i 节点的函数。在注释中，“逻辑块”、“磁盘块”或“区段”的含义相同。

### 8.4.2 代码注释

列表 linux/fs/inode.c 程序

```

1  /*
2   * linux/fs/inode.c
3   *
4   * (C) 1991 Linus Torvalds
5   */
6
7  #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。
8  #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。
9
10 #include <linux/sched.h>      // 调度程序头文件，定义了任务结构 task_struct、初始任务 0 的数据，
    // 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
11 #include <linux/kernel.h>     // 内核头文件。含有一些内核常用函数的原形定义。
12 #include <linux/mm.h>         // 内存管理头文件。含有页面大小定义和一些页面释放函数原型。
13 #include <asm/system.h>       // 系统头文件。定义了设置或修改描述符/中断门等的嵌入式汇编宏。
14
15 struct m_inode inode_table[NR_INODE]={{0},{},}; // 内存中 i 节点表 (NR_INODE=32 项)。
16
17 static void read_inode(struct m_inode * inode);
18 static void write_inode(struct m_inode * inode);
19
    // 等待指定的 i 节点可用。
    // 如果 i 节点已被锁定，则将当前任务置为不可中断的等待状态。直到该 i 节点解锁。
20 static inline void wait_on_inode(struct m_inode * inode)
21 {

```

```

22     cli();
23     while (inode->i_lock)
24         sleep_on(&inode->i_wait);
25     sti();
26 }
27
28     ///// 对指定的 i 节点上锁（锁定指定的 i 节点）。
29     // 如果 i 节点已被锁定，则将当前任务置为不可中断的等待状态。直到该 i 节点解锁，然后对其上锁。
30 static inline void lock_inode(struct m_inode * inode)
31 {
32     cli();
33     while (inode->i_lock)
34         sleep_on(&inode->i_wait);
35     inode->i_lock=1;                // 置锁定标志。
36     sti();
37 }
38
39     ///// 对指定的 i 节点解锁。
40     // 复位 i 节点的锁定标志，并明确地唤醒等待此 i 节点的进程。
41 static inline void unlock_inode(struct m_inode * inode)
42 {
43     inode->i_lock=0;
44     wake_up(&inode->i_wait);
45 }
46
47     ///// 释放内存中设备 dev 的所有 i 节点。
48     // 扫描内存中的 i 节点表数组，如果是指定设备使用的 i 节点就释放之。
49 void invalidate_inodes(int dev)
50 {
51     int i;
52     struct m_inode * inode;
53
54     inode = 0+inode_table;        // 让指针首先指向 i 节点表指针数组首项。
55     for(i=0 ; i<NR_INODE ; i++,inode++) { // 扫描 i 节点表指针数组中的所有 i 节点。
56         wait_on_inode(inode);        // 等待该 i 节点可用（解锁）。
57         if (inode->i_dev == dev) {    // 如果是指定设备的 i 节点，则
58             if (inode->i_count)      // 如果其引用数不为 0，则显示出错警告；
59                 printk("inode in use on removed disk\n\r");
60             inode->i_dev = inode->i_dirt = 0; // 释放该 i 节点(置设备号为 0 等)。
61         }
62     }
63 }
64
65     ///// 同步所有 i 节点。
66     // 同步内存与设备上的所有 i 节点信息。
67 void sync_inodes(void)
68 {
69     int i;
70     struct m_inode * inode;
71
72     inode = 0+inode_table;        // 让指针首先指向 i 节点表指针数组首项。
73     for(i=0 ; i<NR_INODE ; i++,inode++) { // 扫描 i 节点表指针数组。
74         wait_on_inode(inode);        // 等待该 i 节点可用（解锁）。
75     }
76 }

```

```

67         if (inode->i_dirt && !inode->i_pipe) // 如果该 i 节点已修改且不是管道节点,
68             write\_inode(inode); // 则写盘。
69     }
70 }
71
72     /// 块映射处理操作。(block 位图处理函数, bmap - block map)
73     // 参数: inode - i 节点指针; block - 数据块号; create - 创建标志。
74     // 如果创建标志置位, 则在对应逻辑块不存在时就申请新磁盘块。
75     // 返回 block 数据块对应在设备上的逻辑块号。
76 static int bmap(struct m\_inode * inode, int block, int create)
77 {
78     struct buffer head * bh;
79     int i;
80
81     // 如果块号小于 0, 则死机。
82     if (block < 0)
83         panic("_bmap: block<0");
84     // 如果块号大于直接块数 + 间接块数 + 二次间接块数, 超出文件系统表示范围, 则死机。
85     if (block >= 7+512+512*512)
86         panic("_bmap: block>big");
87
88     // 如果该块号小于 7, 则使用直接块表示。
89     if (block < 7) {
90         // 如果创建标志置位, 并且 i 节点中对应该块的逻辑块(区段)字段为 0, 则向相应设备申请一磁盘
91         // 块(逻辑块, 区段), 并将磁盘上的实际逻辑块号填入逻辑块字段中。然后设置 i 节点修改时间,
92         // 置 i 节点已修改标志。最后返回磁盘上实际逻辑块号。
93         if (create && !inode->i_zone[block])
94             if (inode->i_zone[block]=new\_block(inode->i_dev)) {
95                 inode->i_ctime=CURRENT\_TIME;
96                 inode->i_dirt=1;
97             }
98         return inode->i_zone[block];
99     }
100
101     // 如果该块号 >= 7, 则如果该块小于 7+512, 则说明是一次间接块。下面对一次间接块进行处理。
102     block -= 7;
103     if (block < 512) {
104         // 如果是创建, 并且该 i 节点中对应该间接块字段为 0, 表明文件是首次使用间接块, 则需申请
105         // 一磁盘块用于存放间接块信息, 并将此实际磁盘块号填入间接块字段中。然后设置 i 节点
106         // 已修改标志和修改时间。
107         if (create && !inode->i_zone[7])
108             if (inode->i_zone[7]=new\_block(inode->i_dev)) {
109                 inode->i_dirt=1;
110                 inode->i_ctime=CURRENT\_TIME;
111             }
112         // 若此时 i 节点间接块字段中为 0, 表明申请磁盘块失败, 返回 0 退出。
113         if (!inode->i_zone[7])
114             return 0;
115         // 读取设备上的一次间接块。
116         if (!(bh = bread(inode->i_dev, inode->i_zone[7])))
117             return 0;
118         // 取该间接块上第 block 项中的逻辑块号。
119         i = ((unsigned short *) (bh->b_data))[block];

```

```

// 如果是创建并且间接块的第 block 项中的逻辑块号为 0 的话，则申请一磁盘块（逻辑块），并让
// 间接块中的第 block 项等于该新逻辑块块号。然后置位间接块的已修改标志。
101         if (create && !i)
102             if (i=new_block(inode->i_dev)) {
103                 ((unsigned short *) (bh->b_data))[block]=i;
104                 bh->b_dirt=1;
105             }
// 最后释放该间接块，返回磁盘上新申请的对应 block 的逻辑块的块号。
106         brelse(bh);
107         return i;
108     }

// 程序运行到此，表明数据块是二次间接块，处理过程与一次间接块类似。下面是对二次间接块的处理。
// 将 block 再减去间接块所容纳的块数(512)。
109     block -= 512;
// 如果是新创建并且 i 节点的二次间接块字段为 0，则需申请一磁盘块用于存放二次间接块的一级块
// 信息，并将此实际磁盘块号填入二次间接块字段中。之后置 i 节点已修改编制和修改时间。
110     if (create && !inode->i_zone[8])
111         if (inode->i_zone[8]=new_block(inode->i_dev)) {
112             inode->i_dirt=1;
113             inode->i_ctime=CURRENT_TIME;
114         }
// 若此时 i 节点二次间接块字段为 0，表明申请磁盘块失败，返回 0 退出。
115     if (!inode->i_zone[8])
116         return 0;
// 读取该二次间接块的一级块。
117     if (!(bh=bread(inode->i_dev, inode->i_zone[8])))
118         return 0;
// 取该二次间接块的一级块上第(block/512)项中的逻辑块号。
119     i = ((unsigned short *)bh->b_data)[block>>9];
// 如果是创建并且二次间接块的一级块上第(block/512)项中的逻辑块号为 0 的话，则需申请一磁盘
// 块（逻辑块）作为二次间接块的二级块，并让二次间接块的一级块中第(block/512)项等于该二级
// 块的块号。然后置位二次间接块的一级块已修改标志。并释放二次间接块的一级块。
120     if (create && !i)
121         if (i=new_block(inode->i_dev)) {
122             ((unsigned short *) (bh->b_data))[block>>9]=i;
123             bh->b_dirt=1;
124         }
125     brelse(bh);
// 如果二次间接块的二级块块号为 0，表示申请磁盘块失败，返回 0 退出。
126     if (!i)
127         return 0;
// 读取二次间接块的二级块。
128     if (!(bh=bread(inode->i_dev, i)))
129         return 0;
// 取该二级块上第 block 项中的逻辑块号。(与上 511 是为了限定 block 值不超过 511)
130     i = ((unsigned short *)bh->b_data)[block&511];
// 如果是创建并且二级块的第 block 项中的逻辑块号为 0 的话，则申请一磁盘块（逻辑块），作为
// 最终存放数据信息的块。并让二级块中的第 block 项等于该新逻辑块块号(i)。然后置位二级块的
// 已修改标志。
131     if (create && !i)
132         if (i=new_block(inode->i_dev)) {
133             ((unsigned short *) (bh->b_data))[block&511]=i;

```

```

134         bh->b_dirt=1;
135     }
    // 最后释放该二次间接块的二级块，返回磁盘上新申请的对应 block 的逻辑块的块号。
136     brelse(bh);
137     return i;
138 }
139
140 // 根据 i 节点信息取数据块 block 在设备上对应的逻辑块号。
141 int bmap(struct m_inode * inode, int block)
142 {
143     return _bmap(inode, block, 0);
144 }
145
146 // 创建数据块 block 在设备上对应的逻辑块，并返回在设备上的逻辑块号。
147 int create_block(struct m_inode * inode, int block)
148 {
149     return _bmap(inode, block, 1);
150 }
151
152 // 释放一个 i 节点(回写入设备)。
153 void iput(struct m_inode * inode)
154 {
155     if (!inode)
156         return;
157     wait_on_inode(inode); // 等待 inode 节点解锁(如果已上锁的话)。
158     if (!inode->i_count)
159         panic("iput: trying to free free inode");
160     // 如果是管道 i 节点，则唤醒等待该管道的进程，引用次数减 1，如果还有引用则返回。否则释放
161     // 管道占用的内存页面，并复位该节点的引用计数值、已修改标志和管道标志，并返回。
162     if (inode->i_pipe) {
163         wake_up(&inode->i_wait);
164         if (--inode->i_count)
165             return;
166         free_page(inode->i_size);
167         inode->i_count=0;
168         inode->i_dirt=0;
169         inode->i_pipe=0;
170         return;
171     }
172     // 如果 i 节点对应的设备号=0，则将此节点的引用计数递减 1，返回。
173     if (!inode->i_dev) {
174         inode->i_count--;
175         return;
176     }
177     // 如果是块设备文件的 i 节点，此时逻辑块字段 0 中是设备号，则刷新该设备。并等待 i 节点解锁。
178     if (S_ISBLK(inode->i_mode)) {
179         sync_dev(inode->i_zone[0]);
180         wait_on_inode(inode);
181     }
182     repeat:
183     // 如果 i 节点的引用计数大于 1，则递减 1。
184     if (inode->i_count>1) {
185         inode->i_count--;

```

```

178         return;
179     }
    // 如果 i 节点的链接数为 0，则释放该 i 节点的所有逻辑块，并释放该 i 节点。
180     if (!inode->i_nlinks) {
181         truncate(inode);
182         free_inode(inode);
183         return;
184     }
    // 如果该 i 节点已作过修改，则更新该 i 节点，并等待该 i 节点解锁。
185     if (inode->i_dirt) {
186         write_inode(inode);    /* we can sleep - so do again */
187         wait_on_inode(inode);
188         goto repeat;
189     }
    // i 节点引用计数递减 1。
190     inode->i_count--;
191     return;
192 }
193
    //// 从 i 节点表(inode_table)中获取一个空闲 i 节点项。
    // 寻找引用计数 count 为 0 的 i 节点，并将其写盘后清零，返回其指针。
194 struct m_inode * get_empty_inode(void)
195 {
196     struct m_inode * inode;
197     static struct m_inode * last_inode = inode_table; // last_inode 指向 i 节点表第一项。
198     int i;
199
200     do {
    // 扫描 i 节点表。
201         inode = NULL;
202         for (i = NR_INODE; i ; i--) {
    // 如果 last_inode 已经指向 i 节点表的最后 1 项之后，则让其重新指向 i 节点表开始处。
203             if (++last_inode >= inode_table + NR_INODE)
204                 last_inode = inode_table;
    // 如果 last_inode 所指向的 i 节点的计数值为 0，则说明可能找到空闲 i 节点项。让 inode 指向
    // 该 i 节点。如果该 i 节点的已修改标志和锁定标志均为 0，则我们可以使用该 i 节点，于是退出循环。
205             if (!last_inode->i_count) {
206                 inode = last_inode;
207                 if (!inode->i_dirt && !inode->i_lock)
208                     break;
209             }
210         }
    // 如果没有找到空闲 i 节点(inode=NULL)，则将整个 i 节点表打印出来供调试使用，并死机。
211         if (!inode) {
212             for (i=0 ; i<NR_INODE ; i++)
213                 printk("%04x: %6d\t",inode_table[i].i_dev,
214                     inode_table[i].i_num);
215             panic("No free inodes in mem");
216         }
    // 等待该 i 节点解锁（如果又被上锁的话）。
217         wait_on_inode(inode);
    // 如果该 i 节点已修改标志被置位的话，则将该 i 节点刷新，并等待该 i 节点解锁。
218         while (inode->i_dirt) {

```



```

219         write_inode(inode);
220         wait_on_inode(inode);
221     }
222     } while (inode->i_count); // 如果 i 节点又被其它占用的话, 则重新寻找空闲 i 节点。
// 已找到空闲 i 节点项。则将该 i 节点项内容清零, 并置引用标志为 1, 返回该 i 节点指针。
223     memset(inode, 0, sizeof(*inode));
224     inode->i_count = 1;
225     return inode;
226 }
227
//// 获取管道节点。返回为 i 节点指针 (如果是 NULL 则失败)。
// 首先扫描 i 节点表, 寻找一个空闲 i 节点项, 然后取得一页空闲内存供管道使用。
// 然后将得到的 i 节点的引用计数置为 2(读者和写者), 初始化管道头和尾, 置 i 节点的管道类型表示。
228 struct m_inode * get_pipe_inode(void)
229 {
230     struct m_inode * inode;
231
232     if (!(inode = get_empty_inode())) // 如果找不到空闲 i 节点则返回 NULL。
233         return NULL;
234     if (!(inode->i_size=get_free_page())) { // 节点的 i_size 字段指向缓冲区。
235         inode->i_count = 0; // 如果已没有空闲内存, 则
236         return NULL; // 释放该 i 节点, 并返回 NULL。
237     }
238     inode->i_count = 2; /* sum of readers/writers */ /* 读/写两者总计 */
239     PIPE_HEAD(*inode) = PIPE_TAIL(*inode) = 0;
240     inode->i_pipe = 1; // 置节点为管道使用的标志。
241     return inode; // 返回 i 节点指针。
242 }
243
//// 从设备上读取指定节点号的 i 节点。
// nr - i 节点号。
244 struct m_inode * iget(int dev, int nr)
245 {
246     struct m_inode * inode, * empty;
247
248     if (!dev)
249         panic("iget with dev==0");
250     // 从 i 节点表中取一个空闲 i 节点。
251     empty = get_empty_inode();
252     // 扫描 i 节点表。寻找指定节点号的 i 节点。并递增该节点的引用次数。
253     inode = inode_table;
254     while (inode < NR_INODE+inode_table) {
255         // 如果当前扫描的 i 节点的设备号不等于指定的设备号或者节点号不等于指定的节点号, 则继续扫描。
256         if (inode->i_dev != dev || inode->i_num != nr) {
257             inode++;
258             continue;
259         }
260         // 找到指定设备号和节点号的 i 节点, 等待该节点解锁 (如果已上锁的话)。
261         wait_on_inode(inode);
262         // 在等待该节点解锁的阶段, 节点表可能会发生变化, 所以再次判断, 如果发生了变化, 则再次重新
263         // 扫描整个 i 节点表。
264         if (inode->i_dev != dev || inode->i_num != nr) {
265             inode = inode_table;

```

```

260         continue;
261     }
    // 将该 i 节点引用计数增 1。
262     inode->i_count++;
263     if (inode->i_mount) {
264         int i;
265
    // 如果该 i 节点是其它文件系统的安装点，则在超级块表中搜寻安装在此 i 节点的超级块。如果没有
    // 找到，则显示出错信息，并释放函数开始获取的空闲节点，返回该 i 节点指针。
266         for (i = 0 ; i<NR_SUPER ; i++)
267             if (super_block[i].s_imount==inode)
268                 break;
269         if (i >= NR_SUPER) {
270             printk("Mounted inode hasn't got sb\n");
271             if (empty)
272                 iput(empty);
273             return inode;
274         }
    // 将该 i 节点写盘。从安装在此 i 节点文件系统的超级块上取设备号，并令 i 节点号为 1。然后重新
    // 扫描整个 i 节点表，取该被安装文件系统的根节点。
275         iput(inode);
276         dev = super_block[i].s_dev;
277         nr = ROOT_INO;
278         inode = inode_table;
279         continue;
280     }
    // 已经找到相应的 i 节点，因此放弃临时申请的空闲节点，返回该找到的 i 节点。
281     if (empty)
282         iput(empty);
283     return inode;
284 }
    // 如果在 i 节点表中没有找到指定的 i 节点，则利用前面申请的空闲 i 节点在 i 节点表中建立该节点。
    // 并从相应设备上读取该 i 节点信息。返回该 i 节点。
285     if (!empty)
286         return (NULL);
287     inode=empty;
288     inode->i_dev = dev;
289     inode->i_num = nr;
290     read_inode(inode);
291     return inode;
292 }
293
    //// 从设备上读取指定 i 节点的信息到内存中。
294 static void read_inode(struct m_inode * inode)
295 {
296     struct super_block * sb;
297     struct buffer_head * bh;
298     int block;
299
    // 首先锁定该 i 节点，取该节点所在设备的超级块。
300     lock_inode(inode);
301     if (!(sb=get_super(inode->i_dev)))
302         panic("trying to read inode without dev");

```

```

// 该 i 节点所在的逻辑块号 = (启动块+超级块) + i 节点位图占用的块数 + 逻辑块位图占用的块数 +
// (i 节点号-1)/每块含有的 i 节点数。
303     block = 2 + sb->s_imap_blocks + sb->s_zmap_blocks +
304           (inode->i_num-1)/INODES_PER_BLOCK;
// 从设备上读取该 i 节点所在的逻辑块, 并将该 inode 指针指向对应 i 节点信息。
305     if (!(bh=bread(inode->i_dev, block)))
306         panic("unable to read i-node block");
307     *(struct d_inode *)inode =
308         ((struct d_inode *)bh->b_data)
309         [(inode->i_num-1)%INODES_PER_BLOCK];
// 最后释放读入的缓冲区, 并解锁该 i 节点。
310     brelse(bh);
311     unlock_inode(inode);
312 }
313
//// 将指定 i 节点信息写入设备 (写入对应的缓冲区中, 待缓冲区刷新时会写入盘中)。
314 static void write_inode(struct m_inode * inode)
315 {
316     struct super_block * sb;
317     struct buffer_head * bh;
318     int block;
319
// 首先锁定该 i 节点, 如果该 i 节点没有被修改过或者该 i 节点的设备号等于零, 则解锁该 i 节点,
// 并退出。
320     lock_inode(inode);
321     if (!(inode->i_dirt || !inode->i_dev) {
322         unlock_inode(inode);
323         return;
324     }
// 获取该 i 节点的超级块。
325     if (!(sb=get_super(inode->i_dev)))
326         panic("trying to write inode without device");
// 该 i 节点所在的逻辑块号 = (启动块+超级块) + i 节点位图占用的块数 + 逻辑块位图占用的块数 +
// (i 节点号-1)/每块含有的 i 节点数。
327     block = 2 + sb->s_imap_blocks + sb->s_zmap_blocks +
328           (inode->i_num-1)/INODES_PER_BLOCK;
// 从设备上读取该 i 节点所在的逻辑块。
329     if (!(bh=bread(inode->i_dev, block)))
330         panic("unable to read i-node block");
// 将该逻辑块对应该 i 节点的项指向该 i 节点信息。
331     ((struct d_inode *)bh->b_data)
332     [(inode->i_num-1)%INODES_PER_BLOCK] =
333         *(struct d_inode *)inode;
// 置缓冲区已修改标志, 而 i 节点修改标志置零。然后释放该含有 i 节点的缓冲区, 并解锁该 i 节点。
334     bh->b_dirt=1;
335     inode->i_dirt=0;
336     brelse(bh);
337     unlock_inode(inode);
338 }
339

```

### 8.4.3 其它信息

## 8.5 buffer.c 文件

### 8.5.1 功能描述

该文件中的函数主要用于对设备高速缓冲的操作和处理。

空闲缓冲区链表中的缓冲区，并不是都是空闲的。!! 只有当被写盘刷新、解锁且没有其它进程引用时（引用计数=0），才能挪作它用。

### 8.5.2 代码注释

列表 linux/fs/buffer.c 程序

---

```

1  /*
2   *  linux/fs/buffer.c
3   *
4   *  (C) 1991  Linus Torvalds
5   */
6
7  /*
8   *  'buffer.c' implements the buffer-cache functions. Race-conditions have
9   *  been avoided by NEVER letting a interrupt change a buffer (except for the
10  *  data, of course), but instead letting the caller do it. NOTE! As interrupts
11  *  can wake up a caller, some cli-sti sequences are needed to check for
12  *  sleep-on-calls. These should be extremely quick, though (I hope).
13  */
14  /*
15   *  'buffer.c' 用于实现缓冲区高速缓存功能。通过不让中断过程改变缓冲区，而是让调用者
16   *  来执行，避免了竞争条件（当然除改变数据以外）。注意！由于中断可以唤醒一个调用者，
17   *  因此就需要开关中断指令（cli-sti）序列来检测等待调用返回。但需要非常地快（希望是这样）。
18   */
19  /*
20   *  NOTE! There is one discordant note here: checking floppies for
21   *  disk change. This is where it fits best, I think, as it should
22   *  invalidate changed floppy-disk-caches.
23   */
24  /*
25   *  注意！这里有一个程序应不属于这里：检测软盘是否更换。但我想这里是
26   *  放置该程序最好的地方了，因为它需要使已更换软盘缓冲失效。
27   */
28
29  #include <stdarg.h>          // 标准参数头文件。以宏的形式定义变量参数列表。主要说明了一个
30                               // 类型(va_list)和三个宏(va_start, va_arg 和 va_end)，用于
31                               // vsprintf、vprintf、vfprintf 函数。
32
33  #include <linux/config.h>    // 内核配置头文件。定义键盘语言和硬盘类型（HD_TYPE）可选项。
34  #include <linux/sched.h>    // 调度程序头文件，定义了任务结构 task_struct、初始任务 0 的数据，
35                               // 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
36  #include <linux/kernel.h>    // 内核头文件。含有一些内核常用函数的原形定义。
37  #include <asm/system.h>      // 系统头文件。定义了设置或修改描述符/中断门等的嵌入式汇编宏。
38  #include <asm/io.h>          // io 头文件。定义硬件端口输入/输出宏汇编语句。

```

```

28
29 extern int end;           // 由连接程序 ld 生成的表明程序末端的变量。[??]
30 struct buffer_head * start_buffer = (struct buffer_head *) &end;
31 struct buffer_head * hash_table[NR_HASH];           // NR_HASH = 307 项。
32 static struct buffer_head * free_list;
33 static struct task_struct * buffer_wait = NULL;
34 int NR_BUFFERS = 0;
35
36 // 等待指定缓冲区解锁。
37 static inline void wait_on_buffer(struct buffer_head * bh)
38 {
39     cli();           // 关中断。
40     while (bh->b_lock)           // 如果已被上锁，则进程进入睡眠，等待其解锁。
41         sleep_on(&bh->b_wait);
42     sti();           // 开中断。
43 }
44
45 // 系统调用。同步设备和内存高速缓冲中数据。
46 int sys_sync(void)
47 {
48     int i;
49     struct buffer_head * bh;
50
51     sync_inodes();           /* write out inodes into buffers */ /*将 i 节点写入高速缓冲*/
52 // 扫描所有高速缓冲区，对于已被修改的缓冲块产生写盘请求，将缓冲中数据与设备中同步。
53     bh = start_buffer;
54     for (i=0 ; i<NR_BUFFERS ; i++,bh++) {
55         wait_on_buffer(bh);           // 等待缓冲区解锁（如果已上锁的话）。
56         if (bh->b_dirt)
57             ll_rw_block(WRITE, bh);           // 产生写设备块请求。
58     }
59     return 0;
60 }
61
62 // 对指定设备进行高速缓冲数据与设备上数据的同步操作。
63 int sync_dev(int dev)
64 {
65     int i;
66     struct buffer_head * bh;
67
68     bh = start_buffer;
69     for (i=0 ; i<NR_BUFFERS ; i++,bh++) {
70         if (bh->b_dev != dev)
71             continue;
72         wait_on_buffer(bh);
73         if (bh->b_dev == dev && bh->b_dirt)
74             ll_rw_block(WRITE, bh);
75     }
76     sync_inodes();           // 将 i 节点数据写入高速缓冲。
77     bh = start_buffer;
78     for (i=0 ; i<NR_BUFFERS ; i++,bh++) {
79         if (bh->b_dev != dev)
80             continue;

```

```

77         wait_on_buffer(bh);
78         if (bh->b_dev == dev && bh->b_dirt)
79             ll_rw_block(WRITE, bh);
80     }
81     return 0;
82 }
83
84 // 使指定设备在高速缓冲区中的数据无效。
85 // 扫描高速缓冲中的所有缓冲块，对于指定设备的缓冲区，复位其有效(更新)标志和已修改标志。
86 void inline invalidate_buffers(int dev)
87 {
88     int i;
89     struct buffer_head * bh;
90
91     bh = start_buffer;
92     for (i=0 ; i<NR_BUFFERS ; i++,bh++) {
93         if (bh->b_dev != dev) // 如果不是指定设备的缓冲块，则
94             continue; // 继续扫描下一块。
95         wait_on_buffer(bh); // 等待该缓冲区解锁（如果已被上锁）。
96     }
97     // 由于进程执行过睡眠等待，所以需要再判断一下缓冲区是否是指定设备的。
98     if (bh->b_dev == dev)
99         bh->b_uptodate = bh->b_dirt = 0;
100 }
101
102 /*
103  * This routine checks whether a floppy has been changed, and
104  * invalidates all buffer-cache-entries in that case. This
105  * is a relatively slow routine, so we have to try to minimize using
106  * it. Thus it is called only upon a 'mount' or 'open'. This
107  * is the best way of combining speed and utility, I think.
108  * People changing diskettes in the middle of an operation deserve
109  * to loose :-))
110  *
111  * NOTE! Although currently this is only for floppies, the idea is
112  * that any additional removable block-device will use this routine,
113  * and that mount/open needn't know that floppies/whatever are
114  * special.
115  */
116
117 /*
118  * 该子程序检查一个软盘是否已经被更换，如果已经更换就使高速缓冲中与该软驱
119  * 对应的所有缓冲区无效。该子程序相对来说较慢，所以我们要尽量少使用它。
120  * 所以仅在执行'mount'或'open'时才调用它。我想这是将速度和实用性相结合的
121  * 最好方法。若在操作过程当中更换软盘，会导致数据的丢失，这是咎由自取☺。
122  *
123  * 注意！尽管目前该子程序仅用于软盘，以后任何可移动介质的块设备都将使用该
124  * 程序，mount/open 操作是不需要知道是否是软盘或其它什么特殊介质的。
125  */
126
127 // 检查磁盘是否更换，如果已更换就使对应高速缓冲区无效。
128 void check_disk_change(int dev)
129 {
130     int i;
131
132

```

```

// 是软盘设备吗？如果不是则退出。
117     if (MAJOR(dev) != 2)
118         return;
// 测试对应软盘是否已更换，如果没有则退出。
119     if (!floppy_change(dev & 0x03))
120         return;
// 软盘已经更换，所以释放对应设备的 i 节点位图和逻辑块位图所占的高速缓冲区；并使该设备的
// i 节点和数据块信息所占的高速缓冲区无效。
121     for (i=0 ; i<NR_SUPER ; i++)
122         if (super_block[i].s_dev == dev)
123             put_super(super_block[i].s_dev);
124     invalidate_inodes(dev);
125     invalidate_buffers(dev);
126 }
127
// hash 函数和 hash 表项的计算宏定义。
128 #define hashfn(dev,block) (((unsigned)(dev^block))%NR_HASH)
129 #define hash(dev,block) hash_table[hashfn(dev,block)]
130
//// 从 hash 队列和空闲缓冲队列中移走指定的缓冲块。
131 static inline void remove_from_queues(struct buffer_head * bh)
132 {
133     /* remove from hash-queue */
134     /* 从 hash 队列中移除缓冲块 */
135     if (bh->b_next)
136         bh->b_next->b_prev = bh->b_prev;
137     if (bh->b_prev)
138         bh->b_prev->b_next = bh->b_next;
139     // 如果该缓冲区是该队列的头一个块，则让 hash 表的对应项指向本队列中的下一个缓冲区。
140     if (hash(bh->b_dev, bh->b_blocknr) == bh)
141         hash(bh->b_dev, bh->b_blocknr) = bh->b_next;
142     /* remove from free list */
143     /* 从空闲缓冲区表中移除缓冲块 */
144     if (!(bh->b_prev_free) || !(bh->b_next_free))
145         panic("Free block list corrupted");
146     bh->b_prev_free->b_next_free = bh->b_next_free;
147     bh->b_next_free->b_prev_free = bh->b_prev_free;
148     // 如果空闲链表头指向本缓冲区，则让其指向下一缓冲区。
149     if (free_list == bh)
150         free_list = bh->b_next_free;
151 }
152
//// 将指定缓冲区插入空闲链表尾并放入 hash 队列中。
153 static inline void insert_into_queues(struct buffer_head * bh)
154 {
155     /* put at end of free list */
156     /* 放在空闲链表末尾处 */
157     bh->b_next_free = free_list;
158     bh->b_prev_free = free_list->b_prev_free;
159     free_list->b_prev_free->b_next_free = bh;
160     free_list->b_prev_free = bh;
161     /* put the buffer in new hash-queue if it has a device */
162     /* 如果该缓冲块对应一个设备，则将其插入新 hash 队列中 */

```

```

157     bh->b_prev = NULL;
158     bh->b_next = NULL;
159     if (!bh->b_dev)
160         return;
161     bh->b_next = hash(bh->b_dev, bh->b_blocknr);
162     hash(bh->b_dev, bh->b_blocknr) = bh;
163     bh->b_next->b_prev = bh;
164 }
165
166 // 在高速缓冲中寻找给定设备和指定块的缓冲区块。
167 // 如果找到则返回缓冲区块的指针，否则返回 NULL。
168 static struct buffer head * find buffer(int dev, int block)
169 {
170     struct buffer head * tmp;
171     for (tmp = hash(dev, block) ; tmp != NULL ; tmp = tmp->b_next)
172         if (tmp->b_dev==dev && tmp->b_blocknr==block)
173             return tmp;
174     return NULL;
175 }
176 /*
177  * Why like this, I hear you say... The reason is race-conditions.
178  * As we don't lock buffers (unless we are readint them, that is),
179  * something might happen to it while we sleep (ie a read-error
180  * will force it bad). This shouldn't really happen currently, but
181  * the code is ready.
182  */
183 /*
184  * 代码为什么会是这样子的？我听见你问... 原因是竞争条件。由于我们没有对
185  * 缓冲区上锁（除非我们正在读取它们中的数据），那么当我们（进程）睡眠时
186  * 缓冲区可能会发生一些问题（例如一个读错误将导致该缓冲区出错）。目前
187  * 这种情况实际上是不会发生的，但处理的代码已经准备好了。
188  */
189 // 在高速缓冲中寻找给定设备和指定块的缓冲区，如果没有找到则返回 NULL，退出。
190 struct buffer head * get hash table(int dev, int block)
191 {
192     struct buffer head * bh;
193     for (;;) {
194         if (!(bh=find buffer(dev, block)))
195             return NULL;
196         // 对该缓冲区增加引用计数，并等待该缓冲区解锁（如果已被上锁）。
197         bh->b_count++;
198         wait on buffer(bh);
199         // 由于经过了睡眠状态，因此有必要再验证该缓冲区块的正确性，并返回缓冲区头指针。
200         if (bh->b_dev == dev && bh->b_blocknr == block)
201             return bh;
202         // 如果该缓冲区所属的设备号或块号在睡眠时发生了改变，则撤消对它的引用计数，重新寻找。
203         bh->b_count--;
204     }
205 }

```



```

197
198 /*
199  * Ok, this is getblk, and it isn't very clear, again to hinder
200  * race-conditions. Most of the code is seldom used, (ie repeating),
201  * so it should be much more efficient than it looks.
202  *
203  * The algorithm is changed: hopefully better, and an elusive bug removed.
204  */
/*
  * OK, 下面是 getblk 函数, 该函数的逻辑并不是很清晰, 同样也是因为要考虑
  * 竞争条件问题。其中大部分代码很少用到, (例如重复操作语句), 因此它应该
  * 比看上去的样子有效得多。
  *
  * 算法以及作了改变: 希望能更好, 而且一个难以琢磨的错误已经去除。
  */
// 下面宏定义用于同时判断缓冲区的修改标志和锁定标志, 并且定义修改标志的权重要比锁定标志大。
205 #define BADNESS(bh) (((bh)->b_dirt<<1)+(bh)->b_lock)
//// 取高速缓冲中指定的缓冲区。
// 检查所指定的缓冲区是否已经在高速缓冲中, 如果不在, 就需要在高速缓冲中建立一个对应的新项。
// 返回相应缓冲区头指针。
206 struct buffer_head * getblk(int dev, int block)
207 {
208     struct buffer_head * tmp, * bh;
209
210 repeat:
  // 搜索 hash 表, 如果指定块已经在高速缓冲中, 则返回对应缓冲区头指针, 退出。
211     if (bh = get_hash_table(dev, block))
212         return bh;
  // 扫描空闲数据块链表, 寻找空闲缓冲区。
  // 首先让 tmp 指向空闲链表的第一个空闲缓冲区头。
213     tmp = free_list;
214     do {
  // 如果该缓冲区正被使用 (引用计数不等于 0), 则继续扫描下一项。
215         if (tmp->b_count)
216             continue;
  // 如果缓冲区头指针 bh 为空, 或者 tmp 所指缓冲区头的标志 (修改、锁定) 少于 (小于) bh 头的标志,
  // 则让 bh 指向该 tmp 缓冲区头。如果该 tmp 缓冲区头表明缓冲区既没有修改也没有锁定标志置位,
  // 则说明已为指定设备上的块取得对应的高速缓冲区, 则退出循环。
217         if (!bh || BADNESS(tmp) < BADNESS(bh)) {
218             bh = tmp;
219             if (!BADNESS(tmp))
220                 break;
221         }
222 /* and repeat until we find something good */ /* 重复操作直到找到适合的缓冲区 */
223     } while ((tmp = tmp->b_next_free) != free_list);
  // 如果所有缓冲区都正被使用 (所有缓冲区的头部引用计数都>0), 则睡眠, 等待有空闲的缓冲区可用。
224     if (!bh) {
225         sleep_on(&buffer_wait);
226         goto repeat;
227     }
  // 等待该缓冲区解锁 (如果已被上锁的话)。
228     wait_on_buffer(bh);
  // 如果该缓冲区又被其它任务使用的话, 只好重复上述过程。

```

```

229         if (bh->b_count)
230             goto repeat;
// 如果该缓冲区已被修改，则将数据写盘，并再次等待缓冲区解锁。如果该缓冲区又被其它任务使用
// 的话，只好再重复上述过程。
231         while (bh->b_dirt) {
232             sync\_dev(bh->b_dev);
233             wait\_on\_buffer(bh);
234             if (bh->b_count)
235                 goto repeat;
236         }
237 /* NOTE!! While we slept waiting for this block, somebody else might */
238 /* already have added "this" block to the cache. check it */
// 注意！！当进程为了等待该缓冲块而睡眠时，其它进程可能已经将该缓冲块 *
// * 加入进高速缓冲中，所以要对此进行检查。*/
// 在高速缓冲 hash 表中检查指定设备和块的缓冲区是否已经被加入进去。如果是的话，就再次重复
// 上述过程。
239         if (find\_buffer(dev, block))
240             goto repeat;
241 /* OK, FINALLY we know that this buffer is the only one of it's kind, */
242 /* and that it's unused (b_count=0), unlocked (b_lock=0), and clean */
// OK，最终我们知道该缓冲区是指定参数的唯一一块，*/
// 而且还没有被使用(b_count=0)，未被上锁(b_lock=0)，并且是干净的（未被修改的）*/
// 于是让我们占用此缓冲区。置引用计数为 1，复位修改标志和有效(更新)标志。
243         bh->b_count=1;
244         bh->b_dirt=0;
245         bh->b_uptodate=0;
// 从 hash 队列和空闲块链表中移出该缓冲区头，让该缓冲区用于指定设备和其上的指定块。
246         remove\_from\_queues(bh);
247         bh->b_dev=dev;
248         bh->b_blocknr=block;
// 然后根据此新的设备号和块号重新插入空闲链表和 hash 队列新位置处。并最终返回缓冲头指针。
249         insert\_into\_queues(bh);
250         return bh;
251     }
252
//// 释放指定的缓冲区。
// 等待该缓冲区解锁。引用计数递减 1。唤醒等待空闲缓冲区的进程。
253 void brelse(struct buffer head * buf)
254 {
255     if (!buf) // 如果缓冲头指针无效则返回。
256         return;
257     wait\_on\_buffer(buf);
258     if (!(buf->b_count--))
259         panic("Trying to free free buffer");
260     wake\_up(&buffer\_wait);
261 }
262
263 /*
264 * bread() reads a specified block and returns the buffer that contains
265 * it. It returns NULL if the block was unreadable.
266 */
//
// * 从设备上读取指定的数据块并返回含有数据的缓冲区。如果指定的块不存在

```

```

    * 则返回 NULL。
    */
    //// 从指定设备上读取指定的数据块。
267 struct buffer\_head * bread(int dev,int block)
268 {
269     struct buffer\_head * bh;
270
    // 在高速缓冲中申请一块缓冲区。如果返回值是 NULL 指针，表示内核出错，死机。
271     if (!(bh=getblk(dev,block)))
272         panic("bread: getblk returned NULL\n");
    // 如果该缓冲区中的数据是有效的（已更新的）可以直接使用，则返回。
273     if (bh->b_uptodate)
274         return bh;
    // 否则调用 ll_rw_block() 函数，产生读设备块请求。并等待缓冲区解锁。
275     ll\_rw\_block(READ, bh);
276     wait\_on\_buffer(bh);
    // 如果该缓冲区已更新，则返回缓冲区头指针，退出。
277     if (bh->b_uptodate)
278         return bh;
    // 否则表明读设备操作失败，释放该缓冲区，返回 NULL 指针，退出。
279     brelse(bh);
280     return NULL;
281 }
282
    //// 复制内存块。
    // 从 from 地址复制一块数据到 to 位置。
283 #define COPYBLK(from,to) \
284 __asm__ ("cld\n\t" \
285         "rep\n\t" \
286         "movsl\n\t" \
287         :: "c" (BLOCK\_SIZE/4), "S" (from), "D" (to) \
288         : "cx", "di", "si")
289
290 /*
291  * bread_page reads four buffers into memory at the desired address. It's
292  * a function of its own, as there is some speed to be got by reading them
293  * all at the same time, not waiting for one to be read, and then another
294  * etc.
295  */
    /*
    * bread_page 一次读四个缓冲块内容读到内存指定的地址。它是一个完整的函数，
    * 因为同时读取四块可以获得速度上的好处，不用等着读一块，再读一块了。
    */
    //// 读设备上一个页面（4 个缓冲区）的内容到内存指定的地址。
296 void bread\_page(unsigned long address,int dev,int b[4])
297 {
298     struct buffer\_head * bh[4];
299     int i;
300
    // 循环执行 4 次，读一页内容。
301     for (i=0 ; i<4 ; i++)
302         if (b[i]) {
    // 取高速缓冲中指定设备和块号的缓冲区，如果该缓冲区数据无效则产生读设备请求。

```

```

303         if (bh[i] = getblk(dev, b[i]))
304             if (!bh[i]->b_uptodate)
305                 ll\_rw\_block(READ, bh[i]);
306     } else
307         bh[i] = NULL;
308 // 将4块缓冲区上的内容顺序复制到指定地址处。
309     for (i=0 ; i<4 ; i++, address += BLOCK\_SIZE)
310         if (bh[i]) {
311             wait\_on\_buffer(bh[i]);    // 等待缓冲区解锁(如果已被上锁的话)。
312             if (bh[i]->b_uptodate)    // 如果该缓冲区中数据有效的话, 则复制。
313                 COPYBLK((unsigned long) bh[i]->b_data, address);
314             brelse(bh[i]);            // 释放该缓冲区。
315         }
316 }
317 /*
318  * Ok, breada can be used as bread, but additionally to mark other
319  * blocks for reading as well. End the argument list with a negative
320  * number.
321  */
322 /*
323  * OK, breada 可以象 bread 一样使用, 但会另外预读一些块。该函数参数列表
324  * 需要使用一个负数来表明参数列表的结束。
325  */
326 // 从指定设备读取指定的一些块。
327 // 成功时返回第1块的缓冲区头指针, 否则返回 NULL。
328 struct buffer\_head * breada(int dev, int first, ...)
329 {
330     va\_list args;
331     struct buffer\_head * bh, *tmp;
332
333     // 取可变参数表中第1个参数(块号)。
334     va\_start(args, first);
335     // 取高速缓冲中指定设备和块号的缓冲区。如果该缓冲区数据无效, 则发出读设备数据块请求。
336     if (!(bh=getblk(dev, first)))
337         panic("bread: getblk returned NULL\n");
338     if (!bh->b_uptodate)
339         ll\_rw\_block(READ, bh);
340     // 然后顺序取可变参数表中其它预读块号, 并作与上面同样处理, 但不引用。
341     while ((first=va\_arg(args, int))>=0) {
342         tmp=getblk(dev, first);
343         if (tmp) {
344             if (!tmp->b_uptodate)
345                 ll\_rw\_block(READA, bh);
346             tmp->b_count--;
347         }
348     }
349     // 可变参数表中所有参数处理完毕。等待第1个缓冲区解锁(如果已被上锁)。
350     va\_end(args);
351     wait\_on\_buffer(bh);
352     // 如果缓冲区中数据有效, 则返回缓冲区头指针, 退出。否则释放该缓冲区, 返回 NULL, 退出。
353     if (bh->b_uptodate)
354         return bh;

```

```

344     brelse(bh);
345     return (NULL);
346 }
347
348     /// 缓冲区初始化函数。
349     // 参数 buffer_end 是指定的缓冲区内内存的末端。对于系统有 16MB 内存，则缓冲区末端设置为 4MB。
350     // 对于系统有 8MB 内存，缓冲区末端设置为 2MB。
351     void buffer_init(long buffer_end)
352     {
353         struct buffer_head * h = start_buffer;
354         void * b;
355         int i;
356
357         // 如果缓冲区高端等于 1Mb，则由于从 640KB-1MB 被显示内存和 BIOS 占用，因此实际可用缓冲区内内存
358         // 高端应该是 640KB。否则内存高端一定大于 1MB。
359         if (buffer_end == 1<<20)
360             b = (void *) (640*1024);
361         else
362             b = (void *) buffer_end;
363
364         // 这段代码用于初始化缓冲区，建立空闲缓冲区环链表，并获取系统中缓冲块的数目。
365         // h 是指向缓冲头结构的指针，而 h+1 是指向内存地址连续的下一个缓冲头地址，也可以说是指向 h
366         // 缓冲头的末端外。为了保证有足够长度的内存来存储一个缓冲头结构，需要 b 所指向的内存块
367         // 地址 >= h 缓冲头的末端，也即要>=h+1。
368         while ( (b -= BLOCK_SIZE) >= ((void *) (h+1)) ) {
369             h->b_dev = 0;           // 使用该缓冲区的设备号。
370             h->b_dirt = 0;          // 脏标志，也即缓冲区修改标志。
371             h->b_count = 0;         // 该缓冲区引用计数。
372             h->b_lock = 0;          // 缓冲区锁定标志。
373             h->b_uptodate = 0;      // 缓冲区更新标志（或称数据有效标志）。
374             h->b_wait = NULL;       // 指向等待该缓冲区解锁的进程。
375             h->b_next = NULL;       // 指向具有相同 hash 值的下一个缓冲头。
376             h->b_prev = NULL;       // 指向具有相同 hash 值的前一个缓冲头。
377             h->b_data = (char *) b; // 指向对应缓冲区数据块（1024 字节）。
378             h->b_prev_free = h-1;   // 指向链表中前一项。
379             h->b_next_free = h+1;   // 指向链表中下一项。
380             h++;                   // h 指向下一新缓冲头位置。
381             NR_BUFFERS++;          // 缓冲区块数累加。
382             if (b == (void *) 0x100000) // 如果地址 b 递减到等于 1MB，则跳过 384KB，
383                 b = (void *) 0xA0000; // 让 b 指向地址 0xA0000 (640KB) 处。
384         }
385         h--;                       // 让 h 指向最后一个有效缓冲头。
386         free_list = start_buffer;  // 让空闲链表头指向头一个缓冲区头。
387         free_list->b_prev_free = h; // 链表头的 b_prev_free 指向前一项（即最后一项）。
388         h->b_next_free = free_list; // h 的下一项指针指向第一项，形成一个环链。
389
390         // 初始化 hash 表（哈希表、散列表），置表中所有的指针为 NULL。
391         for (i=0; i<NR_HASH; i++)
392             hash_table[i]=NULL;
393     }
394 }
395
396

```

### 8.5.3 其它信息

## 8.6 block\_dev.c 文件

### 8.6.1 功能描述

该文件包括 `block_read()` 和 `block_write()` 两个函数。这两个函数是供系统调用函数 `read()` 和 `write()` 调用的，其它地方没有引用。

块设备读写操作所使用的函数之间的层次关系为：

- `read()`, `write()`
- `block_read()`, `block_write()`, `file_read()`, `file_write()`, `read_pipe()`, `write()`
- `bread()` 或 `breada()`
- `getblk()`
- `ll_rw_block()`

### 8.6.2 代码注释

列表 linux/fs/block\_dev.c 程序

```

1  /*
2   *  linux/fs/block_dev.c
3   *
4   *  (C) 1991  Linus Torvalds
5   */
6
7  #include <errno.h>          // 错误号头文件。包含系统中各种出错号。(Linus 从 minix 中引进的)。
8
9  #include <linux/sched.h>    // 调度程序头文件，定义了任务结构 task_struct、初始任务 0 的数据，
                             // 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
10 #include <linux/kernel.h>   // 内核头文件。含有一些内核常用函数的原形定义。
11 #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。
12 #include <asm/system.h>     // 系统头文件。定义了设置或修改描述符/中断门等的嵌入式汇编宏。
13
14  ///// 数据块写函数 - 向指定设备从给定偏移处写入指定字节数据。
15  // 参数: dev - 设备号; pos - 设备文件中偏移量指针; buf - 用户地址空间中缓冲区地址;
16  //        count - 要传送的字节数。
17  // 对于内核来说，写操作是向高速缓冲区中写入数据，什么时候数据最终写入设备是由高速缓冲管理
18  // 程序决定并处理的。另外，因为设备是以块为单位进行读写的，因此对于写开始位置不处于块起始
19  // 处时，需要先将开始字节所在的整个块读出，然后将需要写的数据从写开始处填写满该块，再将完
20  // 整的一块数据写盘（即交由高速缓冲程序去处理）。
21  int block_write(int dev, long * pos, char * buf, int count)
22  {
23      // 由 pos 地址换算成开始读写块的块序号 block。并求出需读第 1 字节在该块中的偏移位置 offset。
24      int block = *pos >> BLOCK_SIZE_BITS;
25      int offset = *pos & (BLOCK_SIZE-1);
26      int chars;

```

```

19     int written = 0;
20     struct buffer head * bh;
21     register char * p;
22
23     // 针对要写入的字节数 count，循环执行以下操作，直到全部写入。
24     while (count>0) {
25         // 计算在该块中可写入的字节数。如果需要写入的字节数不满一块，则只需写 count 字节。
26         chars = BLOCK\_SIZE - offset;
27         if (chars > count)
28             chars=count;
29         // 如果正好要写 1 块数据，则直接申请 1 块高速缓冲块，否则需要读入将被修改的数据块，并预读
30         // 下两块数据，然后将块号递增 1。
31         if (chars == BLOCK\_SIZE)
32             bh = getblk(dev, block);
33         else
34             bh = breada(dev, block, block+1, block+2, -1);
35         block++;
36         // 如果写出错，则返回已写字节数，如果没有写入任何字节，则返回出错号（负数）。
37         if (!bh)
38             return written?written:-EIO;
39         // p 指向读出数据块中开始写的位置。若最后写入的数据不足一块，则需从块开始填写（修改）所需
40         // 的字节，因此这里需置 offset 为零。
41         p = offset + bh->b_data;
42         offset = 0;
43         // 将文件中偏移指针前移已写字节数。累加已写字节数 chars。传送计数值减去此次已传送字节数。
44         *pos += chars;
45         written += chars;
46         count -= chars;
47         // 从用户缓冲区复制 chars 字节到 p 指向的高速缓冲区中开始写入的位置。
48         while (chars-->0)
49             *(p++) = get\_fs\_byte(buf++);
50         // 置该缓冲区块已修改标志，并释放该缓冲区（也即该缓冲区引用计数递减 1）。
51         bh->b_dirt = 1;
52         brelse(bh);
53     }
54     return written; // 返回已写入的字节数，正常退出。
55 }
56
57 // 数据块读函数 - 从指定设备和位置读入指定字节数的数据到高速缓冲中。
58 int block\_read(int dev, unsigned long * pos, char * buf, int count)
59 {
60     // 由 pos 地址换算成开始读写块的块序号 block。并求出需读第 1 字节在该块中的偏移位置 offset。
61     int block = *pos >> BLOCK\_SIZE\_BITS;
62     int offset = *pos & (BLOCK\_SIZE-1);
63     int chars;
64     int read = 0;
65     struct buffer head * bh;
66     register char * p;
67
68     // 针对要读入的字节数 count，循环执行以下操作，直到全部读入。
69     while (count>0) {
70         // 计算在该块中需读入的字节数。如果需要读入的字节数不满一块，则只需读 count 字节。
71         chars = BLOCK\_SIZE-offset;

```



```

58         if (chars > count)
59             chars = count;
// 读入需要的数据块，并预读下两块数据，如果读操作出错，则返回已读字节数，如果没有读入任何
// 字节，则返回出错号。然后将块号递增 1。
60         if (!(bh = breada(dev, block, block+1, block+2, -1)))
61             return read?read:-EIO;
62         block++;
// p 指向从设备读出数据块中需要读取的开始位置。若最后需要读取的数据不足一块，则需从块开始
// 读取所需的字节，因此这里需将 offset 置零。
63         p = offset + bh->b_data;
64         offset = 0;
// 将文件中偏移指针前移已读出字节数 chars。累加已读字节数。传送计数值减去此次已传送字节数。
65         *pos += chars;
66         read += chars;
67         count -= chars;
// 从高速缓冲区中 p 指向的开始位置复制 chars 字节数据到用户缓冲区，并释放该高速缓冲区。
68         while (chars-->0)
69             put_fs_byte(*(p++), buf++);
70         brelse(bh);
71     }
72     return read; // 返回已读取的字节数，正常退出。
73 }
74

```

### 8.6.3 其它信息

## 8.7 file\_dev.c 文件

### 8.7.1 功能描述

该文件包括 file\_read() 和 file\_write() 两个函数。这两个函数是供系统调用函数 read() 和 write() 调用的，其它地方没有引用。

### 8.7.2 代码注释

列表 linux/fs/file\_dev.c 程序

```

1  /*
2  *  linux/fs/file_dev.c
3  *
4  *  (C) 1991  Linus Torvalds
5  */
6
7  #include <errno.h> // 错误号头文件。包含系统中各种出错号。(Linus 从 minix 中引进的)
8  #include <fcntl.h> // 文件控制头文件。用于文件及其描述符的操作控制常数符号的定义。
9
10 #include <linux/sched.h> // 调度程序头文件，定义了任务结构 task_struct、初始任务 0 的数据，
// 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。

```



```

11 #include <linux/kernel.h> // 内核头文件。含有一些内核常用函数的原形定义。
12 #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。
13
14 #define MIN(a,b) (((a)<(b))?(a):(b)) // 取 a,b 中的最小值。
15 #define MAX(a,b) (((a)>(b))?(a):(b)) // 取 a,b 中的最大值。
16
17 // 文件读函数 - 根据 i 节点和文件结构，读设备数据。
18 // 由 i 节点可以知道设备号，由 filp 结构可以知道文件中当前读写指针位置。buf 指定用户态中
19 // 缓冲区的位置，count 为需要读取的字节数。返回值是实际读取的字节数，或出错号(小于 0)。
20 int file_read(struct m_inode * inode, struct file * filp, char * buf, int count)
21 {
22     int left, chars, nr;
23     struct buffer_head * bh;
24
25     // 若需要读取的字节计数值小于等于零，则返回。
26     if ((left=count)<=0)
27         return 0;
28     // 若还需要读取的字节数不等于 0，就循环执行以下操作，直到全部读出。
29     while (left) {
30         // 根据 i 节点和文件表结构信息，取数据块文件当前读写位置在设备上对应的逻辑块号 nr。若 nr 不
31         // 为 0，则从 i 节点指定的设备上读取该逻辑块，如果读操作失败则退出循环。若 nr 为 0，表示指定
32         // 的数据块不存在，置缓冲块指针为 NULL。
33         if (nr = bmap(inode, (filp->f_pos)/BLOCK_SIZE)) {
34             if (!(bh=bread(inode->i_dev, nr)))
35                 break;
36         } else
37             bh = NULL;
38         // 计算文件读写指针在数据块中的偏移值 nr，则该块中可读字节数为 (BLOCK_SIZE-nr)，然后与还需
39         // 读取的字节数 left 作比较，其中小值即为本次需读的字节数 chars。若 (BLOCK_SIZE-nr) 大则说明
40         // 该块是需要读取的最后一块数据，反之则还需要读取一块数据。
41         nr = filp->f_pos % BLOCK_SIZE;
42         chars = MIN( BLOCK_SIZE-nr, left );
43         // 调整读写文件指针。指针前移此次将读取的字节数 chars。剩余字节计数相应减去 chars。
44         filp->f_pos += chars;
45         left -= chars;
46         // 若从设备上读到了数据，则将 p 指向读出数据块缓冲区中开始读取的位置，并且复制 chars 字节
47         // 到用户缓冲区 buf 中。否则往用户缓冲区中填入 chars 个 0 值字节。
48         if (bh) {
49             char * p = nr + bh->b_data;
50             while (chars-->0)
51                 put_fs_byte(*(p++), buf++);
52             brelse(bh);
53         } else {
54             while (chars-->0)
55                 put_fs_byte(0, buf++);
56         }
57     }
58     // 修改该 i 节点的访问时间为当前时间。返回读取的字节数，若读取字节数为 0，则返回出错号。
59     inode->i_atime = CURRENT_TIME;
60     return (count-left)?(count-left):-ERROR;
61 }
62
63 // 文件写函数 - 根据 i 节点和文件结构信息，将用户数据写入指定设备。

```

```

// 由 i 节点可以知道设备号, 由 filp 结构可以知道文件中当前读写指针位置。buf 指定用户态中
// 缓冲区的位置, count 为需要写入的字节数。返回值是实际写入的字节数, 或出错号(小于 0)。
48 int file_write(struct m_inode * inode, struct file * filp, char * buf, int count)
49 {
50     off_t pos;
51     int block, c;
52     struct buffer_head * bh;
53     char * p;
54     int i=0;
55
56     /*
57      * ok, append may not work when many processes are writing at the same time
58      * but so what. That way leads to madness anyway.
59      */
60     /*
61      * ok, 当许多进程同时写时, append 操作可能不行, 但那又怎样。不管怎样那样做会
62      * 导致混乱一团。
63      */
64     // 如果是要向文件后添加数据, 则将文件读写指针移到文件尾部。否则就将在文件读写指针处写入。
65     if (filp->f_flags & O_APPEND)
66         pos = inode->i_size;
67     else
68         pos = filp->f_pos;
69     // 若已写入字节数 i 小于需要写入的字节数 count, 则循环执行以下操作。
70     while (i < count) {
71         // 创建数据块号(pos/BLOCK_SIZE)在设备上对应的逻辑块, 并返回在设备上的逻辑块号。如果逻辑
72         // 块号=0, 则表示创建失败, 退出循环。
73         if (!(block = create_block(inode, pos/BLOCK_SIZE)))
74             break;
75         // 根据该逻辑块号读取设备上的相应数据块, 若出错则退出循环。
76         if (!(bh=bread(inode->i_dev, block)))
77             break;
78         // 求出文件读写指针在数据块中的偏移值 c, 将 p 指向读出数据块缓冲区中开始读取的位置。置该
79         // 缓冲区已修改标志。
80         c = pos % BLOCK_SIZE;
81         p = c + bh->b_data;
82         bh->b_dirt = 1;
83         // 从开始读写位置到块末共可写入 c=(BLOCK_SIZE-c)个字节。若 c 大于剩余还需写入的字节数
84         // (count-i), 则此次只需再写入 c=(count-i)即可。
85         c = BLOCK_SIZE-c;
86         if (c > count-i) c = count-i;
87         // 文件读写指针前移此次需写入的字节数。如果当前文件读写指针位置值超过了文件的大小, 则
88         // 修改 i 节点中文件大小字段, 并置 i 节点已修改标志。
89         pos += c;
90         if (pos > inode->i_size) {
91             inode->i_size = pos;
92             inode->i_dirt = 1;
93         }
94         // 已写入字节计数累加此次写入的字节数 c。从用户缓冲区 buf 中复制 c 个字节到高速缓冲区中 p
95         // 指向开始的位置处。然后释放该缓冲区。
96         i += c;
97         while (c-->0)
98             *(p++) = get_fs_byte(buf++);

```

---

```

82         brelse(bh);
83     }
    // 更改文件修改时间为当前时间。
84     inode->i_mtime = CURRENT_TIME;
    // 如果此次文件操作是在文件尾添加数据, 则仍将文件读写指针指向文件尾, 并更改 i 节点修改
    // 时间为当前时间。
85     if (!(filp->f_flags & O_APPEND)) {
86         filp->f_pos = pos;
87         inode->i_ctime = CURRENT_TIME;
88     }
    // 返回写入的字节数, 若写入字节数为 0, 则返回出错号-1。
89     return (i?i:-1);
90 }
91

```

---

### 8.7.3 其它信息

## 8.8 file\_table.c 文件

### 8.8.1 功能描述

该程序目前是空的, 仅定义了文件表数组。

### 8.8.2 代码注释

列表 linux/fs/file\_table.c 程序

---

```

1  /*
2   *  linux/fs/file_table.c
3   *
4   *  (C) 1991  Linus Torvalds
5   */
6
7  #include <linux/fs.h>    // 文件系统头文件。定义文件表结构 (file, buffer_head, m_inode 等)。
8
9  struct file file_table[NR_FILE]; // 文件表数组 (64 项)。
10

```

---

### 8.8.3 其它信息

## 8.9 fcntl.c 文件

### 8.9.1 功能描述

### 8.9.2 代码注释

列表 linux/fs/fcntl.c 程序

```

1  /*
2  *  linux/fs/fcntl.c
3  *
4  *  (C) 1991  Linus Torvalds
5  */
6
7  #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。
8  #include <errno.h>          // 错误号头文件。包含系统中各种出错号。(Linus 从 minix 中引进的)。
9  #include <linux/sched.h>    // 调度程序头文件，定义了任务结构 task_struct、初始任务 0 的数据，
                               // 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
10 #include <linux/kernel.h>   // 内核头文件。含有一些内核常用函数的原形定义。
11 #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。
12
13 #include <fcntl.h>          // 文件控制头文件。用于文件及其描述符的操作控制常数符号的定义。
14 #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。
15
16 extern int sys_close(int fd);    // 关闭文件系统调用。(fs/open.c, 192)
17
18     ///// 复制文件句柄(描述符)。
19     // 参数 fd 是欲复制的文件句柄，arg 指定新文件句柄的最小数。
20     // 返回新文件句柄或出错码。
21 static int dupfd(unsigned int fd, unsigned int arg)
22 {
23     // 如果文件句柄值大于一个程序最多打开文件数 NR_OPEN，或者该句柄的文件结构不存在，则出错，
24     // 返回出错码并退出。
25     if (fd >= NR_OPEN || !current->filp[fd])
26         return -EBADF;
27     // 如果指定的新句柄值 arg 大于最多打开文件数，则出错，返回出错码并退出。
28     if (arg >= NR_OPEN)
29         return -EINVAL;
30     // 在当前进程的文件结构指针数组中寻找索引号大于等于 arg 但还没有使用的项。
31     while (arg < NR_OPEN)
32         if (current->filp[arg])
33             arg++;
34     else
35         break;
36     // 如果找到的新句柄值 arg 大于最多打开文件数，则出错，返回出错码并退出。
37     if (arg >= NR_OPEN)
38         return -EMFILE;
39     // 在执行时关闭标志位图中复位该句柄位。也即在运行 exec() 类函数时不关闭该句柄。
40     current->close_on_exec &= ~(1<<arg);
41     // 令该文件结构指针等于原句柄 fd 的指针，并将文件引用计数增 1。
42     (current->filp[arg] = current->filp[fd])->f_count++;

```

```

33         return arg;                // 返回新的文件句柄。
34     }
35
36     ///// 复制文件句柄系统调用函数。
37     // 复制指定文件句柄 oldfd 为句柄值等于 newfd 的句柄。如果 newfd 已经打开，则首先关闭之。
38     int sys_dup2(unsigned int oldfd, unsigned int newfd)
39     {
40         sys_close(newfd);           // 若句柄 newfd 已经打开，则首先关闭之。
41         return dupfd(oldfd, newfd); // 复制并返回新句柄。
42     }
43
44     ///// 复制文件句柄系统调用函数。
45     // 复制指定文件句柄 oldfd 为句柄值最小的未用句柄。
46     int sys_dup(unsigned int fildes)
47     {
48         return dupfd(fildes, 0);
49     }
50
51     ///// 文件句柄(描述符)操作系统调用函数。
52     // 参数 fd 是文件句柄，cmd 是操作命令(参见 include/fcntl.h, 23-30 行)。
53     int sys_fcntl(unsigned int fd, unsigned int cmd, unsigned long arg)
54     {
55         struct file * filp;
56
57         // 如果文件句柄值大于一个程序最多打开文件数 NR_OPEN，或者该句柄的文件结构指针为空，则出错，
58         // 返回出错码并退出。
59         if (fd >= NR_OPEN || !(filp = current->filp[fd]))
60             return -EBADF;
61         // 根据不同命令 cmd 进行分别处理。
62         switch (cmd) {
63             case F_DUPFD:           // 复制文件句柄。
64                 return dupfd(fd, arg);
65             case F_GETFD:           // 取文件句柄的执行时关闭标志。
66                 return (current->close_on_exec >> fd) & 1;
67             case F_SETFD:           // 设置句柄执行时关闭标志。arg 位 0 置位是设置，否则关闭。
68                 if (arg & 1)
69                     current->close_on_exec |= (1 << fd);
70                 else
71                     current->close_on_exec &= ~(1 << fd);
72                 return 0;
73             case F_GETFL:           // 取文件状态标志和访问模式。
74                 return filp->f_flags;
75             case F_SETFL:           // 设置文件状态和访问模式(根据 arg 设置添加、非阻塞标志)。
76                 filp->f_flags &= ~(O_APPEND | O_NONBLOCK);
77                 filp->f_flags |= arg & (O_APPEND | O_NONBLOCK);
78                 return 0;
79             case F_GETLK:           case F_SETLK:           case F_SETLKW: // 未实现。
80                 return -1;
81             default:
82                 return -1;
83         }
84     }
85 }

```

### 8.9.3 其它信息

## 8.10 stat.c 文件

### 8.10.1 功能描述

该程序实现取文件状态信息系统调用函数 `stat()` 和 `fstat()`，并将信息存放在用户的文件状态结构缓冲区中。`stat()` 是利用文件名取信息，而 `fstat()` 是使用文件句柄(描述符)来取信息。

### 8.10.2 代码注释

列表 linux/fs/stat.c 程序

```

1  /*
2   *  linux/fs/stat.c
3   *
4   *  (C) 1991  Linus Torvalds
5   */
6
7  #include <errno.h>          // 错误号头文件。包含系统中各种出错号。(Linus 从 minix 中引进的)。
8  #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。
9
10 #include <linux/fs.h>       // 文件系统头文件。定义文件表结构 (file, buffer_head, m_inode 等)。
11 #include <linux/sched.h>    // 调度程序头文件，定义了任务结构 task_struct、初始任务 0 的数据，
12                             // 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
13 #include <linux/kernel.h>   // 内核头文件。含有一些内核常用函数的原形定义。
14 #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。
15
16     /// 复制文件状态信息。
17     // 参数 inode 是文件对应的 i 节点，statbuf 是 stat 文件状态结构指针，用于存放取得的状态信息。
18 static void cp_stat(struct m_inode * inode, struct stat * statbuf)
19 {
20     struct stat tmp;
21     int i;
22
23     // 首先验证(或分配)存放数据的内存空间。
24     verify_area(statbuf, sizeof (* statbuf));
25     // 然后临时复制相应节点上的信息。
26     tmp.st_dev = inode->i_dev;          // 文件所在的设备号。
27     tmp.st_ino = inode->i_num;          // 文件 i 节点号。
28     tmp.st_mode = inode->i_mode;        // 文件属性。
29     tmp.st_nlink = inode->i_nlinks;     // 文件的连接数。
30     tmp.st_uid = inode->i_uid;          // 文件的用户 id。
31     tmp.st_gid = inode->i_gid;          // 文件的组 id。
32     tmp.st_rdev = inode->i_zone[0];     // 设备号(如果文件是特殊的字符文件或块文件)。
33     tmp.st_size = inode->i_size;        // 文件大小(字节数)(如果文件是常规文件)。
34     tmp.st_atime = inode->i_atime;      // 最后访问时间。
35     tmp.st_mtime = inode->i_mtime;      // 最后修改时间。
36     tmp.st_ctime = inode->i_ctime;      // 最后节点修改时间。

```

```

// 最后将这些状态信息复制到用户缓冲区中。
32     for (i=0 ; i<sizeof (tmp) ; i++)
33         put_fs_byte(((char *) &tmp)[i],&((char *) statbuf)[i]);
34 }
35
//// 文件状态系统调用函数 - 根据文件名获取文件状态信息。
// 参数 filename 是指定的文件名, statbuf 是存放状态信息的缓冲区指针。
// 返回 0, 若出错则返回出错码。
36 int sys_stat(char * filename, struct stat * statbuf)
37 {
38     struct m_inode * inode;
39
// 首先根据文件名找出对应的 i 节点, 若出错则返回错误码。
40     if (!(inode=namei(filename)))
41         return -ENOENT;
// 将 i 节点上的文件状态信息复制到用户缓冲区中, 并释放该 i 节点。
42     cp_stat(inode, statbuf);
43     iput(inode);
44     return 0;
45 }
46
//// 文件状态系统调用 - 根据文件句柄获取文件状态信息。
// 参数 fd 是指定文件的句柄(描述符), statbuf 是存放状态信息的缓冲区指针。
// 返回 0, 若出错则返回出错码。
47 int sys_fstat(unsigned int fd, struct stat * statbuf)
48 {
49     struct file * f;
50     struct m_inode * inode;
51
// 如果文件句柄值大于一个程序最多打开文件数 NR_OPEN, 或者该句柄的文件结构指针为空, 或者
// 对应文件结构的 i 节点字段为空, 则出错, 返回出错码并退出。
52     if (fd >= NR_OPEN || !(f=current->filp[fd]) || !(inode=f->f_inode))
53         return -EBADF;
// 将 i 节点上的文件状态信息复制到用户缓冲区中。
54     cp_stat(inode, statbuf);
55     return 0;
56 }
57

```

### 8.10.3 其它信息

## 8.11 pipe.c 文件

### 8.11.1 功能描述

本程序执行管道文件的读写操作, 并实现了管道系统调用 pipe()。

在初始化管道时, 管道 i 节点的 i\_size 字段中被设置为指向管道缓冲区的指针, 管道数据头部指针存放在 i\_zone[0] 字段中, 而管道数据尾部指针存放在 i\_zone[1] 字段中。对于读管道操作, 数据是从管道尾读出, 并使管道尾指针前移读取字节数个位置; 对于往管道中的写入操作, 数据是向管道头部写入,

并使管道头指针前移写入字节数个位置。参见下面的管道示意图。

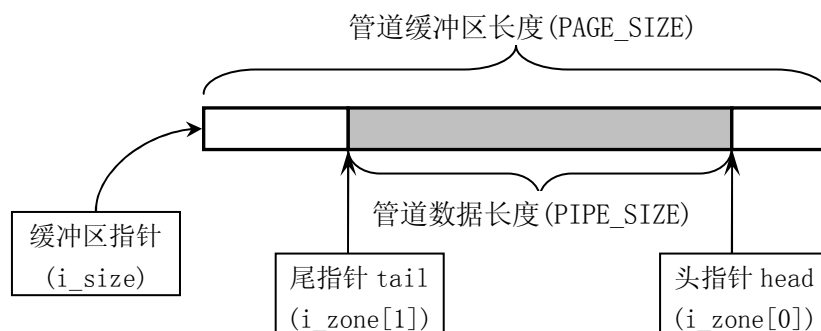


图 5.1 管道缓冲区操作示意图

### 8.11.2 代码注释

列表 linux/fs/pipe.c 程序

```

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，信号结构以及信号操作函数原型。
8
9 #include <linux/sched.h>     // 调度程序头文件，定义了任务结构 task_struct、初始任务 0 的数据，
10                             // 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
11 #include <linux/mm.h>        /* for get_free_page */ /* 使用其中的 get_free_page */
12                             // 内存管理头文件。含有页面大小定义和一些页面释放函数原型。
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 // 管道读操作函数。
16 // 参数 inode 是管道对应的 i 节点，buf 是数据缓冲区指针，count 是读取的字节数。
17 int read_pipe(struct m_inode * inode, char * buf, int count)
18 {
19     int chars, size, read = 0;
20
21     // 若欲读取的字节计数值 count 大于 0，则循环执行以下操作。
22     while (count > 0) {
23         // 若当前管道中没有数据(size=0)，则唤醒等待该节点的进程，如果已没有写管道者，则返回已读
24         // 字节数，退出。否则在该 i 节点上睡眠，等待信息。
25         while (!(size=PIPE_SIZE(*inode))) {
26             wake_up(&inode->i_wait);
27             if (inode->i_count != 2) /* are there any writers? */
28                 return read;
29             sleep_on(&inode->i_wait);
30         }
31         // 取管道尾到缓冲区末端的字节数 chars。如果其大于还需要读取的字节数 count，则令其等于 count。
32         // 如果 chars 大于当前管道中含有数据的长度 size，则令其等于 size。

```



```

24         chars = PAGE_SIZE-PIPE_TAIL(*inode);
25         if (chars > count)
26             chars = count;
27         if (chars > size)
28             chars = size;
// 读字节计数减去此次可读的字节数 chars，并累加已读字节数。
29         count -= chars;
30         read += chars;
// 令 size 指向管道尾部，调整当前管道尾指针（前移 chars 字节）。
31         size = PIPE_TAIL(*inode);
32         PIPE_TAIL(*inode) += chars;
33         PIPE_TAIL(*inode) &= (PAGE_SIZE-1);
// 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 字段中是管道缓冲块指针。
34         while (chars-->0)
35             put_fs_byte(((char *)inode->i_size)[size++], buf++);
36     }
// 唤醒等待该 i 节点的进程，并返回读取的字节数。
37     wake_up(&inode->i_wait);
38     return read;
39 }
40
//// 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。
41 int write_pipe(struct m_inode * inode, char * buf, int count)
42 {
43     int chars, size, written = 0;
44
// 若将写入的字节计数值 count 还大于 0，则循环执行以下操作。
45     while (count>0) {
// 若当前管道中没有已经满了(size=0)，则唤醒等待该节点的进程，如果已没有读管道者，则返回
// 已写入的字节数并退出。若写入 0 字节，则返回-1。否则在该 i 节点上睡眠，等待管道腾出空间。
46         while (!(size=(PAGE_SIZE-1)-PIPE_SIZE(*inode))) {
47             wake_up(&inode->i_wait);
48             if (inode->i_count != 2) { /* no readers */
49                 current->signal |= (1<<(SIGPIPE-1));
50                 return written?written:-1;
51             }
52             sleep_on(&inode->i_wait);
53         }
// 取管道头部到缓冲区末端空间字节数 chars。如果其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。
54         chars = PAGE_SIZE-PIPE_HEAD(*inode);
55         if (chars > count)
56             chars = count;
57         if (chars > size)
58             chars = size;
// 写字节计数减去此次可写入的字节数 chars，并累加已写字节数到 written。
59         count -= chars;
60         written += chars;
// 令 size 指向管道数据头部，调整当前管道数据头部指针（前移 chars 字节）。
61         size = PIPE_HEAD(*inode);
62         PIPE_HEAD(*inode) += chars;
63         PIPE_HEAD(*inode) &= (PAGE_SIZE-1);

```

```

// 从用户缓冲区复制 chars 个字节到管道中。对于管道 i 节点，其 i_size 字段中是管道缓冲块指针。
64         while (chars-->0)
65             ((char *)inode->i_size)[size++]=get\_fs\_byte(buf++);
66     }
// 唤醒等待该 i 节点的进程，返回已写入的字节数，退出。
67     wake\_up(&inode->i_wait);
68     return written;
69 }
70
//// 创建管道系统调用函数。
// 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。fildes[0]
// 用于读管道中数据，fildes[1]用于向管道中写入数据。
// 成功时返回 0，出错时返回-1。
71 int sys\_pipe(unsigned long * fildes)
72 {
73     struct m\_inode * inode;
74     struct file * f[2];
75     int fd[2];
76     int i, j;
77
// 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。
78     j=0;
79     for(i=0; j<2 && i<NR\_FILE; i++)
80         if (!file\_table[i].f_count)
81             (f[j++] = i + file\_table)->f_count++;
// 如果只有一个空闲项，则释放该项(引用计数复位)。
82     if (j==1)
83         f[0]->f_count=0;
// 如果没有找到两个空闲项，则返回-1。
84     if (j<2)
85         return -1;
// 针对上面取得的两个文件结构，分别分配一文件句柄，并使进程的文件结构指针分别指向这两个
// 文件结构。
86     j=0;
87     for(i=0; j<2 && i<NR\_OPEN; i++)
88         if (!current->filp[i]) {
89             current->filp[ fd[j]=i ] = f[j];
90             j++;
91         }
// 如果只有一个空闲文件句柄，则释放该句柄。
92     if (j==1)
93         current->filp[fd[0]]=NULL;
// 如果没有找到两个空闲句柄，则释放上面获取的两个文件结构项(复位引用计数值)，并返回-1。
94     if (j<2) {
95         f[0]->f_count=f[1]->f_count=0;
96         return -1;
97     }
// 申请管道节点。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。
98     if (! (inode=get\_pipe\_inode()) ) {
99         current->filp[fd[0]] =
100             current->filp[fd[1]] = NULL;
101         f[0]->f_count = f[1]->f_count = 0;
102         return -1;

```

```

103     }
    // 初始化两个文件结构，都指向同一个 i 节点，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为写。
104     f[0]->f_inode = f[1]->f_inode = inode;
105     f[0]->f_pos = f[1]->f_pos = 0;
106     f[0]->f_mode = 1;                /* read */
107     f[1]->f_mode = 2;                /* write */
    // 将文件句柄数组复制到对应的用户数组中，并返回 0，退出。
108     put_fs_long(fd[0], 0+fildes);
109     put_fs_long(fd[1], 1+fildes);
110     return 0;
111 }
112

```

### 8.11.3 其它信息

## 8.12 read\_write.c 文件

### 8.12.1 功能描述

### 8.12.2 代码注释

列表 linux/fs/read\_write.c 程序

```

1  /*
2  *  linux/fs/read_write.c
3  *
4  *  (C) 1991  Linus Torvalds
5  */
6
7  #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。
8  #include <errno.h>      // 错误号头文件。包含系统中各种出错号。(Linus 从 minix 中引进的)。
9  #include <sys/types.h>  // 类型头文件。定义了基本的系统数据类型。
10
11 #include <linux/kernel.h> // 内核头文件。含有一些内核常用函数的原形定义。
12 #include <linux/sched.h> // 调度程序头文件，定义了任务结构 task_struct、初始任务 0 的数据，
    // 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
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
    // 字符设备读写函数。
15 extern int rw_char(int rw, int dev, char * buf, int count, off_t * pos);
    // 读管道操作函数。
16 extern int read_pipe(struct m_inode * inode, char * buf, int count);
    // 写管道操作函数。
17 extern int write_pipe(struct m_inode * inode, char * buf, int count);
    // 块设备读操作函数。
18 extern int block_read(int dev, off_t * pos, char * buf, int count);
    // 块设备写操作函数。
19 extern int block_write(int dev, off_t * pos, char * buf, int count);

```

```

// 读文件操作函数。
20 extern int file_read(struct m_inode * inode, struct file * filp,
21                     char * buf, int count);
// 写文件操作函数。
22 extern int file_write(struct m_inode * inode, struct file * filp,
23                      char * buf, int count);
24
//// 重定位文件读写指针系统调用函数。
// 参数 fd 是文件句柄, offset 是新的文件读写指针偏移值, origin 是偏移的起始位置, 是 SEEK_SET
// (0, 从文件开始处)、SEEK_CUR(1, 从当前读写位置)、SEEK_END(2, 从文件尾处)三者之一。
25 int sys_lseek(unsigned int fd, off_t offset, int origin)
26 {
27     struct file * file;
28     int tmp;
29
// 如果文件句柄值大于程序最多打开文件数 NR_OPEN, 或者该句柄的文件结构指针为空, 或者
// 对应文件结构的 i 节点字段为空, 或者指定设备文件指针不可定位, 则返回出错码并退出。
30     if (fd >= NR_OPEN || !(file=current->filp[fd]) || !(file->f_inode)
31         || !IS_SEEKABLE(MAJOR(file->f_inode->i_dev)))
32         return -EBADF;
// 如果文件对应的 i 节点是管道节点, 则返回出错码, 退出。管道头尾指针不可随意移动!
33     if (file->f_inode->i_pipe)
34         return -ESPIPE;
// 根据设置的起始位置, 分别重新定位文件读写指针。
35     switch (origin) {
// origin = SEEK_SET, 要求以文件起始处作为原点设置文件读写指针。若偏移值小于零, 则出错返
// 回错误码。否则设置文件读写指针等于 offset。
36         case 0:
37             if (offset < 0) return -EINVAL;
38             file->f_pos = offset;
39             break;
// origin = SEEK_CUR, 要求以文件当前读写指针处作为原点重定位读写指针。如果文件当前指针加
// 上偏移值小于 0, 则返回出错码退出。否则在当前读写指针上加上偏移值。
40         case 1:
41             if (file->f_pos + offset < 0) return -EINVAL;
42             file->f_pos += offset;
43             break;
// origin = SEEK_END, 要求以文件末尾作为原点重定位读写指针。此时若文件大小加上偏移值小于零
// 则返回出错码退出。否则重定位读写指针为文件长度加上偏移值。
44         case 2:
45             if ((tmp=file->f_inode->i_size+offset) < 0)
46                 return -EINVAL;
47             file->f_pos = tmp;
48             break;
// origin 设置出错, 返回出错码退出。
49         default:
50             return -EINVAL;
51     }
52     return file->f_pos;    // 返回重定位后的文件读写指针值。
53 }
54
//// 读文件系统调用函数。
// 参数 fd 是文件句柄, buf 是缓冲区, count 是欲读字节数。

```

```

55 int sys_read(unsigned int fd, char * buf, int count)
56 {
57     struct file * file;
58     struct m_inode * inode;
59
60     // 如果文件句柄值大于程序最多打开文件数 NR_OPEN, 或者需要读取的字节计数值小于 0, 或者该句柄
61     // 的文件结构指针为空, 则返回出错码并退出。
62     if (fd >= NR_OPEN || count < 0 || !(file = current->filp[fd]))
63         return -EINVAL;
64     // 若需读取的字节数 count 等于 0, 则返回 0, 退出
65     if (!count)
66         return 0;
67     // 验证存放数据的缓冲区内内存限制。
68     verify_area(buf, count);
69     // 取文件对应的 i 节点。若是管道文件, 并且是读管道文件模式, 则进行读管道操作, 若成功则返回
70     // 读取的字节数, 否则返回出错码, 退出。
71     inode = file->f_inode;
72     if (inode->i_pipe)
73         return (file->f_mode & 1) ? read_pipe(inode, buf, count) : -EIO;
74     // 如果是字符型文件, 则进行读字符设备操作, 返回读取的字符数。
75     if (S_ISCHR(inode->i_mode))
76         return rw_char(READ, inode->i_zone[0], buf, count, &file->f_pos);
77     // 如果是块设备文件, 则执行块设备读操作, 并返回读取的字节数。
78     if (S_ISBLK(inode->i_mode))
79         return block_read(inode->i_zone[0], &file->f_pos, buf, count);
80     // 如果是目录文件或者是常规文件, 则首先验证读取数 count 的有效性并进行调整 (若读取字节数加上
81     // 文件当前读写指针值大于文件大小, 则重新设置读取字节数为文件长度-当前读写指针值, 若读取数
82     // 等于 0, 则返回 0 退出), 然后执行文件读操作, 返回读取的字节数并退出。
83     if (S_ISDIR(inode->i_mode) || S_ISREG(inode->i_mode)) {
84         if (count + file->f_pos > inode->i_size)
85             count = inode->i_size - file->f_pos;
86         if (count <= 0)
87             return 0;
88         return file_read(inode, file, buf, count);
89     }
90     // 否则打印节点文件属性, 并返回出错码退出。
91     printk(" (Read) inode->i_mode=%06o\n|r", inode->i_mode);
92     return -EINVAL;
93 }
94
95 int sys_write(unsigned int fd, char * buf, int count)
96 {
97     struct file * file;
98     struct m_inode * inode;
99
100    // 如果文件句柄值大于程序最多打开文件数 NR_OPEN, 或者需要写入的字节计数小于 0, 或者该句柄
101    // 的文件结构指针为空, 则返回出错码并退出。
102    if (fd >= NR_OPEN || count < 0 || !(file = current->filp[fd]))
103        return -EINVAL;
104    // 若需读取的字节数 count 等于 0, 则返回 0, 退出
105    if (!count)
106        return 0;
107    // 取文件对应的 i 节点。若是管道文件, 并且是写管道文件模式, 则进行写管道操作, 若成功则返回

```

```

// 写入的字节数，否则返回出错码，退出。
92     inode=file->f_inode;
93     if (inode->i_pipe)
94         return (file->f_mode&2)?write_pipe(inode, buf, count):-EIO;
// 如果是字符型文件，则进行写字符设备操作，返回写入的字符数，退出。
95     if (S_ISCHR(inode->i_mode))
96         return rw_char(WRITE, inode->i_zone[0], buf, count, &file->f_pos);
// 如果是块设备文件，则进行块设备写操作，并返回写入的字节数，退出。
97     if (S_ISBLK(inode->i_mode))
98         return block_write(inode->i_zone[0], &file->f_pos, buf, count);
// 若是常规文件，则执行文件写操作，并返回写入的字节数，退出。
99     if (S_ISREG(inode->i_mode))
100        return file_write(inode, file, buf, count);
// 否则，显示对应节点的文件模式，返回出错码，退出。
101    printk("(Write)inode->i_mode=%06o\n\r", inode->i_mode);
102    return -EINVAL;
103 }
104

```

### 8.12.3 其它信息

## 8.13 truncate.c 文件

### 8.13.1 功能描述

本程序用于释放指定 i 节点在设备上占用的所有逻辑块，包括直接块、一次间接块和二次间接块。

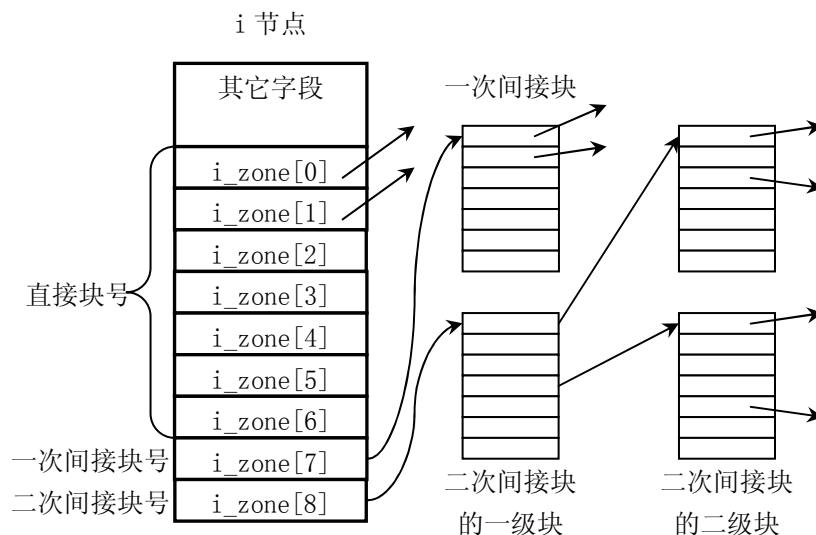


图 索引节点(i 节点)的逻辑块连接方式

### 8.13.2 代码注释

列表 linux/fs/truncate.c 程序

```

2  * linux/fs/truncate.c
3  *
4  * (C) 1991 Linus Torvalds
5  */
6
7 #include <linux/sched.h> // 调度程序头文件，定义了任务结构 task_struct、初始任务 0 的数据，
                        // 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
8
9 #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。
10
11 // 释放一次间接块。
12 static void free_ind(int dev, int block)
13 {
14     struct buffer_head * bh;
15     unsigned short * p;
16     int i;
17
18     // 如果逻辑块号为 0，则返回。
19     if (!block)
20         return;
21
22     // 读取一次间接块，并释放其上表明使用的所有逻辑块，然后释放该一次间接块的缓冲区。
23     if (bh=bread(dev, block)) {
24         p = (unsigned short *) bh->b_data; // 指向数据缓冲区。
25         for (i=0; i<512; i++, p++)        // 每个逻辑块上可有 512 个块号。
26             if (*p)
27                 free_block(dev, *p); // 释放指定的逻辑块。
28         brelse(bh); // 释放缓冲区。
29     }
30
31     // 释放设备上的一次间接块。
32     free_block(dev, block);
33 }
34
35 // 释放二次间接块。
36 static void free_dind(int dev, int block)
37 {
38     struct buffer_head * bh;
39     unsigned short * p;
40     int i;
41
42     // 如果逻辑块号为 0，则返回。
43     if (!block)
44         return;
45
46     // 读取二次间接块的一级块，并释放其上表明使用的所有逻辑块，然后释放该一级块的缓冲区。
47     if (bh=bread(dev, block)) {
48         p = (unsigned short *) bh->b_data; // 指向数据缓冲区。
49         for (i=0; i<512; i++, p++)        // 每个逻辑块上可连接 512 个二级块。
50             if (*p)
51                 free_ind(dev, *p); // 释放所有一次间接块。
52         brelse(bh); // 释放缓冲区。
53     }
54
55     // 最后释放设备上的二次间接块。
56     free_block(dev, block);
57 }

```

```

46  // 将节点对应的文件长度截为 0，并释放占用的设备空间。
47 void truncate(struct m_inode * inode)
48 {
49     int i;
50
51     // 如果不是常规文件或者是目录文件，则返回。
52     if (!(S_ISREG(inode->i_mode) || S_ISDIR(inode->i_mode)))
53         return;
54     // 释放 i 节点的 7 个直接逻辑块，并将这 7 个逻辑块项全置零。
55     for (i=0; i<7; i++)
56         if (inode->i_zone[i]) { // 如果块号不为 0，则释放之。
57             free_block(inode->i_dev, inode->i_zone[i]);
58             inode->i_zone[i]=0;
59         }
60     free_ind(inode->i_dev, inode->i_zone[7]); // 释放一次间接块。
61     free_dind(inode->i_dev, inode->i_zone[8]); // 释放二次间接块。
62     inode->i_zone[7] = inode->i_zone[8] = 0; // 逻辑块项 7、8 置零。
63     inode->i_size = 0; // 文件大小置零。
64     inode->i_dirt = 1; // 置节点已修改标志。
65     inode->i_mtime = inode->i_ctime = CURRENT_TIME; // 重置文件和节点修改时间为当前时间。
66 }

```

### 8.13.3 其它信息

## 8.14 super.c 文件

### 8.14.1 功能描述

### 8.14.2 代码注释

列表 linux/fs/super.c 程序

```

1  /*
2   * linux/fs/super.c
3   *
4   * (C) 1991 Linus Torvalds
5   */
6
7  /*
8   * super.c contains code to handle the super-block tables.
9   */
10 #include <linux/config.h> // 内核配置头文件。定义键盘语言和硬盘类型 (HD_TYPE) 可选项。
11 #include <linux/sched.h> // 调度程序头文件，定义了任务结构 task_struct、初始任务 0 的数据，
12 // 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
13 #include <linux/kernel.h> // 内核头文件。含有一些内核常用函数的原形定义。
14 #include <asm/system.h> // 系统头文件。定义了设置或修改描述符/中断门等的嵌入式汇编宏。

```



```

14
15 #include <errno.h>          // 错误号头文件。包含系统中各种出错号。(Linus 从 minix 中引进的)。
16 #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。
17
18 int sync_dev(int dev);      // 对指定设备执行高速缓冲与设备上数据的同步操作。(fs/buffer.c, 59)
19 void wait_for_keypress(void); // 等待击键。(kernel/chr_drv/tty_io.c, 140)
20
21 /* set_bit uses setb, as gas doesn't recognize setc */
   /* set_bit() 使用了 setb 指令，因为汇编编译器 gas 不能识别指令 setc */
   //// 测试指定偏移处比特位的值(0 或 1)，并返回该比特位值。(应该取名为 test_bit() 更妥帖)
   // 嵌入式汇编宏。参数 bitnr 是比特位偏移值，addr 是测试比特位操作的起始地址。
   // %0 - ax(__res), %1 - 0, %2 - bitnr, %3 - addr
22 #define set_bit(bitnr, addr) ({ \
23     register int __res __asm__("ax"); \
24     __asm__("bt %2, %3; setb %%al": "=a" (__res): "a" (0), "r" (bitnr), "m" (*(addr))); \
25     __res; })
26
27 struct super_block super_block[NR_SUPER]; // 超级块结构数组(共 8 项)。
28 /* this is initialized in init/main.c */
   /* ROOT_DEV 已在 init/main.c 中被初始化 */
29 int ROOT_DEV = 0;
30
   //// 锁定指定的超级块。
31 static void lock_super(struct super_block * sb)
32 {
33     cli(); // 关中断。
34     while (sb->s_lock) // 如果该超级块已经上锁，则睡眠等待。
35         sleep_on(&(sb->s_wait));
36     sb->s_lock = 1; // 给该超级块加锁(置锁定标志)。
37     sti(); // 开中断。
38 }
39
   //// 对指定超级块解锁。(如果使用 unlock_super 这个名称则更妥帖)。
40 static void free_super(struct super_block * sb)
41 {
42     cli(); // 关中断。
43     sb->s_lock = 0; // 复位锁定标志。
44     wake_up(&(sb->s_wait)); // 唤醒等待该超级块的进程。
45     sti(); // 开中断。
46 }
47
   //// 睡眠等待超级块解锁。
48 static void wait_on_super(struct super_block * sb)
49 {
50     cli(); // 关中断。
51     while (sb->s_lock) // 如果超级块已经上锁，则睡眠等待。
52         sleep_on(&(sb->s_wait));
53     sti(); // 开中断。
54 }
55
   //// 取指定设备的超级块。返回该超级块结构指针。
56 struct super_block * get_super(int dev)
57 {

```

```

58     struct super\_block * s;
59
60     // 如果没有指定设备，则返回空指针。
61     if (!dev)
62         return NULL;
63     // s 指向超级块数组开始处。搜索整个超级块数组，寻找指定设备的超级块。
64     s = 0+super\_block;
65     while (s < NR\_SUPER+super\_block)
66     // 如果当前搜索项是指定设备的超级块，则首先等待该超级块解锁（若已经被其它进程上锁的话）。
67     // 在等待期间，该超级块有可能被其它设备使用，因此此时需再判断一次是否是指定设备的超级块，
68     // 如果是则返回该超级块的指针。否则就重新对超级块数组再搜索一遍，因此 s 重又指向超级块数组
69     // 开始处。
70     if (s->s_dev == dev) {
71         wait\_on\_super(s);
72         if (s->s_dev == dev)
73             return s;
74         s = 0+super\_block;
75     // 如果当前搜索项不是，则检查下一项。如果没有找到指定的超级块，则返回空指针。
76     } else
77         s++;
78     return NULL;
79 }
80
81 // 释放指定设备的超级块。
82 // 释放设备所使用的超级块数组项（置 s_dev=0），并释放该设备 i 节点位图和逻辑块位图所占用
83 // 的高速缓冲块。如果超级块对应的文件系统是根文件系统，或者其 i 节点上已经安装有其它的文件
84 // 系统，则不能释放该超级块。
85 void put\_super(int dev)
86 {
87     struct super\_block * sb;
88     struct m\_inode * inode;
89     int i;
90
91     // 如果指定设备是根文件系统设备，则显示警告信息“根系统盘改变了，准备生死决战吧”，并返回。
92     if (dev == ROOT\_DEV) {
93         printk("root diskette changed: prepare for armageddon\n|r");
94         return;
95     }
96     // 如果找不到指定设备的超级块，则返回。
97     if (!(sb = get\_super(dev)))
98         return;
99     // 如果该超级块指明本文件系统 i 节点上安装有其它的文件系统，则显示警告信息，返回。
100    if (sb->s_imount) {
101        printk("Mounted disk changed - tssk, tssk\n|r");
102        return;
103    }
104    // 找到指定设备的超级块后，首先锁定该超级块，然后置该超级块对应的设备号字段为 0，也即将
105    // 放弃该超级块。
106    lock\_super(sb);
107    sb->s_dev = 0;
108    // 然后释放该设备 i 节点位图和逻辑块位图所占用的缓冲块。
109    for(i=0;i<I\_MAP\_SLOTS;i++)
110        brelse(sb->s_imap[i]);

```

```

94         for(i=0;i<Z_MAP_SLOTS;i++)
95             brelse(sb->s_zmap[i]);
// 最后对该超级块解锁，并返回。
96         free_super(sb);
97         return;
98     }
99
// 从设备上读取超级块到内存中。
// 如果该设备的超级块已经在高速缓冲中并且有效，则直接返回该超级块的指针。
100 static struct super_block * read_super(int dev)
101 {
102     struct super_block * s;
103     struct buffer_head * bh;
104     int i,block;
105
// 如果没有指明设备，则返回空指针。
106     if (!dev)
107         return NULL;
// 首先检查该设备是否可更换过盘片（也即是否是软盘设备），如果更换过盘，则高速缓冲区有关该
// 设备的所有缓冲块均失效，需要进行失效处理。
108     check_disk_change(dev);
// 如果该设备的超级块已经在高速缓冲中，则直接返回该超级块的指针。
109     if (s = get_super(dev))
110         return s;
// 否则，首先在超级块数组中找出一个空项（也即其 s_dev=0 的项）。如果数组已经占满则返回空指针。
111     for (s = 0+super_block ;; s++) {
112         if (s >= NR_SUPER+super_block)
113             return NULL;
114         if (!s->s_dev)
115             break;
116     }
// 找到超级块空项后，就将该超级块用于指定设备，对该超级块进行部分初始化。
117     s->s_dev = dev;
118     s->s_isup = NULL;
119     s->s_imount = NULL;
120     s->s_time = 0;
121     s->s_rd_only = 0;
122     s->s_dirt = 0;
// 然后所定该超级块，并从设备上读取超级块信息到 bh 指向的缓冲区中。如果读超级块操作失败，
// 则释放上面选定的超级块数组中的项，并解锁该项，返回空指针退出。
123     lock_super(s);
124     if (!(bh = bread(dev,1))) {
125         s->s_dev=0;
126         free_super(s);
127         return NULL;
128     }
// 将设备上读取的超级块信息复制到内存超级块结构中。并释放存放读取信息的高速缓冲块。
129     *((struct d_super_block *) s) =
130         *((struct d_super_block *) bh->b_data);
131     brelse(bh);
// 如果读取的超级块的文件系统魔数字段内容不对，说明设备上不是正确的文件系统，因此同上面
// 一样，释放上面选定的超级块数组中的项，并解锁该项，返回空指针退出。
// 对于该版 linux 内核，只支持 minix 文件系统版本 1，其魔数是 0x137f。

```

```

132     if (s->s_magic != SUPER\_MAGIC) {
133         s->s_dev = 0;
134         free\_super(s);
135         return NULL;
136     }
137     // 下面开始读取设备上 i 节点位图和逻辑块位图数据。首先初始化内存超级块结构中位图空间。
138     for (i=0;i<I\_MAP\_SLOTS;i++)
139         s->s_imap[i] = NULL;
140     for (i=0;i<Z\_MAP\_SLOTS;i++)
141         s->s_zmap[i] = NULL;
142     // 然后从设备上读取 i 节点位图和逻辑块位图信息，并存放在超级块对应字段中。
143     block=2;
144     for (i=0 ; i < s->s_imap_blocks ; i++)
145         if (s->s_imap[i]=bread(dev,block))
146             block++;
147     else
148         break;
149     for (i=0 ; i < s->s_zmap_blocks ; i++)
150         if (s->s_zmap[i]=bread(dev,block))
151             block++;
152     else
153         break;
154     // 如果读出的位图逻辑块数不等于位图应该占有的逻辑块数，说明文件系统位图信息有问题，超级块
155     // 初始化失败。因此只能释放前面申请的所有资源，返回空指针并退出。
156     if (block != 2+s->s_imap_blocks+s->s_zmap_blocks) {
157         // 释放 i 节点位图和逻辑块位图占用的高速缓冲区。
158         for(i=0;i<I\_MAP\_SLOTS;i++)
159             brelse(s->s_imap[i]);
160         for(i=0;i<Z\_MAP\_SLOTS;i++)
161             brelse(s->s_zmap[i]);
162         //释放上面选定的超级块数组中的项，并解锁该超级块项，返回空指针退出。
163         s->s_dev=0;
164         free\_super(s);
165         return NULL;
166     }
167     // 否则一切成功。对于申请空闲 i 节点的函数来讲，如果设备上所有的 i 节点已经全被使用，则查找
168     // 函数会返回 0 值。因此 0 号 i 节点是不能用的，所以这里将位图中的最低位设置为 1，以防止文件
169     // 系统分配 0 号 i 节点。同样的道理，也将逻辑块位图的最低位设置为 1。
170     s->s_imap[0]->b_data[0] |= 1;
171     s->s_zmap[0]->b_data[0] |= 1;
172     // 解锁该超级块，并返回超级块指针。
173     free\_super(s);
174     return s;
175 }
176
177 // 卸载文件系统调用函数。
178 // 参数 dev_name 是设备文件名。
179 int sys\_umount(char * dev_name)
180 {
181     struct m\_inode * inode;
182     struct super\_block * sb;
183     int dev;

```

```

// 首先根据设备文件名找到对应的 i 节点，并取其中的设备号。
173     if (!(inode=namei(dev_name)))
174         return -ENOENT;
175     dev = inode->i_zone[0];
// 如果不是块设备文件，则释放刚申请的 i 节点 dev_i，返回出错码。
176     if (!S_ISBLK(inode->i_mode)) {
177         iput(inode);
178         return -ENOTBLK;
179     }
// 释放刚申请的 i 节点 dev_i。
180     iput(inode);
// 如果设备是根文件系统，则不能被卸载，返回出错号。
181     if (dev==ROOT_DEV)
182         return -EBUSY;
// 如果取设备的超级块失败，或者该设备文件系统没有安装过，则返回出错码。
183     if (!(sb=get_super(dev)) || !(sb->s_imount))
184         return -ENOENT;
// 如果超级块所指明的被安装到的 i 节点没有置位其安装标志，则显示警告信息。
185     if (!sb->s_imount->i_mount)
186         printk("Mounted inode has i_mount=0\n");
// 查找 i 节点表，看是否有进程在使用该设备上的文件，如果有则返回忙出错码。
187     for (inode=inode_table+0 ; inode<inode_table+NR_INODE ; inode++)
188         if (inode->i_dev==dev && inode->i_count)
189             return -EBUSY;
// 复位被安装到的 i 节点的安装标志，释放该 i 节点。
190     sb->s_imount->i_mount=0;
191     iput(sb->s_imount);
// 置超级块中被安装 i 节点字段为空，并释放设备文件系统的根 i 节点，置超级块中被安装系统
// 根 i 节点指针为空。
192     sb->s_imount = NULL;
193     iput(sb->s_isup);
194     sb->s_isup = NULL;
// 释放该设备的超级块，并对该设备执行高速缓冲与设备上数据的同步操作。
195     put_super(dev);
196     sync_dev(dev);
197     return 0;
198 }
199
//// 安装文件系统调用函数。
// 参数 dev_name 是设备文件名，dir_name 是安装到的目录名，rw_flag 被安装文件的读写标志。
200 int sys_mount(char * dev_name, char * dir_name, int rw_flag)
201 {
202     struct m_inode * dev_i, * dir_i;
203     struct super_block * sb;
204     int dev;
205
// 首先根据设备文件名找到对应的 i 节点，并取其中的设备号。
206     if (!(dev_i=namei(dev_name)))
207         return -ENOENT;
208     dev = dev_i->i_zone[0];
// 如果不是块设备文件，则释放刚申请的 i 节点 dev_i，返回出错码。
209     if (!S_ISBLK(dev_i->i_mode)) {
210         iput(dev_i);

```

```

211         return -EPERM;
212     }
    // 释放刚申请的 i 节点 dev_i。
213     iput(dev_i);
    // 根据给定的目录文件名找到对应的 i 节点 dir_i。
214     if (!(dir_i=namei(dir_name)))
215         return -ENOENT;
    // 如果该 i 节点的引用计数不为 1（仅在这里引用），或者该 i 节点的节点号是根文件系统的节点
    // 号 1，则释放该 i 节点，返回出错码。
216     if (dir_i->i_count != 1 || dir_i->i_num == ROOT_INO) {
217         iput(dir_i);
218         return -EBUSY;
219     }
    // 如果该节点不是一个目录文件节点，则也释放该 i 节点，返回出错码。
220     if (!S_ISDIR(dir_i->i_mode)) {
221         iput(dir_i);
222         return -EPERM;
223     }
    // 读取将安装文件系统的超级块，如果失败则也释放该 i 节点，返回出错码。
224     if (!(sb=read_super(dev))) {
225         iput(dir_i);
226         return -EBUSY;
227     }
    // 如果将要被安装的文件系统已经安装在其它地方，则释放该 i 节点，返回出错码。
228     if (sb->s_imount) {
229         iput(dir_i);
230         return -EBUSY;
231     }
    // 如果将要安装到的 i 节点已经安装了文件系统(安装标志已经置位)，则释放该 i 节点，返回出错码。
232     if (dir_i->i_mount) {
233         iput(dir_i);
234         return -EPERM;
235     }
    // 被安装文件系统超级块的被安装到 i 节点字段指向安装到的目录的 i 节点。
236     sb->s_imount=dir_i;
    // 置安装 i 节点的安装标志和节点已修改标志。/* 注意！这里没有 iput(dir_i) */
237     dir_i->i_mount=1;                /* 这将在 umount 内操作 */
238     dir_i->i_dirt=1;                /* NOTE! we don't iput(dir_i) */
239     return 0;                        /* we do that in umount */
240 }
241
///// 安装根文件系统。
// 该函数是在系统开机初始化设置时(sys_setup())调用的。( kernel/blk_drv/hd.c, 157 )
242 void mount_root(void)
243 {
244     int i, free;
245     struct super_block * p;
246     struct m_inode * mi;
247
    // 如果磁盘 i 节点结构不是 32 个字节，则出错，死机。该判断是用于防止修改源代码时的一致性。
248     if (32 != sizeof (struct d_inode))
249         panic("bad i-node size");
    // 初始化文件表数组（共 64 项，也即系统同时只能打开 64 个文件），将所有文件结构中的引用计数

```

```

// 设置为 0。[??为什么放在这里初始化?]
250     for(i=0;i<NR_FILE;i++)
251         file_table[i].f_count=0;
// 如果根文件系统所在设备是软盘的话,就提示“插入根文件系统盘,并按回车键”,并等待按键。
252     if (MAJOR(ROOT_DEV) == 2) {
253         printk("Insert root floppy and press ENTER");
254         wait_for_keypress();
255     }
// 初始化超级块数组(共8项)。
256     for(p = &super_block[0] ; p < &super_block[NR_SUPER] ; p++) {
257         p->s_dev = 0;
258         p->s_lock = 0;
259         p->s_wait = NULL;
260     }
// 如果读设备上超级块失败,则显示信息,并死机。
261     if (!(p=read_super(ROOT_DEV)))
262         panic("Unable to mount root");
// 从设备上读取文件系统的根 i 节点(1),如果失败则显示出错信息,死机。
263     if (!(mi=iget(ROOT_DEV,ROOT_INO)))
264         panic("Unable to read root i-node");
// 该 i 节点引用次数递增 3 次。[?? Why?]
265     mi->i_count += 3;    /* NOTE! it is logically used 4 times, not 1 */
// 置该超级块的被安装文件系统 i 节点和被安装到的 i 节点为该 i 节点。
266     p->s_isup = p->s_imount = mi;
// 设置当前进程的当前工作目录和根目录 i 节点。
267     current->pwd = mi;
268     current->root = mi;
// 统计该设备上空闲块数。首先令 i 等于超级块中表明的设备逻辑块总数。
269     free=0;
270     i=p->s_nzones;
// 然后根据逻辑块位图中相应比特位的占用情况统计出空闲块数。宏函数 set_bit() 其实只是测试
// 比特位,而非设置比特位。“i&8191”用于取得 i 节点号在当前块中的偏移值。“i>>13”是将 i 除以
// 8192,也即除一个磁盘块包含的比特位数。
271     while (-- i >= 0)
272         if (!set_bit(i&8191,p->s_zmap[i>>13]->b_data))
273             free++;
// 显示设备上空闲逻辑块数/逻辑块总数。
274     printk("%d/%d free blocks\n\r",free,p->s_nzones);
// 统计设备上空闲 i 节点数。首先令 i 等于超级块中表明的设备上 i 节点总数+1。加 1 是将 0 节点
// 也统计进去。[??]
275     free=0;
276     i=p->s_ninodes+1;
// 然后根据 i 节点位图中相应比特位的占用情况计算出空闲 i 节点数。
277     while (-- i >= 0)
278         if (!set_bit(i&8191,p->s_imap[i>>13]->b_data))
279             free++;
// 显示设备上可用的空闲 i 节点数/i 节点总数。
280     printk("%d/%d free inodes\n\r",free,p->s_ninodes);
281 }
282

```

### 8.14.3 其它信息



## 8.15 open.c 文件

### 8.15.1 功能描述

### 8.15.2 代码注释

列表 linux/fs/open.c 程序

```

1  /*
2  *  linux/fs/open.c
3  *
4  *  (C) 1991  Linus Torvalds
5  */
6
7  #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。
8  #include <errno.h>          // 错误号头文件。包含系统中各种出错号。(Linus 从 minix 中引进的)。
9  #include <fcntl.h>          // 文件控制头文件。用于文件及其描述符的操作控制常数符号的定义。
10 #include <sys/types.h>      // 类型头文件。定义了基本的系统数据类型。
11 #include <utime.h>          // 用户时间头文件。定义了访问和修改时间结构以及 utime() 原型。
12 #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。
13
14 #include <linux/sched.h>    // 调度程序头文件，定义了任务结构 task_struct、初始任务 0 的数据，
                            // 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
15 #include <linux/tty.h>     // tty 头文件，定义了有关 tty_io，串行通信方面的参数、常数。
16 #include <linux/kernel.h>  // 内核头文件。含有一些内核常用函数的原形定义。
17 #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。
18
19 // 取文件系统信息函数调用函数。
20 int sys_ustat(int dev, struct ustat * ubuf)
21 {
22     return -ENOSYS;
23 }
24
25 // 设置文件访问和修改时间。
26 // 参数 filename 是文件名，times 是访问和修改时间结构指针。
27 // 如果 times 指针不为 NULL，则取 utimbuf 结构中的时间信息来设置文件的访问和修改时间。如果
28 // times 指针是 NULL，则取系统当前时间来设置指定文件的访问和修改时间域。
29 int sys_utime(char * filename, struct utimbuf * times)
30 {
31     struct m_inode * inode;
32     long actime, modtime;
33
34     // 根据文件名寻找对应的 i 节点，如果没有找到，则返回出错码。
35     if (!(inode=namei(filename)))
36         return -ENOENT;
37
38     // 如果访问和修改时间数据结构指针不为 NULL，则从结构中读取用户设置的时间值。
39     if (times) {
40         actime = get_fs_long((unsigned long *) &times->actime);
41         modtime = get_fs_long((unsigned long *) &times->modtime);
42         // 否则将访问和修改时间置为当前时间。
43     } else

```



```

35         actime = modtime = CURRENT\_TIME;
// 修改 i 节点中的访问时间字段和修改时间字段。
36         inode->i_atime = actime;
37         inode->i_mtime = modtime;
// 置 i 节点已修改标志, 释放该节点, 并返回 0。
38         inode->i_dirt = 1;
39         iput(inode);
40         return 0;
41     }
42
43     /*
44      * XXX should we use the real or effective uid?  BSD uses the real uid,
45      * so as to make this call useful to setuid programs.
46      */
// 文件属性 XXX, 我们该用真实用户 id 还是有效用户 id? BSD 系统使用了真实用户 id,
// 以使该调用可以供 setuid 程序使用。(注: POSIX 标准建议使用真实用户 ID)
// 检查对文件的访问权限。
// 参数 filename 是文件名, mode 是屏蔽码, 由 R_OK(4)、W_OK(2)、X_OK(1)和 F_OK(0)组成。
// 如果请求访问允许的话, 则返回 0, 否则返回出错码。
47 int sys\_access(const char * filename,int mode)
48 {
49     struct m\_inode * inode;
50     int res, i_mode;
51
// 屏蔽码由低 3 位组成, 因此清除所有高比特位。
52     mode &= 0007;
// 如果文件名对应的 i 节点不存在, 则返回出错码。
53     if (!(inode=namei(filename)))
54         return -EACCES;
// 取文件的属性码, 并释放该 i 节点。
55     i_mode = res = inode->i_mode & 0777;
56     iput(inode);
// 如果当前进程是该文件的宿主, 则取文件宿主属性。
57     if (current->uid == inode->i_uid)
58         res >>= 6;
// 否则如果当前进程是与该文件同属一组, 则取文件组属性。
59     else if (current->gid == inode->i_gid)
60         res >>= 6;
// 如果文件属性具有查询的属性位, 则访问许可, 返回 0。
61     if ((res & 0007 & mode) == mode)
62         return 0;
63
64     /*
65      * XXX we are doing this test last because we really should be
66      * swapping the effective with the real user id (temporarily),
67      * and then calling suser() routine.  If we do call the
68      * suser() routine, it needs to be called last.
69      */
// XXX 我们最后才做下面的测试, 因为我们实际上需要交换有效用户 id 和
// 真实用户 id (临时地), 然后才调用 suser() 函数。如果我们确实要调用
// suser() 函数, 则需要最后才被调用。

```

```

        */
// 如果当前用户 id 为 0 (超级用户) 并且屏蔽码执行位是 0 或文件可以被任何人访问, 则返回 0。
69     if ((!current->uid) &&
70         (!(mode & 1) || (i_mode & 0111)))
71         return 0;
// 否则返回出错码。
72     return -EACCES;
73 }
74
///// 改变当前工作目录系统调用函数。
// 参数 filename 是目录名。
// 操作成功则返回 0, 否则返回出错码。
75 int sys_chdir(const char * filename)
76 {
77     struct m_inode * inode;
78
// 如果文件名对应的 i 节点不存在, 则返回出错码。
79     if (!(inode = namei(filename)))
80         return -ENOENT;
// 如果该 i 节点不是目录的 i 节点, 则释放该节点, 返回出错码。
81     if (!S_ISDIR(inode->i_mode)) {
82         iput(inode);
83         return -ENOTDIR;
84     }
// 释放当前进程原工作目录 i 节点, 并指向该新置的工作目录 i 节点。返回 0。
85     iput(current->pwd);
86     current->pwd = inode;
87     return (0);
88 }
89
///// 改变根目录系统调用函数。
// 将指定的路径名改为根目录 '/'。
// 如果操作成功则返回 0, 否则返回出错码。
90 int sys_chroot(const char * filename)
91 {
92     struct m_inode * inode;
93
// 如果文件名对应的 i 节点不存在, 则返回出错码。
94     if (!(inode = namei(filename)))
95         return -ENOENT;
// 如果该 i 节点不是目录的 i 节点, 则释放该节点, 返回出错码。
96     if (!S_ISDIR(inode->i_mode)) {
97         iput(inode);
98         return -ENOTDIR;
99     }
// 释放当前进程的根目录 i 节点, 并重置为这里指定目录名的 i 节点, 返回 0。
100    iput(current->root);
101    current->root = inode;
102    return (0);
103 }
104
///// 修改文件属性系统调用函数。
// 参数 filename 是文件名, mode 是新的文件属性。

```

```

// 若操作成功则返回 0，否则返回出错码。
105 int sys_chmod(const char * filename,int mode)
106 {
107     struct m_inode * inode;
108
109     // 如果文件名对应的 i 节点不存在，则返回出错码。
110     if (!(inode=namei(filename)))
111         return -ENOENT;
112     // 如果当前进程的有效用户 id 不等于文件 i 节点的用户 id，并且当前进程不是超级用户，则释放该
113     // 文件 i 节点，返回出错码。
114     if ((current->euid != inode->i_uid) && !suser()) {
115         iput(inode);
116         return -EACCES;
117     }
118     // 重新设置 i 节点的文件属性，并置该 i 节点已修改标志。释放该 i 节点，返回 0。
119     inode->i_mode = (mode & 07777) | (inode->i_mode & ~07777);
120     inode->i_dirt = 1;
121     iput(inode);
122     return 0;
123 }
124
125 // 修改文件宿主系统调用函数。
126 // 参数 filename 是文件名，uid 是用户标识符(用户 id)，gid 是组 id。
127 // 若操作成功则返回 0，否则返回出错码。
128 int sys_chown(const char * filename,int uid,int gid)
129 {
130     struct m_inode * inode;
131
132     // 如果文件名对应的 i 节点不存在，则返回出错码。
133     if (!(inode=namei(filename)))
134         return -ENOENT;
135     // 若当前进程不是超级用户，则释放该 i 节点，返回出错码。
136     if (!suser()) {
137         iput(inode);
138         return -EACCES;
139     }
140     // 设置文件对应 i 节点的用户 id 和组 id，并置 i 节点已经修改标志，释放该 i 节点，返回 0。
141     inode->i_uid=uid;
142     inode->i_gid=gid;
143     inode->i_dirt=1;
144     iput(inode);
145     return 0;
146 }
147
148 // 打开（或创建）文件系统调用函数。
149 // 参数 filename 是文件名，flag 是打开文件标志：只读 O_RDONLY、只写 O_WRONLY 或读写 O_RDWR，
150 // 以及 O_CREAT、O_EXCL、O_APPEND 等其它一些标志的组合，若本函数创建了一个新文件，则 mode
151 // 用于指定使用文件的许可属性，这些属性有 S_IRWXU(文件宿主具有读、写和执行权限)、S_IRUSR
152 // (用户具有读文件权限)、S_IRWXG(组成员具有读、写和执行权限)等等。对于新创建的文件，这些
153 // 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个可读写的文件句柄。
154 // 若操作成功则返回文件句柄(文件描述符)，否则返回出错码。(参见 sys/stat.h, fcntl.h)
155 int sys_open(const char * filename,int flag,int mode)
156 {

```

```

140     struct m\_inode * inode;
141     struct file * f;
142     int i, fd;
143
144     // 将用户设置的模式与进程的模式屏蔽码相与，产生许可的文件模式。
145     mode &= 0777 & ~current->umask;
146     // 搜索进程结构中文件结构指针数组，查找一个空闲项，若已经没有空闲项，则返回出错码。
147     for(fd=0 ; fd<NR\_OPEN ; fd++)
148         if (!current->filp[fd])
149             break;
150     if (fd>=NR\_OPEN)
151         return -EINVAL;
152     // 设置执行时关闭文件句柄位图，复位对应比特位。
153     current->close_on_exec &= ~(1<<fd);
154     // 令 f 指向文件表数组开始处。搜索空闲文件结构项(句柄引用计数为 0 的项)，若已经没有空闲
155     // 文件表结构项，则返回出错码。
156     f=0+file\_table;
157     for (i=0 ; i<NR\_FILE ; i++, f++)
158         if (!f->f_count) break;
159     if (i>=NR\_FILE)
160         return -EINVAL;
161     // 让进程的对应文件句柄的文件结构指针指向搜索到的文件结构，并令句柄引用计数递增 1。
162     (current->filp[fd]=f)->f_count++;
163     // 调用函数执行打开操作，若返回值小于 0，则说明出错，释放刚申请到的文件结构，返回出错码。
164     if ((i=open\_namei(filename, flag, mode, &inode))<0) {
165         current->filp[fd]=NULL;
166         f->f_count=0;
167         return i;
168     }
169     /* ttys are somewhat special (ttyxx major==4, tty major==5) */
170     /* ttys 有些特殊 (ttyxx 主号==4, tty 主号==5) */
171     // 如果是字符设备文件，那么如果设备号是 4 的话，则设置当前进程的 tty 号为该 i 节点的子设备号。
172     // 并设置当前进程 tty 对应的 tty 表项的父进程组号等于进程的父进程组号。
173     if (S\_ISCHR(inode->i_mode))
174         if (MAJOR(inode->i_zone[0])==4) {
175             if (current->leader && current->tty<0) {
176                 current->tty = MINOR(inode->i_zone[0]);
177                 tty\_table[current->tty].pgrp = current->pgrp;
178             }
179             // 否则如果该字符文件设备号是 5 的话，若当前进程没有 tty，则说明出错，释放 i 节点和申请到的
180             // 文件结构，返回出错码。
181             } else if (MAJOR(inode->i_zone[0])==5)
182                 if (current->tty<0) {
183                     iput(inode);
184                     current->filp[fd]=NULL;
185                     f->f_count=0;
186                     return -EPERM;
187                 }
188     /* Likewise with block-devices: check for floppy_change */
189     /* 同样对于块设备文件：需要检查盘片是否被更换 */
190     // 如果打开的是块设备文件，则检查盘片是否更换，若更换则需要是高速缓冲中对应该设备的所有
191     // 缓冲块失效。
192     if (S\_ISBLK(inode->i_mode))

```

```

178         check\_disk\_change(inode->i_zone[0]);
// 初始化文件结构。置文件结构属性和标志，置句柄引用计数为 1，设置 i 节点字段，文件读写指针
// 初始化为 0。返回文件句柄。
179         f->f_mode = inode->i_mode;
180         f->f_flags = flag;
181         f->f_count = 1;
182         f->f_inode = inode;
183         f->f_pos = 0;
184         return (fd);
185     }
186
// 创建文件系统调用函数。
// 参数 pathname 是路径名，mode 与上面的 sys_open() 函数相同。
// 成功则返回文件句柄，否则返回出错码。
187 int sys\_creat(const char * pathname, int mode)
188 {
189     return sys\_open(pathname, O\_CREAT | O\_TRUNC, mode);
190 }
191
// 关闭文件系统调用函数。
// 参数 fd 是文件句柄。
// 成功则返回 0，否则返回出错码。
192 int sys\_close(unsigned int fd)
193 {
194     struct file * filp;
195
// 若文件句柄值大于程序同时能打开的文件数，则返回出错码。
196     if (fd >= NR\_OPEN)
197         return -EINVAL;
// 复位进程的执行时关闭文件句柄位图对应位。
198     current->close_on_exec &= ~(1<<fd);
// 若该文件句柄对应的文件结构指针是 NULL，则返回出错码。
199     if (! (filp = current->filp[fd]))
200         return -EINVAL;
// 置该文件句柄的文件结构指针为 NULL。
201     current->filp[fd] = NULL;
// 若在关闭文件之前，对应文件结构中的句柄引用计数已经为 0，则说明内核出错，死机。
202     if (filp->f_count == 0)
203         panic("Close: file count is 0");
// 否则将对应文件结构的句柄引用计数减 1，如果还不为 0，则返回 0（成功）。若已等于 0，说明该
// 文件已经没有句柄引用，则释放该文件 i 节点，返回 0。
204     if (--filp->f_count)
205         return (0);
206     iput(filp->f_inode);
207     return (0);
208 }
209

```

### 8.15.3 其它信息

## 8.16 exec.c 程序

### 8.16.1 功能描述

本源程序实现对二进制可执行文件和 shell 脚本文件的加载与执行。其中主要的函数是函数 `do_execve()`，它是系统中断调用 (`int 0x80`) 功能号 `__NR_execve()` 调用的 C 处理函数，是 `exec()` 函数簇的主要实现函数。其主要功能为：

- 执行对参数和环境参数空间页面的初始化操作 —— 设置初始空间起始指针；初始化空间页面指针数组为 (NULL)；根据执行文件名取执行对象的 I 节点；计算参数个数和环境变量个数；检查文件类型，执行权限；
- 根据执行文件开始部分的头数据结构，对其中信息进行处理 —— 根据被执行文件 I 节点读取文件头部信息；若是 Shell 脚本程序（第一行以 `#!` 开始），则分析 Shell 程序名及其参数，并以被执行文件作为参数执行该执行的 Shell 程序；执行根据文件的幻数以及段长度等信息判断是否可执行；
- 对当前调用进程进行运行新文件前初始化操作 —— 指向新执行文件的 I 节点；复位信号处理句柄；根据头结构信息设置局部描述符基址和段长；设置参数和环境参数页面指针；修改进行各执行字段内容；
- 替换堆栈上原调用 `execve()` 程序的返回地址为新执行程序运行地址，运行新加载的程序。

`execve()` 函数有大量对参数和环境空间的处理操作，参数和环境空间共可有 `MAX_ARG_PAGES` 个页面，总长度可达 128kB 字节。在该空间中存放数据的方式类似于堆栈操作，即是从假设的 128kB 空间末端处逆向开始存放参数或环境变量字符串的。在初始时，程序定义了一个指向该空间末端 (128kB-4 字节) 处空间内偏移值 `p`，该偏移值随着存放数据的增多而后退，由图中可以看出，`p` 明确地指出了当前参数环境空间中还剩余多少可用空间。在分析程序中 `copy_string()` 函数时，可参照此图。

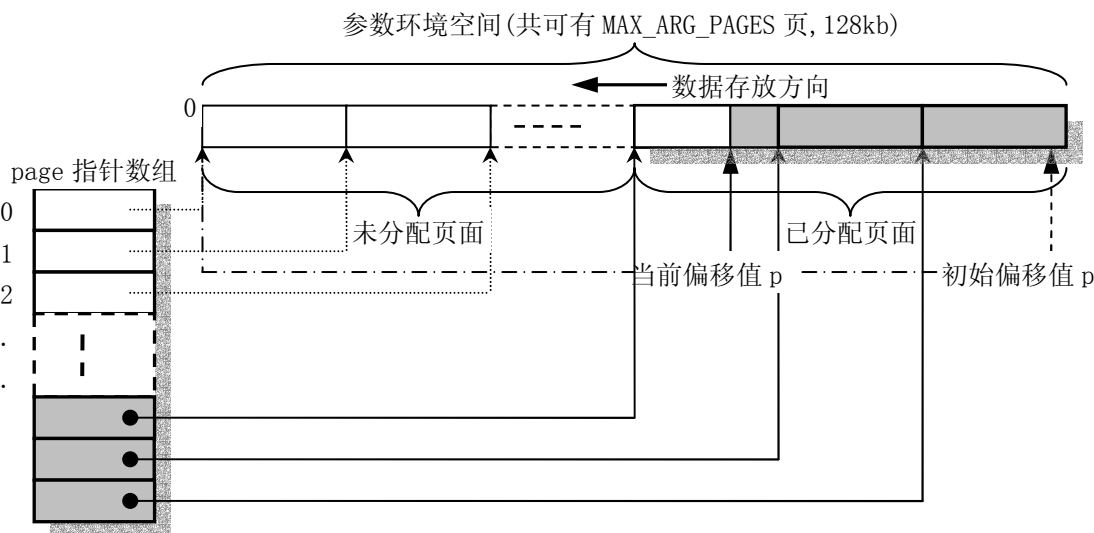


图 参数和环境变量字符串空间

`create_tables()` 函数用于根据给定的当前堆栈指针值 `p` 以及参数变量个数 `argc` 和环境变量个数 `envc`，在新的程序堆栈中创建环境和参数变量指针表，并返回此时的堆栈指针值 `sp`。创建完毕后堆栈指针表的形式见下图所示。

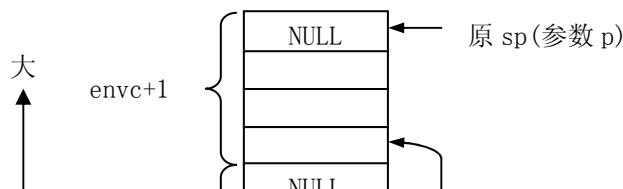


图 新程序堆栈中指针表示意图

## 8.16.2 代码注释

列表 linux/fs/exec.c 程序

---

```

1  /*
2   * linux/fs/exec.c
3   *
4   * (C) 1991 Linus Torvalds
5   */
6
7  /*
8   * #-checking implemented by tytso.
9   */
10 /*
11  * #-开始的程序检测部分是由 tytso 实现的。
12  */
13
14 /*
15  * Demand-loading implemented 01.12.91 - no need to read anything but
16  * the header into memory. The inode of the executable is put into
17  * "current->executable", and page faults do the actual loading. Clean.
18  *
19  * Once more I can proudly say that linux stood up to being changed: it
20  * was less than 2 hours work to get demand-loading completely implemented.
21  */
22 /*
23  * 需求时加载是于 1991.12.1 实现的 - 只需将执行文件头部分读进内存而无须
24  * 将整个执行文件都加载进内存。执行文件的 i 节点被放在当前进程的可执行字段中
25  * ("current->executable"), 而页异常会进行执行文件的实际加载操作以及清理工作。
26  *
27  * 我可以再一次自豪地说, linux 经得起修改: 只用了不到 2 小时的工作时间就完全
28  * 实现了需求加载处理。
29  */
30
31 #include <errno.h>          // 错误号头文件。包含系统中各种出错号。(Linus 从 minix 中引进的)。
32 #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。
33 #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。
34 #include <a.out.h>          // a.out 头文件。定义了 a.out 执行文件格式和一些宏。

```

```

24
25 #include <linux/fs.h>      // 文件系统头文件。定义文件表结构 (file, buffer_head, m_inode 等)。
26 #include <linux/sched.h>   // 调度程序头文件，定义了任务结构 task_struct、初始任务 0 的数据，
                             // 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
27 #include <linux/kernel.h>  // 内核头文件。含有一些内核常用函数的原形定义。
28 #include <linux/mm.h>      // 内存管理头文件。含有页面大小定义和一些页面释放函数原型。
29 #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。
30
31 extern int sys_exit(int exit_code); // 程序退出系统调用。
32 extern int sys_close(int fd);      // 文件关闭系统调用。
33
34 /*
35  * MAX_ARG_PAGES defines the number of pages allocated for arguments
36  * and envelope for the new program. 32 should suffice, this gives
37  * a maximum env+arg of 128kB !
38  */
39 /*
40  * MAX_ARG_PAGES 定义了新程序分配给参数和环境变量使用的内存最大页数。
41  * 32 页内存应该足够了，这使得环境和参数 (env+arg) 空间的总合达到 128kB!
42  */
43 #define MAX_ARG_PAGES 32
44
45 /*
46  * create_tables() parses the env- and arg-strings in new user
47  * memory and creates the pointer tables from them, and puts their
48  * addresses on the "stack", returning the new stack pointer value.
49  */
50 /*
51  * create_tables() 函数在新用户内存中解析环境变量和参数字符串，由此
52  * 创建指针表，并将它们的地址放到“堆栈”上，然后返回新栈的指针值。
53  */
54 // 在新用户堆栈中创建环境和参数变量指针表。
55 // 参数: p - 以数据段为起点的参数和环境信息偏移指针; argc - 参数个数; envc - 环境变量数。
56 // 返回: 堆栈指针。
57 static unsigned long * create_tables(char * p, int argc, int envc)
58 {
59     unsigned long *argv, *envp;
60     unsigned long * sp;
61
62     // 堆栈指针是以 4 字节 (1 节) 为边界寻址的，因此这里让 sp 为 4 的整数倍。
63     sp = (unsigned long *) (0xffffffff & (unsigned long) p);
64     // sp 向下移动，空出环境参数占用的空间个数，并让环境参数指针 envp 指向该处。
65     sp -= envc+1;
66     envp = sp;
67     // sp 向下移动，空出命令行参数指针占用的空间个数，并让 argv 指针指向该处。
68     // 下面指针加 1，sp 将递增指针宽度字节值。
69     sp -= argc+1;
70     argv = sp;
71     // 将环境参数指针 envp 和命令行参数指针以及命令行参数个数压入堆栈。
72     put_fs_long((unsigned long) envp, --sp);
73     put_fs_long((unsigned long) argv, --sp);
74     put_fs_long((unsigned long) argc, --sp);
75     // 将命令行各参数指针放入前面空出来的相应地方，最后放置一个 NULL 指针。

```



```

59     while (argc-->0) {
60         put_fs_long((unsigned long) p, argv++);
61         while (get_fs_byte(p++)) /* nothing */; // p 指针前移 4 字节。
62     }
63     put_fs_long(0, argv);
64 // 将环境变量各指针放入前面空出来的相应地方，最后放置一个 NULL 指针。
65     while (envc-->0) {
66         put_fs_long((unsigned long) p, envp++);
67         while (get_fs_byte(p++)) /* nothing */;
68     }
69     put_fs_long(0, envp);
70     return sp; // 返回构造的当前新堆栈指针。
71 }
72 /*
73  * count() counts the number of arguments/envelopes
74  */
75 /*
76  * count() 函数计算命令行参数/环境变量的个数。
77  */
78 // 计算参数个数。
79 // 参数: argv - 参数指针数组，最后一个指针项是 NULL。
80 // 返回: 参数个数。
81 static int count(char ** argv)
82 {
83     int i=0;
84     char ** tmp;
85
86     if (tmp = argv)
87         while (get_fs_long((unsigned long *) (tmp++)))
88             i++;
89
90     return i;
91 }
92 /*
93  * 'copy_string()' copies argument/envelope strings from user
94  * memory to free pages in kernel mem. These are in a format ready
95  * to be put directly into the top of new user memory.
96  *
97  * Modified by TYT, 11/24/91 to add the from_kmem argument, which specifies
98  * whether the string and the string array are from user or kernel segments:
99  *
100  * from_kmem    argv *    argv **
101  * 0            user space user space
102  * 1            kernel space user space
103  * 2            kernel space kernel space
104  *
105  * We do this by playing games with the fs segment register. Since it
106  * is expensive to load a segment register, we try to avoid calling
107  * set_fs() unless we absolutely have to.
108  */
109 /*

```

```

* 'copy_string()' 函数从用户内存空间拷贝参数和环境字符串到内核空闲页面内存中。
* 这些已具有直接放到新用户内存中的格式。
*
* 由 TYT(Tytso)于 1991.12.24 日修改, 增加了 from_kmem 参数, 该参数指明了字符串或
* 字符串数组是来自用户段还是内核段。
*
* from_kmem    argv *      argv **
*   0          用户空间    用户空间
*   1          内核空间    用户空间
*   2          内核空间    内核空间
*
* 我们是通过巧妙处理 fs 段寄存器来操作的。由于加载一个段寄存器代价太大, 所以
* 我们尽量避免调用 set_fs(), 除非实在必要。
*/
//// 复制指定个数的参数字符串到参数和环境空间。
// 参数: argc - 欲添加的参数个数; argv - 参数指针数组; page - 参数和环境空间页面指针数组。
//      p - 在参数表空间中的偏移指针, 始终指向已复制串的头; from_kmem - 字符串来源标志。
// 在 do_execve() 函数中, p 初始化为指向参数表(128kB)空间的最后一个长字处, 参数字符串
// 是以堆栈操作方式逆向往其中复制存放的, 因此 p 指针会始终指向参数字符串的头。
// 返回: 参数和环境空间当前头部指针。
104 static unsigned long copy_strings(int argc, char ** argv, unsigned long *page,
105                                unsigned long p, int from_kmem)
106 {
107     char *tmp, *pag;
108     int len, offset = 0;
109     unsigned long old_fs, new_fs;
110
111     if (!p)
112         return 0;      /* bullet-proofing */ /* 偏移指针验证 */
// 取 ds 寄存器值到 new_fs, 并保存原 fs 寄存器值到 old_fs。
113     new_fs = get_ds();
114     old_fs = get_fs();
// 如果字符串和字符串数组来自内核空间, 则设置 fs 段寄存器指向内核数据段 (ds)。
115     if (from_kmem==2)
116         set_fs(new_fs);
// 循环处理各个参数, 从最后一个参数逆向开始复制, 复制到指定偏移地址处。
117     while (argc-- > 0) {
// 如果字符串在用户空间而字符串数组在内核空间, 则设置 fs 段寄存器指向内核数据段 (ds)。
118         if (from_kmem == 1)
119             set_fs(new_fs);
// 从最后一个参数开始逆向操作, 取 fs 段中最后一参数指针到 tmp, 如果为空, 则出错死机。
120         if (!(tmp = (char *)get_fs_long(((unsigned long *)argv)+argc)))
121             panic("argc is wrong");
// 如果字符串在用户空间而字符串数组在内核空间, 则恢复 fs 段寄存器原值。
122         if (from_kmem == 1)
123             set_fs(old_fs);
// 计算该参数字符串长度 len, 并使 tmp 指向该参数字符串末端。
124         len=0;      /* remember zero-padding */
125         do {        /* 我们知道串是以 NULL 字节结尾的 */
126             len++;
127         } while (get_fs_byte(tmp++));
// 如果该字符串长度超过此时参数和环境空间中还剩余的空闲长度, 则恢复 fs 段寄存器并返回 0。
128         if (p-len < 0) {      /* this shouldn't happen - 128kB */

```

```

129         set_fs(old_fs); /* 不会发生-因为有 128kB 的空间 */
130         return 0;
131     }
    // 复制 fs 段中当前指定的参数字符串, 是从该字符串尾逆向开始复制。
132     while (len) {
133         --p; --tmp; --len;
    // 函数刚开始执行时, 偏移变量 offset 被初始化为 0, 因此若 offset-1<0, 说明是首次复制字符串,
    // 则令其等于 p 指针在页面内的偏移值, 并申请空闲页面。
134         if (--offset < 0) {
135             offset = p % PAGE_SIZE;
    // 如果字符串和字符串数组在内核空间, 则恢复 fs 段寄存器原值。
136             if (from_kmem==2)
137                 set_fs(old_fs);
    // 如果当前偏移值 p 所在的串空间页面指针数组项 page[p/PAGE_SIZE]==0, 表示相应页面还不存在,
    // 则需申请新的内存空闲页面, 将该页面指针填入指针数组, 并且也使 pag 指向该新页面, 若申请不
    // 到空闲页面则返回 0。
138             if (!(pag = (char *) page[p/PAGE_SIZE])) &&
139                 !(pag = (char *) page[p/PAGE_SIZE] =
140                     (unsigned long *) get_free_page()))
141                 return 0;
    // 如果字符串和字符串数组来自内核空间, 则设置 fs 段寄存器指向内核数据段 (ds)。
142             if (from_kmem==2)
143                 set_fs(new_fs);
144         }
145     }
    // 从 fs 段中复制参数字符串中一字节到 pag+offset 处。
146     *(pag + offset) = get_fs_byte(tmp);
147 }
148 }
    // 如果字符串和字符串数组在内核空间, 则恢复 fs 段寄存器原值。
149 if (from_kmem==2)
150     set_fs(old_fs);
    // 最后, 返回参数和环境空间中已复制参数信息的头部偏移值。
151 return p;
152 }
153
    //// 修改局部描述符表中的描述符基址和段限长, 并将参数和环境空间页面放置在数据段末端。
    // 参数: text_size - 执行文件头部中 a_text 字段给出的代码段长度值;
    //      page - 参数和环境空间页面指针数组。
    // 返回: 数据段限长值 (64MB)。
154 static unsigned long change_ldt(unsigned long text_size, unsigned long * page)
155 {
156     unsigned long code_limit, data_limit, code_base, data_base;
157     int i;
158
    // 根据执行文件头部 a_text 值, 计算以页面长度为边界的代码段限长。并设置数据段长度为 64MB。
159     code_limit = text_size+PAGE_SIZE -1;
160     code_limit &= 0xFFFFF000;
161     data_limit = 0x4000000;
    // 取当前进程中局部描述符表代码段描述符中代码段基址, 代码段基址与数据段基址相同。
162     code_base = get_base(current->ldt[1]);
163     data_base = code_base;
    // 重新设置局部表中代码段和数据段描述符的基址和段限长。

```

```

164     set_base(current->ldt[1], code_base);
165     set_limit(current->ldt[1], code_limit);
166     set_base(current->ldt[2], data_base);
167     set_limit(current->ldt[2], data_limit);
168 /* make sure fs points to the NEW data segment */
    /* 要确信 fs 段寄存器已指向新的数据段 */
    // fs 段寄存器中放入局部表数据段描述符的选择符(0x17)。
169     __asm__ ("pushl $0x17\n\tpop %%fs");
    // 将参数和环境空间已存放数据的页面（共可有 MAX_ARG_PAGES 页，128kB）放到数据段线性地址的
    // 末端。是调用函数 put_page() 进行操作的（mm/memory.c, 197）。
170     data_base += data_limit;
171     for (i=MAX_ARG_PAGES-1 ; i>=0 ; i--) {
172         data_base -= PAGE_SIZE;
173         if (page[i]) // 如果该页面存在，
174             put_page(page[i], data_base); // 就放置该页面。
175     }
176     return data_limit; // 最后返回数据段限长(64MB)。
177 }
178
179 /*
180  * 'do_execve()' executes a new program.
181  */
    /*
    * 'do_execve()' 函数执行一个新程序。
    */
    //// execve() 系统中中断调用函数。加载并执行子进程（其它程序）。
    // 该函数系统中中断调用(int 0x80)功能号 __NR_execve 调用的函数。
    // 参数: eip - 指向堆栈中调用系统中中断的程序代码指针 eip 处，参见 kernel/system_call.s 程序
    // 开始部分的说明; tmp - 系统中中断调用本函数时的返回地址，无用；
    //      filename - 被执行程序文件名; argv - 命令行参数指针数组; envp - 环境变量指针数组。
    // 返回: 如果调用成功，则不返回；否则设置出错号，并返回-1。
182 int do_execve(unsigned long * eip, long tmp, char * filename,
183               char ** argv, char ** envp)
184 {
185     struct m_inode * inode; // 内存中 I 节点指针结构变量。
186     struct buffer_head * bh; // 高速缓存块头指针。
187     struct exec ex; // 执行文件头部数据结构变量。
188     unsigned long page[MAX_ARG_PAGES]; // 参数和环境字符串空间的页面指针数组。
189     int i, argc, envc;
190     int e_uid, e_gid; // 有效用户 id 和有效组 id。
191     int retval; // 返回值。
192     int sh_bang = 0; // 控制是否需要执行脚本处理代码。
    // 参数和环境字符串空间中的偏移指针，初始化为指向该空间的最后一个长字处。
193     unsigned long p=PAGE_SIZE*MAX_ARG_PAGES-4;
194
    // eip[1]中是原代码段寄存器 cs，其中的选择符不可以是内核段选择符，也即内核不能调用本函数。
195     if ((0xffff & eip[1]) != 0x000f)
196         panic("execve called from supervisor mode");
    // 初始化参数和环境串空间的页面指针数组（表）。
197     for (i=0 ; i<MAX_ARG_PAGES ; i++) /* clear page-table */
198         page[i]=0;
    // 取可执行文件的对应 i 节点号。
199     if (!(inode=namei(filename))) /* get executables inode */

```

```

200         return -ENOENT;
// 计算参数个数和环境变量个数。
201     argc = count(argv);
202     envc = count(envp);
203
// 执行文件必须是常规文件。若不是常规文件则置出错返回码, 跳转到 exec_error2(第 347 行)。
204 restart_interp:
205     if (!S_ISREG(inode->i_mode)) { /* must be regular file */
206         retval = -EACCES;
207         goto exec_error2;
208     }
// 检查被执行文件的执行权限。根据其属性(对应 i 节点的 uid 和 gid), 看本进程是否有权执行它。
209     i = inode->i_mode;
210     e_uid = (i & S_ISUID) ? inode->i_uid : current->euid;
211     e_gid = (i & S_ISGID) ? inode->i_gid : current->egid;
212     if (current->euid == inode->i_uid)
213         i >>= 6;
214     else if (current->egid == inode->i_gid)
215         i >>= 3;
216     if (!(i & 1) &&
217         !((inode->i_mode & 0111) && suser())) {
218         retval = -ENOEXEC;
219         goto exec_error2;
220     }
// 读取执行文件的第一块数据到高速缓冲区, 若出错则置出错码, 跳转到 exec_error2 处去处理。
221     if (!(bh = bread(inode->i_dev, inode->i_zone[0]))) {
222         retval = -EACCES;
223         goto exec_error2;
224     }
// 下面对执行文件的头结构数据进行处理, 首先让 ex 指向执行头部分的数据结构。
225     ex = *((struct exec *) bh->b_data); /* read exec-header */ /* 读取执行头部分 */
// 如果执行文件开始的两个字节为 '#!', 并且 sh_bang 标志没有置位, 则处理脚本文件的执行。
226     if ((bh->b_data[0] == '#') && (bh->b_data[1] == '!') && (!sh_bang)) {
227         /*
228          * This section does the #! interpretation.
229          * Sorta complicated, but hopefully it will work. -TYT
230          */
231         /*
232          * 这部分处理对 '#!' 的解释, 有些复杂, 但希望能工作。-TYT
233          */
234
235         char buf[1023], *cp, *interp, *i_name, *i_arg;
236         unsigned long old_fs;
237
// 复制执行程序头一行字符 '#!' 后面的字符串到 buf 中, 其中含有脚本处理程序名。
238         strncpy(buf, bh->b_data+2, 1022);
// 释放高速缓冲块和该执行文件 i 节点。
239         brelse(bh);
240         iput(inode);
// 取第一行内容, 并删除开始的空格、制表符。
241         buf[1022] = '\0';
242         if (cp = strchr(buf, '\n')) {
243             *cp = '\0';

```

```

241         for (cp = buf; (*cp == ' ') || (*cp == '\t'); cp++);
242     }
    // 若该行没有其它内容, 则出错。置出错码, 跳转到 exec_error1 处。
243     if (!cp || *cp == '\0') {
244         retval = -ENOEXEC; /* No interpreter name found */
245         goto exec_error1;
246     }
    // 否则就得到了开头是脚本解释执行程序名称的一行内容。
247     interp = i_name = cp;
    // 下面分析该行。首先取第一个字符串, 其应该是脚本解释程序名, i_name 指向该名称。
248     i_arg = 0;
249     for ( ; *cp && (*cp != ' ') && (*cp != '\t'); cp++) {
250         if (*cp == '/')
251             i_name = cp+1;
252     }
    // 若文件名后还有字符, 则应该是参数串, 令 i_arg 指向该串。
253     if (*cp) {
254         *cp++ = '\0';
255         i_arg = cp;
256     }
257     /*
258      * OK, we've parsed out the interpreter name and
259      * (optional) argument.
260      */
    /*
    * OK, 我们已经解析出解释程序的文件名以及(可选的)参数。
    */
    // 若 sh_bang 标志没有设置, 则设置它, 并复制指定个数的环境变量串和参数串到参数和环境空间中。
261     if (sh_bang++ == 0) {
262         p = copy_strings(envc, envp, page, p, 0);
263         p = copy_strings(--argc, argv+1, page, p, 0);
264     }
265     /*
266      * Splice in (1) the interpreter's name for argv[0]
267      *          (2) (optional) argument to interpreter
268      *          (3) filename of shell script
269      *
270      * This is done in reverse order, because of how the
271      * user environment and arguments are stored.
272      */
    /*
    * 拼接 (1) argv[0] 中放解释程序的名称
    *      (2) (可选的)解释程序的参数
    *      (3) 脚本程序的名称
    *
    * 这是以逆序进行处理的, 是由于用户环境和参数的存放方式造成的。
    */
    // 复制脚本程序文件名到参数和环境空间中。
273     p = copy_strings(1, &filename, page, p, 1);
    // 复制解释程序的参数到参数和环境空间中。
274     argc++;
275     if (i_arg) {
276         p = copy_strings(1, &i_arg, page, p, 2);

```

```

277         argc++;
278     }
    // 复制解释程序文件名到参数和环境空间中。若出错，则置出错码，跳转到 exec_error1。
279     p = copy\_strings(l, &i_name, page, p, 2);
280     argc++;
281     if (!p) {
282         retval = -ENOMEM;
283         goto exec_error1;
284     }
285     /*
286      * OK, now restart the process with the interpreter's inode.
287      */
    /*
    * OK, 现在使用解释程序的 i 节点重启进程。
    */
    // 保留原 fs 段寄存器（原指向用户数据段），现置其指向内核数据段。
288     old_fs = get\_fs();
289     set\_fs(get\_ds());
    // 取解释程序的 i 节点，并跳转到 restart_interp 处重新处理。
290     if (!(inode=namei(interp))) { /* get executables inode */
291         set\_fs(old_fs);
292         retval = -ENOENT;
293         goto exec_error1;
294     }
295     set\_fs(old_fs);
296     goto restart_interp;
297 }
    // 释放该缓冲区。
298     brelse(bh);
    // 下面对执行头信息进行处理。
    // 对于下列情况，将不执行程序：如果执行文件不是需求页可执行文件(ZMAGIC)、或者代码重定位部分
    // 长度 a_trsize 不等于 0、或者数据重定位信息长度不等于 0、或者代码段+数据段+堆段长度超过 50MB、
    // 或者 i 节点表明的该执行文件长度小于代码段+数据段+符号表长度+执行头部分长度的总和。
299     if (N\_MAGIC(ex) != ZMAGIC || ex.a_trsize || ex.a_drsize ||
300         ex.a_text+ex.a_data+ex.a_bss>0x3000000 ||
301         inode->i_size < ex.a_text+ex.a_data+ex.a_syms+N_TXTOFF(ex)) {
302         retval = -ENOEXEC;
303         goto exec_error2;
304     }
    // 如果执行文件执行头部分长度不等于一个内存块大小（1024 字节），也不能执行。转 exec_error2。
305     if (N\_TXTOFF(ex) != BLOCK\_SIZE) {
306         printk("%s: N_TXTOFF != BLOCK_SIZE. See a.out.h.", filename);
307         retval = -ENOEXEC;
308         goto exec_error2;
309     }
    // 如果 sh_bang 标志没有设置，则复制指定个数的环境变量字符串和参数到参数和环境空间中。
    // 若 sh_bang 标志已经设置，则表明是将运行脚本程序，此时环境变量页面已经复制，无须再复制。
310     if (!sh_bang) {
311         p = copy\_strings(envc, envp, page, p, 0);
312         p = copy\_strings(argc, argv, page, p, 0);
    // 如果 p=0，则表示环境变量与参数空间页面已经被占满，容纳不下了。转至出错处理处。
313     if (!p) {
314         retval = -ENOMEM;

```



```

315         goto exec_error2;
316     }
317 }
318 /* OK, This is the point of no return */
319 /* OK, 下面开始就没有返回的地方了 */
320 // 如果原程序也是一个执行程序, 则释放其 i 节点, 并让进程 executable 字段指向新程序 i 节点。
321     if (current->executable)
322         iput(current->executable);
323     current->executable = inode;
324 // 清复位所有信号处理句柄。但对于 SIG_IGN 句柄不能复位, 因此在 322 与 323 行之间需添加一条
325 // if 语句: if (current->sa[I].sa_handler != SIG_IGN)。这是源代码中的一个 bug。
326     for (i=0 ; i<32 ; i++)
327         current->sigaction[i].sa_handler = NULL;
328 // 根据执行时关闭(close_on_exec)文件句柄位图标志, 关闭指定的打开文件, 并复位该标志。
329     for (i=0 ; i<NR\_OPEN ; i++)
330         if ((current->close_on_exec>>i)&1)
331             sys\_close(i);
332     current->close_on_exec = 0;
333 // 根据指定的基地址和限长, 释放原来程序代码段和数据段所对应的内存页表指定的内存块及页表本身。
334     free\_page\_tables(get\_base(current->ldt[1]), get\_limit(0x0f));
335     free\_page\_tables(get\_base(current->ldt[2]), get\_limit(0x17));
336 // 如果“上次任务使用了协处理器”指向的是当前进程, 则将其置空, 并复位使用了协处理器的标志。
337     if (last\_task\_used\_math == current)
338         last\_task\_used\_math = NULL;
339     current->used_math = 0;
340 // 根据 a_text 修改局部表中描述符基址和段限长, 并将参数和环境空间页面放置在数据段末端。
341 // 执行下面语句之后, p 此时是以数据段起始处为原点的偏移值, 仍指向参数和环境空间数据开始处,
342 // 也即转换成为堆栈的指针。
343     p += change\_ldt(ex.a_text, page)-MAX\_ARG\_PAGES*PAGE\_SIZE;
344 // create_tables() 在新用户堆栈中创建环境和参数变量指针表, 并返回该堆栈指针。
345     p = (unsigned long) create\_tables((char *)p, argc, envc);
346 // 修改当前进程各字段为新执行程序的信息。令进程代码段尾值字段 end_code = a_text; 令进程数据
347 // 段尾字段 end_data = a_data + a_text; 令进程堆结尾字段 brk = a_text + a_data + a_bss。
348     current->brk = ex.a_bss +
349         (current->end_data = ex.a_data +
350         (current->end_code = ex.a_text));
351 // 设置进程堆栈开始字段为堆栈指针所在的页面, 并重新设置进程的用户 id 和组 id。
352     current->start_stack = p & 0xfffff000;
353     current->euid = e_uid;
354     current->egid = e_gid;
355 // 初始化一页 bss 段数据, 全为零。
356     i = ex.a_text+ex.a_data;
357     while (i&0xfff)
358         put\_fs\_byte(0, (char *) (i++));
359 // 将原调用系统中中断的程序在堆栈上的代码指针替换为指向新执行程序的入口点, 并将堆栈指针替换
360 // 为新执行程序的堆栈指针。返回指令将弹出这些堆栈数据并使得 CPU 去执行新的执行程序, 因此不会
361 // 返回到原调用系统中中断的程序中去了。
362     eip[0] = ex.a_entry;
363     eip[3] = p;
364     return 0;
365 exec_error2:
366     iput(inode);
367 exec_error1:

```



```

350         for (i=0 ; i<MAX_ARG_PAGES ; i++)
351             free_page(page[i]);
352         return(retval);
353     }
354

```

## 8.16.3 其它信息

### 8.16.3.1 a.out 执行文件格式

Linux 内核 0.11 版仅支持 a.out (Assembly & link editor output) 执行文件格式，虽然这种格式目前已经渐渐不用，而使用功能更为齐全的 ELF (Executable and Link Format) 格式，但是由于其简单性，作为学习入门的材料正好比较适用。下面全面介绍一下 a.out 格式。

在头文件<a.out.h>中声明了三个数据结构以及一些宏函数。这些数据结构描述了系统上可执行的机器码文件（二进制文件）。

一个执行文件共可有七个部分（七节）组成。按照顺序，这些部分是：

执行头部分(exec header)

执行文件头部分。该部分中含有一些参数，内核使用这些参数将执行文件加载到内存中并执行，而链接程序(ld)使用这些参数将一些二进制目标文件组合成一个可执行文件。这是唯一必要的组成部分。

代码段部分(text segment)

含有程序执行使被加载到内存中的指令代码和相关数据。可以以只读形式进行加载。

数据段部分(data segment)

这部分含有已经初始化过的数据，总是被加载到可读写的内存中。

代码重定位部分(text relocations)

这部分含有供链接程序使用的记录数据。在组合二进制目标文件时用于定位代码段中的指针或地址。

数据重定位部分(data relocations)

与代码重定位部分的作用类似，但是是用于数据段中指针的重定位。

符号表部分(symbol table)

这部分同样含有供链接程序使用的记录数据，用于在二进制目标文件之间对命名的变量和函数（符号）进行交叉引用。

字符串表部分(string table)

该部分含有与符号名相对应的字符串。

每个二进制执行文件均以一個执行数据结构（exec structure）开始。该数据结构的形式如下：

```

struct exec {
    unsigned long a_midmag;

```

```

    unsigned long a_text;
    unsigned long a_data;
    unsigned long a_bss;
    unsigned long a_syms;
    unsigned long a_entry;
    unsigned long a_trsize;
    unsigned long a_drsize;
};

```

各个字段的功能如下：

|          |   |
|----------|---|
| a_midmag | 该字段含有被 N_GETFLAG()、N_GETMID 和 N_GETMAGIC() 访问的子部分，是由链接程序在运行时加载到进程地址空间。宏 N_GETMID() 用于返回机器标识符(machine-id)，指示出二进制文件将在什么机器上运行。N_GETMAGIC() 宏指明魔数，它唯一地确定了二进制执行文件与其它加载的文件之间的区别。字段中必须包含以下值之一：<br>OMAGIC 表示代码和数据段紧随在执行头后面并且是连续存放的。内核将代码和数据段都加载到可读写内存中。<br>NMAGIC 同 OMAGIC 一样，代码和数据段紧随在执行头后面并且是连续存放的。然而内核将代码加载到了只读内存中，并把数据段加载到代码段后下一页可读写内存边界开始。<br>ZMAGIC 内核在必要时从二进制执行文件中加载独立的页面。执行头部、代码段和数据段都被链接程序处理成多个页面大小的块。内核加载的代码页面时只读的，而数据段的页面是可写的。 |
| a_text   | 该字段含有代码段的长度值，字节数。   |
| a_data   | 该字段含有数据段的长度值，字节数。   |
| a_bss    | 含有‘bss 段’的长度，内核用其设置在数据段后初始的 break (brk)。内核在加载程序时，这段可写内存显现出处于数据段后面，并且初始时为零。  |
| a_syms   | 含有符号表部分的字节长度值。  |
| a_entry  | 含有内核将执行文件加载到内存中以后，程序执行起始点的内存地址。   |
| a_trsize | 该字段含有代码重定位表的大小，是字节数。  |
| a_drsize | 该字段含有数据重定位表的大小，是字节数。  |

在 a.out.h 头文件中定义了几个宏，这些宏使用 exec 结构来测试一致性或者定位执行文件中各个部分（节）的位置偏移值。这些宏有：

N\_BADMAG(exec) 如果 a\_magic 字段不能被识别，则返回非零值。

N\_TXTOFF(exec) 代码段的起始位置字节偏移值。

N\_DATAOFF(exec) 数据段的起始位置字节偏移值。

N\_DRELOFF(exec) 数据重定位表的起始位置字节偏移值。

N\_TRELOFF(exec) 代码重定位表的起始位置字节偏移值。

N\_SYMOFF(exec) 符号表的起始位置字节偏移值。

N\_STROFF(exec) 字符串表的起始位置字节偏移值。

重定位记录具有标准格式，它使用重定位信息(relocation\_info)结构来描述：

```
struct relocation_info {
    int            r_address;
    unsigned int   r_symbolnum : 24,
                  r_pcrel : 1,
                  r_length : 2,
                  r_extern : 1,
                  r_basereel : 1,
                  r_jmptable : 1,
                  r_relative : 1,
                  r_copy : 1;
};
```

该结构中各字段的含义如下：

**r\_address** 该字段含有需要链接程序处理（编辑）的指针的字节偏移值。代码重定位的偏移值是从代码段开始处计数的，数据重定位的偏移值是从数据段开始处计算的。链接程序会将已经存储在该偏移处的值与使用重定位记录计算出的新值相加。

**r\_symbolnum** 该字段含有符号表中一个符号结构的序号值（不是字节偏移值）。链接程序在算出符号的绝对地址以后，就将该地址加到正在进行重定位的指针上。（如果 r\_extern 比特位是 0，那么情况就不同，见下面。）

**r\_pcrel** 如果设置了该位，链接程序就认为正在更新一个指针，该指针使用 pc 相关寻址方式，是属于机器码指令部分。当运行程序使用这个被重定位的指针时，该指针的地址被隐式地加到该指针上。

**r\_length** 该字段含有指针长度的 2 的次方值：0 表示 1 字节长，1 表示 2 字节长，2 表示 4 字节长。

**r\_extern** 如果被置位，表示该重定位需要一个外部引用；此时链接程序必须使用一个符号地址来更新相应指针。当该位是 0 时，则重定位是“局部”的；链接程序更新指针以反映各个段加载地址中

的变化,而不是反映一个符号值的变化(除非同时设置了 `r_baserel`,见下面)。在这种情况下,`r_symbolnum` 字段的内容是一个 `n_type` 值(见下面);这类字段告诉链接程序被重定位的指针指向那个段。

`r_baserel`      如果设置了该位,则 `r_symbolnum` 字段指定的符号将被重定位成全局偏移表(Global Offset Table)中的一个偏移值。在运行时刻,全局偏移表该偏移处被设置为符号的地址。

`r_jmptable`      如果被置位,则 `r_symbolnum` 字段指定的符号将被重定位成过程链接表(Procedure Linkage Table)中的一个偏移值。

`r_relative`      如果被置位,则说明此重定位与该目标文件将成为其组成部分的映象文件在运行时被加载的地址相关。这类重定位仅在共享目标文件中出现。

`r_copy`          如果被置位,该重定位记录指定了一个符号,该符号的内容将被复制到 `r_address` 指定的地方。该复制操作是通过共享目标模块中一个合适的数项中的运行时刻链接程序完成的。

符号将名称映射为地址(或者更通俗地讲是字符串映射到值)。由于链接程序对地址的调整,一个符号的名称必须用来表示其地址,直到已被赋予一个绝对地址值。符号是由符号表中固定长度的记录以及字符串表中的可变长度名称组成。符号表是 `nlist` 结构的一个数组,如下所示:

```
struct nlist {
    union {
        char    *n_name;
        long    n_strx;
    } n_un;
    unsigned char n_type;
    char          n_other;
    short         n_desc;
    unsigned long n_value;
};
```

其中各字段的含义为:

`n_un.n_strx` 含有本符号的名称在字符串表中的字节偏移值。当程序使用 `nlist()` 函数访问一个符号表时,该字段被替换为 `n_un.n_name` 字段,这是内存中字符串的指针。

`n_type`          用于链接程序确定如何更新符号的值。使用位屏蔽(bitmasks)可以将 `n_type` 字段分割成三个子字段,对于 `N_EXT` 类型位置位的符号,链接程序将它们看作是“外部的”符号,并且允许其它二进制目标文件对它们的引用。`N_TYPE` 屏蔽码用于链接程序感兴趣的比特位:

`N_UNDF`      一个未定义的符号。链接程序必须在其它二进制目标文件中定位一个具有相同名称的外部符号,以确定该符号的绝对数据值。特殊情况下,如果 `n_type` 字段是非零值,并且没有二进制文件定义了这个符号,则链接程序在 BSS 段中将该符号解析为一个地址,保留长度等于 `n_value` 的字节。如果符号在多于一个二进制目标文件中都没有定义并且这些二进制目标文件对其长度值不一致,则链接程序将选择所有二进制目标文件

中最大的长度。

**N\_ABS** 一个绝对符号。链接程序不会更新一个绝对符号。

**N\_TEXT** 一个代码符号。该符号的值是代码地址，链接程序在合并二进制目标文件时会更新其值。

**N\_DATA** 一个数据符号；与 **N\_TEXT** 类似，但是用于数据地址。对应代码和数据符号的值不是文件的偏移值而是地址；为了找出文件的偏移，就有必要确定相关部分开始加载的地址并减去它，然后加上该部分的偏移。

**N\_BSS** 一个 BSS 符号；与代码或数据符号类似，但在二进制目标文件中没有对应的偏移。

**N\_FN** 一个文件名符号。在合并二进制目标文件时，链接程序会将该符号插入在二进制文件中的符号之前。符号的名称就是给予链接程序的文件名，而其值是二进制文件中首个代码段地址。链接和加载时不需要文件名符号，但对于调式程序非常有用。

**N\_STAB** 屏蔽码用于选择符号调式程序（例如 `gdb`）感兴趣的位；其值在 `stab()` 中说明。

**n\_other** 该字段按照 **n\_type** 确定的段，提供有关符号重定位操作的符号独立性信息。目前，**n\_other** 字段的最低 4 位含有两个值之一：**AUX\_FUNC** 和 **AUX\_OBJECT**（有关定义参见 `<link.h>`）。**AUX\_FUNC** 将符号与可调用的函数相关，**AUX\_OBJECT** 将符号与数据相关，而不管它们是位于代码段还是数据段。该字段主要用于链接程序 `ld`，用于动态可执行程序创建。

**n\_desc** 保留给调式程序使用；链接程序不对其进行处理。不同的调试程序将该字段用作不同的用途。

**n\_value** 含有符号的值。对于代码、数据和 BSS 符号，这是一个地址；对于其它符号（例如调式程序符号），值可以是任意的。

字符串表是由长度为 `u_int32_t` 后跟一 `null` 结尾的符号字符串组成。长度代表整个表的字节大小，所以在 32 位的机器上其最小值（或者是第 1 个字符串的偏移）总是 4。

## 8.17 namei.c 文件

### 8.17.1 功能描述

### 8.17.2 代码注释

列表 `linux/fs/namei.c` 程序

```

1  /*
2   * linux/fs/namei.c
3   *
4   * (C) 1991 Linus Torvalds
5   */
6
7  /*
8   * Some corrections by tytso.
9   */
10 /*
11  * tytso 作了一些纠正。
12  */
13
14 #include <linux/sched.h> // 调度程序头文件，定义了任务结构 task_struct、初始任务 0 的数据，
15 // 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
16 #include <linux/kernel.h> // 内核头文件。含有一些内核常用函数的原形定义。
17 #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。
18
19 #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。
20 #include <fcntl.h> // 文件控制头文件。用于文件及其描述符的操作控制常数符号的定义。
21 #include <errno.h> // 错误号头文件。包含系统中各种出错号。(Linus 从 minix 中引进的)。
22 #include <const.h> // 常数符号头文件。目前仅定义了 i 节点中 i_mode 字段的各标志位。
23 #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。
24
25 #define ACC_MODE(x) ("004006377"[(x)&0_ACCMODE]) // 访问模式宏(前面是 8 进制数)。
26
27 /*
28  * comment out this line if you want names > NAME_LEN chars to be
29  * truncated. Else they will be disallowed.
30  */
31 /*
32  * 如果能让文件名长度>NAME_NAME 的字符被截掉，就将下面定义注释掉。
33  */
34 #define NO_TRUNCATE
35
36 #define MAY_EXEC 1 // 可执行(可进入)。
37 #define MAY_WRITE 2 // 可写。
38 #define MAY_READ 4 // 可读。
39
40 /*
41  *
42  * permission()
43  *
44  * is used to check for read/write/execute permissions on a file.
45  * I don't know if we should look at just the euid or both euid and
46  * uid, but that should be easily changed.
47  */
48 /*
49  *
50  * permission()
51  * 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 还是
52  * 需要检查 euid 和 uid 两者，不过这很容易修改。
53  */
54 // 检测文件访问许可权限。
55 // 参数: inode - 文件对应的 i 节点; mask - 访问属性屏蔽码。

```

```

// 返回: 访问许可返回 1, 否则返回 0。
40 static int permission(struct m_inode * inode, int mask)
41 {
42     int mode = inode->i_mode;
43
44     /* special case: not even root can read/write a deleted file */
45     /* 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 */
46     /* 如果 i 节点有对应的设备, 但该 i 节点的连接数等于 0, 则返回。*/
47     if (inode->i_dev && !inode->i_nlinks)
48         return 0;
49     /* 否则, 如果进程的有效用户 id(euid)与 i 节点的用户 id 相同, 则取文件宿主的用户访问权限。*/
50     else if (current->euid==inode->i_uid)
51         mode >>= 6;
52     /* 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。*/
53     else if (current->egid==inode->i_gid)
54         mode >>= 3;
55     /* 如果上面所取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。*/
56     if (((mode & mask & 0007) == mask) || suser())
57         return 1;
58     return 0;
59 }
60
61 /*
62  * ok, we cannot use strncmp, as the name is not in our data space.
63  * Thus we'll have to use match. No big problem. Match also makes
64  * some sanity tests.
65  *
66  * NOTE! unlike strncmp, match returns 1 for success, 0 for failure.
67  */
68 /*
69  * 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(不在内核空间)。
70  * 因而我们只能使用 match()。问题不大。match() 同样也处理一些完整的测试。
71  *
72  * 注意! 与 strncmp 不同的是 match() 成功时返回 1, 失败时返回 0。
73  */
74 //// 指定长度字符串比较函数。
75 // 参数: len - 比较的字符串长度; name - 文件名指针; de - 目录项结构。
76 // 返回: 相同返回 1, 不同返回 0。
77 static int match(int len, const char * name, struct dir_entry * de)
78 {
79     register int same __asm__ ("ax");
80
81     /* 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。*/
82     if (!de || !de->inode || len > NAME_LEN)
83         return 0;
84     /* 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, 则返回 0。*/
85     if (len < NAME_LEN && de->name[len])
86         return 0;
87     /* 下面嵌入汇编语句, 在用户数据空间(fs)执行字符串的比较操作。*/
88     /* %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); */
89     /* %4 - ecx(比较的字节长度值 len)。
90     __asm__ ("cld\n\t" // 清方向位。
91             "fs ; repe ; cmpsb\n\t" // 用户空间执行循环比较[esi++]和[edi++]操作,

```

```

73         "setz %%a1"                                // 若比较结果一样(z=0)则设置 a1=1(same=eax)。
74         : "=a" (same)
75         : "" (0), "S" ((long) name), "D" ((long) de->name), "c" (len)
76         : "cx", "di", "si");
77     return same;                                    // 返回比较结果。
78 }
79
80 /*
81  *      find_entry()
82  *
83  * finds an entry in the specified directory with the wanted name. It
84  * returns the cache buffer in which the entry was found, and the entry
85  * itself (as a parameter - res_dir). It does NOT read the inode of the
86  * entry - you'll have to do that yourself if you want to.
87  *
88  * This also takes care of the few special cases due to '..'-traversal
89  * over a pseudo-root and a mount point.
90  */
91 /*
92  *      find_entry()
93  * 在指定的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
94  * 缓冲区以及目录项本身(作为一个参数 - res_dir)。并不读目录项的 i 节点 - 如
95  * 果需要的话需自己操作。
96  *
97  * '..' 目录项, 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
98  * 及安装点。
99  */
100 // 查找指定目录和文件名的目录项。
101 // 参数: dir - 指定目录 i 节点的指针; name - 文件名; namelen - 文件名长度;
102 // 返回: 高速缓冲区指针; res_dir - 返回的目录项结构指针;
103 static struct buffer_head * find_entry(struct m_inode ** dir,
104     const char * name, int namelen, struct dir_entry ** res_dir)
105 {
106     int entries;
107     int block, i;
108     struct buffer_head * bh;
109     struct dir_entry * de;
110     struct super_block * sb;
111
112     // 如果定义了 NO_TRUNCATE, 则若文件名长度超过最大长度 NAME_LEN, 则返回。
113 #ifdef NO_TRUNCATE
114     if (namelen > NAME_LEN)
115         return NULL;
116     //如果没有定义 NO_TRUNCATE, 则若文件名长度超过最大长度 NAME_LEN, 则截短之。
117 #else
118     if (namelen > NAME_LEN)
119         namelen = NAME_LEN;
120 #endif
121 // 计算本目录中目录项项数 entries。置空返回目录项结构指针。
122     entries = (*dir)->i_size / (sizeof (struct dir_entry));
123     *res_dir = NULL;
124 // 如果文件名长度等于 0, 则返回 NULL, 退出。
125     if (!namelen)

```



```

110         return NULL;
111 /* check for '..', as we might have to do some "magic" for it */
112 /* 检查目录项 '..', 因为可能需要对其特别处理 */
113 if (namelen==2 && get_fs_byte(name)=='.' && get_fs_byte(name+1)=='.') {
114 /* '..' in a pseudo-root results in a faked '.' (just change namelen) */
115 /* 伪根中的 '..' 如同一个假 '.' (只需改变名字长度) */
116 /* 如果当前进程的根节点指针即是指定的目录, 则将文件名修改为 '.',
117     if ((*dir) == current->root)
118         namelen=1;
119 // 否则如果该目录的 i 节点号等于 ROOT_INO(1) 的话, 说明是文件系统根节点。则取文件系统的超级块。

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 '..' 将导致目录交换到安装到文件系统的目录 i 节点。
124     注意! 由于设置了 mounted 标志, 因而我们能够取出该新目录 */
125     sb=get_super((*dir)->i_dev);
126 // 如果被安装到的 i 节点存在, 则先释放原 i 节点, 然后对被安装到的 i 节点进行处理。
127 // 让*dir 指向该被安装到的 i 节点; 该 i 节点的引用数加 1。
128     if (sb->s_imount) {
129         input(*dir);
130         (*dir)=sb->s_imount;
131         (*dir)->i_count++;
132     }
133 }
134 // 如果该 i 节点所指向的第一个直接磁盘块号为 0, 则返回 NULL, 退出。
135 if (!(block = (*dir)->i_zone[0]))
136     return NULL;
137 // 从节点所在设备读取指定的目录项数据块, 如果不成功, 则返回 NULL, 退出。
138 if (!(bh = bread((*dir)->i_dev, block)))
139     return NULL;
140 // 在目录项数据块中搜索匹配指定文件名的目录项, 首先让 de 指向数据块, 并在不超过目录中目录项数
141 // 的条件下, 循环执行搜索。
142 i = 0;
143 de = (struct dir_entry *) bh->b_data;
144 while (i < entries) {
145 // 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。
146     if ((char *)de >= BLOCK_SIZE+bh->b_data) {
147         brelse(bh);
148         bh = NULL;
149 // 在读入下一目录项数据块。若这块为空, 则只要还没有搜索完目录中的所有目录项, 就跳过该块,
150 // 继续读下一目录项数据块。若该块不空, 就让 de 指向该目录项数据块, 继续搜索。
151         if (!(block = bmap(*dir, i/DIR_ENTRIES_PER_BLOCK)) ||
152             !(bh = bread((*dir)->i_dev, block))) {
153             i += DIR_ENTRIES_PER_BLOCK;
154             continue;
155         }
156         de = (struct dir_entry *) bh->b_data;
157     }
158 // 如果找到匹配的目录项的话, 则返回该目录项结构指针和该目录项数据块指针, 退出。
159     if (match(namelen, name, de)) {
160         *res_dir = de;

```

```

146         return bh;
147     }
    // 否则继续在目录项数据块中比较下一个目录项。
148     de++;
149     i++;
150 }
    // 若指定目录中的所有目录项都搜索完还没有找到相应的目录项，则释放目录项数据块，返回 NULL。
151     brelse(bh);
152     return NULL;
153 }
154
155 /*
156  *      add_entry()
157  *
158  * adds a file entry to the specified directory, using the same
159  * semantics as find_entry(). It returns NULL if it failed.
160  *
161  * NOTE!! The inode part of 'de' is left at 0 - which means you
162  * may not sleep between calling this and putting something into
163  * the entry, as someone else might have used it while you slept.
164  */
    /*
    *      add_entry()
    * 使用与 find_entry() 同样的方法，往指定目录中添加一文件目录项。
    * 如果失败则返回 NULL。
    *
    * 注意！！'de' (指定目录项结构指针) 的 i 节点部分被设置为 0 - 这表示
    * 在调用该函数和往目录项中添加信息之间不能睡眠，因为若睡眠那么其它
    * 人(进程)可能会已经使用了该目录项。
    */
    // 根据指定的目录和文件名添加目录项。
    // 参数: dir - 指定目录的 i 节点; name - 文件名; namelen - 文件名长度;
    // 返回: 高速缓冲区指针; res_dir - 返回的目录项结构指针;
165 static struct buffer_head * add_entry(struct m_inode * dir,
166     const char * name, int namelen, struct dir_entry ** res_dir)
167 {
168     int block, i;
169     struct buffer_head * bh;
170     struct dir_entry * de;
171
172     *res_dir = NULL;
    // 如果定义了 NO_TRUNCATE，则若文件名长度超过最大长度 NAME_LEN，则返回。
173 #ifdef NO_TRUNCATE
174     if (namelen > NAME_LEN)
175         return NULL;
    // 如果没有定义 NO_TRUNCATE，则若文件名长度超过最大长度 NAME_LEN，则截短之。
176 #else
177     if (namelen > NAME_LEN)
178         namelen = NAME_LEN;
179 #endif
    // 如果文件名长度等于 0，则返回 NULL，退出。
180     if (!namelen)
181         return NULL;

```

```

// 如果该目录 i 节点所指向的第一个直接磁盘块号为 0，则返回 NULL 退出。
182     if (!(block = dir->i_zone[0]))
183         return NULL;
// 如果读取该磁盘块失败，则返回 NULL 并退出。
184     if (!(bh = bread(dir->i_dev, block)))
185         return NULL;
// 在目录项数据块中循环查找最后未使用的目录项。首先让目录项结构指针 de 指向高速缓冲的数据块
// 开始处，也即第一个目录项。
186     i = 0;
187     de = (struct dir_entry *) bh->b_data;
188     while (1) {
// 如果当前判别的目录项已经超出当前数据块，则释放该数据块，重新申请一块磁盘块 block。如果
// 申请失败，则返回 NULL，退出。
189         if ((char *)de >= BLOCK_SIZE+bh->b_data) {
190             brelse(bh);
191             bh = NULL;
192             block = create_block(dir, i/DIR_ENTRIES_PER_BLOCK);
193             if (!block)
194                 return NULL;
// 如果读取磁盘块返回的指针为空，则跳过该块继续。
195             if (!(bh = bread(dir->i_dev, block))) {
196                 i += DIR_ENTRIES_PER_BLOCK;
197                 continue;
198             }
// 否则，让目录项结构指针 de 指向该块的高速缓冲数据块开始处。
199             de = (struct dir_entry *) bh->b_data;
200         }
// 如果当前所操作的目录项序号 i*目录结构大小已经超过了该目录所指出的大小 i_size，则说明该第 i
// 个目录项还未使用，我们可以使用它。于是对该目录项进行设置(置该目录项的 i 节点指针为空)。并
// 更新该目录的长度值(加上一个目录项的长度，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。
201         if (i*sizeof(struct dir_entry) >= dir->i_size) {
202             de->inode=0;
203             dir->i_size = (i+1)*sizeof(struct dir_entry);
204             dir->i_dirt = 1;
205             dir->i_ctime = CURRENT_TIME;
206         }
// 若该目录项的 i 节点为空，则表示找到一个还未使用的目录项。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，置相应的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区的指针，退出。
207         if (!de->inode) {
208             dir->i_mtime = CURRENT_TIME;
209             for (i=0; i < NAME_LEN; i++)
210                 de->name[i]=(i<namelen)?get_fs_byte(name+i):0;
211             bh->b_dirt = 1;
212             *res_dir = de;
213             return bh;
214         }
// 如果该目录项已经被使用，则继续检测下一个目录项。
215         de++;
216         i++;
217     }
// 执行不到这里。也许 Linus 在写这段代码时是先复制了上面 find_entry() 的代码，而后修改的。

```

```

218     brelse(bh);
219     return NULL;
220 }
221
222 /*
223  *      get_dir()
224  *
225  * Getdir traverses the pathname until it hits the topmost directory.
226  * It returns NULL on failure.
227  */
228 /*
229  *      get_dir()
230  * 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。
231  * 如果失败则返回 NULL。
232  */
233 // 搜寻指定路径名的目录。
234 // 参数: pathname - 路径名。
235 // 返回: 目录的 i 节点指针。失败时返回 NULL。
236 static struct m_inode * get_dir(const char * pathname)
237 {
238     char c;
239     const char * thisname;
240     struct m_inode * inode;
241     struct buffer head * bh;
242     int namelen, inr, idev;
243     struct dir_entry * de;
244
245     // 如果进程没有设定根 i 节点，或者该进程根 i 节点的引用为 0，则系统出错，死机。
246     if (!current->root || !current->root->i_count)
247         panic("No root inode");
248     // 如果进程的当前工作目录指针为空，或者该当前目录 i 节点的引用计数为 0，也是系统有问题，死机。
249     if (!current->pwd || !current->pwd->i_count)
250         panic("No cwd inode");
251     // 如果用户指定的路径名的第 1 个字符是 '/'，则说明路径名是绝对路径名。则从根 i 节点开始操作。
252     if ((c=get_fs_byte(pathname))== '/') {
253         inode = current->root;
254         pathname++;
255     }
256     // 否则若第一个字符是其它字符，则表示给定的是相对路径名。应从进程的当前工作目录开始操作。
257     // 则取进程当前工作目录的 i 节点。
258     } else if (c)
259         inode = current->pwd;
260     // 则表示路径名为空，出错。返回 NULL，退出。
261     else
262         return NULL;    /* empty name is bad */    /* 空的路径名是错误的 */
263     // 将取得的 i 节点引用计数增 1。
264     inode->i_count++;
265     while (1) {
266         // 若该 i 节点不是目录节点，或者没有可进入的访问许可，则释放该 i 节点，返回 NULL，退出。
267         thisname = pathname;
268         if (!S_ISDIR(inode->i_mode) || !permission(inode, MAY_EXEC)) {
269             iput(inode);
270             return NULL;
271         }
272     }

```

```

// 从路径名开始起搜索检测字符，直到字符已是结尾符(NULL)或者是'/'，此时 namelen 正好是当前处理
// 目录名的长度。如果最后也是一个目录名，但其后没有加'/'，则不会返回该最后目录的 i 节点！
// 比如：/var/log/httpd，将只返回 log/目录的 i 节点。
255         for(namelen=0; (c=get_fs_byte(pathname++))&&(c!='/') ;namelen++)
256             /* nothing */;
// 若字符是结尾符 NULL，则表明已经到达指定目录，则返回该 i 节点指针，退出。
257         if (!c)
258             return inode;
// 调用查找指定目录和文件名的目录项函数，在当前处理目录中寻找子目录项。如果没有找到，则释放
// 该 i 节点，并返回 NULL，退出。
259         if (!(bh = find_entry(&inode, thisname, namelen, &de))) {
260             iput(inode);
261             return NULL;
262         }
// 取该子目录项的 i 节点号 inr 和设备号 idev，释放包含该目录项的高速缓冲块和该 i 节点。
263         inr = de->inode;
264         idev = inode->i_dev;
265         brelse(bh);
266         iput(inode);
// 取节点号 inr 的 i 节点信息，若失败，则返回 NULL，退出。否则继续以该子目录的 i 节点进行操作。
267         if (!(inode = iget(idev, inr)))
268             return NULL;
269     }
270 }
271
272 /*
273  *      dir_namei()
274  *
275  * dir_namei() returns the inode of the directory of the
276  * specified name, and the name within that directory.
277  */
/*
 *      dir_namei()
 * dir_namei() 函数返回指定目录名的 i 节点指针，以及在最顶层目录的名称。
 */
// 参数：pathname - 目录路径名；namelen - 路径名长度。
// 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及其长度。
278 static struct m_inode * dir_namei(const char * pathname,
279     int * namelen, const char ** name)
280 {
281     char c;
282     const char * basename;
283     struct m_inode * dir;
284
// 取指定路径名最顶层目录的 i 节点，若出错则返回 NULL，退出。
285     if (!(dir = get_dir(pathname)))
286         return NULL;
// 对路径名 pathname 进行搜索检测，查处最后一个'/'后面的名字字符串，计算其长度，并返回最顶
// 层目录的 i 节点指针。
287     basename = pathname;
288     while (c=get_fs_byte(pathname++))
289         if (c=='/')
290             basename=pathname;

```

```

291     *namelen = pathname-basename-1;
292     *name = basename;
293     return dir;
294 }
295
296 /*
297  *      namei()
298  *
299  * is used by most simple commands to get the inode of a specified name.
300  * Open, link etc use their own routines, but this is enough for things
301  * like 'chmod' etc.
302  */
303
304 /*
305  *      namei()
306  * 该函数被许多简单的命令用于取得指定路径名称的 i 节点。open、link 等则使用它们
307  * 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。
308  */
309
310 // 取指定路径名的 i 节点。
311 // 参数: pathname - 路径名。
312 // 返回: 对应的 i 节点。
313
314 struct m\_inode * namei(const char * pathname)
315 {
316     const char * basename;
317     int inr, dev, namelen;
318     struct m\_inode * dir;
319     struct buffer\_head * bh;
320     struct dir\_entry * de;
321
322     // 首先查找指定路径的最顶层目录的目录名及其 i 节点，若不存在，则返回 NULL，退出。
323     if (!(dir = dir\_namei(pathname, &namelen, &basename)))
324         return NULL;
325     // 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。
326     if (!namelen)
327         /* special case: '/usr/' etc */
328         return dir;
329     /* 对应于'/usr/'等情况 */
330
331     // 在返回的顶层目录中寻找指定文件名的目录项的 i 节点。因为如果最后也是一个目录名，但其后没
332     // 有加'/'，则不会返回该最后目录的 i 节点！比如：/var/log/httpd，将只返回 log/目录的 i 节点。
333     // 因此 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
334     // 情况使用寻找目录项 i 节点函数 find_entry() 进行处理。
335     bh = find\_entry(&dir, basename, namelen, &de);
336     if (!bh) {
337         iput(dir);
338         return NULL;
339     }
340     // 取该目录项的 i 节点号和目录的设备号，并释放包含该目录项的高速缓冲区以及目录 i 节点。
341     inr = de->inode;
342     dev = dir->i_dev;
343     brelse(bh);
344     iput(dir);
345     // 取对应节号的 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。
346     dir = iget(dev, inr);
347     if (dir) {
348         dir->i_atime = CURRENT\_TIME;
349         dir->i_dirt = 1;

```

```

328     }
329     return dir;
330 }
331
332 /*
333  *      open_namei()
334  *
335  * namei for open - this is in fact almost the whole open-routine.
336 */
337 /*
338  *      open_namei()
339  * open() 所使用的 namei 函数 - 这其实几乎是完整的打开文件程序。
340 */
341 // 文件打开 namei 函数。
342 // 参数: pathname - 文件路径名; flag - 文件打开标志; mode - 文件访问许可属性;
343 // 返回: 成功返回 0, 否则返回出错码; res_inode - 返回的对应文件路径名的 i 节点指针。
344 int open_namei(const char * pathname, int flag, int mode,
345                struct m\_inode ** res_inode)
346 {
347     const char * basename;
348     int inr, dev, namelen;
349     struct m\_inode * dir, * inode;
350     struct buffer head * bh;
351     struct dir\_entry * de;
352
353     // 如果文件访问许可模式标志是只读(0), 但文件截 0 标志 O_TRUNC 却置位了, 则改为只写标志。
354     if ((flag & O\_TRUNC) && !(flag & O\_ACCMODE))
355         flag |= O\_WRONLY;
356     // 使用进程的文件访问许可屏蔽码, 屏蔽掉给定模式中的相应位, 并添上普通文件标志。
357     mode &= 0777 & ~current->umask;
358     mode |= I\_REGULAR;
359     // 根据路径名寻找到对应的 i 节点, 以及最顶端文件名及其长度。
360     if (!(dir = dir\_namei(pathname, &namelen, &basename)))
361         return -ENOENT;
362     // 如果最顶端文件名长度为 0(例如 '/usr/' 这种路径名的情况), 那么若打开操作不是创建、截 0,
363     // 则表示打开一个目录名, 直接返回该目录的 i 节点, 并退出。
364     if (!namelen) {
365         /* special case: '/usr/' etc */
366         if (!(flag & (O\_ACCMODE | O\_CREAT | O\_TRUNC))) {
367             *res_inode = dir;
368             return 0;
369         }
370     }
371     // 否则释放该 i 节点, 返回出错码。
372     iput(dir);
373     return -EISDIR;
374 }
375 // 在 dir 节点对应的目录中取文件名对应的目录项结构 de 和该目录项所在的高速缓冲区。
376 bh = find\_entry(&dir, basename, namelen, &de);
377 // 如果该高速缓冲指针为 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。
378 if (!bh) {
379     // 如果不是创建文件, 则释放该目录的 i 节点, 返回出错号退出。
380     if (!(flag & O\_CREAT)) {
381         iput(dir);
382         return -ENOENT;

```



```

365     }
    // 如果用户在该目录没有写的权力, 则释放该目录的 i 节点, 返回出错号退出。
366     if (!permission(dir, MAY\_WRITE)) {
367         iput(dir);
368         return -EACCES;
369     }
    // 在目录节点对应的设备上申请一个新 i 节点, 若失败, 则释放目录的 i 节点, 并返回没有空间出错码。
370     inode = new\_inode(dir->i_dev);
371     if (!inode) {
372         iput(dir);
373         return -ENOSPC;
374     }
    // 否则使用该新 i 节点, 对其进行初始设置: 置节点的用户 id; 对应节点访问模式; 置已修改标志。
375     inode->i_uid = current->euid;
376     inode->i_mode = mode;
377     inode->i_dirt = 1;
    // 然后在指定目录 dir 中添加一新目录项。
378     bh = add\_entry(dir, basename, namelen, &de);
    // 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。于是将该
    // 新 i 节点的引用连接计数减 1; 并释放该 i 节点与目录的 i 节点, 返回出错码, 退出。
379     if (!bh) {
380         inode->i_nlinks--;
381         iput(inode);
382         iput(dir);
383         return -ENOSPC;
384     }
    // 初始设置该新目录项: 置 i 节点号为新申请到的 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, 释放目录的 i 节点。返回新目录项的 i 节点指针, 退出。
385     de->inode = inode->i_num;
386     bh->b_dirt = 1;
387     brelse(bh);
388     iput(dir);
389     *res_inode = inode;
390     return 0;
391 }
    // 若上面在目录中取文件名对应的目录项结构操作成功(也即 bh 不为 NULL), 取出该目录项的 i 节点号
    // 和其所在的设备号, 并释放该高速缓冲区以及目录的 i 节点。
392     inr = de->inode;
393     dev = dir->i_dev;
394     brelse(bh);
395     iput(dir);
    // 如果独占使用标志 O_EXCL 置位, 则返回文件已存在出错码, 退出。
396     if (flag & O\_EXCL)
397         return -EEXIST;
    // 如果取该目录项对应 i 节点的操作失败, 则返回访问出错码, 退出。
398     if (!(inode=iget(dev, inr)))
399         return -EACCES;
    // 若该 i 节点是一个目录的节点并且访问模式是只读或只写, 或者没有访问的许可权限, 则释放该 i
    // 节点, 返回访问权限出错码, 退出。
400     if ((S\_ISDIR(inode->i_mode) && (flag & O\_ACCMODE)) ||
401         !permission(inode, ACC\_MODE(flag))) {
402         iput(inode);
403         return -EPERM;

```



```

404     }
// 更新该 i 节点的访问时间字段为当前时间。
405     inode->i_atime = CURRENT_TIME;
// 如果设立了截 0 标志, 则将该 i 节点的文件长度截为 0。
406     if (flag & O_TRUNC)
407         truncate(inode);
// 最后返回该目录项 i 节点的指针, 并返回 0 (成功)。
408     *res_inode = inode;
409     return 0;
410 }
411
///// 系统调用函数 - 创建一个特殊文件或普通文件节点(node)。
// 创建名称为 filename, 由 mode 和 dev 指定的文件系统节点(普通文件、设备特殊文件或命名管道)。
// 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。
// 返回: 成功则返回 0, 否则返回出错码。
412 int sys_mknod(const char * filename, int mode, int dev)
413 {
414     const char * basename;
415     int namelen;
416     struct m_inode * dir, * inode;
417     struct buffer_head * bh;
418     struct dir_entry * de;
419
// 如果不是超级用户, 则返回访问许可出错码。
420     if (!suser())
421         return -EPERM;
// 如果找不到对应路径名目录的 i 节点, 则返回出错码。
422     if (!(dir = dir_namei(filename, &namelen, &basename)))
423         return -ENOENT;
// 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, 释放该目录 i 节点, 返回
// 出错码, 退出。
424     if (!namelen) {
425         iput(dir);
426         return -ENOENT;
427     }
// 如果在该目录中没有写的权限, 则释放该目录的 i 节点, 返回访问许可出错码, 退出。
428     if (!permission(dir, MAY_WRITE)) {
429         iput(dir);
430         return -EPERM;
431     }
// 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高速缓冲区, 释放目录
// 的 i 节点, 返回文件已经存在出错码, 退出。
432     bh = find_entry(&dir, basename, namelen, &de);
433     if (bh) {
434         brelse(bh);
435         iput(dir);
436         return -EEXIST;
437     }
// 申请一个新的 i 节点, 如果不成功, 则释放目录的 i 节点, 返回无空间出错码, 退出。
438     inode = new_inode(dir->i_dev);
439     if (!inode) {
440         iput(dir);
441         return -ENOSPC;

```

```

442     }
    // 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 节点的直接块
    // 指针 0 等于设备号。
443     inode->i_mode = mode;
444     if (S_ISBLK(mode) || S_ISCHR(mode))
445         inode->i_zone[0] = dev;
    // 设置该 i 节点的修改时间、访问时间为当前时间。
446     inode->i_mtime = inode->i_atime = CURRENT_TIME;
447     inode->i_dirt = 1;
    // 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区指针为 NULL)，则释放目录的
    // i 节点；所申请的 i 节点引用连接计数复位，并释放该 i 节点。返回出错码，退出。
448     bh = add_entry(dir, basename, namelen, &de);
449     if (!bh) {
450         iput(dir);
451         inode->i_nlinks=0;
452         iput(inode);
453         return -ENOSPC;
454     }
    // 令该目录项的 i 节点字段等于新 i 节点号，置高速缓冲区已修改标志，释放目录和新的 i 节点，释放
    // 高速缓冲区，最后返回 0(成功)。
455     de->inode = inode->i_num;
456     bh->b_dirt = 1;
457     iput(dir);
458     iput(inode);
459     brelse(bh);
460     return 0;
461 }
462
    ///// 系统调用函数 - 创建目录。
    // 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。
463 int sys_mkdir(const char * pathname, int mode)
464 {
465     const char * basename;
466     int namelen;
467     struct m_inode * dir, * inode;
468     struct buffer_head * bh, *dir_block;
469     struct dir_entry * de;
470
    // 如果不是超级用户，则返回访问许可出错码。
471     if (!suser())
472         return -EPERM;
    // 如果找不到对应路径名目录的 i 节点，则返回出错码。
473     if (!(dir = dir_namei(pathname, &namelen, &basename)))
474         return -ENOENT;
    // 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，释放该目录 i 节点，返回
    // 出错码，退出。
475     if (!namelen) {
476         iput(dir);
477         return -ENOENT;
478     }
    // 如果在该目录中没有写的权限，则释放该目录的 i 节点，返回访问许可出错码，退出。
479     if (!permission(dir, MAY_WRITE)) {

```

```

480         iput(dir);
481         return -EPERM;
482     }
    // 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高速缓冲区，释放目录
    // 的 i 节点，返回文件已经存在出错码，退出。
483     bh = find_entry(&dir, basename, namelen, &de);
484     if (bh) {
485         brelse(bh);
486         iput(dir);
487         return -EEXIST;
488     }
    // 申请一个新的 i 节点，如果不成功，则释放目录的 i 节点，返回无空间出错码，退出。
489     inode = new_inode(dir->i_dev);
490     if (!inode) {
491         iput(dir);
492         return -ENOSPC;
493     }
    // 置该新 i 节点对应的文件长度为 32(一个目录项的大小)，置节点已修改标志，以及节点的修改时间
    // 和访问时间。
494     inode->i_size = 32;
495     inode->i_dirt = 1;
496     inode->i_mtime = inode->i_atime = CURRENT_TIME;
    // 为该 i 节点申请一磁盘块，并令节点第一个直接块指针等于该块号。如果申请失败，则释放对应目录
    // 的 i 节点；复位新申请的 i 节点连接计数；释放该新的 i 节点，返回没有空间出错码，退出。
497     if (!(inode->i_zone[0]=new_block(inode->i_dev))) {
498         iput(dir);
499         inode->i_nlinks--;
500         iput(inode);
501         return -ENOSPC;
502     }
    // 置该新的 i 节点已修改标志。
503     inode->i_dirt = 1;
    // 读新申请的磁盘块。若出错，则释放对应目录的 i 节点；释放申请的磁盘块；复位新申请的 i 节点
    // 连接计数；释放该新的 i 节点，返回没有空间出错码，退出。
504     if (!(dir_block=bread(inode->i_dev, inode->i_zone[0]))) {
505         iput(dir);
506         free_block(inode->i_dev, inode->i_zone[0]);
507         inode->i_nlinks--;
508         iput(inode);
509         return -ERROR;
510     }
    // 令 de 指向目录项数据块，置该目录项的 i 节点号字段等于新申请的 i 节点号，名字字段等于“.”。
511     de = (struct dir_entry *) dir_block->b_data;
512     de->inode=inode->i_num;
513     strcpy(de->name, ".");
    // 然后 de 指向下一个目录项结构，该结构用于存放上级目录的节点号和名字“..”。
514     de++;
515     de->inode = dir->i_num;
516     strcpy(de->name, "..");
517     inode->i_nlinks = 2;
    // 然后设置该高速缓冲区已修改标志，并释放该缓冲区。
518     dir_block->b_dirt = 1;
519     brelse(dir_block);

```

```

// 初始化设置新 i 节点的模式字段，并置该 i 节点已修改标志。
520     inode->i_mode = I_DIRECTORY | (mode & 0777 & ~current->umask);
521     inode->i_dirt = 1;
// 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区指针为 NULL)，则释放目录的
// i 节点；所申请的 i 节点引用连接计数复位，并释放该 i 节点。返回出错码，退出。
522     bh = add_entry(dir, basename, namelen, &de);
523     if (!bh) {
524         iput(dir);
525         free_block(inode->i_dev, inode->i_zone[0]);
526         inode->i_nlinks=0;
527         iput(inode);
528         return -ENOSPC;
529     }
// 令该目录项的 i 节点字段等于新 i 节点号，置高速缓冲区已修改标志，释放目录和新的 i 节点，释放
// 高速缓冲区，最后返回 0(成功)。
530     de->inode = inode->i_num;
531     bh->b_dirt = 1;
532     dir->i_nlinks++;
533     dir->i_dirt = 1;
534     iput(dir);
535     iput(inode);
536     brelse(bh);
537     return 0;
538 }
539
540 /*
541  * routine to check that the specified directory is empty (for rmdir)
542  */
543 /*
544  * 用于检查指定的目录是否为空的子程序(用于 rmdir 系统调用函数)。
545  */
546 // 检查指定目录是否是空的。
547 // 参数: inode - 指定目录的 i 节点指针。
548 // 返回: 0 - 是空的; 1 - 不空。
549 static int empty_dir(struct m_inode * inode)
550 {
551     int nr, block;
552     int len;
553     struct buffer_head * bh;
554     struct dir_entry * de;
555
556     // 计算指定目录中现有目录项的个数(应该起码有 2 个，即“.”和“..”两个文件目录项)。
557     len = inode->i_size / sizeof (struct dir_entry);
558     // 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，或者相应磁盘
559     // 块读不出，则显示警告信息“设备 dev 上目录错”，返回 0(失败)。
560     if (len<2 || !inode->i_zone[0] ||
561         !(bh=bread(inode->i_dev, inode->i_zone[0]))) {
562         printk("warning - bad directory on dev %04x\n", inode->i_dev);
563         return 0;
564     }
565     // 让 de 指向含有读出磁盘块数据的高速缓冲区中第 1 项目录项。
566     de = (struct dir_entry *) bh->b_data;
567     // 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。
557     if (de[0].inode != inode->i_num || !de[1].inode ||
558         strcmp(".",de[0].name) || strcmp("..",de[1].name)) {
559         printk("warning - bad directory on dev %04x\n",inode->i_dev);
560         return 0;
561     }
// 令 nr 等于目录项序号; de 指向第三个目录项。
562     nr = 2;
563     de += 2;
// 循环检测该目录中所有的目录项(len-2 个),看有没有目录项的 i 节点号字段不为 0(被使用)。
564     while (nr<len) {
// 如果该块磁盘块中的目录项已经检测完,则释放该磁盘块的高速缓冲区,读取下一块含有目录项的
// 磁盘块。若相应块没有使用(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。
565         if ((void *) de >= (void *) (bh->b_data+BLOCK_SIZE)) {
566             brelse(bh);
567             block=bmap(inode,nr/DIR_ENTRIES_PER_BLOCK);
568             if (!block) {
569                 nr += DIR_ENTRIES_PER_BLOCK;
570                 continue;
571             }
572             if (!(bh=bread(inode->i_dev,block)))
573                 return 0;
574             de = (struct dir_entry *) bh->b_data;
575         }
// 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。
576         if (de->inode) {
577             brelse(bh);
578             return 0;
579         }
// 否则,若还没有查询完该目录中的所有目录项,则继续检测。
580         de++;
581         nr++;
582     }
// 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),则返回缓冲区,返回 1。
583     brelse(bh);
584     return 1;
585 }
586
///// 系统调用函数 - 删除指定名称的目录。
// 参数: name - 目录名(路径名)。
// 返回: 返回 0 表示成功,否则返回出错号。
587 int sys_rmdir(const char * name)
588 {
589     const char * basename;
590     int namelen;
591     struct m_inode * dir, * inode;
592     struct buffer_head * bh;
593     struct dir_entry * de;
594
// 如果不是超级用户,则返回访问许可出错码。

```

```

595         if (!suser())
596             return -EPERM;
// 如果找不到对应路径名目录的 i 节点，则返回出错码。
597         if (!(dir = dir_namei(name, &namelen, &basename)))
598             return -ENOENT;
// 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，释放该目录 i 节点，返回
// 出错码，退出。
599         if (!namelen) {
600             iput(dir);
601             return -ENOENT;
602         }
// 如果在该目录中没有写的权限，则释放该目录的 i 节点，返回访问许可出错码，退出。
603         if (!permission(dir, MAY_WRITE)) {
604             iput(dir);
605             return -EPERM;
606         }
// 如果对应路径名上最后的文件名的目录项不存在，则释放包含该目录项的高速缓冲区，释放目录
// 的 i 节点，返回文件已经存在出错码，退出。否则 dir 是包含要被删除目录名的目录 i 节点，de
// 是要被删除目录的目录项结构。
607         bh = find_entry(&dir, basename, namelen, &de);
608         if (!bh) {
609             iput(dir);
610             return -ENOENT;
611         }
// 取该目录项指明的 i 节点。若出错则释放目录的 i 节点，并释放含有目录项的高速缓冲区，返回
// 出错号。
612         if (!(inode = iget(dir->i_dev, de->inode))) {
613             iput(dir);
614             brelse(bh);
615             return -EPERM;
616         }
// 若该目录设置了受限删除标志并且进程的有效用户 id 不等于该 i 节点的用户 id，则表示没有权限删
// 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，释放高速缓冲区，返
// 回出错码。
617         if ((dir->i_mode & S_ISVTX) && current->euid &&
618             inode->i_uid != current->euid) {
619             iput(dir);
620             iput(inode);
621             brelse(bh);
622             return -EPERM;
623         }
// 如果要被删除的目录项的 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，释放高速缓冲区，返回出错码。
624         if (inode->i_dev != dir->i_dev || inode->i_count > 1) {
625             iput(dir);
626             iput(inode);
627             brelse(bh);
628             return -EPERM;
629         }
// 如果要被删除目录的目录项 i 节点的节点号等于包含该需删除目录的 i 节点号，则表示试图删除“.”
// 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，释放高速缓冲区，返回
// 出错码。

```

```

630         if (inode == dir) {           /* we may not delete ".", but "../dir" is ok */
631             iput(inode);               /* 我们不可以删除".", 但可以删除 "../dir" */
632             iput(dir);
633             brelse(bh);
634             return -EPERM;
635         }
        // 若要被删除的目录的 i 节点的属性表明这不是一个目录, 则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, 释放高速缓冲区, 返回出错码。
636         if (!S_ISDIR(inode->i_mode)) {
637             iput(inode);
638             iput(dir);
639             brelse(bh);
640             return -ENOTDIR;
641         }
        // 若该需被删除的目录不空, 则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
        // 高速缓冲区, 返回出错码。
642         if (!empty_dir(inode)) {
643             iput(inode);
644             iput(dir);
645             brelse(bh);
646             return -ENOTEMPTY;
647         }
        // 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。
648         if (inode->i_nlinks != 2)
649             printk("empty directory has nlink!=2 (%d)", inode->i_nlinks);
        // 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
        // 缓冲区已修改标志, 并释放该缓冲区。
650         de->inode = 0;
651         bh->b_dirt = 1;
652         brelse(bh);
        // 置被删除目录的 i 节点的连接数为 0, 并置 i 节点已修改标志。
653         inode->i_nlinks=0;
654         inode->i_dirt=1;
        // 将包含被删除目录名的目录的 i 节点引用计数减 1, 修改其改变时间和修改时间为当前时间, 并置
        // 该节点已修改标志。
655         dir->i_nlinks--;
656         dir->i_ctime = dir->i_mtime = CURRENT_TIME;
657         dir->i_dirt=1;
        // 最后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, 返回 0(成功)。
658         iput(dir);
659         iput(inode);
660         return 0;
661     }
662
        //// 系统调用函数 - 删除文件名以及可能也删除其相关的文件。
        // 从文件系统删除一个名字。如果是一个文件的最后一个连接, 并且没有进程正打开该文件, 则该文件
        // 也将被删除, 并释放所占用的设备空间。
        // 参数: name - 文件名。
        // 返回: 成功则返回 0, 否则返回出错号。
663 int sys_unlink(const char * name)
664 {
665     const char * basename;
666     int namelen;

```

```

667     struct m\_inode * dir, * inode;
668     struct buffer head * bh;
669     struct dir\_entry * de;
670
671     // 如果找不到对应路径名目录的 i 节点，则返回出错码。
672     if (!(dir = dir\_namei(name, &namelen, &basename)))
673         return -ENOENT;
674
675     // 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，释放该目录 i 节点，返回
676     // 出错码，退出。
677     if (!namelen) {
678         iput(dir);
679         return -ENOENT;
680     }
681
682     // 如果在该目录中没有写的权限，则释放该目录的 i 节点，返回访问许可出错码，退出。
683     if (!permission(dir, MAY\_WRITE)) {
684         iput(dir);
685         return -EPERM;
686     }
687
688     // 如果对应路径名上最后的文件名的目录项不存在，则释放包含该目录项的高速缓冲区，释放目录
689     // 的 i 节点，返回文件已经存在出错码，退出。否则 dir 是包含要被删除目录名的目录 i 节点，de
690     // 是要被删除目录的目录项结构。
691     bh = find\_entry(&dir, basename, namelen, &de);
692     if (!bh) {
693         iput(dir);
694         return -ENOENT;
695     }
696
697     // 取该目录项指明的 i 节点。若出错则释放目录的 i 节点，并释放含有目录项的高速缓冲区，返回
698     // 出错号。
699     if (!(inode = iget(dir->i_dev, de->inode))) {
700         iput(dir);
701         brelse(bh);
702         return -ENOENT;
703     }
704
705     // 如果该目录设置了受限删除标志并且用户不是超级用户，并且进程的有效用户 id 不等于被删除文件
706     // 名 i 节点的用户 id，并且进程的有效用户 id 也不等于目录 i 节点的用户 id，则没有权限删除该文件
707     // 名。则释放该目录 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区，返回出错号。
708     if ((dir->i_mode & S\_ISVTX) && !suser() &&
709         current->euid != inode->i_uid &&
710         current->euid != dir->i_uid) {
711         iput(dir);
712         iput(inode);
713         brelse(bh);
714         return -EPERM;
715     }
716
717     // 如果该指定文件名是一个目录，则也不能删除，释放该目录 i 节点和该文件名目录项的 i 节点，释放
718     // 包含该目录项的缓冲区，返回出错号。
719     if (S\_ISDIR(inode->i_mode)) {
720         iput(inode);
721         iput(dir);
722         brelse(bh);
723         return -EPERM;
724     }
725
726     // 如果该 i 节点的连接数已经为 0，则显示警告信息，修正其为 1。

```



```

705     if (!inode->i_nlinks) {
706         printk("Deleting nonexistent file (%04x:%d), %d\n",
707             inode->i_dev, inode->i_num, inode->i_nlinks);
708         inode->i_nlinks=1;
709     }
    // 将该文件名的目录项中的 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区
    // 已修改标志，释放该高速缓冲区。
710     de->inode = 0;
711     bh->b_dirt = 1;
712     brelse(bh);
    // 该 i 节点的连接数减 1，置已修改标志，更新改变时间为当前时间。最后释放该 i 节点和目录的 i 节
    // 点，返回 0(成功)。
713     inode->i_nlinks--;
714     inode->i_dirt = 1;
715     inode->i_ctime = CURRENT_TIME;
716     iput(inode);
717     iput(dir);
718     return 0;
719 }
720
    //// 系统调用函数 - 为文件建立一个文件名。
    // 为一个已经存在的文件创建一个新连接(也称为硬连接 - 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。
721 int sys_link(const char * oldname, const char * newname)
722 {
723     struct dir_entry * de;
724     struct m_inode * oldinode, * dir;
725     struct buffer_head * bh;
726     const char * basename;
727     int namelen;
728
    // 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。
729     oldinode=namei(oldname);
730     if (!oldinode)
731         return -ENOENT;
    // 如果原路径名对应的是一个目录名，则释放该 i 节点，返回出错号。
732     if (S_ISDIR(oldinode->i_mode)) {
733         iput(oldinode);
734         return -EPERM;
735     }
    // 查找新路径名的最顶层目录的 i 节点，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，
    // 则释放原路径名的 i 节点，返回出错号。
736     dir = dir_namei(newname, &namelen, &basename);
737     if (!dir) {
738         iput(oldinode);
739         return -EACCES;
740     }
    // 如果新路径名中不包括文件名，则释放原路径名 i 节点和新路径名目录的 i 节点，返回出错号。
741     if (!namelen) {
742         iput(oldinode);
743         iput(dir);
744         return -EPERM;

```

```

745     }
    // 如果新路径名目录的设备号与原路径名的设备号不一样, 则也不能建立连接, 于是释放新路径名
    // 目录的 i 节点和原路径名的 i 节点, 返回出错号。
746     if (dir->i_dev != oldinode->i_dev) {
747         iput(dir);
748         iput(oldinode);
749         return -EXDEV;
750     }
    // 如果用户没有在新目录中写的权限, 则也不能建立连接, 于是释放新路径名目录的 i 节点和原路径名
    // 的 i 节点, 返回出错号。
751     if (!permission(dir, MAY_WRITE)) {
752         iput(dir);
753         iput(oldinode);
754         return -EACCES;
755     }
    // 查询该新路径名是否已经存在, 如果存在, 则也不能建立连接, 于是释放包含该已存在目录项的高速
    // 缓冲区, 释放新路径名目录的 i 节点和原路径名的 i 节点, 返回出错号。
756     bh = find_entry(&dir, basename, namelen, &de);
757     if (bh) {
758         brelse(bh);
759         iput(dir);
760         iput(oldinode);
761         return -EEXIST;
762     }
    // 在新目录中添加一个目录项。若失败则释放该目录的 i 节点和原路径名的 i 节点, 返回出错号。
763     bh = add_entry(dir, basename, namelen, &de);
764     if (!bh) {
765         iput(dir);
766         iput(oldinode);
767         return -ENOSPC;
768     }
    // 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高速缓冲区
    // 已修改标志, 释放该缓冲区, 释放目录的 i 节点。
769     de->inode = oldinode->i_num;
770     bh->b_dirt = 1;
771     brelse(bh);
772     iput(dir);
    // 将原节点的应用计数加 1, 修改其改变时间为当前时间, 并设置 i 节点已修改标志, 最后释放原
    // 路径名的 i 节点, 并返回 0(成功)。
773     oldinode->i_nlinks++;
774     oldinode->i_ctime = CURRENT_TIME;
775     oldinode->i_dirt = 1;
776     iput(oldinode);
777     return 0;
778 }
779

```

## 8.18 char\_dev.c 文件

### 8.18.1 功能描述

### 8.18.2 代码注释

列表 linux/fs/char\_dev.c 程序

```

1  /*
2   *  linux/fs/char_dev.c
3   *
4   *  (C) 1991  Linus Torvalds
5   */
6
7  #include <errno.h>          // 错误号头文件。包含系统中各种出错号。(Linus 从 minix 中引进的)。
8  #include <sys/types.h>     // 类型头文件。定义了基本的系统数据类型。
9
10 #include <linux/sched.h>    // 调度程序头文件，定义了任务结构 task_struct、初始任务 0 的数据，
    // 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
11 #include <linux/kernel.h>  // 内核头文件。含有一些内核常用函数的原形定义。
12
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 <asm/io.h>        // io 头文件。定义硬件端口输入/输出宏汇编语句。
15
16 extern int tty_read(unsigned minor, char * buf, int count);    // 终端读。
17 extern int tty_write(unsigned minor, char * buf, int count);    // 终端写。
18
    // 定义字符设备读写函数指针类型。
19 typedef (*crw_ptr)(int rw, unsigned minor, char * buf, int count, off_t * pos);
20
    // 终端读写操作函数。
    // 参数: rw - 读写命令; minor - 终端子设备号; buf - 缓冲区; cout - 读写字节数;
    //      pos - 读写操作当前指针，对于终端操作，该指针无用。
    // 返回: 实际读写的字节数。
21 static int rw_ttyx(int rw, unsigned minor, char * buf, int count, off_t * pos)
22 {
23     return ((rw==READ)?tty_read(minor, buf, count):
24             tty_write(minor, buf, count));
25 }
26
    // 终端读写操作函数。
    // 同上 rw_ttyx(), 只是增加了对进程是否有控制终端的检测。
27 static int rw_tty(int rw, unsigned minor, char * buf, int count, off_t * pos)
28 {
    // 若进程没有对应的控制终端，则返回出错号。
29     if (current->tty<0)
30         return -EPERM;
    // 否则调用终端读写函数 rw_ttyx(), 并返回实际读写字节数。
31     return rw_ttyx(rw, current->tty, buf, count, pos);
32 }
33

```

```

    /// 内存数据读写。未实现。
34 static int rw\_ram(int rw, char * buf, int count, off\_t * pos)
35 {
36     return -EIO;
37 }
38
    /// 内存数据读写操作函数。未实现。
39 static int rw\_mem(int rw, char * buf, int count, off\_t * pos)
40 {
41     return -EIO;
42 }
43
    /// 内核数据区读写函数。未实现。
44 static int rw\_kmem(int rw, char * buf, int count, off\_t * pos)
45 {
46     return -EIO;
47 }
48
    // 端口读写操作函数。
    // 参数: rw - 读写命令; buf - 缓冲区; cout - 读写字节数; pos - 端口地址。
    // 返回: 实际读写的字节数。
49 static int rw\_port(int rw, char * buf, int count, off\_t * pos)
50 {
51     int i = *pos;
52
    // 对于所要求读写的字节数, 并且端口地址小于 64k 时, 循环执行单个字节的读写操作。
53     while (count-->0 && i<65536) {
    // 若是读命令, 则从端口 i 中读取一字节内容并放到用户缓冲区中。
54         if (rw==READ)
55             put\_fs\_byte(inb(i), buf++);
    // 若是写命令, 则从用户数据缓冲区中取一字节输出到端口 i。
56         else
57             outb(get\_fs\_byte(buf++), i);
    // 前移一个端口。[??]
58         i++;
59     }
    // 计算读/写的字节数, 并相应调整读写指针。
60     i -= *pos;
61     *pos += i;
    // 返回读/写的字节数。
62     return i;
63 }
64
    /// 内存读写操作函数。
65 static int rw\_memory(int rw, unsigned minor, char * buf, int count, off\_t * pos)
66 {
    // 根据内存设备子设备号, 分别调用不同的内存读写函数。
67     switch(minor) {
68         case 0:
69             return rw\_ram(rw, buf, count, pos);
70         case 1:
71             return rw\_mem(rw, buf, count, pos);
72         case 2:

```

```

73         return rw\_kmem(rw, buf, count, pos);
74     case 3:
75         return (rw==READ)?0:count;        /* rw\_null */
76     case 4:
77         return rw\_port(rw, buf, count, pos);
78     default:
79         return -EIO;
80     }
81 }
82
83 // 定义系统中设备种数。
84 #define NRDEVS ((sizeof (crw\_table))/(sizeof (crw\_ptr)))
85
86 // 字符设备读写函数指针表。
87 static crw\_ptr crw\_table[]={
88     NULL,                /* nodev */                /* 无设备(空设备) */
89     rw\_memory,           /* /dev/mem etc */          /* /dev/mem 等 */
90     NULL,                 /* /dev/fd */              /* /dev/fd 软驱 */
91     NULL,                 /* /dev/hd */              /* /dev/hd 硬盘 */
92     rw\_ttyx,              /* /dev/ttyx */            /* /dev/ttyx 串口终端 */
93     rw\_tty,               /* /dev/tty */             /* /dev/tty 终端 */
94     NULL,                 /* /dev/lp */              /* /dev/lp 打印机 */
95     NULL};                /* unnamed pipes */      /* 未命名管道 */
96
97 // 字符设备读写操作函数。
98 // 参数: rw - 读写命令; dev - 设备号; buf - 缓冲区; count - 读写字节数; pos - 读写指针。
99 // 返回: 实际读/写字节数。
100 int rw\_char(int rw,int dev, char * buf, int count, off\_t * pos)
101 {
102     crw\_ptr call_addr;
103
104     // 如果设备号超出系统设备数, 则返回出错码。
105     if (MAJOR(dev)>=NRDEVS)
106         return -ENODEV;
107     // 若该设备没有对应的读/写函数, 则返回出错码。
108     if (!(call_addr=crw\_table[MAJOR(dev)]))
109         return -ENODEV;
110     // 调用对应设备的读写操作函数, 并返回实际读/写的字节数。
111     return call_addr(rw,MINOR(dev),buf,count,pos);
112 }
113

```

## 8.19 ioctl.c 文件

### 8.19.1 功能描述

## 8.19.2 代码注释

列表 linux/fs/ioctl.c 程序

```

1  /*
2  *  linux/fs/ioctl.c
3  *
4  *  (C) 1991  Linus Torvalds
5  */
6
7  #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。
8  #include <errno.h>          // 错误号头文件。包含系统中各种出错号。(Linus 从 minix 中引进的)。
9  #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。
10
11 #include <linux/sched.h>    // 调度程序头文件，定义了任务结构 task_struct、初始任务 0 的数据，
    // 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
12
13 extern int tty\_ioctl(int dev, int cmd, int arg); // 终端 ioctl(chr_drv/tty_ioctl.c, 115)。
14
    // 定义输入输出控制(ioctl)函数指针。
15 typedef int (\*ioctl\_ptr)(int dev, int cmd, int arg);
16
    // 定义系统中设备种数。
17 #define NRDEVS ((sizeof (ioctl\_table))/(sizeof (ioctl\_ptr)))
18
    // ioctl 操作函数指针表。
19 static ioctl\_ptr ioctl\_table[]={
20     NULL,          /* nodev */
21     NULL,          /* /dev/mem */
22     NULL,          /* /dev/fd */
23     NULL,          /* /dev/hd */
24     tty\_ioctl,      /* /dev/ttyx */
25     tty\_ioctl,      /* /dev/tty */
26     NULL,          /* /dev/lp */
27     NULL};          /* named pipes */
28
29
    //// 系统调用函数 - 输入输出控制函数。
    // 参数: fd - 文件描述符; cmd - 命令码; arg - 参数。
    // 返回: 成功则返回 0, 否则返回出错码。
30 int sys\_ioctl(unsigned int fd, unsigned int cmd, unsigned long arg)
31 {
32     struct file * filp;
33     int dev, mode;
34
    // 如果文件描述符超出可打开的文件数，或者对应描述符的文件结构指针为空，则返回出错码，退出。
35     if (fd >= NR\_OPEN || !(filp = current->filp[fd]))
36         return -EBADF;
    // 取对应文件的属性。如果该文件不是字符文件，也不是块设备文件，则返回出错码，退出。
37     mode=filp->f_inode->i_mode;
38     if (!S\_ISCHR(mode) && !S\_ISBLK(mode))
39         return -EINVAL;
    // 从字符或块设备文件的 i 节点中取设备号。如果设备号大于系统现有的设备数，则返回出错号。
40     dev = filp->f_inode->i_zone[0];

```

---

```
41         if (MAJOR(dev) >= NRDEVS)
42             return -ENODEV;
43         // 如果该设备在 ioctl 函数指针表中没有对应函数，则返回出错码。
44         if (!ioctl_table[MAJOR(dev)])
45             return -ENOTTY;
46         // 否则返回实际 ioctl 函数返回码，成功则返回 0，否则返回出错码。
47         return ioctl_table[MAJOR(dev)](dev, cmd, arg);
48     }
```

---








# 第 9 章 内存管理(mm)

## 9.1 概述

内存管理目录中共有三个程序，如下表所示：

|   | Name                     | Size        | Last modified (GMT) | Description |
|---|--------------------------|-------------|---------------------|-------------|
|  | <a href="#">Makefile</a> | 813 bytes   | 1991-12-02 03:21:45 | m           |
|  | <a href="#">memory.c</a> | 11223 bytes | 1991-12-03 00:48:01 | m           |
|  | <a href="#">page.s</a>   | 508 bytes   | 1991-10-02 14:16:30 | m           |

下面对每一个文件分别描述。

## 9.2 Makefile 文件

### 9.2.1 功能描述

### 9.2.2 代码注释

列表 linux/mm/Makefile 文件

|   |        |   |  |
|---|--------|---|--|
| 1 | CC     | =gcc  | # GNU C 语言编译器。   |
| 2 | CFLAGS | =-O -Wall -fstrength-reduce -fcombine-regs -fomit-frame-pointer \ |  |
| 3 |        | -finline-functions -nostdinc -I../include                         |  |
|   |        |   | # C 编译程序选项。-Wall 显示所有的警告信息；-O 优化选项，优化代码长度和执行时间；                  |
|   |        |   | # -fstrength-reduce 优化循环执行代码，排除重复变量；-fomit-frame-pointer 省略保存不必要 |
|   |        |   | # 的框架指针；-fcombine-regs 合并寄存器，减少寄存器类的使用；-finline-functions 将所有简   |
|   |        |   | # 单短小的函数代码嵌入调用程序中；-nostdinc -I../include 不使用默认路径中的包含文件，而         |
|   |        |   | # 使用这里指定目录中的(../include)。  |
| 4 | AS     | =gas  | # GNU 的汇编程序。   |
| 5 | AR     | =gar  | # GNU 的二进制文件处理程序，用于创建、修改以及从归档文件中抽取文件。                            |
| 6 | LD     | =gld  | # GNU 的连接程序。   |
| 7 | CPP    | =gcc -E -nostdinc -I../include                                    |  |
|   |        |   | # C 前处理选项。-E 只运行 C 前处理，对所有指定的 C 程序进行预处理并将处理结果输出到标准输              |
|   |        |   | # 出设备或指定的输出文件中；-nostdinc -I../include 同前。                        |
| 8 |        |   |  |
|   |        |   | # 下面的规则指示 make 利用下面的命令将所有的.c 文件编译生成.s 汇编程序。该规则的命令                |
|   |        |   | # 指使 gcc 采用 CFLAGS 所指定的选项对 C 代码编译后不进行汇编就停止（-S），从而产生与             |

```

# 输入的各个 C 文件对应的汇编代码文件。默认情况下所产生的汇编程序文件名是原 C 文件名
# 去掉.c 而加上.s 后缀。-o 表示其后是输出文件的名称。其中$.s (或$@) 是自动目标变量,
# $<代表第一个先决条件, 这里即是符合条件*.c 的文件。
9 .c.o:
10     $(CC) $(CFLAGS) \
11     -c -o $.o $<
# 下面规则表示将所有.s 汇编程序文件编译成.o 目标文件。22 行是实现该操作的具体命令。
12 .s.o:
13     $(AS) -o $.o $<
14 .c.s:                                # 类似上面, *.c 文件->*.s 汇编程序文件。不进行连接。
15     $(CC) $(CFLAGS) \
16     -S -o $.s $<
17
18 OBJS    = memory.o page.o    # 定义目标文件变量 OBJS。
19
20 all: mm.o
21
22 mm.o: $(OBJS)                      # 在有了先决条件 OBJS 后使用下面的命令连接成目标 mm.o
23     $(LD) -r -o mm.o $(OBJS)
24
# 下面的规则用于清理工作。当执行'make clean'时, 就会执行 26--27 行上的命令, 去除所有编译
# 连接生成的文件。'rm'是文件删除命令, 选项-f 含义是忽略不存在的文件, 并且不显示删除信息。
25 clean:
26     rm -f core *.o *.a tmp_make
27     for i in *.c;do rm -f `basename $$i .c`.s;done
28
# 下面得目标或规则用于检查各文件之间的依赖关系。方法如下:
# 使用字符串编辑程序 sed 对 Makefile 文件 (这里即是自己) 进行处理, 输出为删除 Makefile
# 文件中'### Dependencies' 行后面的所有行 (下面从 35 开始的行), 并生成 tmp_make
# 临时文件 (30 行的作用)。然后对 mm/目录下的每一个 C 文件执行 gcc 预处理操作。
# -M 标志告诉预处理程序输出描述每个目标文件相关性的规则, 并且这些规则符合 make 语法。
# 对于每一个源文件, 预处理程序输出一个 make 规则, 其结果形式是相应源程序文件的目标
# 文件名加上其依赖关系--该源文件中包含的所有头文件列表。把预处理结果都添加到临时
# 文件 tmp_make 中, 然后将该临时文件复制成新的 Makefile 文件。

29 dep:
30     sed '/\#\#\# Dependencies/q' < Makefile > tmp_make
31     (for i in *.c;do $(CPP) -M $$i;done) >> tmp_make
32     cp tmp_make Makefile
33
34 ### Dependencies:
35 memory.o : memory.c ../include/signal.h ../include/sys/types.h \
36     ../include/asm/system.h ../include/linux/sched.h ../include/linux/head.h \
37     ../include/linux/fs.h ../include/linux/mm.h ../include/linux/kernel.h

```

## 9.3 memory.c 程序

### 9.3.1 功能描述

## 9.3.2 代码注释

列表 linux/mm/memory.c 程序

```

1  /*
2   * linux/mm/memory.c
3   *
4   * (C) 1991 Linus Torvalds
5   */
6
7  /*
8   * demand-loading started 01.12.91 - seems it is high on the list of
9   * things wanted, and it should be easy to implement. - Linus
10  */
11  /*
12   * 需求加载是从 01.12.91 开始编写的 - 在程序编制表中是呼是最重要的程序，
13   * 并且应该是很容易编制的 - linus
14   */
15
16  /*
17   * Ok, demand-loading was easy, shared pages a little bit trickier. Shared
18   * pages started 02.12.91, seems to work. - Linus.
19   *
20   * Tested sharing by executing about 30 /bin/sh: under the old kernel it
21   * would have taken more than the 6M I have free, but it worked well as
22   * far as I could see.
23   *
24   * Also corrected some "invalidate()"s - I wasn't doing enough of them.
25   */
26  /*
27   * OK, 需求加载是比较容易编写的，而共享页面却需要点技巧。共享页面程序是
28   * 02.12.91 开始编写的，好象能够工作 - Linus。
29   *
30   * 通过执行大约 30 个/bin/sh 对共享操作进行了测试：在老内核当中需要占用多于
31   * 6M 的内存，而目前却不用。目前看来工作的很好。
32   *
33   * 对"invalidate()"函数也进行了修正 - 在这方面我还做的不够。
34   */
35
36 #include <signal.h>          // 信号头文件。定义信号符号常量，信号结构以及信号操作函数原型。
37
38 #include <asm/system.h>      // 系统头文件。定义了设置或修改描述符/中断门等的嵌入式汇编宏。
39
40 #include <linux/sched.h>     // 调度程序头文件，定义了任务结构 task_struct、初始任务 0 的数据，
41                               // 还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。
42
43 #include <linux/head.h>      // head 头文件，定义了与 head.s 程序描述符有关的 C 程序。
44 #include <linux/kernel.h>    // 内核头文件。含有一些内核常用函数的原形定义。
45
46 volatile void do_exit(long code);
47
48 //// 内存已用完。
49 static inline volatile void oom(void)
50 {
51     printk("out of memory\n\r");

```

```

36     do_exit(SIGSEGV);
37 }
38
// 刷新页变换高速缓冲。
// 为了提高地址转换的效率，CPU 将最近使用的页表数据存放在芯片中高速缓冲中。在修改过页表
// 信息之后，就需要刷新该缓冲区。这里使用重新加载页目录基址寄存器 cr3 的方法来进行刷新。
39 #define invalidate() \
40 __asm__ ("movl %%eax, %%cr3":: "a" (0))
41
42 /* these are not to be changed without changing head.s etc */
// 下面代码需要与 head.s 等中的相关信息一起改变 */
43 #define LOW_MEM 0x100000 // 内存低端（1MB）。
44 #define PAGING_MEMORY (15*1024*1024) // 分页内存 15MB。
45 #define PAGING_PAGES (PAGING_MEMORY>>12) // 分页后的页数。
46 #define MAP_NR(addr) (((addr)-LOW_MEM)>>12) // 映射页号。
47 #define USED 100
48
// CODE_SPACE(addr) (((addr)+0xfff)&~0xfff) < current->start_code + current->end_code。
49 #define CODE_SPACE(addr) (((addr)+4095)&~4095) < \
50 current->start_code + current->end_code)
51
52 static long HIGH_MEMORY = 0;
53
// 复制页面（4K 字节）。
54 #define copy_page(from,to) \
55 __asm__ ("cld ; rep ; movsl":: "S" (from), "D" (to), "c" (1024): "cx", "di", "si")
56
57 static unsigned char mem_map [ PAGING_PAGES ] = {0,}; // 内存映射字节图(1 字节代表 1 页)。
58
59 /*
60  * Get physical address of first (actually last :- ) free page, and mark it
61  * used. If no free pages left, return 0.
62  */
// 获取首个(实际上是最后 1 个:-)空闲页面，并标记为已使用。如果没有空闲页面了，
// 就返回 0。
// 取空闲页面。如果已经没有内存了，则返回 0。
// 输入：%1(ax=0) - 0；%2(LOW_MEM)；%3(cx=PAGING_PAGES)；%4(di=mem_map+PAGING_PAGES-1)。
// 输出：返回%0(ax=页面号)。
// 从内存映像末端开始向前扫描所有页面标志（页面总数为 PAGING_PAGES），如果有页面空闲（对应
// 内存映像位为 0）则返回页面地址。
63 unsigned long get_free_page(void)
64 {
65     register unsigned long __res asm("ax");
66
67     __asm__ ("std ; repne ; scasb\n\t" // 方向位置位，将 a1(0)与对应每个页面的(di)内容比较，
68             "jne 1f\n\t" // 如果没有等于 0 的字节，则跳转结束（返回 0）。
69             "movb $1, 1(%%edi)\n\t" // 将对应页面的内存映像位置 1。
70             "sall $12, %%ecx\n\t" // 页面数*4K = 相对页面起始地址。
71             "addl %2, %%ecx\n\t" // 再加上低端内存地址，即获得页面实际起始地址。
72             "movl %%ecx, %%edx\n\t" // 将页面实际起始地址→edx 寄存器。
73             "movl $1024, %%ecx\n\t" // 寄存器 ecx 置计数值 1024。

```

```

74     "leal 4092(%%edx), %%edi\n\t" // 将 4092+edx 的位置→edi (该页面的末端)。
75     "rep ; stosl\n\t"           // 将 edi 所指内存清零 (反方向, 也即将该页面清零)。
76     "movl %%edx, %%eax\n\t"    // 将页面起始地址→eax (返回值)。
77     "l:"
78     : "=a" (__res)
79     : "" (0), "i" (LOW_MEM), "c" (PAGING_PAGES),
80     "D" (mem_map+PAGING_PAGES-1)
81     : "di", "cx", "dx";
82 return __res;                  // 返回空闲页面地址 (如果无空闲也则返回 0)。
83 }
84
85 /*
86  * Free a page of memory at physical address 'addr'. Used by
87  * 'free_page_tables()'
88  */
89
90 // 释放物理地址 'addr' 的页面。用于 'free_page_tables()'。
91 // 1MB 以下的内存空间用于操作系统 (内核), 不作为分配页面的内存。
92 void free_page(unsigned long addr)
93 {
94     if (addr < LOW_MEM) return; // 如果物理地址 addr 小于内存低端 (1MB), 则返回。
95     if (addr >= HIGH_MEMORY)    // 如果物理地址 addr>=内存高端, 则显示出错信息。
96         panic("trying to free nonexistent page");
97     addr -= LOW_MEM;             // 物理地址-低端内存位置, 再除以 4KB, 得页面号。
98     addr >>= 12;
99     if (mem_map[addr]--) return; // 如果对应内存页面映射字节不等于 0, 则减 1 返回。
100    mem_map[addr]=0;             // 否则置对应页面映射字节为 0, 并显示出错信息, 死机。
101    panic("trying to free free page");
102 }
103
104 /*
105  * This function frees a continuous block of page tables, as needed
106  * by 'exit()'. As does copy_page_tables(), this handles only 4Mb blocks.
107  */
108
109 // 下面函数释放页表连续的内存块, 'exit()' 需要该函数。与 copy_page_tables()
110 // 类似, 该函数仅处理 4Mb 的内存块。
111
112 // 根据指定的地址和限长, 释放对应内存页表所指定的内存块及页表本身。
113 // 页目录位于物理地址 0 开始处, 共 1024 项, 占 4K 字节。每个目录项指定一个页表。
114 // 页表从物理地址 0x1000 处开始 (紧接着目录空间), 每个页表有 1024 项, 也占 4K 内存。
115 // 每个页表项指定一页内存 (4K)。目录项和页表项的大小均为 4 个字节。
116 // 参数: from - 起始基地址; size - 释放的长度。
117 int free_page_tables(unsigned long from, unsigned long size)
118 {
119     unsigned long *pg_table;
120     unsigned long *dir, nr;
121
122     if (from & 0x3fffff) // 要释放内存块的地址需以 4M 为边界。
123         panic("free_page_tables called with wrong alignment");
124     if (!from)           // 出错, 试图释放内核所占空间。

```

```

113         panic("Trying to free up swapper memory space");
// 计算所占页目录项数(4M 的进位整数倍), 也即所占页表数。
114         size = (size + 0x3ffff) >> 22;
// 下面一句计算起始目录项。对应的目录项号=from>>22, 因每项占 4 字节, 又由于页目录是从
// 物理地址 0 开始, 因此实际的目录项指针=目录项号<<2, 也即(from>>20)。与上 0xffc 确保
// 目录项指针范围有效。
115         dir = (unsigned long *) ((from>>20) & 0xffc); /* _pg_dir = 0 */
116         for ( ; size-->0 ; dir++) { // size 现在是需释放内存的目录项数。
117             if (!(1 & *dir)) // 如果该目录项无效(P 位=0), 则继续。
118                 continue;
119             pg_table = (unsigned long *) (0xfffff000 & *dir); // 取目录项中页表地址。
120             for (nr=0 ; nr<1024 ; nr++) { // 每个页表有 1024 个页项。
121                 if (1 & *pg_table) // 如果该页表项有效(P 位=1), 则释放该页。
122                     free_page(0xfffff000 & *pg_table);
123                 *pg_table = 0; // 该页表项内容清零。
124                 pg_table++; // 指向页表中下一项。
125             }
126             free_page(0xfffff000 & *dir); // 释放该页表所占内存页面。
127             *dir = 0; // 对相应页表的目录项清零。
128         }
129         invalidate(); // 刷新页变换高速缓冲。
130         return 0;
131     }
132
133 /*
134  * Well, here is one of the most complicated functions in mm. It
135  * copies a range of linear addresses by copying only the pages.
136  * Let's hope this is bug-free, 'cause this one I don't want to debug :-)
137  *
138  * Note! We don't copy just any chunks of memory - addresses have to
139  * be divisible by 4Mb (one page-directory entry), as this makes the
140  * function easier. It's used only by fork anyway.
141  *
142  * NOTE 2!! When from==0 we are copying kernel space for the first
143  * fork(). Then we DONT want to copy a full page-directory entry, as
144  * that would lead to some serious memory waste - we just copy the
145  * first 160 pages - 640kB. Even that is more than we need, but it
146  * doesn't take any more memory - we don't copy-on-write in the low
147  * 1 Mb-range, so the pages can be shared with the kernel. Thus the
148  * special case for nr=xxxx.
149  */
150 /*
151  * 好了, 下面是内存管理 mm 中最为复杂的程序之一。它通过只复制内存页面
152  * 来拷贝一定范围内的线性地址中的内容。希望代码中没有错误, 因为我不想
153  * 再调试这块代码了☺。
154  *
155  * 注意! 我们并不是仅复制任何内存块 - 内存块的地址需要是 4Mb 的倍数 (正好
156  * 一个页目录项对应的内存大小), 因为这样处理可使函数很简单。不管怎样,
157  * 它仅被 fork() 使用。
158  *
159  * 注意 2!! 当 from==0 时, 是在为第一次 fork() 调用复制内核空间。此时我们
160  * 不想复制整个页目录项对应的内存, 因为这样做会导致严重的内存浪费 - 我们
161  * 只复制头 160 个页面 - 对应 640kB。即使是复制这些页面也已经超出我们的需求,

```

\* 但这不会占用更多的内存 – 在低 1Mb 内存范围内我们不执行写时复制操作，所以  
 \* 这些页面可以与内核共享。因此这是 nr=xxxx 的特殊情况。

\*/

//// 复制进程的页目录页表。

```

150 int copy_page_tables(unsigned long from,unsigned long to,long size)
151 {
152     unsigned long * from_page_table;
153     unsigned long * to_page_table;
154     unsigned long this_page;
155     unsigned long * from_dir, * to_dir;
156     unsigned long nr;
157
158     // 源地址和目的地址都需要是 4Mb 的倍数。否则出错，死机。
159     if ((from&0x3fffff) || (to&0x3fffff))
160         panic("copy_page_tables called with wrong alignment");
161     // 取得源地址和目的地址的目录项(from_dir 和 to_dir)。
162     from_dir = (unsigned long *) ((from>>20) & 0xffc); /* _pg_dir = 0 */
163     to_dir = (unsigned long *) ((to>>20) & 0xffc);
164     // 计算要复制的内存块占用的页表数（也即目录项数）。
165     size = ((unsigned) (size+0x3fffff)) >> 22;
166     // 下面开始对每个占用的页表依次进行复制操作。
167     for( ; size-->0 ; from_dir++,to_dir++) {
168         // 如果目的目录项指定的页表已经存在(P=1)，则出错，死机。
169         if (1 & *to_dir)
170             panic("copy_page_tables: already exist");
171         // 如果此源目录项未被使用，则不用复制对应页表，跳过。
172         if (!(1 & *from_dir))
173             continue;
174         // 取当前源目录项中页表的地址→from_page_table。
175         from_page_table = (unsigned long *) (0xfffff000 & *from_dir);
176         // 为目的页表取一页空闲内存，如果返回是 0 则说明没有申请到空闲内存页面。返回值=-1，退出。
177         if (!(to_page_table = (unsigned long *) get_free_page()))
178             return -1; /* Out of memory, see freeing */
179         // 设置目的目录项信息。7 是标志信息，表示(Usr, R/W, Present)。
180         *to_dir = ((unsigned long) to_page_table) | 7;
181         // 针对当前处理的页表，设置需复制的页面数。如果是在内核空间，则仅需复制头 160 页，否则需要
182         // 复制 1 个页表中的所有 1024 页面。
183         nr = (from==0)?0xA0:1024;
184         // 对于当前页表，开始复制指定数目 nr 个内存页面。
185         for ( ; nr-- > 0 ; from_page_table++,to_page_table++) {
186             this_page = *from_page_table; // 取源页表项内容。
187             if (!(1 & this_page)) // 如果当前源页面没有使用，则不用复制。
188                 continue;
189             // 复位页表项中 R/W 标志(置 0)。(如果 U/S 位是 0，则 R/W 就没有作用。如果 U/S 是 1，而 R/W 是 0，
190             // 那么运行在用户层的代码就只能读页面。如果 U/S 和 R/W 都置位，则就有写的权限。)
191             this_page &= ~2;
192             *to_page_table = this_page; // 将该页表项复制到目的页表中。
193             // 如果该页表项所指页面的地址在 1M 以上，则需要设置内存页面映射数组 mem_map[]，于是计算
194             // 页面号，并以它为索引在页面映射数组相应项中增加引用次数。
195             if (this_page > LOW_MEM) {
196                 *from_page_table = this_page; // 令源页表项也只读[??]。
197                 this_page -= LOW_MEM;
198                 this_page >>= 12;

```



```

183             mem_map[this_page]++;
184         }
185     }
186 }
187 invalidate();           // 刷新页变换高速缓冲。
188 return 0;
189 }
190
191 /*
192  * This function puts a page in memory at the wanted address.
193  * It returns the physical address of the page gotten, 0 if
194  * out of memory (either when trying to access page-table or
195  * page.)
196  */
197 /*
198  * 下面函数将一页面放置在内存中指定地址处。它返回页面的物理地址，如果
199  * 内存不够(在访问页表或页面时)，则返回 0。
200  */
201 // 在指定物理地址处放置一页面。
202 // 主要工作是在页目录和页表中设置指定页面的信息。
203 // 若成功则返回页面地址。
204 unsigned long put_page(unsigned long page, unsigned long address)
205 {
206     unsigned long tmp, *page_table;
207
208     /* NOTE !!! This uses the fact that _pg_dir=0 */
209     /* 注意!!!这里使用了页目录基址_pg_dir=0 的条件 */
210
211     // 如果申请的页面位置低于 LOW_MEM(1Mb)或超出系统实际含有内存高端 HIGH_MEMORY，则发出警告。
212     if (page < LOW_MEM || page >= HIGH_MEMORY)
213         printk("Trying to put page %p at %p\n", page, address);
214     // 如果申请的页面在内存页面映射字节图中没有置位，则显示警告信息。
215     if (mem_map[(page-LOW_MEM)>>12] != 1)
216         printk("mem_map disagrees with %p at %p\n", page, address);
217     // 计算指定地址在页目录中对应的目录项指针。
218     page_table = (unsigned long *) ((address>>20) & 0xffc);
219     // 如果该目录项有效(P=1)(也即指定的页表在内存中)，则从中取得指定页表的地址→page_table。
220     if ((*page_table)&1)
221         page_table = (unsigned long *) (0xfffff000 & *page_table);
222     else {
223         // 否则，申请空闲页面给页表使用，并在对应目录项中置相应标志 7 (User, U/S, R/W)。然后将
224         // 该页表的地址→page_table。
225         if (!(tmp=get_free_page()))
226             return 0;
227         *page_table = tmp|7;
228         page_table = (unsigned long *) tmp;
229     }
230     // 在页表中设置指定地址页面的页表项内容。每个页表共可有 1024 项(0x3ff)。
231     page_table[(address>>12) & 0x3ff] = page | 7;
232     /* no need for invalidate */
233     /* 不需要刷新页变换高速缓冲 */
234     return page;           // 返回页面地址。
235 }

```



```

220 220
221 221 // 取消写保护页面函数（写时复制）。
222 222 // 输入参数为页表项指针。
223 223 // [?? un_wp_page 意思是否是取消页面的写保护?? wp -- Write Protected.]
224 224 void un_wp_page(unsigned long * table_entry)
225 225 {
226 226     unsigned long old_page, new_page;
227 227
228 228     old_page = 0xfffff000 & *table_entry; // 取原页面地址。
229 229 // 如果原页面地址大于内存低端 LOW_MEM(1Mb) 并且其在页面映射字节数组中值为 1（表示仅被
230 230 // 引用 1 次，页面没有被共享），则在该页面的页表项中置 R/W 标志（可写），并刷新页变换
231 231 // 高速缓冲，然后返回。
232 232 if (old_page >= LOW_MEM && mem_map[MAP_NR(old_page)]==1) {
233 233     *table_entry |= 2;
234 234     invalidate();
235 235     return;
236 236 }
237 237 // 否则，申请一页空闲页面。
238 238 if (!(new_page=get_free_page()))
239 239     oom(); // Out of Memory。内存不够处理。
240 240 // 如果原页面大于内存低端（则意味着 mem_map[]>1，页面是共享的），则将原页面的页面映射
241 241 // 数组值递减 1。然后将指定页表项内容更新为新页面的地址，并置可读写等标志 (U/S, R/W, P)。
242 242 // 刷新页变换高速缓冲。最后将原页面内容复制到新页面。
243 243 if (old_page >= LOW_MEM)
244 244     mem_map[MAP_NR(old_page)]--;
245 245 *table_entry = new_page | 7;
246 246 invalidate();
247 247 copy_page(old_page, new_page);
248 248 }
249 249
250 250 /*
251 251  * This routine handles present pages, when users try to write
252 252  * to a shared page. It is done by copying the page to a new address
253 253  * and decrementing the shared-page counter for the old page.
254 254  *
255 255  * If it's in code space we exit with a segment error.
256 256  */
257 257 /*
258 258  * 当用户试图往一个共享页面上写时，该函数处理已存在的内存页面，（写时复制）
259 259  * 它是通过将页面复制到一个新地址上并递减原页面的共享页面计数值实现的。
260 260  *
261 261  * 如果它在代码空间，我们就以段错误信息退出。
262 262  */
263 263 // 页面异常中断处理调用函数。写共享页面处理函数。
264 264 // 参数 error_code 是由 CPU 自动产生，address 是页面线性地址。
265 264 // 写共享页面时，需复制页面（写时复制）。
266 266 void do_wp_page(unsigned long error_code, unsigned long address)
267 267 {
268 268 #if 0
269 269 /* we cannot do this yet: the estdio library writes to code space */
270 270 /* stupid, stupid. I really want the libc.a from GNU */
271 271 /* 我们现在还不能这样做：因为 estdio 库会在代码空间执行写操作 */
272 271 /* 真是太愚蠢了。我真想从 GNU 得到 libc.a 库。*/

```

```

252         if (CODE_SPACE(address))          // 如果地址位于代码空间，则终止执行程序。
253             do_exit(SIGSEGV);
254 #endif
// 处理取消页面保护。参数指定页面在页表中的页表项指针，其计算方法是：
// ((address>>10) & 0xffc)：计算指定地址的页面在页表中的偏移地址；
// (0xfffff000 & ((address>>20) & 0xffc))：取目录项中页表的地址值，
// 其中((address>>20) & 0xffc)计算页面所在页表的目录项指针；
// 两者相加即得指定地址对应页面的页表项指针。
255         un_wp_page((unsigned long *)
256                     (((address>>10) & 0xffc) + (0xfffff000 &
257                                                     *((unsigned long *) ((address>>20) & 0xffc)))));
258
259     }
260
261     // 写页面验证。
262 void write_verify(unsigned long address)
263 {
264     unsigned long page;
265
266     // 判断指定地址所对应的页表页面是否存在(P)，若不存在(P=0)则返回。
267     if (!(page = *((unsigned long *) ((address>>20) & 0xffc)) & 1))
268         return;
269     // 取页表的地址，加上指定地址的页面在页表中的页表项偏移值，得对应页面的页表项指针。
270     page &= 0xfffff000;
271     page += ((address>>10) & 0xffc);
272     // 如果该页面不可写(标志 R/W 没有置位)，则执行共享检验和复制页面操作(写时复制)。
273     if ((3 & *(unsigned long *) page) == 1) /* non-writeable, present */
274         un_wp_page((unsigned long *) page);
275     return;
276 }
277
278 // 取指定地址的空页面。
279 // 与 get_free_page() 不同。get_free_page() 仅置内存页面映射数组 mem_map[] 中的引用标志。
280 // 而这里 get_empty_page() 不仅是取指定地址处的页面，还进一步调用 put_page()，将页面信息
281 // 添加到页目录和页表中。
282 void get_empty_page(unsigned long address)
283 {
284     unsigned long tmp;
285
286     if (!(tmp=get_free_page()) || !put_page(tmp,address)) {
287         free_page(tmp);          /* 0 is ok - ignored */
288         oom();
289     }
290 }
291
292 /*
293  * try_to_share() checks the page at address "address" in the task "p",
294  * to see if it exists, and if it is clean. If so, share it with the current
295  * task.
296  *
297  * NOTE! This assumes we have checked that p != current, and that they
298  * share the same executable.
299  */

```

```

/*
 * try_to_share() 在任务“p”中检查位于地址“address”处的页面，看页面是否存在，是否干净。
 * 如果是干净的话，就与当前任务共享。
 *
 * 注意！这里我们已假定 p != 当前任务，并且它们共享同一个执行程序。
 */
///// 尝试共享。
// 返回 1-成功，0-失败。
292 static int try\_to\_share(unsigned long address, struct task\_struct * p)
293 {
294     unsigned long from;
295     unsigned long to;
296     unsigned long from_page;
297     unsigned long to_page;
298     unsigned long phys_addr;
299
    // 求指定内存地址的页目录项。
300     from_page = to_page = ((address>>20) & 0xffc);
    // 计算地址在指定进程 p 中对应的页目录项。
301     from_page += ((p->start_code>>20) & 0xffc);
    // 计算地址在当前进程中对应的页目录项。
302     to_page += ((current->start_code>>20) & 0xffc);
303 /* is there a page-directory at from? */
    // 在 from 处是否存在页目录？*/
    // *** 对 p 进程页面进行操作。
    // 取页目录项内容。如果该目录项无效(P=0)，则返回。否则取该目录项对应页表地址→from。
304     from = *(unsigned long *) from_page;
305     if (!(from & 1))
306         return 0;
307     from &= 0xfffff000;
    // 计算地址对应的页表项指针值，并取出该页表项内容→phys_addr。
308     from_page = from + ((address>>10) & 0xffc);
309     phys_addr = *(unsigned long *) from_page;
310 /* is the page clean and present? */
    // 页面干净并且存在吗？*/
    // 0x41 对应页表项中的 Dirty 和 Present 标志。如果页面不干净或无效则返回。
311     if ((phys_addr & 0x41) != 0x01)
312         return 0;
    // 取页面的地址→phys_addr。如果该页面地址不存在或小于内存低端(1M)也返回退出。
313     phys_addr &= 0xfffff000;
314     if (phys_addr >= HIGH MEMORY || phys_addr < LOW MEM)
315         return 0;
    // *** 对当前进程页面进行操作。
    // 取页目录项内容→to。如果该目录项无效(P=0)，则取空闲页面，并让更新 to_page 所指的目录项。
316     to = *(unsigned long *) to_page;
317     if (!(to & 1))
318         if (to = get\_free\_page())
319             *(unsigned long *) to_page = to | 7;
320     else
321         oom();
    // 取对应页表地址→to，页表项地址→to_page。如果对应的页面已经存在，则出错，死机。
322     to &= 0xfffff000;
323     to_page = to + ((address>>10) & 0xffc);

```

```

324         if (1 & *(unsigned long *) to_page)
325             panic("try_to_share: to_page already exists");
326 /* share them: write-protect */
327 /* 对它们进行共享处理：写保护 */
328 // 对 p 进程中页面置写保护标志(置 R/W=0)。并且当前进程中的对应页表项指向它。
329     *(unsigned long *) from_page &= ~2;
330     *(unsigned long *) to_page = *(unsigned long *) from_page;
331 // 刷新页变换高速缓冲。
332     invalidate();
333 // 计算所操作页面的页面号，并将对应页面映射数组项中的引用递增 1。
334     phys_addr -= LOW\_MEM;
335     phys_addr >>= 12;
336     mem\_map[phys_addr]++;
337     return 1;
338 }
339
340 /*
341 * share_page() tries to find a process that could share a page with
342 * the current one. Address is the address of the wanted page relative
343 * to the current data space.
344 * We first check if it is at all feasible by checking executable->i_count.
345 * It should be >1 if there are other tasks sharing this inode.
346 */
347 /*
348 * share_page() 试图找到一个进程，它可以与当前进程共享页面。参数 address 是
349 * 当前数据空间中期望共享的某页面地址。
350 * 首先我们通过检测 executable->i_count 来查证是否可行。如果有其它任务已共享
351 * 该 inode，则它应该大于 1。
352 */
353 //// 共享页面。
354 // 返回 1 - 成功，0 - 失败。
355 static int share\_page(unsigned long address)
356 {
357     struct task\_struct ** p;
358
359     // 如果是不可执行的，则返回。
360     if (!current->executable)
361         return 0;
362     // 如果只能单独执行(executable->i_count=1)，也退出。
363     if (current->executable->i_count < 2)
364         return 0;
365     // 搜索任务数组中所有任务。寻找与当前进程可共享页面的进程，并尝试对指定地址的页面进行共享。
366     for (p = &LAST\_TASK ; p > &FIRST\_TASK ; --p) {
367         if (!*p) // 如果该任务项空闲，则继续寻找。
368             continue;
369         if (current == *p) // 如果就是当前任务，也继续寻找。
370             continue;
371         if ((*p)->executable != current->executable) // 如果 executable 不等，也继续。
372             continue;
373         if (try\_to\_share(address, *p)) // 尝试共享页面。
374             return 1;
375     }

```

```

361     }
362     return 0;
363 }
364
365 // 页面异常中断处理调用函数。缺页异常处理函数。
366 // 参数 error_code 是由 CPU 自动产生，address 是页面线性地址。
367 void do_no_page(unsigned long error_code, unsigned long address)
368 {
369     int nr[4];
370     unsigned long tmp;
371     unsigned long page;
372     int block, i;
373
374     address &= 0xfffff000; // 页面地址。
375     // 如果当前进程的 executable 空，或者指定地址超出进程的数据末端，则取指定地址空页面并退出。
376     tmp = address - current->start_code;
377     if (!current->executable || tmp >= current->end_data) {
378         get_empty_page(address);
379         return;
380     }
381     // 如果在所有进程中尝试共享页面失败，则退出。
382     if (share_page(tmp))
383         return;
384     // 取空闲页面，如果内存不够了，则进行内存不够处理。
385     if (!(page = get_free_page()))
386         oom();
387     /* remember that 1 block is used for header */
388     /* 记住，（任务）头要使用 1 个数据块 */
389     //
390     block = 1 + tmp/BLOCK_SIZE;
391     for (i=0 ; i<4 ; block++, i++)
392         nr[i] = bmap(current->executable, block);
393     bread_page(page, current->executable->i_dev, nr);
394     i = tmp + 4096 - current->end_data;
395     tmp = page + 4096;
396     while (i-- > 0) {
397         tmp--;
398         *(char *)tmp = 0;
399     }
400     if (put_page(page, address))
401         return;
402     free_page(page);
403     oom();
404 }
405
406 // 内存初始化。
407 // 参数：start_mem - 可用作分页处理的物理内存起始位置（已去除 RAMDISK 所占内存空间等）。
408 //       end_mem   - 实际物理内存最大地址。
409 // 在该版的 linux 内核中，最多能使用 16Mb 的内存，大于 16Mb 的内存将不予考虑，弃置不用。
410 // 0 - 1Mb 内存空间用于内核系统（其实是 0-640Kb）。
411 void mem_init(long start_mem, long end_mem)
412 {
413     int i;

```

```

402
403     HIGH_MEMORY = end_mem;           // 设置内存最高端。
404     for (i=0 ; i<PAGING_PAGES ; i++) // 首先置所有页面为已占用 (USED=100) 状态,
405         mem_map[i] = USED;           // 即将页面映射数组全置成 USED。
406     i = MAP_NR(start_mem);           // 然后计算可使用起始内存的页面号。
407     end_mem -= start_mem;             // 再计算可分页处理的内存块大小。
408     end_mem >>= 12;                   // 从而计算出可用于分页处理的页面数。
409     while (end_mem-->0)               // 最后将这些可用页面对应的页面映射数组清零。
410         mem_map[i++]=0;
411 }
412
413 // 计算内存空闲页面数并显示。
414 void calc_mem(void)
415 {
416     int i, j, k, free=0;
417     long * pg_tbl;
418
419     // 扫描内存页面映射数组 mem_map[], 获取空闲页面数并显示。
420     for(i=0 ; i<PAGING_PAGES ; i++)
421         if (!mem_map[i]) free++;
422     printk("%d pages free (of %d)\n", free, PAGING_PAGES);
423 // 扫描所有页目录项 (除 0, 1 项), 如果页目录项有效, 则统计对应页表中有效页面数, 并显示。
424     for(i=2 ; i<1024 ; i++) {
425         if (l&pg_dir[i]) {
426             pg_tbl=(long *) (0xffffffff & pg_dir[i]);
427             for(j=k=0 ; j<1024 ; j++)
428                 if (pg_tbl[j]&1)
429                     k++;
430             printk("Pg-dir[%d] uses %d pages\n", i, k);
431         }
432     }

```

## 9.4 page.s 程序

### 9.4.1 功能描述

### 9.4.2 代码注释

列表 linux/mm/page.s 程序

```

1 /*
2  * linux/mm/page.s
3  *
4  * (C) 1991 Linus Torvalds
5  */
6
7 /*

```

```

8  * page.s contains the low-level page-exception code.
9  * the real work is done in mm.c
10 */
/*
   * page.s 程序包含底层页异常处理代码。实际的工作在 memory.c 中完成。
   */
11
// 该文件包括页异常中断处理程序（中断 14），主要分两种情况处理。一是由于缺页引起的页异常中断，
// 通过调用 do_no_page(error_code, address) 来处理；二是其它页保护引起的页异常，此时调用页写
// 保护处理函数 do_wp_page(error_code, address) 进行处理。其中的出错码 (error_code) 是由 CPU 自
// 动产生并压入堆栈的，出现异常时访问的线性地址是从控制寄存器 CR2 中取得的。
12 .globl _page_fault
13
14 _page_fault:
15     xchgl %eax, (%esp)    # 取出错码到 eax。
16     pushl %ecx
17     pushl %edx
18     push %ds
19     push %es
20     push %fs
21     movl $0x10, %edx     # 置内核数据段选择符。
22     mov %dx, %ds
23     mov %dx, %es
24     mov %dx, %fs
25     movl %cr2, %edx      # 取引起页面异常的线性地址
26     pushl %edx           # 将该线性地址和出错码压入堆栈，作为调用函数的参数。
27     pushl %eax
28     testl $1, %eax       # 测试标志 P，如果不是缺页引起的异常则跳转。
29     jne 1f
30     call _do_no_page     # 调用缺页处理函数（mm/memory.c, 365 行）。
31     jmp 2f
32 1:    call _do_wp_page    # 调用写保护处理函数（mm/memory.c, 247 行）。
33 2:    addl $8, %esp       # 丢弃压入栈的两个参数。
34     pop %fs
35     pop %es
36     pop %ds
37     popl %edx
38     popl %ecx
39     popl %eax
40     iret

```

## 9.4.3 其它信息

### 9.4.3.1 页异常的处理

当处理器在转换线性地址到物理地址的过程中检测到以下两种条件时，就会发生页异常中断，中断 14。

- o 当 CPU 发现对应页目录项或页表项的存在位（Present）标志为 0。
- o 当前进程没有访问指定页面的权限。

对于页异常处理中断，CPU 提供了两项信息用来诊断页异常和从中恢复运行。

- (1) 放在堆栈上的出错码。该出错码指出了异常是由于页不存在引起的还是违反了访问权限引起的；在发生异常时 CPU 的当前特权层；以及是读操作还是写操作。出错码的格式是一个 32 位的长字。但只用了最后的 3 个比特位。分别说明导致异常发生时的原因：

位 2(U/S) - 0 表示在超级用户模式下执行, 1 表示在用户模式下执行;

位 1(W/R) - 0 表示读操作, 1 表示写操作;

位 0(P) - 0 表示页不存在, 1 表示页级保护。

- (2) CR2(控制寄存器 2)。CPU 将造成异常的用于访问的线性地址存放在 CR2 中。异常处理程序可以使用这个地址来定位相应的页目录和页表项。如果在页异常处理程序执行期间允许发生另一个页异常, 那么处理程序应该将 CR2 压入堆栈中。



## 第 10 章 包含文件(include)

















### 10.1 概述

Linux 0.11 版内核中共有 32 个头文件(\*.h)，其中 asm/子目录中含有 4 个，linux/子目录中含有 10 个，sys/子目录中含有 5 个。从下一节开始我们首先描述 include/目录下的 13 个头文件，然后依次说明每个子目录中的文件。说明顺序按照文件名称排序进行。

### 10.2 include/目录下的文件

该目录下的文件列表如下（列表 10.1）：

列表 linux/include/目录下的文件

|   | Name                      | Size       | Last modified (GMT) | Description |
|---|---------------------------|------------|---------------------|-------------|
|    | <a href="#">asm/</a>      |            | 1991-09-17 13:08:31 |             |
|    | <a href="#">linux/</a>    |            | 1991-11-02 13:35:49 |             |
|  | <a href="#">sys/</a>      |            | 1991-09-17 15:06:07 |             |
|  | <a href="#">a.out.h</a>   | 6047 bytes | 1991-09-17 15:10:49 | m           |
|  | <a href="#">const.h</a>   | 321 bytes  | 1991-09-17 15:12:39 | m           |
|  | <a href="#">ctype.h</a>   | 1049 bytes | 1991-11-07 17:30:47 | m           |
|  | <a href="#">errno.h</a>   | 1268 bytes | 1991-09-17 15:04:15 | m           |
|  | <a href="#">fcntl.h</a>   | 1374 bytes | 1991-09-17 15:12:39 | m           |
|  | <a href="#">signal.h</a>  | 1762 bytes | 1991-09-22 19:58:04 | m           |
|  | <a href="#">stdarg.h</a>  | 780 bytes  | 1991-09-17 15:02:23 | m           |
|  | <a href="#">stddef.h</a>  | 286 bytes  | 1991-09-17 15:02:17 | m           |
|  | <a href="#">string.h</a>  | 7881 bytes | 1991-09-17 15:04:09 | m           |
|  | <a href="#">termios.h</a> | 5325 bytes | 1991-11-25 20:02:08 | m           |
|  | <a href="#">time.h</a>    | 734 bytes  | 1991-09-17 15:02:02 | m           |
|  | <a href="#">unistd.h</a>  | 6410 bytes | 1991-11-25 20:18:55 | m           |
|  | <a href="#">utime.h</a>   | 225 bytes  | 1991-09-17 15:03:38 | m           |

### 10.3 a.out.h 文件

#### 10.3.1 功能描述

a.out.h 头文件主要定义了二进制执行文件 a.out (Assembly out) 的格式。其中包括三个数据结构和一些宏函数。

## 10.3.2 代码注释

列表 linux/include/a.out.h 文件

```

1 #ifndef A_OUT_H
2 #define A_OUT_H
3
4 #define GNU_EXEC_MACROS
5
6 // 执行文件结构。
7 // =====
8 // unsigned long a_magic      // 执行文件魔数。使用 N_MAGIC 等宏访问。
9 // unsigned a_text           // 代码长度，字节数。
10 // unsigned a_data           // 数据长度，字节数。
11 // unsigned a_bss            // 文件中的未初始化数据区长度，字节数。
12 // unsigned a_syms           // 文件中的符号表长度，字节数。
13 // unsigned a_entry          // 执行开始地址。
14 // unsigned a_trsize         // 代码重定位信息长度，字节数。
15 // unsigned a_drsize         // 数据重定位信息长度，字节数。
16 // -----
17 struct exec {
18     unsigned long a_magic;      /* Use macros N_MAGIC, etc for access */
19     unsigned a_text;           /* length of text, in bytes */
20     unsigned a_data;           /* length of data, in bytes */
21     unsigned a_bss;            /* length of uninitialized data area for file, in bytes */
22     unsigned a_syms;           /* length of symbol table data in file, in bytes */
23     unsigned a_entry;          /* start address */
24     unsigned a_trsize;         /* length of relocation info for text, in bytes */
25     unsigned a_drsize;         /* length of relocation info for data, in bytes */
26 };
27
28 // 用于取执行结构中的魔数。
29 #ifndef N_MAGIC
30 #define N_MAGIC(exec) ((exec).a_magic)
31 #endif
32
33 #ifndef OMAGIC
34 /* Code indicating object file or impure executable. */
35 /* 指明为目标文件或者不纯的可执行文件的代号 */
36 #define OMAGIC 0407
37 /* Code indicating pure executable. */
38 /* 指明为纯可执行文件的代号 */
39 #define NMAGIC 0410
40 /* Code indicating demand-paged executable. */
41 /* 指明为需求分页处理的可执行文件 */
42 #define ZMAGIC 0413
43 #endif /* not OMAGIC */
44
45 // 如果魔数不能被识别，则返回真。
46 #ifndef N_BADMAG
47 #define N_BADMAG(x) \
48     (N_MAGIC(x) != OMAGIC && N_MAGIC(x) != NMAGIC \
49     && N_MAGIC(x) != ZMAGIC)
50 #endif

```

```

35
36 #define _N_BADMAG(x) \
37 (N_MAGIC(x) != OMAGIC && N_MAGIC(x) != NMAGIC \
38  && N_MAGIC(x) != ZMAGIC)
39
// 程序头在内存中的偏移位置。
40 #define N_HDROFF(x) (SEGMENT_SIZE - sizeof (struct exec))
41
// 代码起始偏移值。
42 #ifndef N_TXTOFF
43 #define N_TXTOFF(x) \
44 (N_MAGIC(x) == ZMAGIC ? N_HDROFF(x) + sizeof (struct exec) : sizeof (struct exec))
45 #endif
46
// 数据起始偏移值。
47 #ifndef N_DATOFF
48 #define N_DATOFF(x) (N_TXTOFF(x) + (x).a_text)
49 #endif
50
// 代码重定位信息偏移值。
51 #ifndef N_TRELOFF
52 #define N_TRELOFF(x) (N_DATOFF(x) + (x).a_data)
53 #endif
54
// 数据重定位信息偏移值。
55 #ifndef N_DRELOFF
56 #define N_DRELOFF(x) (N_TRELOFF(x) + (x).a_trsize)
57 #endif
58
// 符号表偏移值。
59 #ifndef N_SYMOFF
60 #define N_SYMOFF(x) (N_DRELOFF(x) + (x).a_drsize)
61 #endif
62
// 字符串信息偏移值。
63 #ifndef N_STROFF
64 #define N_STROFF(x) (N_SYMOFF(x) + (x).a_syms)
65 #endif
66
67 /* Address of text segment in memory after it is loaded. */
68 /* 代码段加载到内存中后的地址 */
69 #ifndef N_TXTADDR
70 #define N_TXTADDR(x) 0
71 #endif
72 /* Address of data segment in memory after it is loaded.
73  Note that it is up to you to define SEGMENT_SIZE
74  on machines not listed here. */
75 /* 数据段加载到内存中后的地址。
76  注意，对于下面没有列出名称的机器，需要你自己来定义
77  对应的 SEGMENT_SIZE */
78 #if defined(vax) || defined(hp300) || defined(pyr)
79 #define SEGMENT_SIZE PAGE_SIZE

```

```

77 #endif
78 #ifdef hp300
79 #define PAGE_SIZE 4096
80 #endif
81 #ifdef sony
82 #define SEGMENT_SIZE 0x2000
83 #endif /* Sony. */
84 #ifdef is68k
85 #define SEGMENT_SIZE 0x20000
86 #endif
87 #if defined(m68k) && defined(PORTAR)
88 #define PAGE_SIZE 0x400
89 #define SEGMENT_SIZE PAGE_SIZE
90 #endif
91
92 #define PAGE_SIZE 4096
93 #define SEGMENT_SIZE 1024
94
95 // 以段为界的大小。
96 #define N_SEGMENT_ROUND(x) (((x) + SEGMENT_SIZE - 1) & ~(SEGMENT_SIZE - 1))
97
98 // 代码段尾地址。
99 #define N_TXTENDADDR(x) (N_TXTADDR(x) + (x).a_text)
100
101 // 数据开始地址。
102 #ifndef N_DATADDR
103 #define N_DATADDR(x) \
104     (N_MAGIC(x) == OMAGIC ? (N_TXTENDADDR(x)) \
105      : (N_SEGMENT_ROUND(N_TXTENDADDR(x))))
106 #endif
107
108 /* Address of bss segment in memory after it is loaded. */
109 /* bss 段加载到内存以后的地址 */
110 #ifndef N_BSSADDR
111 #define N_BSSADDR(x) (N_DATADDR(x) + (x).a_data)
112 #endif
113
114 // nlist 结构。
115 #ifndef N_NLIST_DECLARED
116 struct nlist {
117     union {
118         char *n_name;
119         struct nlist *n_next;
120         long n_strx;
121     } n_un;
122     unsigned char n_type;
123     char n_other;
124     short n_desc;
125     unsigned long n_value;
126 };
127 #endif
128
129 // 下面定义 exec 结构中的变量偏移值。

```

```

124 #ifndef N\_UNDF
125 #define N\_UNDF 0
126 #endif
127 #ifndef N\_ABS
128 #define N\_ABS 2
129 #endif
130 #ifndef N\_TEXT
131 #define N\_TEXT 4
132 #endif
133 #ifndef N\_DATA
134 #define N\_DATA 6
135 #endif
136 #ifndef N\_BSS
137 #define N\_BSS 8
138 #endif
139 #ifndef N\_COMM
140 #define N\_COMM 18
141 #endif
142 #ifndef N\_FN
143 #define N\_FN 15
144 #endif
145
146 #ifndef N\_EXT
147 #define N\_EXT 1
148 #endif
149 #ifndef N\_TYPE
150 #define N\_TYPE 036
151 #endif
152 #ifndef N\_STAB
153 #define N\_STAB 0340
154 #endif
155
156 /* The following type indicates the definition of a symbol as being
157    an indirect reference to another symbol. The other symbol
158    appears as an undefined reference, immediately following this symbol.
159
160    Indirection is asymmetrical. The other symbol's value will be used
161    to satisfy requests for the indirect symbol, but not vice versa.
162    If the other symbol does not have a definition, libraries will
163    be searched to find a definition. */
164 /* 下面的类型指明了符号的定义作为对另一个符号的间接引用。紧接该符号的其它
165    * 的符号呈现为未定义的引用。
166    *
167    * 间接性是不对称的。其它符号的值将被用于满足间接符号的请求，但反之不然。
168    * 如果其它符号并没有定义，则将搜索库来寻找一个定义 */
164 #define N\_INDR 0xa
165
166 /* The following symbols refer to set elements.
167    All the N\_SET\[ATDB\] symbols with the same name form one set.
168    Space is allocated for the set in the text section, and each set
169    element's value is stored into one word of the space.
170    The first word of the space is the length of the set (number of elements).
171

```

```

172 The address of the set is made into an N_SETV symbol
173 whose name is the same as the name of the set.
174 This symbol acts like a N_DATA global symbol
175 in that it can satisfy undefined external references. */
/* 下面的符号与集合元素有关。所有具有相同名称 N_SET[ATDB] 的符号
   形成一个集合。在代码部分中已为集合分配了空间，并且每个集合元素
   的值存放在一个字（word）的空间。空间的第一个字存有集合的长度（集合元素数目）。

```

集合的地址被放入一个 N\_SETV 符号，它的名称与集合同名。

在满足未定义的外部引用方面，该符号的行为象一个 N\_DATA 全局符号。\*/

```

176
177 /* These appear as input to LD, in a .o file. */
/* 以下这些符号在目标文件中是作为链接程序 LD 的输入。*/
178 #define N_SETA 0x14 /* Absolute set element symbol */
                        /* 绝对集合元素符号 */
179 #define N_SETT 0x16 /* Text set element symbol */
                        /* 代码集合元素符号 */
180 #define N_SETD 0x18 /* Data set element symbol */
                        /* 数据集合元素符号 */
181 #define N_SETB 0x1A /* Bss set element symbol */
                        /* Bss 集合元素符号 */
182
183 /* This is output from LD. */
/* 下面是 LD 的输出。*/
184 #define N_SETV 0x1C /* Pointer to set vector in data area. */
                        /* 指向数据区中集合向量。*/
185
186 #ifndef N_RELOCATION_INFO_DECLARED
187
188 /* This structure describes a single relocation to be performed.
189 The text-relocation section of the file is a vector of these structures,
190 all of which apply to the text section.
191 Likewise, the data-relocation section applies to the data section. */
/* 下面的结构描述执行一个重定位的操作。
   文件的代码重定位部分是这些结构的一个向量，所有这些适用于代码部分。
   类似地，数据重定位部分适用于数据部分。*/
192
193 // 重定位信息结构。
194 struct relocation_info
195 {
196     /* Address (within segment) to be relocated. */
        /* 需要重定位的地址（在段内）。*/
197     int r_address;
198     /* The meaning of r_symbolnum depends on r_extern. */
        /* r_symbolnum 的含义与 r_extern 有关。*/
199     unsigned int r_symbolnum:24;
200     /* Nonzero means value is a pc-relative offset
        and it should be relocated for changes in its own address
201     as well as for changes in the symbol or section specified. */
        /* 非零意味着值是一个 pc 相关的偏移值，因而需要被重定位到自己
           的地址处以及符号或节指定的改变。*/ [??]
202     unsigned int r_pcrel:1;
203     /* Length (as exponent of 2) of the field to be relocated.

```

```

204      Thus, a value of 2 indicates 1<<2 bytes.  */
      /* 需要被重定位的字段长度（是 2 的次方）。
         因此，若值是 2 则表示 1<<2 字节数。*/
205 unsigned int r_length:2;
206 /* 1 => relocate with value of symbol.
207      r_symbolnum is the index of the symbol
208      in file's the symbol table.
209      0 => relocate with the address of a segment.
210      r_symbolnum is N_TEXT, N_DATA, N_BSS or N_ABS
211      (the N_EXT bit may be set also, but signifies nothing).  */
      /* 1 => 以符号的值重定位。
         r_symbolnum 是文件符号表中符号的索引。
         0 => 以段的地址进行重定位。
         r_symbolnum 是 N_TEXT、N_DATA、N_BSS 或 N_ABS
         (N_EXT 比特位也可以被设置，但是毫无意义)。*/
212 unsigned int r_extern:1;
213 /* Four bits that aren't used, but when writing an object file
214      it is desirable to clear them.  */
      /* 没有使用的 4 个比特位，但是当进行写一个目标文件时
         最好将它们复位掉。*/
215 unsigned int r_pad:4;
216 };
217 #endif /* no N_RELOCATION_INFO_DECLARED.  */
218
219
220 #endif /* __A_OUT_GNU_H__  */
221

```

### 10.3.3 其它信息

#### a.out 执行文件格式

Linux 内核 0.11 版仅支持 a.out (Assembly out) 执行文件格式，虽然这种格式目前已经渐渐不用，而使用功能更为齐全的 ELF (Executable and Link Format) 格式，但是由于其简单性，作为学习入门的材料正好比较适用。下面全面介绍一下 a.out 格式。

在头文件<a.out.h>中声明了三个数据结构以及一些宏函数。这些数据结构描述了系统上可执行的机器码文件（二进制文件）。

一个执行文件共可有七个部分（七节）组成。按照顺序，这些部分是：

执行头部分(exec header)

执行文件头部分。该部分中含有一些参数，内核使用这些参数将执行文件加载到内存中并执行，而链接程序(ld)使用这些参数将一些二进制目标文件组合成一个可执行文件。这是唯一必要的组成部分。

代码段部分(text segment)

含有程序执行使被加载到内存中的指令代码和相关数据。可以以只读形式进行加载。

数据段部分(data segment)

这部分含有已经初始化过的数据，总是被加载到可读写的内存中。

代码重定位部分(text relocations)

这部分含有供链接程序使用的记录数据。在组合二进制目标文件时用于定位代码段中的指针或地址。

数据重定位部分(data relocations)

与代码重定位部分的作用类似，但是是用于数据段中指针的重定位。

符号表部分(symbol table)

这部分同样含有供链接程序使用的记录数据，用于在二进制目标文件之间对命名的变量和函数（符号）进行交叉引用。

字符串表部分(string table)

该部分含有与符号名相对应的字符串。

每个二进制执行文件均以执行数据结构（exec structure）开始。该数据结构的形式如下：

```
struct exec {
    unsigned long a_midmag;
    unsigned long a_text;
    unsigned long a_data;
    unsigned long a_bss;
    unsigned long a_syms;
    unsigned long a_entry;
    unsigned long a_trsize;
    unsigned long a_drsize;
};
```

各个字段的功能如下：

- |          |   |
|----------|---|
| a_midmag | 该字段含有被 N_GETFLAG()、N_GETMID 和 N_GETMAGIC() 访问的子部分，是由链接程序在运行时加载到进程地址空间。宏 N_GETMID() 用于返回机器标识符(machine-id)，指示出二进制文件将在什么机器上运行。N_GETMAGIC() 宏指明魔数，它唯一地确定了二进制执行文件与其它加载的文件之间的区别。字段中必须包含以下值之一：<br>OMAGIC 表示代码和数据段紧随在执行头后面并且是连续存放的。内核将代码和数据段都加载到可读写内存中。<br>NMAGIC 同 OMAGIC 一样，代码和数据段紧随在执行头后面并且是连续存放的。然而内核将代码加载到了只读内存中，并把数据段加载到代码段后下一页可读写内存边界开始。<br>ZMAGIC 内核在必要时从二进制执行文件中加载独立的页面。执行头部、代码段和数据段都被链接程序处理成多个页面大小的块。内核加载的代码页面时只读的，而数据段的页面是可写的。 |
| a_text   | 该字段含有代码段的长度值，字节数。   |
| a_data   | 该字段含有数据段的长度值，字节数。   |



|          |   |
|----------|---|
| a_bss    | 含有‘bss 段’的长度，内核用其设置在数据段后初始的 break (brk)。内核在加载程序时，这段可写内存显现出处于数据段后面，并且初始时为全零。 |
| a_syms   | 含有符号表部分的字节长度值。  |
| a_entry  | 含有内核将执行文件加载到内存中以后，程序执行起始点的内存地址。   |
| a_trsize | 该字段含有代码重定位表的大小，是字节数。  |
| a_drsize | 该字段含有数据重定位表的大小，是字节数。  |

在 a.out.h 头文件中定义了几个宏，这些宏使用 exec 结构来测试一致性或者定位执行文件中各个部分（节）的位置偏移值。这些宏有：

N\_BADMAG(exec) 如果 a\_magic 字段不能被识别，则返回非零值。

N\_TXTOFF(exec) 代码段的起始位置字节偏移值。

N\_DATOFF(exec) 数据段的起始位置字节偏移值。

N\_DRELOFF(exec) 数据重定位表的起始位置字节偏移值。

N\_TRELOFF(exec) 代码重定位表的起始位置字节偏移值。

N\_SYMOFF(exec) 符号表的起始位置字节偏移值。

N\_STROFF(exec) 字符串表的起始位置字节偏移值。

重定位记录具有标准格式，它使用重定位信息(relocation\_info)结构来描述：

```
struct relocation_info {
    int            r_address;
    unsigned int   r_symbolnum : 24,
                  r_pcrel : 1,
                  r_length : 2,
                  r_extern : 1,
                  r_baserel : 1,
                  r_jmptable : 1,
                  r_relative : 1,
                  r_copy : 1;
};
```

该结构中各字段的含义如下：

**r\_address** 该字段含有需要链接程序处理（编辑）的指针的字节偏移值。代码重定位的偏移值是从代码段开始处计数的，数据重定位的偏移值是从数据段开始处计算的。链接程序会将已经存储在该偏移处的值与使用重定位记录计算出的新值相加。

**r\_symbolnum** 该字段含有符号表中一个符号结构的序号值（不是字节偏移值）。链接程序在算出符号的绝对地址以后，就将该地址加到正在进行重定位的指针上。（如果 **r\_extern** 比特位是 0，那么情况就不同，见下面。）

**r\_pcrel** 如果设置了该位，链接程序就认为正在更新一个指针，该指针使用 **pc** 相关寻址方式，是属于机器码指令部分。当运行程序使用这个被重定位的指针时，该指针的地址被隐式地加到该指针上。

**r\_length** 该字段含有指针长度的 2 的次方值：0 表示 1 字节长，1 表示 2 字节长，2 表示 4 字节长。

**r\_extern** 如果被置位，表示该重定位需要一个外部引用；此时链接程序必须使用一个符号地址来更新相应指针。当该位是 0 时，则重定位是“局部”的；链接程序更新指针以反映各个段加载地址中的变化，而不是反映一个符号值的变化（除非同时设置了 **r\_baserel**，见下面）。在这种情况下，**r\_symbolnum** 字段的内容是一个 **n\_type** 值（见下面）；这类字段告诉链接程序被重定位的指针指向那个段。

**r\_baserel** 如果设置了该位，则 **r\_symbolnum** 字段指定的符号将被重定位成全局偏移表 (Global Offset Table) 中的一个偏移值。在运行时刻，全局偏移表该偏移处被设置为符号的地址。

**r\_jmptable** 如果被置位，则 **r\_symbolnum** 字段指定的符号将被重定位成过程链接表 (Procedure Linkage Table) 中的一个偏移值。

**r\_relative** 如果被置位，则说明此重定位与该目标文件将成为其组成部分的映象文件在运行时被加载的地址相关。这类重定位仅在共享目标文件中出现。

**r\_copy** 如果被置位，该重定位记录指定了一个符号，该符号的内容将被复制到 **r\_address** 指定的地方。该复制操作是通过共享目标模块中一个合适的数据项中的运行时刻链接程序完成的。

符号将名称映射为地址（或者更通俗地讲是字符串映射到值）。由于链接程序对地址的调整，一个符号的名称必须用来表示其地址，直到已被赋予一个绝对地址值。符号是由符号表中固定长度的记录以及字符串表中的可变长度名称组成。符号表是 **nlist** 结构的一个数组，如下所示：

```
struct nlist {
    union {
        char    *n_name;
        long    n_strx;
    } n_un;
    unsigned char n_type;
    char          n_other;
    short         n_desc;
```

```
        unsigned long    n_value;
    };
```

其中各字段的含义为：

`n_un.n_strx` 含有本符号的名称在字符串表中的字节偏移值。当程序使用 `nlist()` 函数访问一个符号表时，该字段被替换为 `n_un.n_name` 字段，这是内存中字符串的指针。

`n_type` 用于链接程序确定如何更新符号的值。使用位屏蔽(bitmasks)可以将 `n_type` 字段分割成三个子字段，对于 `N_EXT` 类型位置位的符号，链接程序将它们看作是“外部的”符号，并且允许其它二进制目标文件对它们的引用。`N_TYPE` 屏蔽码用于链接程序感兴趣的比特位：

`N_UNDF` 一个未定义的符号。链接程序必须在其它二进制目标文件中定位一个具有相同名称的外部符号，以确定该符号的绝对数据值。特殊情况下，如果 `n_type` 字段是非零值，并且没有二进制文件定义了这个符号，则链接程序在 BSS 段中将该符号解析为一个地址，保留长度等于 `n_value` 的字节。如果符号在多于一个二进制目标文件中都没有定义并且这些二进制目标文件对其长度值不一致，则链接程序将选择所有二进制目标文件中最大的长度。

`N_ABS` 一个绝对符号。链接程序不会更新一个绝对符号。

`N_TEXT` 一个代码符号。该符号的值是代码地址，链接程序在合并二进制目标文件时会更新其值。

`N_DATA` 一个数据符号；与 `N_TEXT` 类似，但是用于数据地址。对应代码和数据符号的值不是文件的偏移值而是地址；为了找出文件的偏移，就有必要确定相关部分开始加载的地址并减去它，然后加上该部分的偏移。

`N_BSS` 一个 BSS 符号；与代码或数据符号类似，但在二进制目标文件中没有对应的偏移。

`N_FN` 一个文件名符号。在合并二进制目标文件时，链接程序会将该符号插入在二进制文件中的符号之前。符号的名称就是给予链接程序的文件名，而其值是二进制文件中首个代码段地址。链接和加载时不需要文件名符号，但对于调式程序非常有用。

`N_STAB` 屏蔽码用于选择符号调式程序(例如 `gdb`)感兴趣的位；其值在 `stab()` 中说明。

`n_other` 该字段按照 `n_type` 确定的段，提供有关符号重定位操作的符号独立性信息。目前，`n_other` 字段的最低 4 位含有两个值之一：`AUX_FUNC` 和 `AUX_OBJECT`（有关定义参见 `<link.h>`）。`AUX_FUNC` 将符号与可调用的函数相关，`AUX_OBJECT` 将符号与数据相关，而不管它们是位于代码段还是数据段。该字段主要用于链接程序 `ld`，用于动态可执行程序创建。

`n_desc` 保留给调式程序使用；链接程序不对其进行处理。不同的调试程序将该字段用作不同的用途。

`n_value` 含有符号的值。对于代码、数据和 BSS 符号，这是一个地址；对于其它符号（例如调式程序符号），值可以是任意的。

字符串表是由长度为 `u_int32_t` 后跟一 `null` 结尾的符号字符串组成。长度代表整个表的字节大小，所以在 32 位的机器上其最小值（或者是第 1 个字符串的偏移）总是 4。

## 10.4 const.h 文件

### 10.4.1 功能描述

该文件中定义了 `i` 节点中文件属性和类型 `i_mode` 字段所用的一些标志位常量符号。

### 10.4.2 代码注释

列表 linux/include/const.h 文件

---

```

1 #ifndef CONST\_H
2 #define CONST\_H
3
4 #define BUFFER\_END 0x200000 // 定义缓冲使用内存的末端(代码中没有使用该常量)。
5
6 // i 节点数据结构中 i_mode 字段的各标志位。
7 #define I\_TYPE 0170000 // 指明 i 节点类型。
8 #define I\_DIRECTORY 0040000 // 是目录文件。
9 #define I\_REGULAR 0100000 // 常规文件，不是目录文件或特殊文件。
10 #define I\_BLOCK\_SPECIAL 0060000 // 块设备特殊文件。
11 #define I\_CHAR\_SPECIAL 0020000 // 字符设备特殊文件。
12 #define I\_NAMED\_PIPE 0010000 // 命名管道。
13 #define I\_SET\_UID\_BIT 0004000 // 在执行时设置有效用户 id 类型。
14 #define I\_SET\_GID\_BIT 0002000 // 在执行时设置有效组 id 类型。
15 #endif
16

```

---

## 10.5 ctype.h 文件

### 10.5.1 功能描述

该文件定义了一些有关字符类型判断和转换的宏，是使用数组（表）进行操作的。当使用宏时，字符是作为一个表（`__ctype`）中的索引，从表中获取一个字节，于是可得到相关的比特位。

### 10.5.2 代码注释

列表 linux/include/ctype.h 文件

---

```

1 #ifndef CTYPE_H
2 #define CTYPE_H
3
4 #define U      0x01    /* upper */           // 该比特位用于大写字符[A-Z]。
5 #define L      0x02    /* lower */           // 该比特位用于小写字符[a-z]。
6 #define D      0x04    /* digit */           // 该比特位用于数字[0-9]。
7 #define C      0x08    /* cntrl */           // 该比特位用于控制字符。
8 #define P      0x10    /* punct */           // 该比特位用于标点字符。
9 #define S      0x20    /* white space (space/lf/tab) */ // 用于空白字符，如空格、\t、\n 等。
10 #define X      0x40    /* hex digit */        // 该比特位用于十六进制数字。
11 #define SP     0x80    /* hard space (0x20) */   // 该比特位用于空格字符(0x20)。
12
13 extern unsigned char ctype[];           // 字符特性数组(表)，定义了各个字符对应上面的属性。
14 extern char ctmp;                       // 一个临时字符变量(在 fs/ctype.c 中定义)。
15
16 // 下面是一些确定字符类型的宏。
17 #define isalnum(c) ((ctype+1)[c]&(U|L|D)) // 是字符或数字[A-Z]、[a-z]或[0-9]。
18 #define isalpha(c) ((ctype+1)[c]&(U|L))      // 是字符。
19 #define iscntrl(c) ((ctype+1)[c]&(C))          // 是控制字符。
20 #define isdigit(c) ((ctype+1)[c]&(D))          // 是数字。
21 #define isgraph(c) ((ctype+1)[c]&(P|U|L|D)) // 是图形字符。
22 #define islower(c) ((ctype+1)[c]&(L))          // 是小写字符。
23 #define isprint(c) ((ctype+1)[c]&(P|U|L|D|SP)) // 是可打印字符。
24 #define ispunct(c) ((ctype+1)[c]&(P))          // 是标点符号。
25 #define isspace(c) ((ctype+1)[c]&(S))          // 是空白字符如空格、\f、\n、\r、\t、\v。
26 #define isupper(c) ((ctype+1)[c]&(U))          // 是大写字符。
27 #define isxdigit(c) ((ctype+1)[c]&(D|X))       // 是十六进制数字。
28
29 #define isascii(c) (((unsigned) c)<=0x7f)         // 是 ASCII 字符。
30 #define toascii(c) (((unsigned) c)&0x7f)          // 转换成 ASCII 字符。
31
32 #define tolower(c) (ctmp=c, isupper(ctmp)? ctmp-'A'+ctmp-'a': ctmp) // 转换成对应小写字符。
33 #define toupper(c) (ctmp=c, islower(ctmp)? ctmp-'a'+ctmp-'A': ctmp) // 转换成对应大写字符。
34
35 #endif

```

---

## 10.6 errno.h 文件

### 10.6.1 功能描述

### 10.6.2 代码注释

列表 linux/include/errno.h 文件

---

```

1 #ifndef ERRNO_H
2 #define ERRNO_H
3
4 /*

```

```

5  * ok, as I hadn't got any other source of information about
6  * possible error numbers, I was forced to use the same numbers
7  * as minix.
8  * Hopefully these are posix or something. I wouldn't know (and posix
9  * isn't telling me - they want $$$ for their f***ing standard).
10 *
11 * We don't use the _SIGN cludge of minix, so kernel returns must
12 * see to the sign by themselves.
13 *
14 * NOTE! Remember to change strerror() if you change this file!
15 */
/*
* ok, 由于我没有得到任何其它有关出错号的资料, 我只能使用与 minix 系统
* 相同的出错号了。
* 希望这些是 POSIX 兼容的或者在一定程度上是这样的, 我不知道 (而且 POSIX
* 没有告诉我 - 要获得他们的混蛋标准需要出钱)。
*
* 我们没有使用 minix 那样的 _SIGN 簇, 所以内核的返回值必须自己辨别正负号。
*
* 注意! 如果你改变该文件的话, 记着也要修改 strerror() 函数。
*/
16
17 extern int errno;
18
19 #define ERROR          99          // 一般错误。
20 #define EPERM          1          // 操作没有许可。
21 #define ENOENT         2          // 文件或目录不存在。
22 #define ESRCH          3          // 指定的进程不存在。
23 #define EINTR          4          // 中断的函数调用。
24 #define EIO            5          // 输入/输出错。
25 #define ENXIO          6          // 指定设备或地址不存在。
26 #define E2BIG          7          // 参数列表太长。
27 #define ENOEXEC        8          // 执行程序格式错误。
28 #define EBADF          9          // 文件句柄(描述符)错误。
29 #define ECHILD        10         // 子进程不存在。
30 #define EAGAIN        11         // 资源暂时不可用。
31 #define ENOMEM        12         // 内存不足。
32 #define EACCES        13         // 没有许可权限。
33 #define EFAULT        14         // 地址错。
34 #define ENOTBLK       15         // 不是块设备文件。
35 #define EBUSY         16         // 资源正忙。
36 #define EEXIST        17         // 文件已存在。
37 #define EXDEV         18         // 非法连接。
38 #define ENODEV        19         // 设备不存在。
39 #define ENOTDIR       20         // 不是目录文件。
40 #define EISDIR        21         // 是目录文件。
41 #define EINVAL       22         // 参数无效。
42 #define ENFILE        23         // 系统打开文件数太多。
43 #define EMFILE        24         // 打开文件数太多。
44 #define ENOTTY        25         // 不恰当的 IO 控制操作(没有 tty 终端)。
45 #define ETXTBSY       26         // 不再使用。
46 #define EFBIG         27         // 文件太大。
47 #define ENOSPC        28         // 设备已满(设备已经没有空间)。

```

---

```

48 #define ESPIPE          29          // 无效的文件指针重定位。
49 #define EROFS           30          // 文件系统只读。
50 #define EMLINK          31          // 连接太多。
51 #define EPIPE           32          // 管道错。
52 #define EDOM            33          // 域(domain)出错。
53 #define ERANGE          34          // 结果太大。
54 #define EDEADLK         35          // 避免资源死锁。
55 #define ENAMETOOLONG    36          // 文件名太长。
56 #define ENOLCK          37          // 没有锁定可用。
57 #define ENOSYS          38          // 功能还没有实现。
58 #define ENOTEMPTY       39          // 目录不空。
59
60 #endif
61

```

---

## 10.7 fcntl.h 文件

### 10.10.1 功能描述

### 10.7.2 代码注释

列表 linux/include/fcntl.h 文件

---

```

1 #ifndef FCNTL\_H
2 #define FCNTL\_H
3
4 #include <sys/types.h>    // 类型头文件。定义了基本的系统数据类型。
5
6 /* open/fcntl - NOCTTY, NDELAY isn't implemented yet */
7 /* open/fcntl - NOCTTY 和 NDELAY 现在还没有实现 */
8 #define O\_ACCMODE        00003      // 文件访问模式屏蔽码。
9 // 打开文件 open() 和文件控制 fcntl() 函数使用的文件访问模式。同时只能使用三者之一。
10 #define O\_RDONLY         00          // 以只读方式打开文件。
11 #define O\_WRONLY         01          // 以只写方式打开文件。
12 #define O\_RDWR          02          // 以读写方式打开文件。
13 // 下面是文件创建标志，用于 open()。可与上面访问模式用'位或'的方式一起使用。
14 #define O\_CREAT           00100      /* not fcntl */ // 如果文件不存在就创建。
15 #define O\_EXCL            00200      /* not fcntl */ // 独占使用文件标志。
16 #define O\_NOCTTY          00400      /* not fcntl */ // 不分配控制终端。
17 #define O\_TRUNC           01000      /* not fcntl */ // 若文件已存在且是写操作，则长度截为 0。
18 #define O\_APPEND          02000      // 以添加方式打开，文件指针置为文件尾。
19 #define O\_NONBLOCK        04000      /* not fcntl */ // 非阻塞方式打开和操作文件。
20 #define O\_NDELAY          O\_NONBLOCK // 非阻塞方式打开和操作文件。
21
22 /* Defines for fcntl-commands. Note that currently
23  * locking isn't supported, and other things aren't really
24  * tested.
25  */
26 /* 下面定义了 fcntl 的命令。注意目前锁定命令还没有支持，而其它

```

```

* 命令实际上还没有测试过。
*/
// 文件句柄(描述符)操作函数 fcntl() 的命令。
23 #define F_DUPFD      0      /* dup */           // 拷贝文件句柄为最小数值的句柄。
24 #define F_GETFD      1      /* get f_flags */  // 取文件句柄标志。
25 #define F_SETFD      2      /* set f_flags */  // 设置文件句柄标志。
26 #define F_GETFL      3      /* more flags (cloexec) */ // 取文件状态标志和访问模式。
27 #define F_SETFL      4      // 设置文件状态标志和访问模式。
// 下面是文件锁定命令。fcntl() 的第三个参数 lock 是指向 flock 结构的指针。
28 #define F_GETLK      5      /* not implemented */ // 返回阻止锁定的 flock 结构。
29 #define F_SETLK      6      // 设置(F_RDLCK 或 F_WRLCK)或清除(F_UNLCK)锁定。
30 #define F_SETLKW     7      // 等待设置或清除锁定。
31
32 /* for F_[GET/SET]FL */
/* 用于 F_GETFL 或 F_SETFL */
// 在执行 exec() 簇函数时关闭文件句柄。(执行时关闭 - Close On EXECution)
33 #define FD_CLOEXEC    1      /* actually anything with low bit set goes */
/* 实际上只要低位为 1 即可 */
34
35 /* Ok, these are locking features, and aren't implemented at any
36  * level. POSIX wants them.
37  */
/* OK, 以下是锁定类型, 任何函数中都还没有实现。POSIX 标准要求这些类型。
*/
38 #define F_RDLCK       0      // 共享或读文件锁定。
39 #define F_WRLCK       1      // 独占或写文件锁定。
40 #define F_UNLCK       2      // 文件解锁。
41
42 /* Once again - not implemented, but ... */
/* 同样 - 也还没有实现, 但是... */
// 文件锁定操作数据结构。描述了受影响文件段的类型(l_type)、开始偏移(l_whence)、
// 相对偏移(l_start)、锁定长度(l_len)和实施锁定的进程 id。
43 struct flock {
44     short l_type;          // 锁定类型 (F_RDLCK, F_WRLCK, F_UNLCK)。
45     short l_whence;        // 开始偏移 (SEEK_SET, SEEK_CUR 或 SEEK_END)。
46     off_t l_start;         // 阻塞锁定的开始处。相对偏移 (字节数)。
47     off_t l_len;          // 阻塞锁定的大小; 如果是 0 则为到文件末尾。
48     pid_t l_pid;          // 加锁的进程 id。
49 };
50
// 以下是使用上述标志或命令的函数原型。
// 创建新文件或重写一个已存在文件。
// 参数 filename 是欲创建文件的文件名, mode 是创建文件的属性 (参见 include/sys/stat.h)。
51 extern int creat(const char * filename, mode_t mode);
// 文件句柄操作, 会影响文件的打开。
// 参数 fildes 是文件句柄, cmd 是操作命令, 见上面 23-30 行。
52 extern int fcntl(int fildes, int cmd, ...);
// 打开文件。在文件与文件句柄之间建立联系。
// 参数 filename 是欲打开文件的文件名, flags 是上面 7-17 行上的标志的组合。
53 extern int open(const char * filename, int flags, ...);
54
55 #endif
56

```



## 10.8 signal.h 文件

### 10.8.1 功能描述

### 10.8.2 文件注释

列表 linux/include/signal.h 文件

```

1 #ifndef \_SIGNAL\_H
2 #define \_SIGNAL\_H
3
4 #include <sys/types.h> // 类型头文件。定义了基本的系统数据类型。
5
6 typedef int sig\_atomic\_t; // 定义信号原子操作类型。
7 typedef unsigned int sigset\_t; /* 32 bits */ // 定义信号集类型。
8
9 #define NSIG 32 // 定义信号种类 -- 32 种。
10 #define NSIG \_NSIG // NSIG = _NSIG
11
12 // 以下这些是 Linux 0.11 内核中定义的信号。
13 #define SIGHUP 1 // Hang Up -- 挂断控制终端或进程。
14 #define SIGINT 2 // Interrupt -- 来自键盘的中断。
15 #define SIGQUIT 3 // Quit -- 来自键盘的退出。
16 #define SIGILL 4 // Illegal -- 非法指令。
17 #define SIGTRAP 5 // Trap -- 跟踪断点。
18 #define SIGABRT 6 // Abort -- 异常结束。
19 #define SIGIOT 6 // IO Trap -- 同上。
20 #define SIGUNUSED 7 // Unused -- 没有使用。
21 #define SIGFPE 8 // FPE -- 协处理器出错。
22 #define SIGKILL 9 // Kill -- 强迫进程终止。
23 #define SIGUSR1 10 // User1 -- 用户信号 1，进程可使用。
24 #define SIGSEGV 11 // Segment Violation -- 无效内存引用。
25 #define SIGUSR2 12 // User2 -- 用户信号 2，进程可使用。
26 #define SIGPIPE 13 // Pipe -- 管道写出错，无读者。
27 #define SIGALRM 14 // Alarm -- 实时定时器报警。
28 #define SIGTERM 15 // Terminate -- 进程终止。
29 #define SIGSTKFLT 16 // Stack Fault -- 栈出错（协处理器）。
30 #define SIGCHLD 17 // Child -- 子进程停止或被终止。
31 #define SIGCONT 18 // Continue -- 恢复进程继续执行。
32 #define SIGSTOP 19 // Stop -- 停止进程的执行。
33 #define SIGTSTP 20 // TTY Stop -- tty 发出停止进程，可忽略。
34 #define SIGTTIN 21 // TTY In -- 后台进程请求输入。
35 #define SIGTTOU 22 // TTY Out -- 后台进程请求输出。
36
37 /* Ok, I haven't implemented sigactions, but trying to keep headers POSIX */
38 /* OK, 我还没有实现 sigactions 的编制，但在头文件中仍希望遵守 POSIX 标准 */
39 #define SA\_NOCLDSTOP 1 // 当子进程处于停止状态，就不对 SIGCHLD 处理。
40 #define SA\_NOMASK 0x40000000 // 不阻止在指定的信号处理程序(信号句柄)中再收到该信号。

```

```

39 #define SA_ONESHOT      0x80000000 // 信号句柄一旦被调用过就恢复到默认处理句柄。
40 // 以下参数用于 sigprocmask() -- 改变阻塞信号集(屏蔽码)。这些参数可以改变该函数的行为。
41 #define SIG_BLOCK        0      /* for blocking signals */
42 // 在阻塞信号集中加上给定的信号集。
43 #define SIG_UNBLOCK      1      /* for unblocking signals */
44 // 从阻塞信号集中删除指定的信号集。
45 #define SIG_SETMASK      2      /* for setting the signal mask */
46 // 设置阻塞信号集(信号屏蔽码)。
47
48 #define SIG_DFL          ((void (*)(int))0) /* default signal handling */
49 // 默认的信号处理程序(信号句柄)。
50 #define SIG_IGN          ((void (*)(int))1) /* ignore signal */
51 // 忽略信号的处理程序。
52
53 // 下面是 sigaction 的数据结构。
54 // sa_handler 是对应某信号指定要采取的行动。可以是上面的 SIG_DFL, 或者是 SIG_IGN 来忽略
55 // 该信号, 也可以是指向处理该信号的一个指针。
56 // sa_mask 给出了对信号的屏蔽码, 在信号程序执行时将阻塞对这些信号的处理。
57 // sa_flags 指定改变信号处理过程的信号集。它是由 37-39 行的位标志定义的。
58 // 另外, 引起触发信号处理的信号也将被阻塞, 除非使用了 SA_NOMASK 标志。
59 struct sigaction {
60     void (*sa_handler)(int);
61     sigset_t sa_mask;
62     int sa_flags;
63     void (*sa_restorer)(void);
64 };
65
66 // 为信号_sig 安装一个新的信号处理程序(信号句柄), 与 sigaction() 类似。
67 void (*signal(int _sig, void (*_func)(int)))(int);
68 // 向当前进程发送一个信号。其作用等价于 kill(getpid(), sig)。
69 int raise(int sig);
70 // 可用于向任何进程组或进程发送任何信号。
71 int kill(pid_t pid, int sig);
72 // 向信号集中添加信号。
73 int sigaddset(sigset_t *mask, int signo);
74 // 从信号集中去除指定的信号。
75 int sigdelset(sigset_t *mask, int signo);
76 // 从信号集中清除指定信号集。
77 int sigemptyset(sigset_t *mask);
78 // 向信号集中置入所有信号。
79 int sigfillset(sigset_t *mask);
80 // 判断一个信号是否是信号集中的。1 -- 是, 0 -- 不是, -1 -- 出错。
81 int sigismember(sigset_t *mask, int signo); /* 1 - is, 0 - not, -1 error */
82 // 对 set 中的信号进行检测, 看是否有挂起的信号。
83 int sigpending(sigset_t *set);
84 // 改变目前的被阻塞信号集(信号屏蔽码)。
85 int sigprocmask(int how, sigset_t *set, sigset_t *oldset);
86 // 用 sigmask 临时替换进程的信号屏蔽码, 然后暂停该进程直到收到一个信号。
87 int sigsuspend(sigset_t *sigmask);
88 // 用于改变进程在收到指定信号时所采取的行动。
89 int sigaction(int sig, struct sigaction *act, struct sigaction *oldact);
90

```

```
68 #endif /* _SIGNAL_H */
69
```

## 10.9 stdarg.h 文件

### 10.9.1 功能描述

stdarg.h 是标准参数头文件。它以宏的形式定义变量参数列表。主要说明了一个类型(va\_list)和三个宏(va\_start, va\_arg 和 va\_end)，用于 vsprintf、vprintf、vfprintf 函数。

### 10.9.2 代码注释

列表 linux/include/stdarg.h 文件

```
1 #ifndef STDARG_H
2 #define STDARG_H
3
4 typedef char *va_list; // 定义 va_list 是一个字符指针类型。
5
6 /* Amount of space required in an argument list for an arg of type TYPE.
7    TYPE may alternatively be an expression whose type is used. */
8 /* 下面给出了类型为 TYPE 的 arg 参数列表所要求的空间容量。
9    TYPE 也可以是使用该类型的一个表达式 */
10
11 // 下面这句定义了取整后的 TYPE 类型的字节长度值。是 int 长度(4)的倍数。
12 #define va_rounded_size(TYPE) \
13 (((sizeof (TYPE) + sizeof (int) - 1) / sizeof (int)) * sizeof (int))
14
15 // 下面这个函数（用宏实现）使 AP 指向传给函数的可变参数表的第一个参数。
16 // 在第一次调用 va_arg 或 va_end 之前，必须首先调用该函数。
17 #ifndef __sparc__
18 #define va_start(AP, LASTARG) \
19 (AP = ((char *) &(LASTARG) + va_rounded_size (LASTARG)))
20 #else
21 #define va_start(AP, LASTARG) \
22 (__builtin_saveregs (), \
23 AP = ((char *) &(LASTARG) + va_rounded_size (LASTARG)))
24 #endif
25
26 // 下面该宏用于被调用函数完成一次正常返回。va_end 可以修改 AP 使其在重新调用
27 // va_start 之前不能被使用。va_end 必须在 va_arg 读完所有的参数后再被调用。
28 void va_end (va_list); /* Defined in gnu lib */ /* 在 gnu lib 中定义 */
29 #define va_end(AP)
30
31 // 下面该宏用于扩展表达式使其与下一个被传递参数具有相同的类型和值。
32 // 对于缺省值，va_arg 可以用字符、无符号字符和浮点类型。
33 // 在第一次使用 va_arg 时，它返回表中的第一个参数，后续的每次调用都将返回表中的
34 // 下一个参数。这是通过先访问 AP，然后把它增加以指向下一项来实现的。
35 // va_arg 使用 TYPE 来完成访问和定位下一项，每调用一次 va_arg，它就修改 AP 以指示
36 // 表中的下一参数。
37 #define va_arg(AP, TYPE) \
```

---

```

25 (AP += \_\_va\_rounded\_size (TYPE),
26 *((TYPE *) (AP - \_\_va\_rounded\_size (TYPE))))
27
28 #endif /* \_STDARG\_H */
29

```

---

## 10.10 stddef.h 文件

### 10.10.1 功能描述

该文件定义了一些常用的类型和宏。内核中很少使用该文件。

### 10.10.2 代码注释

列表 linux/include/stddef.h 文件

---

```

1 #ifndef STDEDEF\_H
2 #define STDEDEF\_H
3
4 #ifndef PTRDIFF\_T
5 #define PTRDIFF\_T
6 typedef long ptrdiff\_t;           // 两个指针相减结果的类型。
7 #endif
8
9 #ifndef SIZE\_T
10 #define SIZE\_T
11 typedef unsigned long size\_t;    // sizeof 返回的类型。
12 #endif
13
14 #undef NULL
15 #define NULL ((void *)0)          // 空指针。
16
17 #define offsetof(TYPE, MEMBER) ((size\_t) &((TYPE *)0)->MEMBER) // 成员在类型中的偏移位置。
18
19 #endif
20

```

---

## 10.11 string.h 文件

### 10.11.1 功能描述

该头文件中以内嵌函数的形式定义了所有字符串操作函数，为了提高执行速度使用了内嵌汇编程序。

### 10.11.2 代码注释

列表 linux/include/string.h 文件

```

1 #ifndef STRING\_H
2 #define STRING\_H
3
4 #ifndef NULL
5 #define NULL ((void *) 0)
6 #endif
7
8 #ifndef SIZE\_T
9 #define SIZE\_T
10 typedef unsigned int size\_t;
11 #endif
12
13 extern char * strerror(int errno);
14
15 /*
16  * This string-include defines all string functions as inline
17  * functions. Use gcc. It also assumes ds=es=data space, this should be
18  * normal. Most of the string-functions are rather heavily hand-optimized,
19  * see especially strtok, strstr, str[c]spn. They should work, but are not
20  * very easy to understand. Everything is done entirely within the register
21  * set, making the functions fast and clean. String instructions have been
22  * used through-out, making for "slightly" unclear code :-)
23  *
24  * (C) 1991 Linus Torvalds
25  */
26
27 /*
28  * 这个字符串头文件以内嵌函数的形式定义了所有字符串操作函数。使用 gcc 时，同时
29  * 假定了 ds=es=数据空间，这应该是常规的。绝大多数字符串函数都是经手工进行大量
30  * 优化的，尤其是函数 strtok、strstr、str[c]spn。它们应该能正常工作，但却不是那
31  * 么容易理解。所有的操作基本上都是使用寄存器集来完成的，这使得函数即快有整洁。
32  * 所有地方都使用了字符串指令，这又使得代码“稍微”难以理解☺
33  *
34  * (C) 1991 Linus Torvalds
35  */
36
37
38 // 将一个字符串(src)拷贝到另一个字符串(dest)，直到遇到 NULL 字符后停止。
39 // 参数: dest - 目的字符串指针, src - 源字符串指针。
40 // %0 - esi(src), %1 - edi(dest)。
41 extern inline char * strcpy(char * dest, const char *src)
42 {
43     __asm__(
44         "cld\n" // 清方向位。
45         "1:|t lodsb\n|t" // 加载 DS:[esi]处 1 字节→al, 并更新 esi。
46         "stosb\n|t" // 存储字节 al→ES:[edi], 并更新 edi。
47         "testb %%al, %%al\n|t" // 刚存储的字节是 0?
48         "jne 1b" // 不是则向后跳转到标号 1 处, 否则结束。
49         :: "S" (src), "D" (dest): "si", "di", "ax");
50     return dest; // 返回目的字符串指针。
51 }
52
53 // 拷贝源字符串 count 个字节到目的字符串。
54 // 如果源串长度小于 count 个字节, 就附加空字符(NULL)到目的字符串。
55 // 参数: dest - 目的字符串指针, src - 源字符串指针, count - 拷贝字节数。
56 // %0 - esi(src), %1 - edi(dest), %2 - ecx(count)。

```

```

38 extern inline char * strncpy(char * dest, const char *src, int count)
39 {
40     __asm__( "cld\n" // 清方向位。
41             "l:|tdecl %2\n|t" // 寄存器 ecx-- (count--)。
42             "js 2f\n|t" // 如果 count<0 则向前跳转到标号 2，结束。
43             "lods b\n|t" // 取 ds:[esi]处 1 字节→al，并且 esi++。
44             "stos b\n|t" // 存储该字节→es:[edi]，并且 edi++。
45             "testb %%al, %%al\n|t" // 该字节是 0？
46             "jne 1b\n|t" // 不是，则向前跳转到标号 1 处继续拷贝。
47             "rep\n|t" // 否则，在目的串中存放剩余个数的空字符。
48             "stos b\n"
49             "2:"
50             :: "S" (src), "D" (dest), "c" (count): "si", "di", "ax", "cx");
51 return dest; // 返回目的字符串指针。
52 }
53
54 // 将源字符串拷贝到目的字符串的末尾处。
55 // 参数：dest - 目的字符串指针，src - 源字符串指针。
56 // %0 - esi(src)，%1 - edi(dest)，%2 - eax(0)，%3 - ecx(-1)。
57 extern inline char * strcat(char * dest, const char * src)
58 {
59     __asm__( "cld\n|t" // 清方向位。
60             "repne\n|t" // 比较 al 与 es:[edi]字节，并更新 edi++，
61             "scas b\n|t" // 直到找到目的串中是 0 的字节，此时 edi 已经指向后 1 字节。
62             "decl %1\n|t" // 让 es:[edi]指向 0 值字节。
63             "l:|tlods b\n|t" // 取源字符串字节 ds:[esi]→al，并 esi++。
64             "stos b\n|t" // 将该字节存到 es:[edi]，并 edi++。
65             "testb %%al, %%al\n|t" // 该字节是 0？
66             "jne 1b" // 不是，则向后跳转到标号 1 处继续拷贝，否则结束。
67             :: "S" (src), "D" (dest), "a" (0), "c" (0xffffffff): "si", "di", "ax", "cx");
68 return dest; // 返回目的字符串指针。
69 }
70
71 // 将源字符串的 count 个字节复制到目的字符串的末尾处，最后添一空字符。
72 // 参数：dest - 目的字符串，src - 源字符串，count - 欲复制的字节数。
73 // %0 - esi(src)，%1 - edi(dest)，%2 - eax(0)，%3 - ecx(-1)，%4 - (count)。
74 extern inline char * strncat(char * dest, const char * src, int count)
75 {
76     __asm__( "cld\n|t" // 清方向位。
77             "repne\n|t" // 比较 al 与 es:[edi]字节，edi++。
78             "scas b\n|t" // 直到找到目的串的末端 0 值字节。
79             "decl %1\n|t" // edi 指向该 0 值字节。
80             "movl %4, %3\n" // 欲复制字节数→ecx。
81             "l:|tdecl %3\n|t" // ecx-- (从 0 开始计数)。
82             "js 2f\n|t" // ecx < 0 ?，是则向前跳转到标号 2 处。
83             "lods b\n|t" // 否则取 ds:[esi]处的字节→al，esi++。
84             "stos b\n|t" // 存储到 es:[edi]处，edi++。
85             "testb %%al, %%al\n|t" // 该字节值为 0？
86             "jne 1b\n" // 不是则向后跳转到标号 1 处，继续复制。
87             "2:|txorl %2, %2\n|t" // 将 al 清零。
88             "stos b" // 存到 es:[edi]处。
89             :: "S" (src), "D" (dest), "a" (0), "c" (0xffffffff), "g" (count)
90             : "si", "di", "ax", "cx");

```

```

85 return dest;                // 返回目的字符串指针。
86 }
87
88 // 将一个字符串与另一个字符串进行比较。
89 // 参数: cs - 字符串 1, ct - 字符串 2。
90 // %0 - eax(__res)返回值, %1 - edi(cs)字符串 1 指针, %2 - esi(ct)字符串 2 指针。
91 // 返回: 如果串 1 > 串 2, 则返回 1; 串 1 = 串 2, 则返回 0; 串 1 < 串 2, 则返回-1。
88 extern inline int strcmp(const char * cs, const char * ct)
89 {
90     register int __res __asm__("ax");    // __res 是寄存器变量(eax)。
91     __asm__("cld\n"                    // 清方向位。
92             "1:\tlodsb\n\t"            // 取字符串 2 的字节 ds:[esi]→al, 并且 esi++。
93             "scasb\n\t"                // al 与字符串 1 的字节 es:[edi]作比较, 并且 edi++。
94             "jne 2f\n\t"                // 如果不相等, 则向前跳转到标号 2。
95             "testb %%al, %%al\n\t"      // 该字节是 0 值字节吗(字符串结尾)?
96             "jne 1b\n\t"                // 不是, 则向后跳转到标号 1, 继续比较。
97             "xorl %%eax, %%eax\n\t"      // 是, 则返回值 eax 清零,
98             "jmp 3f\n\t"                // 向前跳转到标号 3, 结束。
99             "2:\tmovl $1, %%eax\n\t"      // eax 中置 1。
100            "jl 3f\n\t"                  // 若前面比较中串 2 字符<串 1 字符, 则返回正值, 结束。
101            "negl %%eax\n\t"              // 否则 eax = -eax, 返回负值, 结束。
102            "3:")
103     : "=a" (__res): "D" (cs), "S" (ct): "si", "di";
104     return __res;                      // 返回比较结果。
105 }
106
107 // 字符串 1 与字符串 2 的前 count 个字符进行比较。
108 // 参数: cs - 字符串 1, ct - 字符串 2, count - 比较的字符数。
109 // %0 - eax(__res)返回值, %1 - edi(cs)串 1 指针, %2 - esi(ct)串 2 指针, %3 - ecx(count)。
110 // 返回: 如果串 1 > 串 2, 则返回 1; 串 1 = 串 2, 则返回 0; 串 1 < 串 2, 则返回-1。
107 extern inline int strncmp(const char * cs, const char * ct, int count)
108 {
109     register int __res __asm__("ax");    // __res 是寄存器变量(eax)。
110     __asm__("cld\n"                    // 清方向位。
111             "1:\tdecl %3\n\t"            // count--。
112             "js 2f\n\t"                  // 如果 count<0, 则向前跳转到标号 2。
113             "lodsb\n\t"                  // 取串 2 的字符 ds:[esi]→al, 并且 esi++。
114             "scasb\n\t"                  // 比较 al 与串 1 的字符 es:[edi], 并且 edi++。
115             "jne 3f\n\t"                  // 如果不相等, 则向前跳转到标号 3。
116             "testb %%al, %%al\n\t"      // 该字符是 NULL 字符吗?
117             "jne 1b\n\t"                  // 不是, 则向后跳转到标号 1, 继续比较。
118             "2:\txorl %%eax, %%eax\n\t"  // 是 NULL 字符, 则 eax 清零(返回值)。
119             "jmp 4f\n\t"                  // 向前跳转到标号 4, 结束。
120             "3:\tmovl $1, %%eax\n\t"      // eax 中置 1。
121             "jl 4f\n\t"                  // 如果前面比较中串 2 字符<串 2 字符, 则返回 1, 结束。
122             "negl %%eax\n\t"              // 否则 eax = -eax, 返回负值, 结束。
123             "4:")
124     : "=a" (__res): "D" (cs), "S" (ct), "c" (count): "si", "di", "cx";
125     return __res;                      // 返回比较结果。
126 }
127
128 // 在字符串中寻找第一个匹配的字符。
129 // 参数: s - 字符串, c - 欲寻找的字符。

```



```

// %0 - eax(__res), %1 - esi(字符串指针 s), %2 - eax(字符 c)。
// 返回: 返回字符串中第一次出现匹配字符的指针。若没有找到匹配的字符, 则返回空指针。
128 extern inline char * strchr(const char * s, char c)
129 {
130     register char * __res __asm__( "ax" ); // __res 是寄存器变量(eax)。
131     __asm__( "cld\n\t" // 清方向位。
132             "movb %%al, %%ah\n\t" // 将欲比较字符移到 ah。
133             "1:\t lodsb\n\t" // 取字符串中字符 ds:[esi]→al, 并且 esi++。
134             "cmpb %%ah, %%al\n\t" // 字符串中字符 al 与指定字符 ah 相比较。
135             "je 2f\n\t" // 若相等, 则向前跳转到标号 2 处。
136             "testb %%al, %%al\n\t" // al 中字符是 NULL 字符吗? (字符串结尾?)
137             "jne 1b\n\t" // 若不是, 则向后跳转到标号 1, 继续比较。
138             "movl $1, %1\n\t" // 是, 则说明没有找到匹配字符, esi 置 1。
139             "2:\t movl %1, %0\n\t" // 将指向匹配字符后一个字节处的指针值放入 eax
140             "decl %0" // 将指针调整为指向匹配的字符。
141             : "=a" (__res): "S" (s), "" (c): "si" );
142     return __res; // 返回指针。
143 }
144
///// 寻找字符串中指定字符最后一次出现的地方。(反向搜索字符串)
// 参数: s - 字符串, c - 欲寻找的字符。
// %0 - edx(__res), %1 - edx(0), %2 - esi(字符串指针 s), %3 - eax(字符 c)。
// 返回: 返回字符串中最后一次出现匹配字符的指针。若没有找到匹配的字符, 则返回空指针。
145 extern inline char * strrchr(const char * s, char c)
146 {
147     register char * __res __asm__( "dx" ); // __res 是寄存器变量(edx)。
148     __asm__( "cld\n\t" // 清方向位。
149             "movb %%al, %%ah\n\t" // 将欲寻找的字符移到 ah。
150             "1:\t lodsb\n\t" // 取字符串中字符 ds:[esi]→al, 并且 esi++。
151             "cmpb %%ah, %%al\n\t" // 字符串中字符 al 与指定字符 ah 作比较。
152             "jne 2f\n\t" // 若不相等, 则向前跳转到标号 2 处。
153             "movl %%esi, %0\n\t" // 将字符指针保存到 edx 中。
154             "decl %0\n\t" // 指针后退一位, 指向字符串中匹配字符处。
155             "2:\t testb %%al, %%al\n\t" // 比较的字符是 0 吗(到字符串尾)?
156             "jne 1b" // 不是则向后跳转到标号 1 处, 继续比较。
157             : "=d" (__res): "" (0), "S" (s), "a" (c): "ax", "si" );
158     return __res; // 返回指针。
159 }
160
///// 在字符串 1 中寻找第 1 个字符序列, 该字符序列中的任何字符都包含在字符串 2 中。
// 参数: cs - 字符串 1 指针, ct - 字符串 2 指针。
// %0 - esi(__res), %1 - eax(0), %2 - ecx(-1), %3 - esi(串 1 指针 cs), %4 - (串 2 指针 ct)。
// 返回字符串 1 中包含字符串 2 中任何字符的首个字符序列的长度值。
161 extern inline int strspn(const char * cs, const char * ct)
162 {
163     register char * __res __asm__( "si" ); // __res 是寄存器变量(esi)。
164     __asm__( "cld\n\t" // 清方向位。
165             "movl %4, %%edi\n\t" // 首先计算串 2 的长度。串 2 指针放入 edi 中。
166             "repne\n\t" // 比较 al(0) 与串 2 中的字符(es:[edi]), 并 edi++。
167             "scasb\n\t" // 如果不相等就继续比较(ecx 逐步递减)。
168             "notl %%ecx\n\t" // ecx 中每位取反。
169             "decl %%ecx\n\t" // ecx--, 得串 2 的长度值。
170             "movl %%ecx, %%edx\n\t" // 将串 2 的长度值暂放入 edx 中。

```



```

171     "l:\tlodsb\n\t"           // 取串 1 字符 ds:[esi]→al, 并且 esi++.
172     "testb %%al, %%al\n\t"     // 该字符等于 0 值吗 (串 1 结尾)?
173     "je 2f\n\t"               // 如果是, 则向前跳转到标号 2 处。
174     "movl %4, %%edi\n\t"       // 取串 2 头指针放入 edi 中。
175     "movl %%edx, %%ecx\n\t"    // 再将串 2 的长度值放入 ecx 中。
176     "repne\n\t"               // 比较 al 与串 2 中字符 es:[edi], 并且 edi++.
177     "scasb\n\t"               // 如果不相等就继续比较。
178     "je 1b\n\t"               // 如果相等, 则向后跳转到标号 1 处。
179     "2:\tdecl %0"              // esi--, 指向最后一个包含在串 2 中的字符。
180     : "=S" (__res): "a" (0), "c" (0xffffffff), "" (cs), "g" (ct)
181     : "ax", "cx", "dx", "di");
182 return __res-cs;              // 返回字符序列的长度值。
183 }
184
185 // 寻找字符串 1 中不包含字符串 2 中任何字符的首个字符序列。
186 // 参数: cs - 字符串 1 指针, ct - 字符串 2 指针。
187 // %0 - esi(__res), %1 - eax(0), %2 - ecx(-1), %3 - esi(串 1 指针 cs), %4 - (串 2 指针 ct)。
188 // 返回字符串 1 中不包含字符串 2 中任何字符的首个字符序列的长度值。
189 extern inline int strcspn(const char * cs, const char * ct)
190 {
191     register char * __res __asm__("si"); // __res 是寄存器变量(esi)。
192     __asm__("cld\n\t"           // 清方向位。
193             "movl %4, %%edi\n\t" // 首先计算串 2 的长度。串 2 指针放入 edi 中。
194             "repne\n\t"         // 比较 al(0) 与串 2 中的字符 (es:[edi]), 并 edi++.
195             "scasb\n\t"         // 如果不相等就继续比较(ecx 逐步递减)。
196             "notl %%ecx\n\t"     // ecx 中每位取反。
197             "decl %%ecx\n\t"     // ecx--, 得串 2 的长度值。
198             "movl %%ecx, %%edx\n\t" // 将串 2 的长度值暂放入 edx 中。
199             "l:\tlodsb\n\t"     // 取串 1 字符 ds:[esi]→al, 并且 esi++.
200             "testb %%al, %%al\n\t" // 该字符等于 0 值吗 (串 1 结尾)?
201             "je 2f\n\t"         // 如果是, 则向前跳转到标号 2 处。
202             "movl %4, %%edi\n\t" // 取串 2 头指针放入 edi 中。
203             "movl %%edx, %%ecx\n\t" // 再将串 2 的长度值放入 ecx 中。
204             "repne\n\t"         // 比较 al 与串 2 中字符 es:[edi], 并且 edi++.
205             "scasb\n\t"         // 如果不相等就继续比较。
206             "jne 1b\n\t"        // 如果不相等, 则向后跳转到标号 1 处。
207             "2:\tdecl %0"       // esi--, 指向最后一个包含在串 2 中的字符。
208             : "=S" (__res): "a" (0), "c" (0xffffffff), "" (cs), "g" (ct)
209             : "ax", "cx", "dx", "di");
210 return __res-cs;              // 返回字符序列的长度值。
211 }
212
213 // 在字符串 1 中寻找首个包含在字符串 2 中的任何字符。
214 // 参数: cs - 字符串 1 的指针, ct - 字符串 2 的指针。
215 // %0 -esi(__res), %1 -eax(0), %2 -ecx(0xffffffff), %3 -esi(串 1 指针 cs), %4 -(串 2 指针 ct)。
216 // 返回字符串 1 中首个包含字符串 2 中字符的指针。
217 extern inline char * strpbrk(const char * cs, const char * ct)
218 {
219     register char * __res __asm__("si"); // __res 是寄存器变量(esi)。
220     __asm__("cld\n\t"           // 清方向位。
221             "movl %4, %%edi\n\t" // 首先计算串 2 的长度。串 2 指针放入 edi 中。
222             "repne\n\t"         // 比较 al(0) 与串 2 中的字符 (es:[edi]), 并 edi++.
223             "scasb\n\t"         // 如果不相等就继续比较(ecx 逐步递减)。

```

```

216     "notl %%ecx\n\t"           // ecx 中每位取反。
217     "decl %%ecx\n\t"           // ecx--, 得串 2 的长度值。
218     "movl %%ecx, %%edx\n\t"     // 将串 2 的长度值暂放入 edx 中。
219     "1:\t\tlodsb\n\t"           // 取串 1 字符 ds:[esi]→al, 并且 esi++。
220     "testb %%al, %%al\n\t"      // 该字符等于 0 值吗 (串 1 结尾)?
221     "je 2f\n\t"                 // 如果是, 则向前跳转到标号 2 处。
222     "movl %4, %%edi\n\t"        // 取串 2 头指针放入 edi 中。
223     "movl %%edx, %%ecx\n\t"     // 再将串 2 的长度值放入 ecx 中。
224     "repne\n\t"                 // 比较 al 与串 2 中字符 es:[edi], 并且 edi++。
225     "scasb\n\t"                 // 如果不相等就继续比较。
226     "jne 1b\n\t"                 // 如果不相等, 则向后跳转到标号 1 处。
227     "decl %0\n\t"               // esi--, 指向一个包含在串 2 中的字符。
228     "jmp 3f\n\t"                 // 向前跳转到标号 3 处。
229     "2:\txorl %0, %0\n\t"       // 没有找到符合条件的, 将返回值为 NULL。
230     "3:"
231     : "=S" (__res): "a" (0), "c" (0xffffffff), "" (cs), "g" (ct)
232     : "ax", "cx", "dx", "di");
233 return __res;                   // 返回指针值。
234 }
235
236     /// 在字符串 1 中寻找首个匹配整个字符串 2 的字符串。
237     // 参数: cs - 字符串 1 的指针, ct - 字符串 2 的指针。
238     // %0 - eax(__res), %1 - eax(0), %2 - ecx(0xffffffff), %3 - esi(串 1 指针 cs), %4 - (串 2 指针 ct)。
239     // 返回: 返回字符串 1 中首个匹配字符串 2 的字符串指针。
240 extern inline char * strstr(const char * cs, const char * ct)
241 {
242     register char * __res __asm__ ("ax"); // __res 是寄存器变量(eax)。
243     __asm__ ("cld\n\t" \
244             "movl %4, %%edi\n\t"          // 清方向位。
245             "repne\n\t"                   // 首先计算串 2 的长度。串 2 指针放入 edi 中。
246             "scasb\n\t"                   // 比较 al(0) 与串 2 中的字符 (es:[edi]), 并 edi++。
247             "notl %%ecx\n\t"              // 如果不相等就继续比较 (ecx 逐步递减)。
248             "decl %%ecx\n\t"              // ecx 中每位取反。
249             /* NOTE! This also sets Z if searchstring='' */
250             /* 注意! 如果搜索串为空, 将设置 Z 标志 */ // 得串 2 的长度值。
251             "movl %%ecx, %%edx\n\t"        // 将串 2 的长度值暂放入 edx 中。
252             "1:\tmovl %4, %%edi\n\t"      // 取串 2 头指针放入 edi 中。
253             "movl %%esi, %%eax\n\t"        // 将串 1 的指针复制到 eax 中。
254             "movl %%edx, %%ecx\n\t"        // 再将串 2 的长度值放入 ecx 中。
255             "repe\n\t"                    // 比较串 1 和串 2 字符 (ds:[esi], es:[edi]), esi++, edi++。
256             "cmpsb\n\t"                    // 若对应字符相等就一直比较下去。
257             "je 2f\n\t"                    /* also works for empty string, see above */
258             /* 对空串同样有效, 见上面 */ // 若全相等, 则转到标号 2。
259             "xchgl %%eax, %%esi\n\t"       // 串 1 头指针→esi, 比较结果的串 1 指针→eax。
260             "incl %%esi\n\t"               // 串 1 头指针指向下一个字符。
261             "cmpb $0, -1(%%eax)\n\t"       // 串 1 指针 (eax-1) 所指字节是 0 吗?
262             "jne 1b\n\t"                   // 不是则跳转到标号 1, 继续从串 1 的第 2 个字符开始比较。
263             "xorl %%eax, %%eax\n\t"        // 清 eax, 表示没有找到匹配。
264             "2:"
265             : "=a" (__res): "" (0), "c" (0xffffffff), "S" (cs), "g" (ct)
266             : "cx", "dx", "di", "si");
267 return __res;                   // 返回比较结果。
268 }

```

```

//// 计算字符串长度。
// 参数: s - 字符串。
// %0 - ecx(__res), %1 - edi(字符串指针 s), %2 - eax(0), %3 - ecx(0xffffffff)。
// 返回: 返回字符串的长度。
263 extern inline int strlen(const char * s)
264 {
265     register int __res __asm__( "cx" );    // __res 是寄存器变量(ecx)。
266     __asm__( "cld\n\t"                    // 清方向位。
267             "repne\n\t"                  // al(0)与字符串中字符 es:[edi]比较,
268             "scasb\n\t"                  // 若不相等就一直比较。
269             "notl %0\n\t"                // ecx 取反。
270             "decl %0"                    // ecx--, 得字符串得长度值。
271             : "=c" (__res) : "D" (s), "a" (0), "" (0xffffffff) : "di" );
272     return __res;                        // 返回字符串长度值。
273 }
274
275 extern char * __strtok;    // 用于临时存放指向下面被分析字符串 1(s)的指针。
276
//// 利用字符串 2 中的字符将字符串 1 分割成标记(token)序列。
// 将串 1 看作是包含零个或多个单词(token)的序列, 并由分割符字符串 2 中的一个或多个字符分开。
// 第一次调用 strtok()时, 将返回指向字符串 1 中第 1 个 token 首字符的指针, 并在返回 token 时将
// 一 null 字符写到分割符处。后续使用 null 作为字符串 1 的调用, 将用这种方法继续扫描字符串 1,
// 直到没有 token 为止。在不同的调用过程中, 分割串 2 可以不同。
// 参数: s - 待处理的字符串 1, ct - 包含各个分割符的字符串 2。
// 汇编输出: %0 - ebx(__res), %1 - esi(__strtok);
// 汇编输入: %2 - ebx(__strtok), %3 - esi(字符串 1 指针 s), %4 - (字符串 2 指针 ct)。
// 返回: 返回字符串 s 中第 1 个 token, 如果没有找到 token, 则返回一个 null 指针。
// 后续使用字符串 s 指针为 null 的调用, 将在原字符串 s 中搜索下一个 token。
277 extern inline char * strtok(char * s, const char * ct)
278 {
279     register char * __res __asm__( "si" );
280     __asm__( "testl %1, %1\n\t"          // 首先测试 esi(字符串 1 指针 s)是否是 NULL。
281             "jne 1f\n\t"                // 如果不是, 则表明是首次调用本函数, 跳转标号 1。
282             "testl %0, %0\n\t"          // 如果是 NULL, 则表示此次是后续调用, 测 ebx(__strtok)。
283             "je 8f\n\t"                 // 如果 ebx 指针是 NULL, 则不能处理, 跳转结束。
284             "movl %0, %1\n\t"            // 将 ebx 指针复制到 esi。
285             "l:\txorl %0, %0\n\t"        // 清 ebx 指针。
286             "movl $-1, %%ecx\n\t"        // 置 ecx = 0xffffffff。
287             "xorl %%eax, %%eax\n\t"      // 清零 eax。
288             "cld\n\t"                   // 清方向位。
289             "movl %4, %%edi\n\t"         // 下面求字符串 2 的长度。edi 指向字符串 2。
290             "repne\n\t"                  // 将 al(0)与 es:[edi]比较, 并且 edi++。
291             "scasb\n\t"                  // 直到找到字符串 2 的结束 null 字符, 或计数 ecx==0。
292             "notl %%ecx\n\t"             // 将 ecx 取反,
293             "decl %%ecx\n\t"             // ecx--, 得到字符串 2 的长度值。
294             "je 7f\n\t"                  /* empty delimiter-string */
                                         /* 分割符字符串空 */ // 若串 2 长度为 0, 则转标号 7。
295             "movl %%ecx, %%edx\n\t"      // 将串 2 长度暂存入 edx。
296             "2:\tlodsb\n\t"              // 取串 1 的字符 ds:[esi]→al, 并且 esi++。
297             "testb %%al, %%al\n\t"       // 该字符为 0 值吗(串 1 结束)?
298             "je 7f\n\t"                 // 如果是, 则跳转标号 7。
299             "movl %4, %%edi\n\t"         // edi 再次指向串 2 首。
300             "movl %%edx, %%ecx\n\t"      // 取串 2 的长度值置入计数器 ecx。

```

```

301     "repne|n|t" // 将 a1 中串 1 的字符与串 2 中所有字符比较，
302     "scasb|n|t" // 判断该字符是否为分割符。
303     "je 2b|n|t" // 若能在串 2 中找到相同字符（分割符），则跳转标号 2。
304     "decl %l|n|t" // 若不是分割符，则串 1 指针 esi 指向此时的该字符。
305     "cmpb $0, (%l)|n|t" // 该字符是 NULL 字符吗？
306     "je 7f|n|t" // 若是，则跳转标号 7 处。
307     "movl %l, %0|n" // 将该字符的指针 esi 存放在 ebx。
308     "3:|t|lods|n|t" // 取串 1 下一个字符 ds:[esi] → al，并且 esi++。
309     "testb %%al, %%al|n|t" // 该字符是 NULL 字符吗？
310     "je 5f|n|t" // 若是，表示串 1 结束，跳转到标号 5。
311     "movl %4, %%edi|n|t" // edi 再次指向串 2 首。
312     "movl %%edx, %%ecx|n|t" // 串 2 长度值置入计数器 ecx。
313     "repne|n|t" // 将 a1 中串 1 的字符与串 2 中每个字符比较，
314     "scasb|n|t" // 测试 al 字符是否是分割符。
315     "jne 3b|n|t" // 若不是分割符则跳转标号 3，检测串 1 中下一个字符。
316     "decl %l|n|t" // 若是分割符，则 esi--，指向该分割符字符。
317     "cmpb $0, (%l)|n|t" // 该分割符是 NULL 字符吗？
318     "je 5f|n|t" // 若是，则跳转到标号 5。
319     "movb $0, (%l)|n|t" // 若不是，则将该分割符用 NULL 字符替换掉。
320     "incl %l|n|t" // esi 指向串 1 中下一个字符，也即剩余串首。
321     "jmp 6f|n" // 跳转标号 6 处。
322     "5:|t|xorl %l, %l|n" // esi 清零。
323     "6:|t|cmpb $0, (%0)|n|t" // ebx 指针指向 NULL 字符吗？
324     "jne 7f|n|t" // 若不是，则跳转标号 7。
325     "xorl %0, %0|n" // 若是，则让 ebx=NULL。
326     "7:|t|testl %0, %0|n|t" // ebx 指针为 NULL 吗？
327     "jne 8f|n|t" // 若不是则跳转 8，结束汇编代码。
328     "movl %0, %l|n" // 将 esi 置为 NULL。
329     "8:"
330     : "=b" (__res), "=S" (__strtok)
331     : "" (__strtok), "l" (s), "g" (ct)
332     : "ax", "cx", "dx", "di");
333 return __res; // 返回指向新 token 的指针。
334 }
335
336 // 内存块复制。从源地址 src 处开始复制 n 个字节到目的地址 dest 处。
337 // 参数：dest - 复制的目的地址，src - 复制的源地址，n - 复制字节数。
338 // %0 - ecx(n)，%1 - esi(src)，%2 - edi(dest)。
339 extern inline void * memcpy(void * dest, const void * src, int n)
340 {
341     __asm__ ("cld|n|t" // 清方向位。
342             "rep|n|t" // 重复执行复制 ecx 个字节，
343             "movsb" // 从 ds:[esi] 到 es:[edi]，esi++，edi++。
344             :: "c" (n), "S" (src), "D" (dest)
345             : "cx", "si", "di");
346 return dest; // 返回目的地址。
347 }
348
349 // 内存块移动。同内存块复制，但考虑移动的方向。
350 // 参数：dest - 复制的目的地址，src - 复制的源地址，n - 复制字节数。
351 // 若 dest < src 则：%0 - ecx(n)，%1 - esi(src)，%2 - edi(dest)。
352 // 否则：%0 - ecx(n)，%1 - esi(src+n-1)，%2 - edi(dest+n-1)。
353 // 这样操作是为了防止在复制时错误地重叠覆盖。

```

```

346 extern inline void * memmove(void * dest, const void * src, int n)
347 {
348     if (dest < src)
349         __asm__( "cld\n|t"                // 清方向位。
350                 "rep\n|t"                // 从 ds:[esi] 到 es:[edi], 并且 esi++, edi++,
351                 "movsb"                  // 重复执行复制 ecx 字节。
352                 :: "c" (n), "S" (src), "D" (dest)
353                 : "cx", "si", "di");
354     else
355         __asm__( "std\n|t"                // 置方向位, 从末端开始复制。
356                 "rep\n|t"                // 从 ds:[esi] 到 es:[edi], 并且 esi--, edi--,
357                 "movsb"                  // 复制 ecx 个字节。
358                 :: "c" (n), "S" (src+n-1), "D" (dest+n-1)
359                 : "cx", "si", "di");
360     return dest;
361 }
362
363 // 比较 n 个字节的内存块 (两个字符串), 即使遇上 NULL 字节也不停止比较。
364 // 参数: cs - 内存块 1 地址, ct - 内存块 2 地址, count - 比较的字节数。
365 // %0 - eax(__res), %1 - eax(0), %2 - edi(内存块 1), %3 - esi(内存块 2), %4 - ecx(count)。
366 // 返回: 若块 1 > 块 2 返回 1; 块 1 < 块 2, 返回 -1; 块 1 == 块 2, 则返回 0。
367 extern inline int memcmp(const void * cs, const void * ct, int count)
368 {
369     register int __res __asm__( "ax"); // __res 是寄存器变量。
370     __asm__( "cld\n|t"                // 清方向位。
371             "repe\n|t"                // 如果相等则重复,
372             "cmpsb\n|t"                // 比较 ds:[esi] 与 es:[edi] 的内容, 并且 esi++, edi++.
373             "je 1f\n|t"                // 如果都相同, 则跳转到标号 1, 返回 0(eax) 值
374             "movl $1, %%eax\n|t"      // 否则 eax 置 1,
375             "j1 1f\n|t"                // 若内存块 2 内容的值 < 内存块 1, 则跳转标号 1。
376             "negl %%eax\n|t"          // 否则 eax = -eax。
377             "1:"
378             : "=a" (__res) : "" (0), "D" (cs), "S" (ct), "c" (count)
379             : "si", "di", "cx");
380     return __res; // 返回比较结果。
381 }
382
383 // 在 n 字节大小的内存块 (字符串) 中寻找指定字符。
384 // 参数: cs - 指定内存块地址, c - 指定的字符, count - 内存块长度。
385 // %0 - edi(__res), %1 - eax(字符 c), %2 - edi(内存块地址 cs), %3 - ecx(字节数 count)。
386 // 返回第一个匹配字符的指针, 如果没有找到, 则返回 NULL 字符。
387 extern inline void * memchr(const void * cs, char c, int count)
388 {
389     register void * __res __asm__( "di"); // __res 是寄存器变量。
390     if (!count) // 如果内存块长度 == 0, 则返回 NULL, 没有找到。
391         return NULL;
392     __asm__( "cld\n|t"                // 清方向位。
393             "repne\n|t"               // 如果不相等则重复执行下面语句,
394             "scasb\n|t"               // al 中字符与 es:[edi] 字符作比较, 并且 edi++,
395             "je 1f\n|t"               // 如果相等则向前跳转到标号 1 处。
396             "movl $1, %0\n|t"         // 否则 edi 中置 1。
397             "1:\tdecl %0"              // 让 edi 指向找到的字符 (或是 NULL)。
398             : "=D" (__res) : "a" (c), "D" (cs), "c" (count)

```

```

391         : "cx");
392 return __res;           // 返回字符指针。
393 }
394
395 // 用字符 c 填充指定长度内存块。
396 // 用字符 c 填充 s 指向的内存区域，共填 count 字节。
397 // %0 - eax(字符 c)，%1 - edi(内存地址)，%2 - ecx(字节数 count)。
398 extern inline void * memset(void * s, char c, int count)
399 {
400     __asm__( "cld\n\t"           // 清方向位。
401             "rep\n\t"           // 重复 ecx 指定的次数，执行
402             "stosb"             // 将 al 中字符存入 es:[edi]中，并且 edi++。
403             :: "a" (c), "D" (s), "c" (count)
404             : "cx", "di");
405 return s;
406 }
407 #endif
408

```

## 10.12 termios.h 文件

### 10.12.1 功能描述

该文件含有终端 I/O 接口定义。包括 termios 数据结构和一些对通用终端接口设置的函数原型。这些函数用来读取或设置终端的属性、线路控制、读取或设置波特率以及读取或设置终端前端进程的组 id。虽然这是 linux 早期的头文件，但已完全符合目前的 POSIX 标准，并作了适当的扩展。

在该文件中定义的两个终端数据结构 termio 和 termios 是分别属于两类 UNIX 系列(或克隆)，termio 是在 AT&T 系统 V 中定义的，而 termios 是 POSIX 标准指定的。两个结构基本一样，只是 termio 使用短整数类型定义模式标志集而 termios 使用长整数定义模式标志集。由于目前这两种结构都在使用，因此为了兼容性，大多数系统都同时支持它们。另外，以前使用的是一类似的 sgtty 结构，目前已基本不用。

### 10.12.2 代码注释

列表 linux/include/termios.h 文件

```

1 #ifndef TERMIOS\_H
2 #define TERMIOS\_H
3
4 #define TTY\_BUF\_SIZE 1024           // tty 中的缓冲区长度。
5
6 /* 0x54 is just a magic number to make these relatively unique ('T') */
7 /* 0x54 只是一个魔数，目的是为了这些常数唯一('T') */
8
9 // tty 设备的 ioctl 调用命令集。ioctl 将命令编码在低位字中。
10 // 下面名称 TC[*]的含义是 tty 控制命令。
11 // 取相应终端 termios 结构中的信息(参见 tcgetattr())。
12 #define TCGETS 0x5401

```



```

// 设置相应终端 termios 结构中的信息(参见 tcsetattr(), TCSANOW)。
9 #define TCSETS          0x5402
// 在设置终端 termios 的信息之前, 需要先等待输出队列中所有数据处理完(耗尽)。对于修改参数
// 会影响输出的情况, 就需要使用这种形式(参见 tcsetattr(), TCSADRAIN 选项)。
10 #define TCSETSW        0x5403
// 在设置 termios 的信息之前, 需要先等待输出队列中所有数据处理完, 并且刷新(清空)输入队列。
// 再设置(参见 tcsetattr(), TCSAFLUSH 选项)。
11 #define TCSETSF        0x5404
// 取相应终端 termio 结构中的信息(参见 tcgetattr())。
12 #define TCGETA         0x5405
// 设置相应终端 termio 结构中的信息(参见 tcsetattr(), TCSANOW 选项)。
13 #define TCSETA         0x5406
// 在设置终端 termio 的信息之前, 需要先等待输出队列中所有数据处理完(耗尽)。对于修改参数
// 会影响输出的情况, 就需要使用这种形式(参见 tcsetattr(), TCSADRAIN 选项)。
14 #define TCSETAW        0x5407
// 在设置 termio 的信息之前, 需要先等待输出队列中所有数据处理完, 并且刷新(清空)输入队列。
// 再设置(参见 tcsetattr(), TCSAFLUSH 选项)。
15 #define TCSETAF        0x5408
// 等待输出队列处理完毕(空), 如果参数值是 0, 则发送一个 break(参见 tcsendbreak(), tcdrain())。
16 #define TCSBRK         0x5409
// 开始/停止控制。如果参数值是 0, 则挂起输出; 如果是 1, 则重新开启挂起的输出; 如果是 2, 则挂起
// 输入; 如果是 3, 则重新开启挂起的输入(参见 tcflow())。
17 #define TCXONC         0x540A
// 刷新已写输出但还没发送或已收但还没有读数据。如果参数是 0, 则刷新(清空)输入队列; 如果是 1,
// 则刷新输出队列; 如果是 2, 则刷新输入和输出队列(参见 tcflush())。
18 #define TCFLSH         0x540B
// 下面名称 TIOC[*]的含义是 tty 输入输出控制命令。
// 设置终端串行线路专用模式。
19 #define TIOCEXCL        0x540C
// 复位终端串行线路专用模式。
20 #define TIOCNXCL        0x540D
// 设置 tty 为控制终端。(TIOCNOTTY - 禁止 tty 为控制终端)。
21 #define TIOCSCTTY       0x540E
// 读取指定终端设备进程的组 id(参见 tcgetpgrp())。
22 #define TIOCGPRGP       0x540F
// 设置指定终端设备进程的组 id(参见 tcsetpgrp())。
23 #define TIOCSPGRP       0x5410
// 返回输出队列中还未送出的字符数。
24 #define TIOCOUTQ        0x5411
// 模拟终端输入。该命令以一个指向字符的指针作为参数, 并假装该字符是在终端上键入的。用户必须
// 在该控制终端上具有超级用户权限或具有读许可权限。
25 #define TIOCSTI         0x5412
// 读取终端设备窗口大小信息(参见 winsize 结构)。
26 #define TIOCGWINSZ       0x5413
// 设置终端设备窗口大小信息(参见 winsize 结构)。
27 #define TIOCSWINSZ       0x5414
// 返回 modem 状态控制引线的当前状态比特位标志集(参见下面 185-196 行)。
28 #define TIOCMGET         0x5415
// 设置单个 modem 状态控制引线的状态(true 或 false)(Individual control line Set)。
29 #define TIOCMBS         0x5416
// 复位单个 modem 状态控制引线的状态(Individual control line clear)。
30 #define TIOCMBS         0x5417
// 设置 modem 状态引线的状态。如果某一比特位置位, 则 modem 对应的状态引线将置为有效。

```

```

31 #define TIOCMSET          0x5418
    // 读取软件载波检测标志(1 - 开启; 0 - 关闭)。
    // 对于本地连接的终端或其它设备, 软件载波标志是开启的, 对于使用 modem 线路的终端或设备则
    // 是关闭的。为了能使用这两个 ioctl 调用, tty 线路应该是以 O_NDELAY 方式打开的, 这样 open()
    // 就不会等待载波。
32 #define TIOCGSOFTCAR      0x5419
    // 设置软件载波检测标志(1 - 开启; 0 - 关闭)。
33 #define TIOCSSOFTCAR      0x541A
    // 返回输入队列中还未取走字符的数目。
34 #define TIOCIQ            0x541B
35
    // 窗口大小(Window size)属性结构。在窗口环境中可用于基于屏幕的应用程序。
    // ioctl 中的 TIOCGWINSZ 和 TIOCSWINSZ 可用来读取或设置这些信息。
36 struct winsize {
37     unsigned short ws_row;        // 窗口字符行数。
38     unsigned short ws_col;        // 窗口字符列数。
39     unsigned short ws_xpixel;     // 窗口宽度, 像素值。
40     unsigned short ws_ypixel;     // 窗口高度, 像素值。
41 };
42
    // AT&T 系统 V 的 termio 结构。
43 #define NCC 8                    // termio 结构中控制字符数组的长度。
44 struct termio {
45     unsigned short c_iflag;        /* input mode flags */ // 输入模式标志。
46     unsigned short c_oflag;        /* output mode flags */ // 输出模式标志。
47     unsigned short c_cflag;        /* control mode flags */ // 控制模式标志。
48     unsigned short c_lflag;        /* local mode flags */ // 本地模式标志。
49     unsigned char c_line;          /* line discipline */ // 线路规程(速率)。
50     unsigned char c_cc[NCC];       /* control characters */ // 控制字符数组。
51 };
52
    // POSIX 的 termios 结构。
53 #define NCCS 17                  // termios 结构中控制字符数组的长度。
54 struct termios {
55     unsigned long c_iflag;          /* input mode flags */ // 输入模式标志。
56     unsigned long c_oflag;          /* output mode flags */ // 输出模式标志。
57     unsigned long c_cflag;          /* control mode flags */ // 控制模式标志。
58     unsigned long c_lflag;          /* local mode flags */ // 本地模式标志。
59     unsigned char c_line;           /* line discipline */ // 线路规程(速率)。
60     unsigned char c_cc[NCCS];       /* control characters */ // 控制字符数组。
61 };
62
63 /* c_cc characters */ /* c_cc 数组中的字符 */
    // 以下是 c_cc 数组对应字符的索引值。
64 #define VINTR 0                  // c_cc[VINTR] = INTR (^C), \003, 中断字符。
65 #define VQUIT 1                  // c_cc[VQUIT] = QUIT (^), \034, 退出字符。
66 #define VERASE 2                 // c_cc[VERASE] = ERASE (^H), \177, 擦出字符。
67 #define VKILL 3                  // c_cc[VKILL] = KILL (^U), \025, 终止字符。
68 #define VEOF 4                   // c_cc[VEOF] = EOF (^D), \004, 文件结束字符。
69 #define VTIME 5                  // c_cc[VTIME] = TIME (\0), \0, 定时器值(参见后面说明)。
70 #define VMIN 6                   // c_cc[VMIN] = MIN (\1), \1, 定时器值。
71 #define VSWTC 7                  // c_cc[VSWTC] = SWTC (\0), \0, 交换字符。
72 #define VSTART 8                 // c_cc[VSTART] = START (^Q), \021, 开始字符。

```



```

73 #define VSTOP 9          // c_cc[VSTOP] = STOP (^S), \023, 停止字符。
74 #define VSUSP 10         // c_cc[VSUSP] = SUSP (^Z), \032, 挂起字符。
75 #define VEOL 11          // c_cc[VEOL] = EOL (\0), \0, 行结束字符。
76 #define VREPRINT 12      // c_cc[VREPRINT] = REPRINT (^R), \022, 重显示字符。
77 #define VDISCARD 13      // c_cc[VDISCARD] = DISCARD (^O), \017, 丢弃字符。
78 #define VWERASE 14       // c_cc[VWERASE] = WERASE (^W), \027, 单词擦除字符。
79 #define VLNEXT 15        // c_cc[VLNEXT] = LNEXT (^V), \026, 下一行字符。
80 #define VEOL2 16         // c_cc[VEOL2] = EOL2 (\0), \0, 行结束 2。
81
82 /* c_iflag bits */      /* c_iflag 比特位 */
// termios 结构输入模式字段 c_iflag 各种标志的符号常数。
83 #define IGNBRK 0000001   // 输入时忽略 BREAK 条件。
84 #define BRKINT 0000002   // 在 BREAK 时产生 SIGINT 信号。
85 #define IGNPAR 0000004   // 忽略奇偶校验出错的字符。
86 #define PARMRK 0000010   // 标记奇偶校验错。
87 #define INPCK 0000020    // 允许输入奇偶校验。
88 #define ISTRIP 0000040    // 屏蔽字符第 8 位。
89 #define INLCR 0000100    // 输入时将换行符 NL 映射成回车符 CR。
90 #define IGNCR 0000200    // 忽略回车符 CR。
91 #define ICRNL 0000400    // 在输入时将回车符 CR 映射成换行符 NL。
92 #define IUCLC 0001000    // 在输入时将大写字母转换成小写字母。
93 #define IXON 0002000     // 允许开始/停止 (XON/XOFF) 输出控制。
94 #define IXANY 0004000    // 允许任何字符重启输出。
95 #define IXOFF 0010000    // 允许开始/停止 (XON/XOFF) 输入控制。
96 #define IMAXBEL 0020000  // 输入队列满时响铃。
97
98 /* c_oflag bits */      /* c_oflag 比特位 */
// termios 结构中输出模式字段 c_oflag 各种标志的符号常数。
99 #define OPOST 0000001    // 执行输出处理。
100 #define OLCUC 0000002    // 在输出时将小写字母转换成大写字母。
101 #define ONLCR 0000004    // 在输出时将换行符 NL 映射成回车-换行符 CR-NL。
102 #define OCRNL 0000010    // 在输出时将回车符 CR 映射成换行符 NL。
103 #define ONOCR 0000020    // 在 0 列不输出回车符 CR。
104 #define ONLRET 0000040   // 换行符 NL 执行回车符的功能。
105 #define OFILL 0000100    // 延迟时使用填充字符而不使用时间延迟。
106 #define OFDEL 0000200    // 填充字符是 ASCII 码 DEL。如果未设置, 则使用 ASCII NULL。
107 #define NLDLY 0000400    // 选择换行延迟。
108 #define NL0 0000000      // 换行延迟类型 0。
109 #define NL1 0000400      // 换行延迟类型 1。
110 #define CRDLY 0003000    // 选择回车延迟。
111 #define CR0 0000000      // 回车延迟类型 0。
112 #define CR1 0001000      // 回车延迟类型 1。
113 #define CR2 0002000      // 回车延迟类型 2。
114 #define CR3 0003000      // 回车延迟类型 3。
115 #define TABDLY 0014000   // 选择水平制表延迟。
116 #define TAB0 0000000     // 水平制表延迟类型 0。
117 #define TAB1 0004000     // 水平制表延迟类型 1。
118 #define TAB2 0010000     // 水平制表延迟类型 2。
119 #define TAB3 0014000     // 水平制表延迟类型 3。
120 #define XTABS 0014000    // 将制表符 TAB 换成空格, 该值表示空格数。
121 #define BSDLY 0020000    // 选择退格延迟。
122 #define BS0 0000000      // 退格延迟类型 0。
123 #define BS1 0020000      // 退格延迟类型 1。

```

```

124 #define VTDLY 0040000 // 纵向制表延迟。
125 #define VTO 0000000 // 纵向制表延迟类型 0。
126 #define VT1 0040000 // 纵向制表延迟类型 1。
127 #define FFDLY 0040000 // 选择换页延迟。
128 #define FFO 0000000 // 换页延迟类型 0。
129 #define FFI 0040000 // 换页延迟类型 1。
130
131 /* c_cflag bit meaning */ /* c_cflag 比特位的含义 */
// termios 结构中的控制模式标志字段 c_cflag 标志的符号常数 (8 进制数)。
132 #define CBAUD 0000017 // 传输速率位屏蔽码。
133 #define B0 0000000 /* hang up */ /* 挂断线路 */
134 #define B50 0000001 // 波特率 50。
135 #define B75 0000002 // 波特率 75。
136 #define B110 0000003 // 波特率 110。
137 #define B134 0000004 // 波特率 134。
138 #define B150 0000005 // 波特率 150。
139 #define B200 0000006 // 波特率 200。
140 #define B300 0000007 // 波特率 300。
141 #define B600 0000010 // 波特率 600。
142 #define B1200 0000011 // 波特率 1200。
143 #define B1800 0000012 // 波特率 1800。
144 #define B2400 0000013 // 波特率 2400。
145 #define B4800 0000014 // 波特率 4800。
146 #define B9600 0000015 // 波特率 9600。
147 #define B19200 0000016 // 波特率 19200。
148 #define B38400 0000017 // 波特率 38400。
149 #define EXTA B19200 // 扩展波特率 A。
150 #define EXTB B38400 // 扩展波特率 B。

151 #define CSIZE 0000060 // 字符位宽度屏蔽码。
152 #define CS5 0000000 // 每字符 5 比特位。
153 #define CS6 0000020 // 每字符 6 比特位。
154 #define CS7 0000040 // 每字符 7 比特位。
155 #define CS8 0000060 // 每字符 8 比特位。
156 #define CSTOPB 0000100 // 设置两个停止位，而不是 1 个。
157 #define CREAD 0000200 // 允许接收。
158 #define CPARENB 0000400 // 开启输出时产生奇偶位、输入时进行奇偶校验。
159 #define CPARODD 0001000 // 输入/输入校验是奇校验。
160 #define HUPCL 0002000 // 最后进程关闭后挂断。
161 #define CLOCAL 0004000 // 忽略调制解调器(modem)控制线路。
162 #define CIBAUD 03600000 /* input baud rate (not used) */ /* 输入波特率(未使用) */
163 #define CRTSCTS 020000000000 /* flow control */ /* 流控制 */
164
165 #define PARENB CPARENB // 开启输出时产生奇偶位、输入时进行奇偶校验。
166 #define PARODD CPARODD // 输入/输入校验是奇校验。
167
168 /* c_lflag bits */ /* c_lflag 比特位 */
// termios 结构中本地模式标志字段 c_lflag 的符号常数。
169 #define ISIG 0000001 // 当收到字符 INTR、QUIT、SUSP 或 DSUSP，产生相应的信号。
170 #define ICANON 0000002 // 开启规范模式(熟模式)。
171 #define XCASE 0000004 // 若设置了 ICANON，则终端是大写字符的。
172 #define ECHO 0000010 // 回显输入字符。
173 #define ECHOE 0000020 // 若设置了 ICANON，则 ERASE/WERASE 将擦除前一字符/单词。

```

```

174 #define ECHOK 0000040 // 若设置了 ICANON, 则 KILL 字符将擦除当前行。
175 #define ECHONL 0000100 // 如设置了 ICANON, 则即使 ECHO 没有开启也回显 NL 字符。
176 #define NOFLSH 0000200 // 当生成 SIGINT 和 SIGQUIT 信号时不刷新输入输出队列, 当
// 生成 SIGSUSP 信号时, 刷新输入队列。
177 #define TOSTOP 0000400 // 发送 SIGTTOU 信号到后台进程的进程组, 该后台进程试图写
// 自己的控制终端。
178 #define ECHOCTL 0001000 // 若设置了 ECHO, 则除 TAB、NL、START 和 STOP 以外的 ASCII
// 控制信号将被回显成象 ^X 式样, X 值是控制符+0x40。
179 #define ECHOPRT 0002000 // 若设置了 ICANON 和 IECHO, 则字符在擦除时将显示。
180 #define ECHOKE 0004000 // 若设置了 ICANON, 则 KILL 通过擦除行上的所有字符被回显。
181 #define FLUSHO 0010000 // 输出被刷新。通过键入 DISCARD 字符, 该标志被翻转。
182 #define PENDIN 0040000 // 当下一个字符是读时, 输入队列中的所有字符将被重显。
183 #define TEXTEN 0100000 // 开启实现时定义的输入处理。
184
185 /* modem lines */ /* modem 线路信号符号常数 */
186 #define TIOCM_LE 0x001 // 线路允许(Line Enable)。
187 #define TIOCM_DTR 0x002 // 数据终端就绪(Data Terminal Ready)。
188 #define TIOCM_RTS 0x004 // 请求发送(Request to Send)。
189 #define TIOCM_ST 0x008 // 串行数据发送(Serial Transfer)。[??]
190 #define TIOCM_SR 0x010 // 串行数据接收(Serial Receive)。[??]
191 #define TIOCM_CTS 0x020 // 清除发送(Clear To Send)。
192 #define TIOCM_CAR 0x040 // 载波监测(Carrier Detect)。
193 #define TIOCM_RNG 0x080 // 响铃指示(Ring indicate)。
194 #define TIOCM_DSR 0x100 // 数据设备就绪(Data Set Ready)。
195 #define TIOCM_CD TIOCM_CAR
196 #define TIOCM_RI TIOCM_RNG
197
198 /* tcflow() and TCXONC use these */ /* tcflow() 和 TCXONC 使用这些符号常数 */
199 #define TCOOFF 0 // 挂起输出。
200 #define TCOON 1 // 重启被挂起的输出。
201 #define TCIOFF 2 // 系统传输一个 STOP 字符, 使设备停止向系统传输数据。
202 #define TCION 3 // 系统传输一个 START 字符, 使设备开始向系统传输数据。
203
204 /* tcflush() and TCFLSH use these */ /* tcflush() 和 TCFLSH 使用这些符号常数 */
205 #define TCIFLUSH 0 // 清接收到的数据但不读。
206 #define TCOFLUSH 1 // 清已写的数据但不传送。
207 #define TCIOFLUSH 2 // 清接收到的数据但不读。清已写的数据但不传送。
208
209 /* tcsetattr uses these */ /* tcsetattr() 使用这些符号常数 */
210 #define TCSANOW 0 // 改变立即发生。
211 #define TCSADRAIN 1 // 改变在所有已写的输出被传输之后发生。
212 #define TCSAFLUSH 2 // 改变在所有已写的输出被传输之后并且在所有接收到但
// 还没有读取的数据被丢弃之后发生。
213
214 typedef int speed_t; // 波特率数值类型。
215
// 返回 termios_p 所指 termios 结构中的接收波特率。
216 extern speed_t cfgetispeed(struct termios *termios_p);
// 返回 termios_p 所指 termios 结构中的发送波特率。
217 extern speed_t cfgetospeed(struct termios *termios_p);
// 将 termios_p 所指 termios 结构中的接收波特率设置为 speed。
218 extern int cfsetispeed(struct termios *termios_p, speed_t speed);
// 将 termios_p 所指 termios 结构中的发送波特率设置为 speed。

```

```

219 extern int cfsetospeed(struct termios *termios_p, speed\_t speed);
    // 等待 fildes 所指对象已写输出数据被传送出去。
220 extern int tcdrain(int fildes);
    // 挂起/重启 fildes 所指对象数据的接收和发送。
221 extern int tcflow(int fildes, int action);
    // 丢弃 fildes 指定对象所有已写但还没传送以及所有已收到但还没有读取的数据。
222 extern int tcflush(int fildes, int queue_selector);
    // 获取与句柄 fildes 对应对象的参数, 并将其保存在 termios_p 所指的地方。
223 extern int tcgetattr(int fildes, struct termios *termios_p);
    // 如果终端使用异步串行数据传输, 则在一定时间内连续传输一系列 0 值比特位。
224 extern int tcsendbreak(int fildes, int duration);
    // 使用 termios 结构指针 termios_p 所指的数据, 设置与终端相关的参数。
225 extern int tcsetattr(int fildes, int optional_actions,
226                     struct termios *termios_p);
227
228 #endif
229

```

### 10.12.3 其它信息

#### 控制字符 TIME、MIN

在非规范模式输入处理中, 输入字符没有被处理成行, 因此擦除和终止处理也就不会发生。MIN 和 TIMEDE 的值即用于确定如何处理接收到的字符。

MIN 表示当满足读操作时 (也即, 当字符返给用户时) 需要读取的最少字符数。TIME 是以 1/10 秒计数的定时值, 用于超时定时和短期数据传输。这两个字符的四种组合情况及其相互作用描述如下:

MIN > 0, TIME > 0 的情况:

在这种情况下, TIME 起字符与字符间的定时器作用, 并在接收到第 1 个字符后开始起作用。由于它是字符与字符间的定时器, 所以在每收到一个字符就会被复位重启。MIN 与 TIME 之间的相互作用如下: 一旦收到一个字符, 字符间定时器就开始工作。如果在定时器超时 (注意定时器每收到一个字符就会重新开始计时) 之前收到了 MIN 个字符, 则读操作即被满足。如果在 MIN 个字符被收到之前定时器超时了, 就将此时已收到的字符返回给用户。注意, 如果 TIME 超时, 则起码有一个接收到的字符将被返回, 因为定时器只有在接收到了一个字符之后才开始起作用 (计时)。在这种情况下 (MIN > 0, TIME > 0), 读操作将会睡眠, 直到接收到第 1 个字符激活 MIN 与 TIME 机制。如果读到字符数少于已有的字符数, 那么定时器将不会被重新激活, 因而随后的读操作将被立刻满足。

MIN > 0, TIME = 0 的情况:

在这种情况下, 由于 TIME 的值是 0, 因此定时器不起作用, 只有 MIN 是有意义的。等待的读操作只有当接收到 MIN 个字符时才会被满足 (等待着的操作将睡眠直到收到 MIN 个字符)。使用这种情况去读基于记录的终端 I/O 的程序将会在读操作中被不确定地 (随意地) 阻塞。

MIN = 0, TIME > 0 的情况:

在这种情况下, 由于 MIN=0, 则 TIME 不再起字符间的定时器作用, 而是一个读操作定时器, 并在读操作一开始就起作用。只要接收到一个字符或者定时器超时就已满足读操作。注意, 在这种情况下, 如果定时器超时了, 将读不到一个字符。如果定时器没有超时, 那么只有在读到一个字符之后读操作才会满足。因此在这种情况下, 读操作不会无限地 (不确定地) 被阻塞, 以等待字符。在读操作开始后, 如果在 TIME\*0.10 秒的时间内没有收到字符, 读操作将以收到 0 个字符而返回。

MIN = 0, TIME = 0 的情况:

在这种情况下, 读操作会立刻返回。所请求读的字符数或缓冲队列中现有字符数中的最小值将被返回, 而不会等待更多的字符被输入缓冲中。

总的来说, 在非规范模式下, 这两个值是超时定时值和字符计数值。MIN 表示为了满足读操作, 需要读取的最少字符数。TIME 是一个十分之一秒计数的计时值。当这两个都设置的话, 读操作将等待, 直到至少读到一个字符, 然后在以读取 MIN 个字符或者时间 TIME 在读取最后一个字符后超时。如果仅设置了 MIN, 那么在读取 MIN 个字符之前读操作将不返回。如果仅设置了 TIME, 那么在读到至少一个字符或者定时超时后读操作将立刻返回。如果两个都没有设置, 则读操作将立刻返回, 仅给出目前已读的字节数。

## 10.13 time.h 文件

### 10.13.1 功能描述

该头文件用于涉及处理时间的函数。在 MINIX 中有一段对时间的描述很有趣: 时间的处理较为复杂, 比如什么是 GMT (格林威治标准时间)、本地时间或其它时间等。尽管主教 Ussher (1581-1656 年) 曾经计算过, 根据圣经, 世界开始之日是公元前 4004 年 10 月 12 日上午 9 点, 但在 UNIX 世界里, 时间是从 GMT 1970 年 1 月 1 日午夜开始的, 在这之前, 所有均是空无的和(无效的)。

### 10.13.2 代码注释

列表 linux/include/time.h 文件

---

```

1 #ifndef TIME\_H
2 #define TIME\_H
3
4 #ifndef TIME\_T
5 #define TIME\_T
6 typedef long time\_t;           // 从 GMT 1970 年 1 月 1 日开始的以秒计数的时间 (日历时间)。
7 #endif
8
9 #ifndef SIZE\_T
10 #define SIZE\_T
11 typedef unsigned int size\_t;
12 #endif
13
14 #define CLOCKS\_PER\_SEC 100      // 系统时钟滴答频率, 100HZ。
15
16 typedef long clock\_t;          // 从进程开始系统经过的时钟滴答数。
17
18 struct tm {
19     int tm_sec;                 // 秒数 [0, 59]。
20     int tm_min;                 // 分钟数 [0, 59]。
21     int tm_hour;                // 小时数 [0, 59]。
22     int tm_mday;                // 1 个月的天数 [0, 31]。

```



```

23     int tm_mon;                // 1 年中月份 [0, 11]。
24     int tm_year;              // 从 1900 年开始的年数。
25     int tm_wday;              // 1 星期中的某天 [0, 6] (星期天 =0)。
26     int tm_yday;              // 1 年中的某天 [0, 365]。
27     int tm_isdst;             // 夏令时标志。
28 };
29
    // 以下是有关时间操作的函数原型。
    // 确定处理器使用时间。返回程序所用处理器时间（滴答数）的近似值。
30 clock_t clock(void);
    // 取时间（秒数）。返回从 1970.1.1:0:0:0 开始的秒数（称为日历时间）。
31 time_t time(time_t * tp);
    // 计算时间差。返回时间 time2 与 time1 之间经过的秒数。
32 double difftime(time_t time2, time_t time1);
    // 将 tm 结构表示的时间转换成日历时间。
33 time_t mktime(struct tm * tp);
34
    // 将 tm 结构表示的时间转换成一个字符串。返回指向该串的指针。
35 char * asctime(const struct tm * tp);
    // 将日历时间转换成一个字符串形式，如 “Wed Jun 30 21:49:08:1993\n”。
36 char * ctime(const time_t * tp);
    // 将日历时间转换成 tm 结构表示的 UTC 时间（UTC - 世界时间代码 Universal Time Code）。
37 struct tm * gmtime(const time_t * tp);
    // 将日历时间转换成 tm 结构表示的指定时间区 (timezone) 的时间。
38 struct tm * localtime(const time_t * tp);
    // 将 tm 结构表示的时间利用格式字符串 fmt 转换成最大长度为 smax 的字符串并将结果存储在 s 中。
39 size_t strftime(char * s, size_t smax, const char * fmt, const struct tm * tp);
    // 初始化时间转换信息，使用环境变量 TZ，对 zname 变量进行初始化。
    // 在与时间区相关的时间转换函数中将自动调用该函数。
40 void tzset(void);
41
42 #endif
43

```

## 10.14 unistd.h 文件

### 10.14.1 功能描述

### 10.14.2 代码注释

列表 linux/include/unistd.h 文件

```

1 #ifndef UNISTD_H
2 #define UNISTD_H
3
4 /* ok, this may be a joke, but I'm working on it */
    /* ok, 这也许是个玩笑，但我正在着手处理 */
    // 下面符号常数指出符合 IEEE 标准 1003.1 实现的版本号，是一个整数值。
5 #define POSIX_VERSION 198808L

```

```

6 // chown() 和 fchown() 的使用受限于进程的权限。/* 只有超级用户可以执行 chown (我想..) */
7 #define POSIX_CHOWN_RESTRICTED /* only root can do a chown (I think..) */
8 // 长于 (NAME_MAX) 的路径名将产生错误，而不会自动截断。/* 路径名不截断 (但是请看内核代码) */
9 #define POSIX_NO_TRUNC /* no pathname truncation (but see in kernel) */
10 // 下面这个符号将定义成字符值，该值将禁止终端对其的处理。/* 禁止象 ^C 这样的字符 */
11 #define POSIX_VDISABLE '\0' /* character to disable things like ^C */
12 // 每个进程都有一保存的 set-user-ID 和一保存的 set-group-ID。/* 我们将着手对此进行处理 */
13 /*#define _POSIX_SAVED_IDS */ /* we'll get to this yet */
14 // 系统实现支持作业控制。/* 我们还没有支持这项标准，希望很快就行 */
15 /*#define _POSIX_JOB_CONTROL */ /* we aren't there quite yet. Soon hopefully */
16
17 #define STDIN_FILENO 0 // 标准输入文件句柄 (描述符) 号。
18 #define STDOUT_FILENO 1 // 标准输出文件句柄号。
19 #define STDERR_FILENO 2 // 标准出错文件句柄号。
20
21 #ifndef NULL
22 #define NULL ((void *)0) // 定义空指针。
23 #endif
24
25 /* access */ /* 文件访问 */
26 // 以下定义的符号常数用于 access() 函数。
27 #define F_OK 0 // 检测文件是否存在。
28 #define X_OK 1 // 检测是否可执行 (搜索)。
29 #define W_OK 2 // 检测是否可写。
30 #define R_OK 4 // 检测是否可读。
31
32 /* lseek */ /* 文件指针重定位 */
33 // 以下符号常数用于 lseek() 和 fcntl() 函数。
34 #define SEEK_SET 0 // 将文件读写指针设置为偏移值。
35 #define SEEK_CUR 1 // 将文件读写指针设置为当前值加上偏移值。
36 #define SEEK_END 2 // 将文件读写指针设置为文件长度加上偏移值。
37
38 /* _SC stands for System Configuration. We don't use them much */
39 /* _SC 表示系统配置。我们很少使用 */
40 // 下面的符号常数用于 sysconf() 函数。
41 #define SC_ARG_MAX 1 // 最大变量数。
42 #define SC_CHILD_MAX 2 // 子进程最大数。
43 #define SC_CLOCKS_PER_SEC 3 // 每秒滴答数。
44 #define SC_NGROUPS_MAX 4 // 最大组数。
45 #define SC_OPEN_MAX 5 // 最大打开文件数。
46 #define SC_JOB_CONTROL 6 // 作业控制。
47 #define SC_SAVED_IDS 7 // 保存的标识符。
48 #define SC_VERSION 8 // 版本。
49
50 /* more (possibly) configurable things - now pathnames */
51 /* 更多的 (可能的) 可配置参数 - 现在用于路径名 */
52 // 下面的符号常数用于 pathconf() 函数。
53 #define PC_LINK_MAX 1 // 连接最大数。
54 #define PC_MAX_CANON 2 // 最大常规文件数。
55 #define PC_MAX_INPUT 3 // 最大输入长度。
56 #define PC_NAME_MAX 4 // 名称最大长度。
57 #define PC_PATH_MAX 5 // 路径最大长度。

```

```

48 #define PC_PIPE_BUF          6    // 管道缓冲大小。
49 #define PC_NO_TRUNC          7    // 文件名不截断。
50 #define PC_VDISABLE          8    //
51 #define PC_CHOWN_RESTRICTED  9    // 改变宿主受限。
52
53 #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。
54 #include <sys/times.h>        // 定义了进程中运行时间结构 tms 以及 times() 函数原型。
55 #include <sys/utsname.h>      // 系统名称结构头文件。
56 #include <utime.h>            // 用户时间头文件。定义了访问和修改时间结构以及 utime() 原型。
57
58 #ifdef LIBRARY
59
// 以下是内核实现的系统调用符号常数，用于作为系统调用函数表中的索引值。(include/linux/sys.h)
60 #define NR_setup              0    /* used only by init, to get system going */
/* __NR_setup 仅用于初始化，以启动系统 */
61 #define NR_exit               1
62 #define NR_fork               2
63 #define NR_read               3
64 #define NR_write              4
65 #define NR_open               5
66 #define NR_close              6
67 #define NR_waitpid            7
68 #define NR_creat              8
69 #define NR_link               9
70 #define NR_unlink             10
71 #define NR_execve             11
72 #define NR_chdir              12
73 #define NR_time               13
74 #define NR_mknod              14
75 #define NR_chmod              15
76 #define NR_chown              16
77 #define NR_break              17
78 #define NR_stat               18
79 #define NR_lseek              19
80 #define NR_getpid             20
81 #define NR_mount              21
82 #define NR_umount             22
83 #define NR_setuid             23
84 #define NR_getuid             24
85 #define NR_stime              25
86 #define NR_ptrace             26
87 #define NR_alarm              27
88 #define NR_fstat              28
89 #define NR_pause              29
90 #define NR_utime              30
91 #define NR_stty               31
92 #define NR_gtty               32
93 #define NR_access             33
94 #define NR_nice               34
95 #define NR_ftime              35
96 #define NR_sync               36
97 #define NR_kill               37
98 #define NR_rename             38

```



```

99 #define  __NR_mkdir      39
100 #define  __NR_rmdir     40
101 #define  __NR_dup       41
102 #define  __NR_pipe      42
103 #define  __NR_times     43
104 #define  __NR_prof      44
105 #define  __NR_brk       45
106 #define  __NR_setgid    46
107 #define  __NR_getgid    47
108 #define  __NR_signal    48
109 #define  __NR_geteuid    49
110 #define  __NR_getegid    50
111 #define  __NR_acct      51
112 #define  __NR_phys      52
113 #define  __NR_lock      53
114 #define  __NR_ioctl     54
115 #define  __NR_fcntl     55
116 #define  __NR_mpx       56
117 #define  __NR_setpgid   57
118 #define  __NR_ulimit    58
119 #define  __NR_uname      59
120 #define  __NR_umask     60
121 #define  __NR_chroot    61
122 #define  __NR_ustat     62
123 #define  __NR_dup2      63
124 #define  __NR_getppid   64
125 #define  __NR_getpgrp   65
126 #define  __NR_setsid    66
127 #define  __NR_sigaction 67
128 #define  __NR_sgetmask  68
129 #define  __NR_ssetmask  69
130 #define  __NR_setreuid   70
131 #define  __NR_setregid   71
132
// 以下定义系统调用嵌入式汇编宏函数。
// 不带参数的系统调用宏函数。type name(void)。
// %0 - eax(__res), %1 - eax(__NR_###name)。其中 name 是系统调用的名称，与 __NR_ 组合形成上面
// 的系统调用符号常数，从而用来对系统调用表中函数指针寻址。
// 返回：如果返回值大于等于 0，则返回该值，否则置出错号 errno，并返回-1。
133 #define  __syscall0(type,name) \
134 type name(void) \
135 { \
136     long __res; \
137     __asm__ volatile ("int $0x80" \           // 调用系统中断 0x80。
138         : "=a" (__res) \                     // 返回值→eax(__res)。
139         : "i" (__NR_###name)); \           // 输入为系统中断调用号__NR_name。
140     if (__res >= 0) \                         // 如果返回值>=0，则直接返回该值。
141         return (type) __res; \
142     errno = -__res; \                         // 否则置出错号，并返回-1。
143     return -1; \
144 }
145
// 有 1 个参数的系统调用宏函数。type name(atype a)

```

```

    // %0 - eax(__res), %1 - eax(__NR_name), %2 - ebx(a)。
146 #define syscall1(type,name,atype,a) \
147 type name(atype a) \
148 { \
149 long __res; \
150 __asm__ volatile ("int $0x80" \
151 : "=a" (__res) \
152 : "" (__NR_##name), "b" ((long)(a))); \
153 if (__res >= 0) \
154     return (type) __res; \
155 errno = -__res; \
156 return -1; \
157 }
158
    // 有 2 个参数的系统调用宏函数。type name(atype a, btype b)
    // %0 - eax(__res), %1 - eax(__NR_name), %2 - ebx(a), %3 - ecx(b)。
159 #define syscall2(type,name,atype,a,btype,b) \
160 type name(atype a,btype b) \
161 { \
162 long __res; \
163 __asm__ volatile ("int $0x80" \
164 : "=a" (__res) \
165 : "" (__NR_##name), "b" ((long)(a)), "c" ((long)(b))); \
166 if (__res >= 0) \
167     return (type) __res; \
168 errno = -__res; \
169 return -1; \
170 }
171
    // 有 3 个参数的系统调用宏函数。type name(atype a, btype b, ctype c)
    // %0 - eax(__res), %1 - eax(__NR_name), %2 - ebx(a), %3 - ecx(b), %4 - edx(c)。
172 #define syscall3(type,name,atype,a,btype,b,ctype,c) \
173 type name(atype a,btype b,ctype c) \
174 { \
175 long __res; \
176 __asm__ volatile ("int $0x80" \
177 : "=a" (__res) \
178 : "" (__NR_##name), "b" ((long)(a)), "c" ((long)(b)), "d" ((long)(c))); \
179 if (__res>=0) \
180     return (type) __res; \
181 errno=-__res; \
182 return -1; \
183 }
184
185 #endif /* __LIBRARY__ */
186
187 extern int errno; // 出错号, 全局变量。
188
    // 对应各系统调用的函数原型定义。
189 int access(const char * filename, mode_t mode);
190 int acct(const char * filename);
191 int alarm(int sec);
192 int brk(void * end_data_segment);

```

```

193 void * sbrk(ptrdiff_t increment);
194 int chdir(const char * filename);
195 int chmod(const char * filename, mode\_t mode);
196 int chown(const char * filename, uid\_t owner, gid\_t group);
197 int chroot(const char * filename);
198 int close(int fildes);
199 int creat(const char * filename, mode\_t mode);
200 int dup(int fildes);
201 int execve(const char * filename, char ** argv, char ** envp);
202 int execv(const char * pathname, char ** argv);
203 int execvp(const char * file, char ** argv);
204 int execl(const char * pathname, char * arg0, ...);
205 int execlp(const char * file, char * arg0, ...);
206 int execle(const char * pathname, char * arg0, ...);
207 volatile void exit(int status);
208 volatile void \_exit(int status);
209 int fcntl(int fildes, int cmd, ...);
210 int fork(void);
211 int getpid(void);
212 int getuid(void);
213 int geteuid(void);
214 int getgid(void);
215 int getegid(void);
216 int ioctl(int fildes, int cmd, ...);
217 int kill(pid\_t pid, int signal);
218 int link(const char * filename1, const char * filename2);
219 int lseek(int fildes, off\_t offset, int origin);
220 int mknod(const char * filename, mode\_t mode, dev\_t dev);
221 int mount(const char * specialfile, const char * dir, int rwflag);
222 int nice(int val);
223 int open(const char * filename, int flag, ...);
224 int pause(void);
225 int pipe(int * fildes);
226 int read(int fildes, char * buf, off\_t count);
227 int setpgrp(void);
228 int setpgid(pid\_t pid, pid\_t pgid);
229 int setuid(uid\_t uid);
230 int setgid(gid\_t gid);
231 void (*signal(int sig, void (*fn)(int)))(int);
232 int stat(const char * filename, struct stat * stat_buf);
233 int fstat(int fildes, struct stat * stat_buf);
234 int stime(time\_t * tptr);
235 int sync(void);
236 time\_t time(time\_t * tloc);
237 time\_t times(struct tms * tbuf);
238 int ulimit(int cmd, long limit);
239 mode\_t umask(mode\_t mask);
240 int umount(const char * specialfile);
241 int uname(struct utsname * name);
242 int unlink(const char * filename);
243 int ustat(dev\_t dev, struct ustat * ubuf);
244 int utime(const char * filename, struct utimbuf * times);
245 pid\_t waitpid(pid\_t pid, int * wait_stat, int options);

```

---

```

246 pid_t wait(int * wait_stat);
247 int write(int fildes, const char * buf, off_t count);
248 int dup2(int oldfd, int newfd);
249 int getppid(void);
250 pid_t getpgrp(void);
251 pid_t setsid(void);
252
253 #endif
254

```

---

## 10.15 utime.h 文件

### 10.15.1 功能描述

该文件定义了文件访问和修改时间结构 `utimbuf {}` 以及 `utime()` 函数原型。

### 10.15.2 代码注释

列表 linux/include/utime.h 文件

---

```





1 #ifndef _UTIME_H
2 #define _UTIME_H
3
4 #include <sys/types.h> /* I know - shouldn't do this, but .. */
5                        /* 我知道 - 不应该这样做, 但是.. */
6 struct utimbuf {
7     time_t actime;      // 文件访问时间。从 1970.1.1:0:0:0 开始的秒数。
8     time_t modtime;     // 文件修改时间。从 1970.1.1:0:0:0 开始的秒数。
9 };
10
11 // 设置文件访问和修改时间函数。
12 extern int utime(const char *filename, struct utimbuf *times);
13 #endif
14

```

---

## 10.16 include/asm/目录下的文件

列表 linux/include/asm/目录下的文件

| Name  | Size       | Last modified (GMT) | Description |
|---|------------|---------------------|-------------|
|  <a href="#">io.h</a>      | 477 bytes  | 1991-08-07 10:17:51 | m           |
|  <a href="#">memory.h</a>  | 507 bytes  | 1991-06-15 20:54:44 | m           |
|  <a href="#">segment.h</a> | 1366 bytes | 1991-11-25 18:48:24 | m           |
|  <a href="#">system.h</a>  | 1711 bytes | 1991-09-17 13:08:31 | m           |

## 10.17 io.h 文件

### 10.17.1 功能描述

该文件中定义了对硬件 I/O 端口访问的嵌入式汇编宏函数：outb()、inb() 以及 outb\_p() 和 inb\_p()。前面两个函数与后面两个的主要区别在于后者代码中使用了 jmp 指令进行了时间延迟。

### 10.17.2 代码注释

列表 linux/include/asm/io.h 文件

```

1  // 硬件端口字节输出函数。
2  // 参数：value - 欲输出字节；port - 端口。
3  1 #define outb(value, port) \
4  2 __asm__ ("outb %%al, %%dx"::"a" (value), "d" (port))
5  3
6  4
7  // 硬件端口字节输入函数。
8  // 参数：port - 端口。返回读取的字节。
9  5 #define inb(port) ({ \
10 6 unsigned char _v; \
11 7 __asm__ volatile ("inb %%dx, %%al"::"a" (_v): "d" (port)); \
12 8 _v; \
13 9 })
14 10
15  // 带延迟的硬件端口字节输出函数。
16  // 参数：value - 欲输出字节；port - 端口。
17 11 #define outb_p(value, port) \
18 12 __asm__ ("outb %%al, %%dx\n" \

```

```

13         "\tjmp 1f\n" \
14         "1:\tjmp 1f\n" \
15         "1:": "a" (value), "d" (port))
16
17     /// 带延迟的硬件端口字节输入函数。
18     // 参数: port - 端口。返回读取的字节。
19 #define inb_p(port) ({ \
20     unsigned char _v; \
21     __asm__ volatile ("inb %%dx, %%al\n" \
22         "\tjmp 1f\n" \
23         "1:\tjmp 1f\n" \
24         "1:": "=a" (_v): "d" (port)); \
25     _v; \
26 })

```

## 10.18 memory.h 文件

### 10.18.1 功能描述

该文件含有一个内存复制嵌入式汇编宏 `memcpy()`。与 `string.h` 中定义的 `memcpy()` 相同，只是后者采用的是嵌入式汇编 C 函数形式定义的。

### 10.18.2 代码注释

列表 linux/include/asm/memory.h 文件

```

1  /*
2  * NOTE!!! memcpy(dest, src, n) assumes ds=es=normal data segment. This
3  * goes for all kernel functions (ds=es=kernel space, fs=local data,
4  * gs=null), as well as for all well-behaving user programs (ds=es=
5  * user data space). This is NOT a bug, as any user program that changes
6  * es deserves to die if it isn't careful.
7  */
8  /*
9  * 注意!!!memcpy(dest, src, n)假设段寄存器 ds=es=通常数据段。在内核中使用的
10 * 所有函数都基于该假设 (ds=es=内核空间, fs=局部数据空间, gs=null), 具有良好
11 * 行为的应用程序也是这样 (ds=es=用户数据空间)。如果任何用户程序随意改动了
12 * es 寄存器而出错, 则并不是由于系统程序错误造成的。
13 */
14 /// 内存块复制。从源地址 src 处开始复制 n 个字节到目的地址 dest 处。
15 // 参数: dest - 复制的目的地址, src - 复制的源地址, n - 复制字节数。
16 // %0 - edi(目的地址 dest), %1 - esi(源地址 src), %2 - ecx(字节数 n),
17 #define memcpy(dest, src, n) ({ \
18     void * _res = dest; \
19     __asm__ ("cld;rep;movsb" \
20         // 从 ds:[esi]复制到 es:[edi], 并且 esi++, edi++.
21         // 共复制 ecx(n) 字节。
22         :: "D" ((long) (_res)), "S" ((long) (src)), "c" ((long) (n)) \
23         : "di", "si", "cx"); \
24     _res; \
25 })

```

14 })  
15

## 10.19 segment.h 文件

### 10.19.1 功能描述

该文件中定义了一些访问段寄存器或与段寄存器有关的内存操作函数。

### 10.19.2 代码注释

列表 linux/include/asm/segment.h 文件

```

1  // 读取 fs 段中指定地址处的字节。
2  // 参数: addr - 指定的内存地址。
3  // %0 - (返回的字节_v); %1 - (内存地址 addr)。
4  // 返回: 返回内存 fs:[addr]处的字节。
5  extern inline unsigned char get_fs_byte(const char * addr)
6  {
7      unsigned register char _v;
8
9      __asm__ ("movb %%fs:%1, %0": "=r" (_v): "m" (*addr));
10     return _v;
11 }
12
13 // 读取 fs 段中指定地址处的字。
14 // 参数: addr - 指定的内存地址。
15 // %0 - (返回的字_v); %1 - (内存地址 addr)。
16 // 返回: 返回内存 fs:[addr]处的字。
17 extern inline unsigned short get_fs_word(const unsigned short *addr)
18 {
19     unsigned short _v;
20
21     __asm__ ("movw %%fs:%1, %0": "=r" (_v): "m" (*addr));
22     return _v;
23 }
24
25 // 读取 fs 段中指定地址处的长字(4 字节)。
26 // 参数: addr - 指定的内存地址。
27 // %0 - (返回的长字_v); %1 - (内存地址 addr)。
28 // 返回: 返回内存 fs:[addr]处的长字。
29 extern inline unsigned long get_fs_long(const unsigned long *addr)
30 {
31     unsigned long _v;
32
33     __asm__ ("movl %%fs:%1, %0": "=r" (_v): "m" (*addr)); \
34     return _v;
35 }
36
37 // 将一字节存放在 fs 段中指定内存地址处。

```

```

// 参数: val - 字节值; addr - 内存地址。
// %0 - 寄存器(字节值 val); %1 - (内存地址 addr)。
25 extern inline void put\_fs\_byte(char val, char *addr)
26 {
27     __asm__ ("movb %0, %%fs:%1":: "r" (val), "m" (*addr));
28 }
29
//// 将一字存放在 fs 段中指定内存地址处。
// 参数: val - 字值; addr - 内存地址。
// %0 - 寄存器(字值 val); %1 - (内存地址 addr)。
30 extern inline void put\_fs\_word(short val, short * addr)
31 {
32     __asm__ ("movw %0, %%fs:%1":: "r" (val), "m" (*addr));
33 }
34
//// 将一长字存放在 fs 段中指定内存地址处。
// 参数: val - 长字值; addr - 内存地址。
// %0 - 寄存器(长字值 val); %1 - (内存地址 addr)。
35 extern inline void put\_fs\_long(unsigned long val, unsigned long * addr)
36 {
37     __asm__ ("movl %0, %%fs:%1":: "r" (val), "m" (*addr));
38 }
39
40 /*
41  * Someone who knows GNU asm better than I should double check the followig.
42  * It seems to work, but I don't know if I'm doing something subtly wrong.
43  * --- TYT, 11/24/91
44  * [ nothing wrong here, Linus ]
45  */
46
/*
 * 比我更懂 GNU 汇编的人应该仔细检查下面的代码。这些代码能使用，但我不知道是否
 * 含有一些小错误。
 * --- TYT, 1991 年 11 月 24 日
 * [ 这些代码没有错误, Linus ]
 */
47
//// 取 fs 段寄存器值(选择符)。
// 返回: fs 段寄存器值。
48 extern inline unsigned long get\_fs()
49 {
50     unsigned short _v;
51     __asm__ ("mov %%fs, %%ax":: "=a" (_v));
52     return _v;
53 }
54
//// 取 ds 段寄存器值。
// 返回: ds 段寄存器值。
55 extern inline unsigned long get\_ds()
56 {
57     unsigned short _v;
58     __asm__ ("mov %%ds, %%ax":: "=a" (_v));
59     return _v;
60 }

```



```

60  // 设置 fs 段寄存器。
61  // 参数: val - 段值 (选择符)。
62  extern inline void set_fs(unsigned long val)
63  {
64      __asm__ ("mov %0, %%fs::%a" ((unsigned short) val));
65  }
66

```

## 10.20 system.h 文件

### 10.20.1 功能描述

该文件中定义了设置或修改描述符/中断门等的嵌入式汇编宏。

其中, 函数 `move_to_user_mode()` 是用于内核在初始化结束时切换到初始进程 (任务 0)。所使用的方法是模拟中断调用返回过程, 也即利用指令 `iret` 运行初始任务 0。

在切换到任务 0 之前, 首先设置堆栈, 模拟具有特权层切换的刚进入中断调用过程时堆栈的内容布置情况, 见下图所示。然后执行 `iret` 指令, 从而引起系统切换到任务 0 去执行。

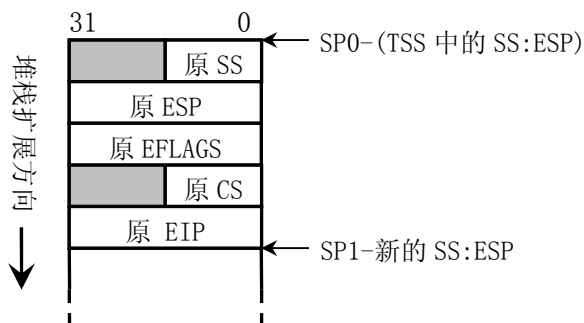


图 中断调用层间切换时堆栈内容

任务 0 是一个特殊进程, 它的数据段和代码段直接映射到内核代码和数据空间, 即从物理地址 0 开始的 640K 内存空间, 其堆栈地址也即内核代码所使用的堆栈。因此图中堆栈中的原 SS 和原 ESP 是直接将现有内核的堆栈指针压入堆栈的。

### 10.20.2 代码注释

列表 linux/include/asm/system.h 文件

```

// 切换到用户模式运行。
// 该函数利用 iret 指令实现从内核模式切换到用户模式 (初始任务 0)。
1 #define move_to_user_mode() \
2 __asm__ ("movl %%esp, %%eax\n\t" \           // 保存堆栈指针 esp 到 eax 寄存器中。
3         "pushl $0x17\n\t" \                 // 首先将堆栈段选择符(SS)入栈。

```

```

4      "pushl %%eax|n|t" \           // 然后将保存的堆栈指针值(esp)入栈。
5      "pushfl|n|t" \               // 将标志寄存器(eflags)内容入栈。
6      "pushl $0x0f|n|t" \          // 将内核代码段选择符(cs)入栈。
7      "pushl $1f|n|t" \            // 将下面标号 1 的偏移地址(eip)入栈。
8      "iret|n" \                   // 执行中断返回指令, 则会跳转到下面标号 1 处。
9      "l:\tmovl $0x17, %%eax|n|t" \ // 此时开始执行任务 0,
10     "movw %%ax, %%ds|n|t" \        // 初始化段寄存器指向本局部表的数据段。
11     "movw %%ax, %%es|n|t" \
12     "movw %%ax, %%fs|n|t" \
13     "movw %%ax, %%gs" \
14     ::: "ax")
15
16 #define sti() __asm__ ("sti"::)    // 开中断嵌入汇编宏函数。
17 #define cli() __asm__ ("cli"::)    // 关中断。
18 #define nop() __asm__ ("nop"::)    // 空操作。
19
20 #define iret() __asm__ ("iret"::)  // 中断返回。
21
22     //// 设置门描述符宏函数。
23     // 参数: gate_addr -描述符地址; type -描述符中类型域值; dpl -描述符特权层值; addr -偏移地址。
24     // %0 - (由 dpl, type 组合成的类型标志字); %1 - (描述符低 4 字节地址);
25     // %2 - (描述符高 4 字节地址); %3 - edx(程序偏移地址 addr); %4 - eax(高字中含有段选择符)。
26 #define set_gate(gate_addr, type, dpl, addr) \
27     __asm__ ("movw %%dx, %%ax|n|t" \           // 将偏移地址低字与选择符组合成描述符低 4 字节(eax)。
28     "movw %0, %%dx|n|t" \                     // 将类型标志字与偏移高字组合成描述符高 4 字节(edx)。
29     "movl %%eax, %1|n|t" \                   // 分别设置门描述符的低 4 字节和高 4 字节。
30     "movl %%edx, %2" \
31     : \
32     : "i" ((short) (0x8000+(dpl<<13)+(type<<8))), \
33     "o" (*(char *) (gate_addr)), \
34     "o" (*(4+(char *) (gate_addr))), \
35     "d" ((char *) (addr)), "a" (0x00080000))
36
37     //// 设置中断门函数。
38     // 参数: n - 中断号; addr - 中断程序偏移地址。
39     // &idt[n]对应中断号在中断描述符表中的偏移值; 中断描述符的类型是 14, 特权级是 0。
40 #define set_intr_gate(n, addr) \
41     set_gate(&idt[n], 14, 0, addr)
42
43     //// 设置陷阱门函数。
44     // 参数: n - 中断号; addr - 中断程序偏移地址。
45     // &idt[n]对应中断号在中断描述符表中的偏移值; 中断描述符的类型是 15, 特权级是 0。
46 #define set_trap_gate(n, addr) \
47     set_gate(&idt[n], 15, 0, addr)
48
49     //// 设置系统调用门函数。
50     // 参数: n - 中断号; addr - 中断程序偏移地址。
51     // &idt[n]对应中断号在中断描述符表中的偏移值; 中断描述符的类型是 15, 特权级是 3。
52 #define set_system_gate(n, addr) \
53     set_gate(&idt[n], 15, 3, addr)
54
55     //// 设置段描述符函数。
56     // 参数: gate_addr -描述符地址; type -描述符中类型域值; dpl -描述符特权层值;

```











```

// base - 段的基地址; limit - 段限长。(参见段描述符的格式)
42 #define set_seg_desc(gate_addr, type, dpl, base, limit) {\
43     *(gate_addr) = ((base) & 0xff000000) | \           // 描述符低 4 字节。
44                     (((base) & 0x00ff0000)>>16) | \
45                     ((limit) & 0xf0000) | \
46                     ((dpl)<<13) | \
47                     (0x00408000) | \
48                     ((type)<<8); \
49     *((gate_addr)+1) = (((base) & 0x0000ffff)<<16) | \ // 描述符高 4 字节。
50                     ((limit) & 0x0ffff); }
51
///// 在全局表中设置任务状态段/局部表描述符。
// 参数: n - 在全局表中描述符项 n 所对应的地址; addr - 状态段/局部表所在内存的基地址。
//      type - 描述符中的标志类型字节。
// %0 - eax(地址 addr); %1 - (描述符项 n 的地址); %2 - (描述符项 n 的地址偏移 2 处);
// %3 - (描述符项 n 的地址偏移 4 处); %4 - (描述符项 n 的地址偏移 5 处);
// %5 - (描述符项 n 的地址偏移 6 处); %6 - (描述符项 n 的地址偏移 7 处);
52 #define set_tssldt_desc(n, addr, type) \
53 __asm_ ("movw $104, %1|n|t" \           // 将 TSS 长度放入描述符长度域(第 0-1 字节)。
54         "movw %%ax, %2|n|t" \           // 将基地址的低字放入描述符第 2-3 字节。
55         "rorl $16, %%eax|n|t" \         // 将基地址高字移入 ax 中。
56         "movb %%al, %3|n|t" \           // 将基地址高字中低字节移入描述符第 4 字节。
57         "movb $" type ", %4|n|t" \      // 将标志类型字节移入描述符的第 5 字节。
58         "movb $0x00, %5|n|t" \         // 描述符的第 6 字节置 0。
59         "movb %%ah, %6|n|t" \           // 将基地址高字中高字节移入描述符第 7 字节。
60         "rorl $16, %%eax" \             // eax 清零。
61         :: "a" (addr), "m" (*(n)), "m" (*(n+2)), "m" (*(n+4)), \
62           "m" (*(n+5)), "m" (*(n+6)), "m" (*(n+7)) \
63         )
64
///// 在全局表中设置任务状态段描述符。
// n - 是该描述符的指针; addr - 是描述符中的基地址值。任务状态段描述符的类型是 0x89。
65 #define set_tss_desc(n, addr) set_tssldt_desc((char *) (n), addr, "0x89")
///// 在全局表中设置局部表描述符。
// n - 是该描述符的指针; addr - 是描述符中的基地址值。局部表描述符的类型是 0x82。
66 #define set_ldt_desc(n, addr) set_tssldt_desc((char *) (n), addr, "0x82")
67

```

## 10.21 include/linux/目录下的文件

列表 linux/include/linux/目录下的文件

| Name  | Size       | Last modified (GMT) | Description |
|---|------------|---------------------|-------------|
|  <a href="#">config.h</a>  | 1289 bytes | 1991-12-08 18:37:16 | m           |
|  <a href="#">fdreg.h</a>   | 2466 bytes | 1991-11-02 10:48:44 | m           |
|  <a href="#">fs.h</a>      | 5474 bytes | 1991-12-01 19:48:26 | m           |
|  <a href="#">hdreg.h</a>   | 1968 bytes | 1991-10-13 15:32:15 | m           |
|  <a href="#">head.h</a>    | 304 bytes  | 1991-06-19 19:24:13 | m           |
|  <a href="#">kernel.h</a>  | 734 bytes  | 1991-12-02 03:19:07 | m           |
|  <a href="#">mm.h</a>    | 219 bytes  | 1991-07-29 17:51:12 | m           |
|  <a href="#">sched.h</a> | 5838 bytes | 1991-11-20 14:40:46 | m           |
|  <a href="#">sys.h</a>   | 2588 bytes | 1991-11-25 20:15:35 | m           |
|  <a href="#">tty.h</a>   | 2173 bytes | 1991-09-21 11:58:05 | m           |

## 10.22 config.h 文件

### 10.22.1 功能描述

内核配置头文件。定义使用的键盘语言类型和硬盘类型（HD\_TYPE）可选项。

### 10.22.2 代码注释

列表 linux/include/linux/config.h 文件

```

1 #ifndef CONFIG\_H
2 #define CONFIG\_H
3
4 /*
```

```

5  * The root-device is no longer hard-coded. You can change the default
6  * root-device by changing the line ROOT_DEV = XXX in boot/bootsect.s
7  */
/*
* 根文件系统设备已不再是硬编码的了。通过修改 boot/bootsect.s 文件中行
* ROOT_DEV = XXX，你可以改变根设备的默认设置值。
*/

8
9 /*
10 * define your keyboard here -
11 * KBD_FINNISH for Finnish keyboards
12 * KBD_US for US-type
13 * KBD_GR for German keyboards
14 * KBD_FR for Frech keyboard
15 */
/*
* 在这里定义你的键盘类型 -
* KBD_FINNISH 是芬兰键盘。
* KBD_US 是美式键盘。
* KBD_GR 是德式键盘。
* KBD_FR 是法式键盘。
*/
16 /*#define KBD_US */
17 /*#define KBD_GR */
18 /*#define KBD_FR */
19 #define KBD_FINNISH
20
21 /*
22 * Normally, Linux can get the drive parameters from the BIOS at
23 * startup, but if this for some unfathomable reason fails, you'd
24 * be left stranded. For this case, you can define HD_TYPE, which
25 * contains all necessary info on your harddisk.
26 *
27 * The HD_TYPE macro should look like this:
28 *
29 * #define HD_TYPE { head, sect, cyl, wpcom, lzone, ctl}
30 *
31 * In case of two harddisks, the info should be sepatated by
32 * commas:
33 *
34 * #define HD_TYPE { h, s, c, wpcom, lz, ctl }, { h, s, c, wpcom, lz, ctl }
35 */
/*
* 通常，Linux 能够在启动时从 BIOS 中获取驱动器德参数，但是若由于未知原因
* 而没有得到这些参数时，会使程序束手无策。对于这种情况，你可以定义 HD_TYPE，
* 其中包括硬盘的所有信息。
*
* HD_TYPE 宏应该象下面这样的形式：
*
* #define HD_TYPE { head, sect, cyl, wpcom, lzone, ctl}
*
* 对于有两个硬盘的情况，参数信息需用逗号分开：
*

```

```

    * #define HD_TYPE { h, s, c, wpcom, lz, ctl }, {h, s, c, wpcom, lz, ctl }
    */
36 /*
37 This is an example, two drives, first is type 2, second is type 3:
38
39 #define HD_TYPE { 4, 17, 615, 300, 615, 8 }, { 6, 17, 615, 300, 615, 0 }
40
41 NOTE: ctl is 0 for all drives with heads<=8, and ctl=8 for drives
42 with more than 8 heads.
43
44 If you want the BIOS to tell what kind of drive you have, just
45 leave HD_TYPE undefined. This is the normal thing to do.
46 */
    /*
    * 下面是一个例子，两个硬盘，第 1 个是类型 2，第 2 个是类型 3：
    *
    * #define HD_TYPE { 4, 17, 615, 300, 615, 8 }, {6, 17, 615, 300, 615, 0 }
    *
    * 注意：对应所有硬盘，若其磁头数<=8，则 ctl 等于 0，若磁头数多于 8 个，
    * 则 ctl=8。
    *
    * 如果你想让 BIOS 给出硬盘的类型，那么只需不定义 HD_TYPE。这是默认操作。
    */
47
48 #endif
49

```

## 10.23 fdreg.h 头文件

### 10.23.1 功能描述

该头文件用以说明软盘系统常用到的一些参数以及所使用的 I/O 端口。由于软盘驱动器的控制比较烦琐，命令也多，因此在阅读代码之前，最好先参考有关微型计算机控制接口原理的书籍，了解软盘控制器 (FDC) 的工作原理，然后你就会觉得这里的定义还是比较合理有序的。

在编程时需要访问 4 个端口，分别对应一个或多个寄存器。对于 1.2M 的软盘控制器有以下一些端口。

表 软盘控制器端口

| I/O 端口 | 读写性 | 寄存器名称            |
|--------|-----|------------------|
| 0x3f2  | 只写  | 数字输出寄存器(数字控制寄存器) |
| 0x3f4  | 只读  | FDC 主状态寄存器       |
| 0x3f5  | 读/写 | FDC 数据寄存器        |
| 0x3f7  | 只读  | 数字输入寄存器          |
| 0x3f7  | 只写  | 磁盘控制寄存器(传输率控制)   |

数字输出端口（数字控制端口）是一个 8 位寄存器，它控制驱动器马达开启、驱动器选择、启动/复位 FDC 以及允许/禁止 DMA 及中断请求。

FDC 的主状态寄存器也是一个 8 位寄存器，用于反映软盘控制器 FDC 和软盘驱动器 FDD 的基本状态。通常，在 CPU 向 FDC 发送命令之前或从 FDC 获取操作结果之前，都要读取主状态寄存器的状态位，以判别当前 FDC 数据寄存器是否就绪，以及确定数据传送的方向。

FDC 的数据端口对应多个寄存器（只写型命令寄存器和参数寄存器、只读型结果寄存器），但任一时刻只能有一个寄存器出现在数据端口 0x3f5。在访问只写型寄存器时，主状态控制的 DIO 方向位必须为 0（CPU → FDC），访问只读型寄存器时则反之。在读取结果时只有在 FDC 不忙之后才算读完结果，通常结果数据最多有 7 个字节。

软盘控制器共可以接受 15 条命令。每个命令均经历三个阶段：命令阶段、执行阶段和结果阶段。

命令阶段是 CPU 向 FDC 发送命令字节和参数字节。每条命令的第一个字节总是命令字节（命令码）。其后跟着 0—8 字节的参数。执行阶段是 FDC 执行命令规定的操作。在执行阶段 CPU 是不加干预的，一般是通过 FDC 发出中断请求获知命令执行的结束。如果 CPU 发出的 FDC 命令是传送数据，则 FDC 可以以中断方式或 DMA 方式进行。中断方式每次传送 1 字节。DMA 方式是在 DMA 控制器管理下，FDC 与内存进行数据的传输直至全部数据传送完。此时 DMA 控制器会将传输字节计数终止信号通知 FDC，最后由 FDC 发出中断请求信号告知 CPU 执行阶段结束。结果阶段是由 CPU 读取 FDC 数据寄存器返回值，从而获得 FDC 命令执行的结果。返回结果数据的长度为 0—7 字节。对于没有返回结果数据的命令，则应向 FDC 发送检测中断状态命令获得操作的状态。

## 10.23.2 文件注释

列表 linux/include/linux/fdreg.h 文件

```

1 /*
2  * This file contains some defines for the floppy disk controller.
3  * Various sources. Mostly "IBM Microcomputers: A Programmers
4  * Handbook", Sanches and Canton.
5  */
6 /*
7  * 该文件中含有一些软盘控制器的一些定义。这些信息有多处来源，大多数取自 Sanches 和 Canton
8  * 编著的"IBM 微型计算机：程序员手册"一书。
9  */
10 #ifndef FDREG_H // 该定义用来排除代码中重复包含此头文件。
11 #define FDREG_H
12
13 // 一些软盘类型函数的原型说明。
14 extern int ticks_to_floppy_on(unsigned int nr);
15 extern void floppy_on(unsigned int nr);
16 extern void floppy_off(unsigned int nr);
17 extern void floppy_select(unsigned int nr);
18 extern void floppy_deselect(unsigned int nr);
19
20 // 下面是有关软盘控制器一些端口和符号的定义。
21 /* Fd controller regs. S&C, about page 340 */
22 /* 软盘控制器(FDC)寄存器端口。摘自 S&C 书中约 340 页 */
23 #define FD_STATUS 0x3f4 // 主状态寄存器端口。
24 #define FD_DATA 0x3f5 // 数据端口。
25 #define FD_DOR 0x3f2 /* Digital Output Register */
26 // 数字输出寄存器（也称为数字控制寄存器）。
27 #define FD_DIR 0x3f7 /* Digital Input Register (read) */
28 // 数字输入寄存器。
29 #define FD_DCR 0x3f7 /* Diskette Control Register (write)*/

```



```

21                                     // 数据传输率控制寄存器。
22 /* Bits of main status register */
23 /* 主状态寄存器各比特位的含义 */
24 #define STATUS_BUSYMASK 0x0F /* drive busy mask */
25                                     // 驱动器忙位（每位对应一个驱动器）。
26                                     /* FDC busy */
27                                     // 软盘控制器忙。
28                                     /* 0- DMA mode */
29                                     // 0 - 为 DMA 数据传输模式，1 - 为非 DMA 模式。
30                                     /* 0- cpu->fdc */
31                                     // 传输方向：0 - CPU → fdc，1 - 相反。
32                                     /* Data reg ready */
33                                     // 数据寄存器就绪位。
34
35 /* Bits of FD_ST0 */
36 /* 状态字节 0 (ST0) 各比特位的含义 */
37 #define ST0_DS 0x03 /* drive select mask */
38                                     // 驱动器选择号（发生中断时驱动器号）。
39 #define ST0_HA 0x04 /* Head (Address) */
40                                     // 磁头号。
41 #define ST0_NR 0x08 /* Not Ready */
42                                     // 磁盘驱动器未准备好。
43 #define ST0_ECE 0x10 /* Equipment chech error */
44                                     // 设备检测出错（零磁道校准出错）。
45 #define ST0_SE 0x20 /* Seek end */
46                                     // 寻道或重新校正操作执行结束。
47 #define ST0_INTR 0xC0 /* Interrupt code mask */
48                                     // 中断代码位（中断原因），00 - 命令正常结束；
49                                     // 01 - 命令异常结束；10 - 命令无效；11 - FDD 就绪状态改变。
50
51 /* Bits of FD_ST1 */
52 /* 状态字节 1 (ST1) 各比特位的含义 */
53 #define ST1_MAM 0x01 /* Missing Address Mark */
54                                     // 未找到地址标志 (ID AM)。
55 #define ST1_WP 0x02 /* Write Protect */
56                                     // 写保护。
57 #define ST1_ND 0x04 /* No Data - unreadable */
58                                     // 未找到指定的扇区。
59 #define ST1_OR 0x10 /* OverRun */
60                                     // 数据传输超时（DMA 控制器故障）。
61 #define ST1_CRC 0x20 /* CRC error in data or addr */
62                                     // CRC 检验出错。
63 #define ST1_EOC 0x80 /* End Of Cylinder */
64                                     // 访问超过一个磁道上的最大扇区号。
65
66 /* Bits of FD_ST2 */
67 /* 状态字节 2 (ST2) 各比特位的含义 */
68 #define ST2_MAM 0x01 /* Missing Address Mark (again) */
69                                     // 未找到数据地址标志。
70 #define ST2_BC 0x02 /* Bad Cylinder */
71                                     // 磁道坏。
72 #define ST2_SNS 0x04 /* Scan Not Satisfied */

```



---

```

49 #define ST2_SEH          0x08          // 检索（扫描）条件不满足。
                                         /* Scan Equal Hit */
                                         // 检索条件满足。
50 #define ST2_WC           0x10          /* Wrong Cylinder */
                                         // 磁道（柱面）号不符。
51 #define ST2_CRC           0x20          /* CRC error in data field */
                                         // 数据场 CRC 校验错。
52 #define ST2_CM           0x40          /* Control Mark = deleted */
                                         // 读数据遇到删除标志。

53
54 /* Bits of FD_ST3 */
   /* 状态字节 3 (ST3) 各比特位的含义 */
55 #define ST3_HA           0x04          /* Head (Address) */
                                         // 磁头号。
56 #define ST3_TZ           0x10          /* Track Zero signal (1=track 0) */
                                         // 零磁道信号。
57 #define ST3_WP           0x40          /* Write Protect */
                                         // 写保护。

58
59 /* Values for FD_COMMAND */
   /* 软盘命令码 */
60 #define FD_RECALIBRATE   0x07          /* move to track 0 */
                                         // 重新校正(磁头退到零磁道)。
61 #define FD_SEEK          0x0F          /* seek track */
                                         // 磁头寻道。
62 #define FD_READ          0xE6          /* read with MT, MFM, SKip deleted */
                                         // 读数据 (MT 多磁道操作, MFM 格式, 跳过删除数据)。
63 #define FD_WRITE         0xC5          /* write with MT, MFM */
                                         // 写数据 (MT, MFM)。
64 #define FD_SENSEI        0x08          /* Sense Interrupt Status */
                                         // 检测中断状态。
65 #define FD_SPECIFY       0x03          /* specify HUT etc */
                                         // 设定驱动器参数 (步进速率、磁头卸载时间等)。

66
67 /* DMA commands */
   /* DMA 命令 */
68 #define DMA_READ         0x46          // DMA 读盘, DMA 方式字 (送 DMA 端口 12, 11)。
69 #define DMA_WRITE        0x4A          // DMA 写盘, DMA 方式字。
70
71 #endif
72

```

---

## 10.24 fs.h 文件

### 10.24.1 功能描述

## 10.24.2 代码注释

列表 linux/include/linux/fs.h 文件

```

1  /*
2  * This file has definitions for some important file table
3  * structures etc.
4  */
5  /*
6  * 本文件含有某些重要文件表结构的定义等。
7  */
8
9  #ifndef \_FS\_H
10 #define \_FS\_H
11
12 #include <sys/types.h>    // 类型头文件。定义了基本的系统数据类型。
13
14 /* devices are as follows: (same as minix, so we can use the minix
15 * file system. These are major numbers.)
16 *
17 * 0 - unused (nodev)
18 * 1 - /dev/mem
19 * 2 - /dev/fd
20 * 3 - /dev/hd
21 * 4 - /dev/ttyx
22 * 5 - /dev/tty
23 * 6 - /dev/lp
24 * 7 - unnamed pipes
25 */
26 /*
27 * 系统所含的设备如下：（与 minix 系统的一样，所以我们可以使用 minix 的
28 * 文件系统。以下这些是主设备号。）
29 *
30 * 0 - 没有用到（nodev）
31 * 1 - /dev/mem      内存设备。
32 * 2 - /dev/fd       软盘设备。
33 * 3 - /dev/hd       硬盘设备。
34 * 4 - /dev/ttyx     tty 串行终端设备。
35 * 5 - /dev/tty      tty 终端设备。
36 * 6 - /dev/lp       打印设备。
37 * 7 - unnamed pipes 没有命名的管道。
38 */
39
40 #define IS\_SEEKABLE(x) ((x)>=1 && (x)<=3)    // 是否是寻找定位的设备。
41
42 #define READ 0
43 #define WRITE 1
44 #define READA 2    /* read-ahead - don't pause */
45 #define WRITEA 3    /* "write-ahead" - silly, but somewhat useful */
46
47 void buffer\_init(long buffer_end);
48
49 #define MAJOR(a) (((unsigned)(a))>>8)    // 取高字节（主设备号）。
50 #define MINOR(a) ((a)&0xff)            // 取低字节（次设备号）。

```

```

35
36 #define NAME_LEN 14 // 名字长度值。
37 #define ROOT_INO 1 // 根 i 节点。
38
39 #define I_MAP_SLOTS 8 // i 节点位图槽数。
40 #define Z_MAP_SLOTS 8 // 逻辑块（区段块）位图槽数。
41 #define SUPER_MAGIC 0x137F // 文件系统魔数。
42
43 #define NR_OPEN 20 // 打开文件数。
44 #define NR_INODE 32
45 #define NR_FILE 64
46 #define NR_SUPER 8
47 #define NR_HASH 307
48 #define NR_BUFFERS nr_buffers
49 #define BLOCK_SIZE 1024 // 数据块长度。
50 #define BLOCK_SIZE_BITS 10 // 数据块长度所占比特位数。
51 #ifndef NULL
52 #define NULL ((void *) 0)
53 #endif
54
55 // 每个逻辑块可存放的 i 节点数。
56 #define INODES_PER_BLOCK ((BLOCK_SIZE)/(sizeof (struct d_inode)))
57 // 每个逻辑块可存放的目录项数。
58 #define DIR_ENTRIES_PER_BLOCK ((BLOCK_SIZE)/(sizeof (struct dir_entry)))
59
60 // 管道头、管道尾、管道大小、管道空?、管道满?、管道头指针递增。
61 #define PIPE_HEAD(inode) ((inode).i_zone[0])
62 #define PIPE_TAIL(inode) ((inode).i_zone[1])
63 #define PIPE_SIZE(inode) ((PIPE_HEAD(inode)-PIPE_TAIL(inode))&(PAGE_SIZE-1))
64 #define PIPE_EMPTY(inode) (PIPE_HEAD(inode)==PIPE_TAIL(inode))
65 #define PIPE_FULL(inode) (PIPE_SIZE(inode)==(PAGE_SIZE-1))
66 #define INC_PIPE(head) \
67 __asm__("incl %0\n\tandl $4095, %0": "m" (head))
68
69 typedef char buffer_block[BLOCK_SIZE]; // 块缓冲区。
70
71 // 缓冲区头数据结构。（极为重要!!!）
72 // 在程序中常用 bh 来表示 buffer_head 类型的缩写。
73 struct buffer_head {
74     char * b_data; // /* pointer to data block (1024 bytes) */ // 数据块。
75     unsigned long b_blocknr; // /* block number */ // 块号。
76     unsigned short b_dev; // /* device (0 = free) */ // 数据源的设备号。
77     unsigned char b_uptodate; // 更新标志：表示数据是否已更新。
78     unsigned char b_dirty; // /* 0-clean, 1-dirty */ // 修改标志：0-未修改，1-已修改。
79     unsigned char b_count; // /* users using this block */ // 使用的用户数。
80     unsigned char b_lock; // /* 0-ok, 1-locked */ // 缓冲区是否被锁定。
81     struct task_struct * b_wait; // 指向等待该缓冲区解锁的任务。
82     struct buffer_head * b_prev; // 前一块（这四个指针用于缓冲区的管理）。
83     struct buffer_head * b_next; // 下一块。
84     struct buffer_head * b_prev_free; // 前一空闲块。
85     struct buffer_head * b_next_free; // 下一空闲块。
86 };
87
88

```

```

// 磁盘上的索引节点(i 节点)数据结构。
83 struct d_inode {
84     unsigned short i_mode;           // 文件类型和属性(rwx 位)。
85     unsigned short i_uid;           // 用户 id (文件拥有者标识符)。
86     unsigned long i_size;           // 文件大小 (字节数)。
87     unsigned long i_time;           // 修改时间 (自 1970.1.1:0 算起, 秒)。
88     unsigned char i_gid;           // 组 id (文件拥有者所在的组)。
89     unsigned char i_nlinks;         // 链接数 (多少个文件目录项指向该 i 节点)。
90     unsigned short i_zone[9];       // 直接 (0-6)、间接 (7) 或双重间接 (8) 逻辑块号。
                                     // zone 是区的意思, 可译成区段, 或逻辑块。
91 };
92
// 这是在内存中的 i 节点结构。前 7 项与 d_inode 完全一样。
93 struct m_inode {
94     unsigned short i_mode;           // 文件类型和属性(rwx 位)。
95     unsigned short i_uid;           // 用户 id (文件拥有者标识符)。
96     unsigned long i_size;           // 文件大小 (字节数)。
97     unsigned long i_mtime;          // 修改时间 (自 1970.1.1:0 算起, 秒)。
98     unsigned char i_gid;           // 组 id (文件拥有者所在的组)。
99     unsigned char i_nlinks;         // 文件目录项链接数。
100    unsigned short i_zone[9];        // 直接 (0-6)、间接 (7) 或双重间接 (8) 逻辑块号。
101    /* these are in memory also */
102    struct task_struct * i_wait;    // 等待该 i 节点的进程。
103    unsigned long i_atime;           // 最后访问时间。
104    unsigned long i_ctime;           // i 节点自身修改时间。
105    unsigned short i_dev;            // i 节点所在的设备号。
106    unsigned short i_num;            // i 节点号。
107    unsigned short i_count;          // i 节点被使用的次数, 0 表示该 i 节点空闲。
108    unsigned char i_lock;            // 锁定标志。
109    unsigned char i_dirt;            // 已修改(脏)标志。
110    unsigned char i_pipe;            // 管道标志。
111    unsigned char i_mount;           // 安装标志。
112    unsigned char i_seek;            // 搜寻标志(lseek 时)。
113    unsigned char i_update;          // 更新标志。
114 };
115
// 文件结构 (用于在文件句柄与 i 节点之间建立关系)
116 struct file {
117     unsigned short f_mode;           // 文件操作模式 (RW 位)
118     unsigned short f_flags;          // 文件打开和控制的标志。
119     unsigned short f_count;          // 对应文件句柄 (文件描述符) 数。
120     struct m_inode * f_inode;      // 指向对应 i 节点。
121     off_t f_pos;                   // 文件位置 (读写偏移值)。
122 };
123
// 内存中磁盘超级块结构。
124 struct super_block {
125     unsigned short s_ninodes;        // 节点数。
126     unsigned short s_nzones;         // 逻辑块数。
127     unsigned short s_imap_blocks;    // i 节点位图所占用的数据块数。
128     unsigned short s_zmap_blocks;    // 逻辑块位图所占用的数据块数。
129     unsigned short s_firstdatazone; // 第一个数据逻辑块号。
130     unsigned short s_log_zone_size;  // log(数据块数/逻辑块)。(以 2 为底)。

```

```

131     unsigned long s_max_size;           // 文件最大长度。
132     unsigned short s_magic;             // 文件系统魔数。
133     /* These are only in memory */
134     struct buffer head * s_imap[8];    // i 节点位图缓冲块指针数组(占用 8 块, 可表示 64M)。
135     struct buffer head * s_zmap[8];    // 逻辑块位图缓冲块指针数组(占用 8 块)。
136     unsigned short s_dev;               // 超级块所在的设备号。
137     struct m\_inode * s_isup;            // 被安装的文件系统根目录的 i 节点。(isup=super i)
138     struct m\_inode * s_imount;         // 被安装到的 i 节点。
139     unsigned long s_time;               // 修改时间。
140     struct task\_struct * s_wait;       // 等待该超级块的进程。
141     unsigned char s_lock;               // 被锁定标志。
142     unsigned char s_rd_only;            // 只读标志。
143     unsigned char s_dirt;               // 已修改(脏)标志。
144 };
145
146 // 磁盘上超级块结构。上面 125-132 行完全一样。
147 struct d\_super\_block {
148     unsigned short s_ninodes;           // 节点数。
149     unsigned short s_nzones;           // 逻辑块数。
150     unsigned short s_imap_blocks;      // i 节点位图所占用的数据块数。
151     unsigned short s_zmap_blocks;      // 逻辑块位图所占用的数据块数。
152     unsigned short s_firstdatazone;    // 第一个数据逻辑块。
153     unsigned short s_log_zone_size;    // log(数据块数/逻辑块)。(以 2 为底)。
154     unsigned long s_max_size;          // 文件最大长度。
155     unsigned short s_magic;            // 文件系统魔数。
156 };
157
158 // 文件目录项结构。
159 struct dir\_entry {
160     unsigned short inode;               // i 节点。
161     char name[NAME_LEN];               // 文件名。
162 };
163
164 extern struct m\_inode inode\_table[NR_INODE]; // 定义 i 节点表数组(32 项)。
165 extern struct file file\_table[NR_FILE];      // 文件表数组(64 项)。
166 extern struct super\_block super\_block[NR_SUPER]; // 超级块数组(8 项)。
167 extern struct buffer head * start\_buffer;      // 缓冲区起始内存位置。
168 extern int nr\_buffers;                       // 缓冲块数。
169
170 // 磁盘操作函数原型。
171 // 检测驱动器中软盘是否改变。
172 extern void check\_disk\_change(int dev);
173 // 检测指定软驱中软盘更换情况。如果软盘更换了则返回 1, 否则返回 0。
174 extern int floppy\_change(unsigned int nr);
175 // 设置启动指定驱动器所需等待的时间(设置等待定时器)。
176 extern int ticks\_to\_floppy\_on(unsigned int dev);
177 // 启动指定驱动器。
178 extern void floppy\_on(unsigned int dev);
179 // 关闭指定的软盘驱动器。
180 extern void floppy\_off(unsigned int dev);
181
182 // 以下是文件系统操作管理用的函数原型。
183 // 将 i 节点指定的文件截为 0。

```

```

173 extern void truncate(struct m\_inode * inode);
    // 刷新 i 节点信息。
174 extern void sync\_inodes(void);
    // 等待指定的 i 节点。
175 extern void wait\_on(struct m\_inode * inode);
    // 逻辑块(区段, 磁盘块)位图操作。取数据块 block 在设备上对应的逻辑块号。
176 extern int bmap(struct m\_inode * inode, int block);
    // 创建数据块 block 在设备上对应的逻辑块, 并返回在设备上的逻辑块号。
177 extern int create\_block(struct m\_inode * inode, int block);
    // 获取指定路径名的 i 节点号。
178 extern struct m\_inode * namei(const char * pathname);
    // 根据路径名为打开文件操作作准备。
179 extern int open\_namei(const char * pathname, int flag, int mode,
180     struct m\_inode ** res_inode);
    // 释放一个 i 节点(回写入设备)。
181 extern void iput(struct m\_inode * inode);
    // 从设备读取指定节点号的一个 i 节点。
182 extern struct m\_inode * iget(int dev, int nr);
    // 从 i 节点表(inode_table)中获取一个空闲 i 节点项。
183 extern struct m\_inode * get\_empty\_inode(void);
    // 获取(申请)管道节点。返回为 i 节点指针(如果是 NULL 则失败)。
184 extern struct m\_inode * get\_pipe\_inode(void);
    // 在哈希表中查找指定的数据块。返回找到块的缓冲头指针。
185 extern struct buffer\_head * get\_hash\_table(int dev, int block);
    // 从设备读取指定块(首先会在 hash 表中查找)。
186 extern struct buffer\_head * getblk(int dev, int block);
    // 读/写数据块。
187 extern void ll\_rw\_block(int rw, struct buffer\_head * bh);
    // 释放指定缓冲块。
188 extern void brelse(struct buffer\_head * buf);
    // 读取指定的数据块。
189 extern struct buffer\_head * bread(int dev, int block);
    // 读 4 块缓冲区到指定地址的内存中。
190 extern void bread\_page(unsigned long addr, int dev, int b[4]);
    // 读取头一个指定的数据块, 并标记后续将要读的块。
191 extern struct buffer\_head * breada(int dev, int block, ...);
    // 向设备 dev 申请一个磁盘块(区段, 逻辑块)。返回逻辑块号
192 extern int new\_block(int dev);
    // 释放设备数据区中的逻辑块(区段, 磁盘块)block。复位指定逻辑块 block 的逻辑块位图比特位。
193 extern void free\_block(int dev, int block);
    // 为设备 dev 建立一个新 i 节点, 返回 i 节点号。
194 extern struct m\_inode * new\_inode(int dev);
    // 释放一个 i 节点(删除文件时)。
195 extern void free\_inode(struct m\_inode * inode);
    // 刷新指定设备缓冲区。
196 extern int sync\_dev(int dev);
    // 读取指定设备的超级块。
197 extern struct super\_block * get\_super(int dev);
198 extern int ROOT\_DEV;
199
    // 安装根文件系统。
200 extern void mount\_root(void);
201

```

202 #endif  
203

## 10.25 hdreg.h 文件

### 10.25.1 功能描述

### 10.25.2 代码注释

列表 linux/include/linux/hdreg.h 文件

```

1  /*
2  * This file contains some defines for the AT-hd-controller.
3  * Various sources. Check out some definitions (see comments with
4  * a ques).
5  */
6  /*
7  * 本文件含有一些 AT 硬盘控制器的定义。来自各种资料。请查证某些
8  * 定义（带有问号的注释）。
9  */
10 #ifndef HDREG_H
11 #define HDREG_H
12
13 /* Hd controller regs. Ref: IBM AT Bios-listing */
14 /* 硬盘控制器寄存器端口。参见：IBM AT Bios 程序 */
15 #define HD_DATA 0x1f0 /* _CTL when writing */
16 #define HD_ERROR 0x1f1 /* see err-bits */
17 #define HD_NSECTOR 0x1f2 /* nr of sectors to read/write */
18 #define HD_SECTOR 0x1f3 /* starting sector */
19 #define HD_LCYL 0x1f4 /* starting cylinder */
20 #define HD_HCYL 0x1f5 /* high byte of starting cyl */
21 #define HD_CURRENT 0x1f6 /* 101dhhh , d=drive, hhhh=head */
22 #define HD_STATUS 0x1f7 /* see status-bits */
23 #define HD_PRECOMP HD_ERROR /* same io address, read=error, write=precomp */
24 #define HD_COMMAND HD_STATUS /* same io address, read=status, write=cmd */
25
26 #define HD_CMD 0x3f6 // 控制寄存器端口。
27
28 /* Bits of HD_STATUS */
29 /* 硬盘状态寄存器各位的定义(HD_STATUS) */
30 #define ERR_STAT 0x01 // 命令执行错误。
31 #define INDEX_STAT 0x02 // 收到索引。
32 #define ECC_STAT 0x04 /* Corrected error */ // ECC 校验错。
33 #define DRQ_STAT 0x08 // 请求服务。
34 #define SEEK_STAT 0x10 // 寻道结束。
35 #define WRERR_STAT 0x20 // 驱动器故障。
36 #define READY_STAT 0x40 // 驱动器准备好（就绪）。
37 #define BUSY_STAT 0x80 // 控制器忙碌。

```



---

```

33 /* Values for HD_COMMAND */
   /* 硬盘命令值 (HD_CMD) */
34 #define WIN_RESTORE          0x10    // 驱动器重新校正 (驱动器复位)。
35 #define WIN_READ              0x20    // 读扇区。
36 #define WIN_WRITE             0x30    // 写扇区。
37 #define WIN_VERIFY            0x40    // 扇区检验。
38 #define WIN_FORMAT            0x50    // 格式化磁道。
39 #define WIN_INIT               0x60    // 控制器初始化。
40 #define WIN_SEEK               0x70    // 寻道。
41 #define WIN_DIAGNOSE          0x90    // 控制器诊断。
42 #define WIN_SPECIFY            0x91    // 建立驱动器参数。
43
44 /* Bits for HD_ERROR */
   /* 错误寄存器各比特位的含义 (HD_ERROR) */
   // 执行控制器诊断命令时含义与其它命令时的不同。下面分别列出:
   // =====
   //          诊断命令时          其它命令时
   // -----
   // 0x01      无错误              数据标志丢失
   // 0x02      控制器出错          磁道 0 错
   // 0x03      扇区缓冲区错
   // 0x04      ECC 部件错          命令放弃
   // 0x05      控制处理器错
   // 0x10
   // 0x40
   // 0x80
   //          ID 未找到
   //          ECC 错误
   //          坏扇区
   // -----
45 #define MARK_ERR              0x01    /* Bad address mark ? */
46 #define TRKO_ERR              0x02    /* couldn't find track 0 */
47 #define ABRT_ERR              0x04    /* ? */
48 #define ID_ERR                0x10    /* ? */
49 #define ECC_ERR                0x40    /* ? */
50 #define BBD_ERR                0x80    /* ? */
51
   // 硬盘分区表结构。参见下面列表后信息。
52 struct partition {
53     unsigned char boot_ind;        /* 0x80 - active (unused) */
54     unsigned char head;            /* ? */
55     unsigned char sector;          /* ? */
56     unsigned char cyl;             /* ? */
57     unsigned char sys_ind;         /* ? */
58     unsigned char end_head;        /* ? */
59     unsigned char end_sector;      /* ? */
60     unsigned char end_cyl;         /* ? */
61     unsigned int start_sect;       /* starting sector counting from 0 */
62     unsigned int nr_sects;         /* nr of sectors in partition */
63 };
64
65 #endif
66

```

---



### 10.25.3 其它信息

#### 硬盘分区表

为了实现多个操作系统共享硬盘资源，硬盘可以在逻辑上分为 1—4 个分区。每个分区之间的扇区号是邻接的。分区表由 4 个表项组成，每个表项由 16 字节组成，对应一个分区的信息，存放有分区的大小和起止的柱面号、磁道号和扇区号，见下表所示。分区表存放在硬盘的 0 柱面 0 头第 1 个扇区的 0x1BE—0x1FD 处。

表 硬盘分区表结构

| 位置        | 名称         | 大小 | 说明  |
|-----------|------------|----|---|
| 0x00      | boot_ind   | 字节 | 引导标志。4 个分区中同时只能有一个分区是可引导的。0x00—不从该分区引导操作系统；0x80—从该分区引导操作系统。 |
| 0x01      | head       | 字节 | 分区起始磁头号。  |
| 0x02      | sector     | 字节 | 分区起始扇区号(位 0–5)和起始柱面号高 2 位(位 6–7)。                           |
| 0x03      | cyl        | 字节 | 分区起始柱面号低 8 位。   |
| 0x04      | sys_ind    | 字节 | 分区类型字节。0x0b—DOS；0x80—Old Minix；0x83—Linux …                 |
| 0x05      | end_head   | 字节 | 分区的结束磁头号。   |
| 0x06      | end_sector | 字节 | 结束扇区号(位 0–5)和结束柱面号高 2 位(位 6–7)。                             |
| 0x07      | end_cyl    | 字节 | 结束柱面号低 8 位。   |
| 0x08—0x0b | start_sect | 长字 | 分区起始物理扇区号。  |
| 0x0c—0x0f | nr_sects   | 长字 | 分区占用的扇区数。   |

### 10.26 head.h 文件

#### 10.26.1 功能描述

head 头文件，定义了与 head.s 程序描述符表有关的 C 程序。

#### 10.26.2 代码注释

列表 linux/include/linux/head.h 文件

```
1 #ifndef HEAD\_H
2 #define HEAD\_H
3
4 typedef struct desc\_struct {           // 定义了段描述符的数据结构。该结构仅说明每个描述
5     unsigned long a,b;                 // 符是由 8 个字节构成，每个描述符表共有 256 项。
6 } desc\_table[256];
7
8 extern unsigned long pg\_dir[1024];    // 内存页目录数组。每个目录项为 4 字节。从物理地址 0 开始。
9 extern desc\_table idt,gdt;           // 中断描述符表，全局描述符表。
10
11 #define GDT\_NUL 0                     // 全局描述符表的第 0 项，不用。
```

```

12 #define GDT_CODE 1 // 第 1 项，是内核代码段描述符项。
13 #define GDT_DATA 2 // 第 2 项，是内核数据段描述符项。
14 #define GDT_TMP 3 // 第 3 项，系统段描述符，Linux 没有使用。
15
16 #define LDT_NUL 0 // 每个局部描述符表的第 0 项，不用。
17 #define LDT_CODE 1 // 第 1 项，是用户程序代码段描述符项。
18 #define LDT_DATA 2 // 第 2 项，是用户程序数据段描述符项。
19
20 #endif
21

```

## 10.27 kernel.h 文件

### 10.27.1 功能描述

定义了一些内核常用的函数原型等。

### 10.27.2 代码注释

列表 linux/include/linux/kernel.h

```

1 /*
2  * 'kernel.h' contains some often-used function prototypes etc
3  */
4 /*
5  * 'kernel.h' 定义了一些常用函数的原型等。
6  */
7 // 验证给定地址开始的内存块是否超限。若超限则追加内存。( kernel/fork.c, 24 )。
8 void verify_area(void * addr, int count);
9 // 显示内核出错信息，然后进入死循环。( kernel/panic.c, 16 )。
10 volatile void panic(const char * str);
11 // 标准打印（显示）函数。( init/main.c, 151)。
12 int printf(const char * fmt, ...);
13 // 内核专用的打印信息函数，功能与 printf() 相同。( kernel/printk.c, 21 )。
14 int printk(const char * fmt, ...);
15 // 往 tty 上写指定长度的字符串。( kernel/chr_drv/tty_io.c, 290 )。
16 int tty_write(unsigned ch, char * buf, int count);
17 // 通用内核内存分配函数。( lib/malloc.c, 117)。
18 void * malloc(unsigned int size);
19 // 释放指定对象占用的内存。( lib/malloc.c, 182)。
20 void free_s(void * obj, int size);
21
22 #define free(x) free_s((x), 0)
23
24 /*
25  * This is defined as a macro, but at some point this might become a
26  * real subroutine that sets a flag if it returns true (to do
27  * BSD-style accounting where the process is flagged if it uses root
28  * privs). The implication of this is that you should do normal
29  * permissions checks first, and check suser() last.
30  */

```

---

```

20 */
   /*
   * 下面函数是以宏的形式定义的，但是在某方面来看它可以成为一个真正的子程序，
   * 如果返回是 true 时它将设置标志（如果使用 root 用户权限的进程设置了标志，则用
   * 于执行 BSD 方式的计帐处理）。这意味着你应该首先执行常规权限检查，最后再
   * 检测 suser()。
   */
21 #define suser() (current->euid == 0)           // 检测是否是超级用户。
22
23

```

---

## 10.28 mm.h 文件

### 10.28.1 功能描述

mm.h 是内存管理头文件。其中主要定义了内存页面的大小和几个页面释放函数原型。

### 10.28.2 代码注释

列表 linux/include/linux/mm.h 文件

---

```

1 #ifndef MM_H
2 #define MM_H
3
4 #define PAGE_SIZE 4096           // 定义内存页面的大小(字节数)。
5
6 // 取空闲页面函数。返回页面地址。扫描页面映射数组 mem_map[]取空闲页面。
6 extern unsigned long get_free_page(void);
// 在指定物理地址处放置一页面。在页目录和页表中放置指定页面信息。
7 extern unsigned long put_page(unsigned long page,unsigned long address);
// 释放物理地址 addr 开始的一页面内存。修改页面映射数组 mem_map[]中引用次数信息。
8 extern void free_page(unsigned long addr);
9
10 #endif
11

```

---

## 10.29 sched.h 文件

### 10.29.1 功能描述

调度程序头文件，定义了任务结构 task\_struct、初始任务 0 的数据，还有一些有关描述符参数设置和获取的嵌入式汇编函数宏语句。

### 10.29.2 代码注释

列表 linux/include/linux/sched.h 文件

---

```

1 #ifndef SCHED_H

```

---

```

2 #define _SCHED_H
3
4 #define NR_TASKS 64 // 系统中同时最多任务（进程）数。
5 #define HZ 100 // 定义系统时钟滴答频率(1 百赫兹，每个滴答 10ms)
6
7 #define FIRST_TASK task[0] // 任务 0 比较特殊，所以特意给它单独定义一个符号。
8 #define LAST_TASK task[NR_TASKS-1] // 任务数组中的最后一项任务。
9
10 #include <linux/head.h> // head 头文件，定义了与 head.s 程序描述符表有关的 C 程序。
11 #include <linux/fs.h> // 文件系统头文件。定义文件表结构（file,buffer_head,m_inode 等）。
12 #include <linux/mm.h> // 内存管理头文件。含有页面大小定义和一些页面释放函数原型。
13 #include <signal.h> // 信号头文件。定义信号符号常量，信号结构以及信号操作函数原型。
14
15 #if (NR_OPEN > 32)
16 #error "Currently the close-on-exec-flags are in one word, max 32 files/proc"
17 #endif
18
19 // 这里定义了进程运行可能处的状态。
20 #define TASK_RUNNING 0 // 进程正在运行或已准备就绪。
21 #define TASK_INTERRUPTIBLE 1 // 进程处于可中断等待状态。
22 #define TASK_UNINTERRUPTIBLE 2 // 进程处于不可中断等待状态，主要用于 I/O 操作等待。
23 #define TASK_ZOMBIE 3 // 进程处于僵死状态，已经停止运行，但父进程还没发信号。
24 #define TASK_STOPPED 4 // 进程已停止。
25
26 #ifndef NULL
27 #define NULL ((void *) 0) // 定义 NULL 为空指针。
28 #endif
29
30 // 复制进程的页目录页表。Linus 认为这是内核中最复杂的函数之一。（mm/memory.c, 105）
31 extern int copy_page_tables(unsigned long from, unsigned long to, long size);
32 // 释放页表所指定的内存块及页表本身。（mm/memory.c, 150）
33 extern int free_page_tables(unsigned long from, unsigned long size);
34
35 // 调度程序的初始化函数。（kernel/sched.c, 385）
36 extern void sched_init(void);
37 // 进程调度函数。（kernel/sched.c, 104）
38 extern void schedule(void);
39 // 异常(陷阱)中断处理初始化函数，设置中断调用门并允许中断请求信号。（kernel/traps.c, 181）
40 extern void trap_init(void);
41 // 显示内核出错信息，然后进入死循环。（kernel/panic.c, 16）。
42 extern void panic(const char *str);
43 // 往 tty 上写指定长度的字符串。（kernel/chr_drv/tty_io.c, 290）。
44 extern int tty_write(unsigned minor, char *buf, int count);
45
46 typedef int (*fn_ptr)(); // 定义函数指针类型。
47
48 // 下面是数学协处理器使用的结构，主要用于保存进程切换时 i387 的执行状态信息。
49 struct i387_struct {
50     long cwd; // 控制字(Control word)。
51     long swd; // 状态字(Status word)。
52     long twd; // 标记字(Tag word)。
53     long fip; // 协处理器代码指针。
54     long fcs; // 协处理器代码段寄存器。
55 }

```

```

46     long    foo;
47     long    fos;
48     long    st_space[20];    /* 8*10 bytes for each FP-reg = 80 bytes */
49 };
50
    // 任务状态段数据结构（参见列表后的信息）。
51 struct tss_struct {
52     long    back_link;    /* 16 high bits zero */
53     long    esp0;
54     long    ss0;    /* 16 high bits zero */
55     long    esp1;
56     long    ss1;    /* 16 high bits zero */
57     long    esp2;
58     long    ss2;    /* 16 high bits zero */
59     long    cr3;
60     long    eip;
61     long    eflags;
62     long    eax, ecx, edx, ebx;
63     long    esp;
64     long    ebp;
65     long    esi;
66     long    edi;
67     long    es;    /* 16 high bits zero */
68     long    cs;    /* 16 high bits zero */
69     long    ss;    /* 16 high bits zero */
70     long    ds;    /* 16 high bits zero */
71     long    fs;    /* 16 high bits zero */
72     long    gs;    /* 16 high bits zero */
73     long    ldt;    /* 16 high bits zero */
74     long    trace_bitmap;    /* bits: trace 0, bitmap 16-31 */
75     struct i387_struct i387;
76 };
77
    // 任务（进程）数据结构。
    // =====
    // long state        任务的运行状态（-1 不可运行，0 可运行（就绪），>0 已停止）。
    // long counter      任务运行时间计数（递减）（滴答数）。
    // long priority     运行优先数。任务开始运行时 counter = priority，越大运行越长。
    // long signal       信号。是位图，每个比特位代表一种信号，信号值=位偏移值+1。
    // struct sigaction sigaction[32] 信号执行属性结构，对应信号将要执行的操作和标志信息。
    // long blocked      进程信号屏蔽码（对应信号位图）。
    // -----
    // int exit_code      任务执行停止的退出码，其父进程会取。
    // unsigned long start_code 代码段地址。
    // unsigned long end_code  代码长度（字节数）。
    // unsigned long end_data  数据长度（字节数）。
    // unsigned long brk       总长度（字节数）。
    // unsigned long start_stack 堆栈段地址。
    // long pid              进程标识号（进程号）。
    // long father           父进程号。
    // long pgrp             父进程组号。
    // long session          会话号。
    // long leader           会话首领。

```

```

// unsigned short uid      用户标识号（用户 id）。
// unsigned short euid     有效用户 id。
// unsigned short suid     保存的用户 id。
// unsigned short gid      组标识号（组 id）。
// unsigned short egid     有效组 id。
// unsigned short sgid     保存的组 id。
// long alarm              报警定时值（滴答数）。
// long utime              用户态运行时间（滴答数）。
// long stime              系统态运行时间（滴答数）。
// long cutime             子进程用户态运行时间。
// long cstime             子进程系统态运行时间。
// long start_time         进程开始运行时刻。
// unsigned short used_math 标志：是否使用了协处理器。
// -----
// int tty                 进程使用 tty 的子设备号。-1 表示没有使用。
// unsigned short umask     文件创建属性屏蔽位。
// struct m_inode * pwd     当前工作目录 i 节点结构。
// struct m_inode * root    根目录 i 节点结构。
// struct m_inode * executable 执行文件 i 节点结构。
// unsigned long close_on_exec 执行时关闭文件句柄位图标志。（参见 include/fcntl.h）
// struct file * filp[NR_OPEN] 进程使用的文件表结构。
// -----
// struct desc_struct ldt[3] 本任务的局部表描述符。0-空，1-代码段 cs，2-数据和堆栈段 ds&ss。
// -----
// struct tss_struct tss    本进程的任务状态段信息结构。
// =====
78 struct task_struct {
79  /* these are hardcoded - don't touch */
80      long state;          /* -1 unrunnable, 0 runnable, >0 stopped */
81      long counter;
82      long priority;
83      long signal;
84      struct sigaction sigaction[32];
85      long blocked;        /* bitmap of masked signals */
86  /* various fields */
87      int exit_code;
88      unsigned long start_code, end_code, end_data, brk, start_stack;
89      long pid, father, pgrp, session, leader;
90      unsigned short uid, euid, suid;
91      unsigned short gid, egid, sgid;
92      long alarm;
93      long utime, stime, cutime, cstime, start_time;
94      unsigned short used_math;
95  /* file system info */
96      int tty;              /* -1 if no tty, so it must be signed */
97      unsigned short umask;
98      struct m_inode * pwd;
99      struct m_inode * root;
100     struct m_inode * executable;
101     unsigned long close_on_exec;
102     struct file * filp[NR_OPEN];
103  /* ldt for this task 0 - zero 1 - cs 2 - ds&ss */
104     struct desc_struct ldt[3];

```

```

105 /* tss for this task */
106     struct tss\_struct tss;
107 };
108
109 /*
110  * INIT_TASK is used to set up the first task table, touch at
111  * your own risk!. Base=0, limit=0x9ffff (=640kB)
112  */
113 /*
114  * INIT_TASK 用于设置第 1 个任务表，若想修改，责任自负☺！
115  * 基址 Base = 0，段长 limit = 0x9ffff (=640kB) 。
116  */
117 // 对应上面任务结构的第 1 个任务的信息。
118 #define INIT\_TASK \
119 /* state etc */ { 0, 15, 15, \      // state, counter, priority
120 /* signals */ 0, {}, 0, \      // signal, sigaction[32], blocked
121 /* ec, brk... */ 0, 0, 0, 0, 0, \  // exit_code, start_code, end_code, end_data, brk, start_stack
122 /* pid etc.. */ 0, -1, 0, 0, 0, \  // pid, father, pgrp, session, leader
123 /* uid etc */ 0, 0, 0, 0, 0, \  // uid, euid, suid, gid, egid, sgid
124 /* alarm */ 0, 0, 0, 0, 0, \  // alarm, utime, stime, cutime, cstime, start_time
125 /* math */ 0, \  // used_math
126 /* fs info */ -1, 0022, NULL, NULL, NULL, 0, \  // tty, umask, pwd, root, executable, close_on_exec
127 /* filp */ {NULL, }, \  // filp[20]
128 { \  // ldt[3]
129 {0, 0}, \
130 /* ldt */ {0x9f, 0xc0fa00}, \  // 代码长 640K, 基址 0x0, G=1, D=1, DPL=3, P=1 TYPE=0x0a
131 {0x9f, 0xc0f200}, \  // 数据长 640K, 基址 0x0, G=1, D=1, DPL=3, P=1 TYPE=0x02
132 }, \
133 /*tss*/ {0, PAGE\_SIZE+(long)&init\_task, 0x10, 0, 0, 0, 0, (long)&pg\_dir, \  // tss
134 0, 0, 0, 0, 0, 0, 0, \
135 0, 0, 0x17, 0x17, 0x17, 0x17, 0x17, \
136 LDT(0), 0x80000000, \
137 {} \
138 }, \
139 }
140
141 extern struct task\_struct *task[NR\_TASKS]; // 任务数组。
142 extern struct task\_struct *last\_task\_used\_math; // 上一个使用过协处理器的进程。
143 extern struct task\_struct *current; // 当前进程变量。
144 extern long volatile jiffies; // 从开机开始算起的滴答数（10ms/滴答）。
145 extern long startup\_time; // 开机时间。从 1970:0:0:0 开始计时的秒数。
146
147 #define CURRENT\_TIME (startup\_time+jiffies/HZ) // 当前时间（秒数）。
148
149 // 添加定时器函数（定时时间 jiffies 滴答数，定时到时调用函数*fn()）。（ kernel/sched.c, 272）
150 extern void add\_timer(long jiffies, void (*fn)(void));
151 // 不可中断的等待睡眠。（ kernel/sched.c, 151 ）
152 extern void sleep\_on(struct task\_struct ** p);
153 // 可中断的等待睡眠。（ kernel/sched.c, 167 ）
154 extern void interruptible\_sleep\_on(struct task\_struct ** p);
155 // 明确唤醒睡眠的进程。（ kernel/sched.c, 188 ）
156 extern void wake\_up(struct task\_struct ** p);
157
158

```



```

149 /*
150  * Entry into gdt where to find first TSS. 0-nul, 1-cs, 2-ds, 3-syscall
151  * 4-TSS0, 5-LDT0, 6-TSS1 etc ...
152 */
/*
* 寻找第 1 个 TSS 在全局表中的入口。0-没有用 nul, 1-代码段 cs, 2-数据段 ds, 3-系统段 syscall
* 4-任务状态段 TSS0, 5-局部表 LTD0, 6-任务状态段 TSS1, 等。
*/
// 全局表中第 1 个任务状态段(TSS)描述符的选择符索引号。
153 #define FIRST_TSS_ENTRY 4
// 全局表中第 1 个局部描述符表(LDT)描述符的选择符索引号。
154 #define FIRST_LDT_ENTRY (FIRST_TSS_ENTRY+1)
// 宏定义, 计算在全局表中第 n 个任务的 TSS 描述符的索引号 (选择符)。
155 #define TSS(n) (((unsigned long) n)<<4)+(FIRST_TSS_ENTRY<<3))
// 宏定义, 计算在全局表中第 n 个任务的 LDT 描述符的索引号。
156 #define LDT(n) (((unsigned long) n)<<4)+(FIRST_LDT_ENTRY<<3))
// 宏定义, 加载第 n 个任务的寄存器 tr。
157 #define ltr(n) __asm__ ("ltr %%ax::\"a\" (TSS(n))"
// 宏定义, 加载第 n 个任务的局部描述符表寄存器 ldtr。
158 #define lldt(n) __asm__ ("lldt %%ax::\"a\" (LDT(n))"
// 取当前运行任务的任务号 (是任务数组中的索引值, 与进程号 pid 不同)。
// 返回: n - 当前任务号。用于 (kernel/traps.c, 79)。
159 #define str(n) \
160 __asm__ ("str %%ax|n|t\" \           // 将任务寄存器中 TSS 段的有效地址→ax
161         "subl %2, %%eax|n|t\" \       // (eax - FIRST_TSS_ENTRY*8)→eax
162         "shrl $4, %%eax\" \           // (eax/16)→eax = 当前任务号。
163         : "=a" (n) \
164         : "a" (0), "i" (FIRST_TSS_ENTRY<<3))
165 /*
166  *      switch_to(n) should switch tasks to task nr n, first
167  * checking that n isn't the current task, in which case it does nothing.
168  * This also clears the TS-flag if the task we switched to has used
169  * the math co-processor latest.
170 */
/*
* switch_to(n)将切换当前任务到任务 nr, 即 n。首先检测任务 n 不是当前任务,
* 如果是则什么也不做退出。如果我们切换到任务最近 (上次运行) 使用过数学
* 协处理器的话, 则还需复位控制寄存器 cr0 中的 TS 标志。
*/
// 输入: %0 - 新 TSS 的偏移地址(*&__tmp.a); %1 - 存放新 TSS 的选择符值(*&__tmp.b);
// dx - 新任务 n 的选择符; ecx - 新任务指针 task[n]。
// 其中临时数据结构 __tmp 中, a 的值是 32 位偏移值, b 为新 TSS 的选择符。在任务切换时, a 值
// 没有用 (忽略)。
171 #define switch_to(n) {\
172 struct {long a,b;} __tmp; \
173 __asm__ ("cmpl %%ecx, _current|n|t\" \ // 任务 n 是当前任务吗?(current ==task[n]?)
174         "je 1f|n|t\" \               // 是, 则什么都不做, 退出。
175         "movw %%dx, %1|n|t\" \       // 将新任务的选择符→*&__tmp.b。
176         "xchgl %%ecx, _current|n|t\" \ // current = task[n]; ecx = 被切换出的任务。
177         "ljmp %0|n|t\" \             // 长跳转, 造成任务切换。
178         "cmpl %%ecx, _last_task_used_math|n|t\" \ // 新任务上次使用过协处理器吗?
179         "jne 1f|n|t\" \             // 没有则跳转, 退出。
180         "clts|n\" \                 // 新任务上次使用过协处理器, 则清 cr0 的 TS 标志。

```



```

181     "l:" \
182     :: "m" (&__tmp.a), "m" (&__tmp.b), \
183     "d" ( TSS(n)), "c" ((long) task[n])); \
184 }
185
186 // 页面地址对准。(在内核代码中没有任何地方引用!!)
187 #define PAGE_ALIGN(n) (((n)+0xfff)&0xfffff000)
188
189 // 设置位于地址 addr 处描述符中的各基址字段(基址是 base)，参见列表后说明。
190 // %0 - 地址 addr 偏移 2; %1 - 地址 addr 偏移 4; %2 - 地址 addr 偏移 7; edx - 基址 base。
191 #define set_base(addr, base) \
192 __asm__ ( "movw %%dx, %0|n|t" \           // 基址 base 低 16 位(位 15-0) → [addr+2]。
193          "rorl $16, %%edx|n|t" \         // edx 中基址高 16 位(位 31-16) → dx。
194          "movb %%dl, %1|n|t" \           // 基址高 16 位中的低 8 位(位 23-16) → [addr+4]。
195          "movb %%dh, %2" \               // 基址高 16 位中的高 8 位(位 31-24) → [addr+7]。
196          :: "m" (*(addr+2)), \
197          "m" (*(addr+4)), \
198          "m" (*(addr+7)), \
199          "d" (base) \
200          : "dx")
201
202 // 设置位于地址 addr 处描述符中的段限长字段(段长是 limit)。
203 // %0 - 地址 addr; %1 - 地址 addr 偏移 6 处; edx - 段长值 limit。
204 #define set_limit(addr, limit) \
205 __asm__ ( "movw %%dx, %0|n|t" \           // 段长 limit 低 16 位(位 15-0) → [addr]。
206          "rorl $16, %%edx|n|t" \         // edx 中的段长高 4 位(位 19-16) → dl。
207          "movb %1, %%dh|n|t" \           // 取原[addr+6]字节 → dh，其中高 4 位是些标志。
208          "andb $0xf0, %%dh|n|t" \        // 清 dh 的低 4 位(将存放段长的位 19-16)。
209          "orb %%dh, %%dl|n|t" \          // 将原高 4 位标志和段长的高 4 位(位 19-16)合成 1 字节，
210          "movb %%dl, %1" \               // 并放会[addr+6]处。
211          :: "m" (*(addr)), \
212          "m" (*(addr+6)), \
213          "d" (limit) \
214          : "dx")
215
216 // 设置局部描述符表中 ldt 描述符的基址字段。
217 #define set_base(ldt, base) set_base( ((char *)&(ldt)) , base )
218 // 设置局部描述符表中 ldt 描述符的段长字段。
219 #define set_limit(ldt, limit) set_limit( ((char *)&(ldt)) , (limit-1)>>12 )
220
221 // 从地址 addr 处描述符中取段基址。功能与_set_base() 正好相反。
222 // edx - 存放基址(__base); %1 - 地址 addr 偏移 2; %2 - 地址 addr 偏移 4; %3 - addr 偏移 7。
223 #define get_base(addr) ({\
224 unsigned long __base; \
225 __asm__ ( "movb %3, %%dh|n|t" \           // 取[addr+7]处基址高 16 位的高 8 位(位 31-24) → dh。
226          "movb %2, %%dl|n|t" \           // 取[addr+4]处基址高 16 位的低 8 位(位 23-16) → dl。
227          "shll $16, %%edx|n|t" \         // 基址高 16 位移到 edx 中高 16 位处。
228          "movw %1, %%dx" \               // 取[addr+2]处基址低 16 位(位 15-0) → dx。
229          : "=d" (__base) \               // 从而 edx 中含有 32 位的段基址。
230          : "m" (*(addr+2)), \
231          "m" (*(addr+4)), \
232          "m" (*(addr+7))) ; \
233 __base; })

```

---

```
225 // 取局部描述符表中 ldt 所指段描述符中的基地址。
226 #define get\_base(ldt) \_get\_base( ((char *)&(ldt)) )
227
228 // 取段选择符 segment 的段长值。
229 // %0 - 存放段长值(字节数); %1 - 段选择符 segment。
230 #define get\_limit(segment) ({ \
231 unsigned long __limit; \
232 __asm__("lsll %1,%0|n|tincl %0": "=r" (__limit): "r" (segment); \
233 __limit;})
234 #endif
```

---

## 10.29.3 其它信息

### 10.29.3.1 任务状态段信息

Figure 7-1. 80386 32-Bit Task State Segment

|                                 |    |                                 |   |   |  |  |  |  |  |  |  |    |    |
|---------------------------------|----|---------------------------------|---|---|--|--|--|--|--|--|--|----|----|
| 31                              | 23 | 15                              | 7 | 0 |  |  |  |  |  |  |  |    |    |
| I/O MAP BASE                    |    | 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 |   |   |  |  |  |  |  |  |  | T  | 64 |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 |    | LDT                             |   |   |  |  |  |  |  |  |  | 60 |    |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 |    | GS                              |   |   |  |  |  |  |  |  |  | 5C |    |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 |    | FS                              |   |   |  |  |  |  |  |  |  | 58 |    |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 |    | DS                              |   |   |  |  |  |  |  |  |  | 54 |    |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 |    | SS                              |   |   |  |  |  |  |  |  |  | 50 |    |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 |    | CS                              |   |   |  |  |  |  |  |  |  | 4C |    |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 |    | ES                              |   |   |  |  |  |  |  |  |  | 48 |    |
|                                 |    | EDI                             |   |   |  |  |  |  |  |  |  | 44 |    |
|                                 |    | ESI                             |   |   |  |  |  |  |  |  |  | 40 |    |
|                                 |    | EBP                             |   |   |  |  |  |  |  |  |  | 3C |    |
|                                 |    | ESP                             |   |   |  |  |  |  |  |  |  | 38 |    |
|                                 |    | EBX                             |   |   |  |  |  |  |  |  |  | 34 |    |
|                                 |    | EDX                             |   |   |  |  |  |  |  |  |  | 30 |    |
|                                 |    | ECX                             |   |   |  |  |  |  |  |  |  | 2C |    |
|                                 |    | EAX                             |   |   |  |  |  |  |  |  |  | 28 |    |
|                                 |    | EFLAGS                          |   |   |  |  |  |  |  |  |  | 24 |    |
|                                 |    | INSTRUCTION POINTER <EIP>       |   |   |  |  |  |  |  |  |  | 20 |    |
|                                 |    | CR3 <PDPR>                      |   |   |  |  |  |  |  |  |  | 1C |    |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 |    | SS2                             |   |   |  |  |  |  |  |  |  | 18 |    |
|                                 |    | ESP2                            |   |   |  |  |  |  |  |  |  | 14 |    |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 |    | SS1                             |   |   |  |  |  |  |  |  |  | 10 |    |
|                                 |    | ESP1                            |   |   |  |  |  |  |  |  |  | 0C |    |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 |    | SS0                             |   |   |  |  |  |  |  |  |  | 8  |    |
|                                 |    | ESP0                            |   |   |  |  |  |  |  |  |  | 4  |    |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 |    | BACK LINK TO PREVIOUS TSS       |   |   |  |  |  |  |  |  |  | 0  |    |

## NOTE

0 MEANS INTEL RESERVED. DO NOT DEFINE.

CPU 管理任务需要的所有信息被存储于一个特殊类型的段中，任务状态段 (task state segment - TSS)。图中显示出执行 80386 任务的 TSS 格式。

TSS 中的字段可以分为两类：

1. CPU 在进行任务切换时更新的动态信息集。这些字段有：
  - o 通用寄存器 (EAX, ECX, EDX, EBX, ESP, EBP, ESI, EDI);
  - o 段寄存器 (ES, CS, SS, DS, FS, GS);
  - o 标志寄存器 (EIP);

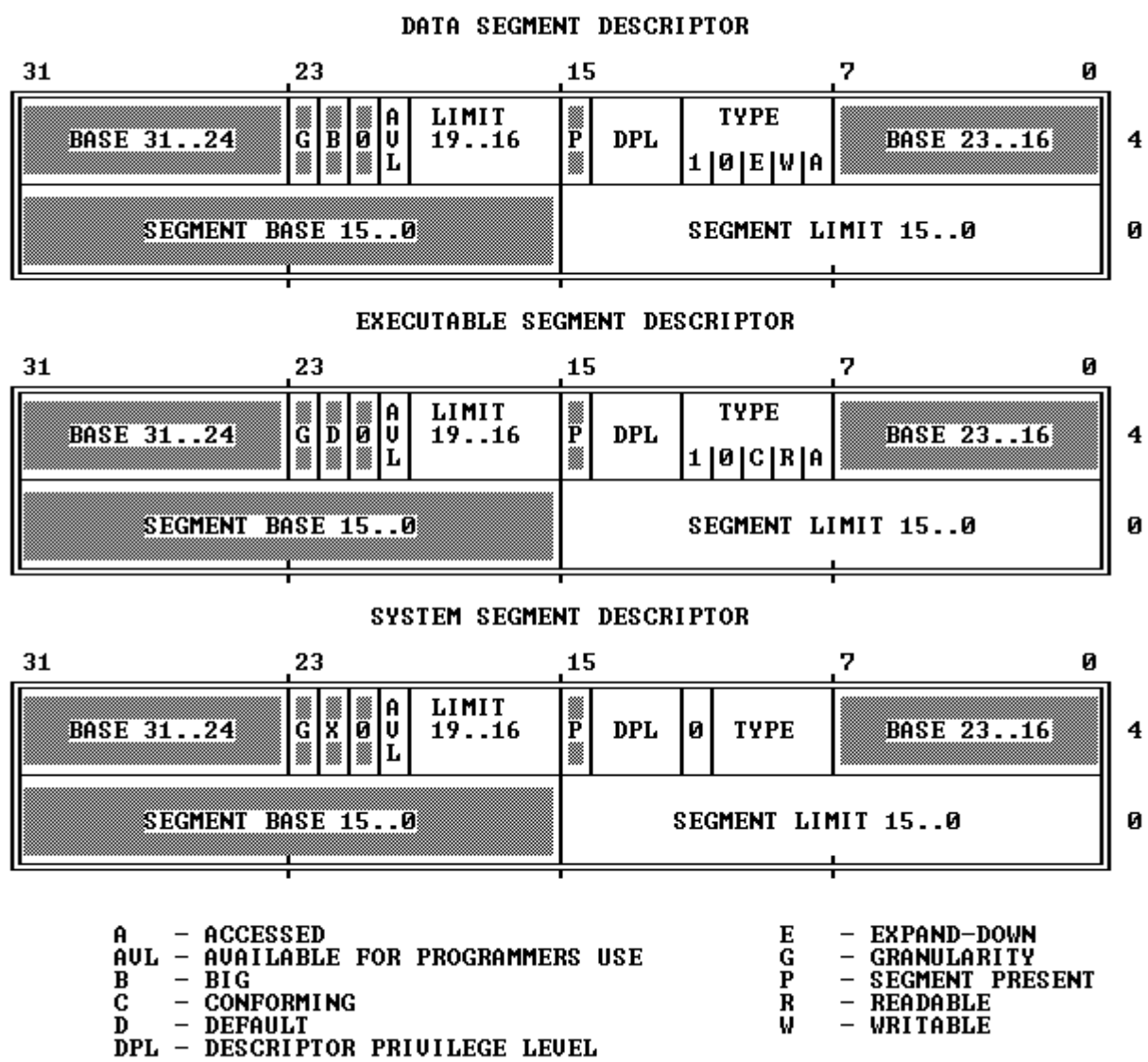
- o 指令指针 (EIP);  
前一个执行任务的 TSS 的选择符 (仅当返回时才更新)。
- 2. CPU 读取但不会更改的静态信息集。这些字段有:
  - o 任务的 LDT 的选择符;
  - o 含有任务页目录基地址的寄存器 (PDBR);
  - o 特权级 0-2 的堆栈指针;
  - o 当任务进行切换时导致 CPU 产生一个调试(debug)异常的 T-比特位 (调试跟踪位);
  - o I/O 比特位图基地址 (其长度上限就是 TSS 的长度上限, 在 TSS 描述符中说明)。

任务状态段可以存放在线形空间的任何地方。与其它各类段相似, 任务状态段也是由描述符来定义的。当前正在执行任务的 TSS 是由任务寄存器 (TR) 来指示的。指令 LTR 和 STR 用来修改和读取任务寄存器中的选择符 (任务寄存器的可见部分)。

I/O 比特位图中的每 1 比特对应 1 个 I/O 端口。比如端口 41 的比特位就是 I/O 位图基地址+5, 位偏移 1 处。在保护模式中, 当遇到 1 个 I/O 指令时 (IN, INS, OUT, OUTS), CPU 首先就会检查当前特权级是否小于标志寄存器的 IOPL, 如果这个条件满足, 就执行该 I/O 操作。如果不满足, 那么 CPU 就会检查 TSS 中的 I/O 比特位图。如果相应比特位是置位的, 就会产生一般保护性异常, 否则就会执行该 I/O 操作。

10.29.3.2 段描述符

Figure 6-1. Protection Fields of Segment Descriptors



10.30 sys.h 文件

10.30.1 功能描述

sys.h 头文件列出了内核中所有系统调用函数的原型，以及系统调用函数指针表。

10.30.2 代码注释

列表 linux/include/linux/sys.h 文件

```
1 extern int sys_setup();           // 系统启动初始化设置函数。 (kernel/blk_drv/hd.c, 71)
2 extern int sys_exit();            // 程序退出。 (kernel/exit.c, 137)
```

|    |  |                    |                             |
|----|--|--------------------|-----------------------------|
| 3  | extern int <a href="#">sys fork()</a> ;    | // 创建进程。           | (kernel/system_call.s, 208) |
| 4  | extern int <a href="#">sys read()</a> ;    | // 读文件。            | (fs/read_write.c, 55)       |
| 5  | extern int <a href="#">sys write()</a> ;   | // 写文件。            | (fs/read_write.c, 83)       |
| 6  | extern int <a href="#">sys open()</a> ;    | // 打开文件。           | (fs/open.c, 138)            |
| 7  | extern int <a href="#">sys close()</a> ;   | // 关闭文件。           | (fs/open.c, 192)            |
| 8  | extern int <a href="#">sys waitpid()</a> ; | // 等待进程终止。         | (kernel/exit.c, 142)        |
| 9  | extern int <a href="#">sys creat()</a> ;   | // 创建文件。           | (fs/open.c, 187)            |
| 10 | extern int <a href="#">sys link()</a> ;    | // 创建一个文件的硬连接。     | (fs/namei.c, 721)           |
| 11 | extern int <a href="#">sys unlink()</a> ;  | // 删除一个文件名(或删除文件)。 | (fs/namei.c, 663)           |
| 12 | extern int <a href="#">sys execve()</a> ;  | // 执行程序。           | (kernel/system_call.s, 200) |
| 13 | extern int <a href="#">sys chdir()</a> ;   | // 更改当前目录。         | (fs/open.c, 75)             |
| 14 | extern int <a href="#">sys time()</a> ;    | // 取当前时间。          | (kernel/sys.c, 102)         |
| 15 | extern int <a href="#">sys mknod()</a> ;   | // 建立块/字符特殊文件。     | (fs/namei.c, 412)           |
| 16 | extern int <a href="#">sys chmod()</a> ;   | // 修改文件属性。         | (fs/open.c, 105)            |
| 17 | extern int <a href="#">sys chown()</a> ;   | // 修改文件宿主和所属组。     | (fs/open.c, 121)            |
| 18 | extern int <a href="#">sys break()</a> ;   | //                 | (-kernel/sys.c, 21)         |
| 19 | extern int <a href="#">sys stat()</a> ;    | // 使用路径名取文件的状态信息。  | (fs/stat.c, 36)             |
| 20 | extern int <a href="#">sys lseek()</a> ;   | // 重新定位读/写文件偏移。    | (fs/read_write.c, 25)       |
| 21 | extern int <a href="#">sys getpid()</a> ;  | // 取进程 id。         | (kernel/sched.c, 348)       |
| 22 | extern int <a href="#">sys mount()</a> ;   | // 安装文件系统。         | (fs/super.c, 200)           |
| 23 | extern int <a href="#">sys umount()</a> ;  | // 卸载文件系统。         | (fs/super.c, 167)           |
| 24 | extern int <a href="#">sys setuid()</a> ;  | // 设置进程用户 id。      | (kernel/sys.c, 143)         |
| 25 | extern int <a href="#">sys getuid()</a> ;  | // 取进程用户 id。       | (kernel/sched.c, 358)       |
| 26 | extern int <a href="#">sys stime()</a> ;   | // 设置系统时间日期。       | (-kernel/sys.c, 148)        |
| 27 | extern int <a href="#">sys ptrace()</a> ;  | // 程序调试。           | (-kernel/sys.c, 26)         |
| 28 | extern int <a href="#">sys alarm()</a> ;   | // 设置报警。           | (kernel/sched.c, 338)       |
| 29 | extern int <a href="#">sys fstat()</a> ;   | // 使用文件句柄取文件的状态信息。 | (fs/stat.c, 47)             |
| 30 | extern int <a href="#">sys pause()</a> ;   | // 暂停进程运行。         | (kernel/sched.c, 144)       |
| 31 | extern int <a href="#">sys utime()</a> ;   | // 改变文件的访问和修改时间。   | (fs/open.c, 24)             |
| 32 | extern int <a href="#">sys stty()</a> ;    | // 修改终端行设置。        | (-kernel/sys.c, 31)         |
| 33 | extern int <a href="#">sys gtty()</a> ;    | // 取终端行设置信息。       | (-kernel/sys.c, 36)         |
| 34 | extern int <a href="#">sys access()</a> ;  | // 检查用户对一个文件的访问权限。 | (fs/open.c, 47)             |
| 35 | extern int <a href="#">sys nice()</a> ;    | // 设置进程执行优先权。      | (kernel/sched.c, 378)       |
| 36 | extern int <a href="#">sys ftime()</a> ;   | // 取日期和时间。         | (-kernel/sys.c, 16)         |
| 37 | extern int <a href="#">sys sync()</a> ;    | // 同步高速缓冲与设备中数据。   | (fs/buffer.c, 44)           |
| 38 | extern int <a href="#">sys kill()</a> ;    | // 终止一个进程。         | (kernel/exit.c, 60)         |
| 39 | extern int <a href="#">sys rename()</a> ;  | // 更改文件名。          | (-kernel/sys.c, 41)         |
| 40 | extern int <a href="#">sys mkdir()</a> ;   | // 创建目录。           | (fs/namei.c, 463)           |
| 41 | extern int <a href="#">sys rmdir()</a> ;   | // 删除目录。           | (fs/namei.c, 587)           |
| 42 | extern int <a href="#">sys dup()</a> ;     | // 复制文件句柄。         | (fs/fcntl.c, 42)            |
| 43 | extern int <a href="#">sys pipe()</a> ;    | // 创建管道。           | (fs/pipe.c, 71)             |
| 44 | extern int <a href="#">sys times()</a> ;   | // 取运行时间。          | (kernel/sys.c, 156)         |
| 45 | extern int <a href="#">sys prof()</a> ;    | // 程序执行时间区域。       | (-kernel/sys.c, 46)         |
| 46 | extern int <a href="#">sys brk()</a> ;     | // 修改数据段长度。        | (kernel/sys.c, 168)         |
| 47 | extern int <a href="#">sys setgid()</a> ;  | // 设置进程组 id。       | (kernel/sys.c, 72)          |
| 48 | extern int <a href="#">sys getgid()</a> ;  | // 取进程组 id。        | (kernel/sched.c, 368)       |
| 49 | extern int <a href="#">sys signal()</a> ;  | // 信号处理。           | (kernel/signal.c, 48)       |
| 50 | extern int <a href="#">sys geteuid()</a> ; | // 取进程有效用户 id。     | (kernel/sched.c, 363)       |
| 51 | extern int <a href="#">sys getegid()</a> ; | // 取进程有效组 id。      | (kernel/sched.c, 373)       |
| 52 | extern int <a href="#">sys acct()</a> ;    | // 进程记帐。           | (-kernel/sys.c, 77)         |
| 53 | extern int <a href="#">sys phys()</a> ;    | //                 | (-kernel/sys.c, 82)         |
| 54 | extern int <a href="#">sys lock()</a> ;    | //                 | (-kernel/sys.c, 87)         |
| 55 | extern int <a href="#">sys ioctl()</a> ;   | // 设备控制。           | (fs/ioctl.c, 30)            |

```

56 extern int sys\_fcntl();           // 文件句柄操作。           (fs/fcntl.c, 47)
57 extern int sys\_mpx();             //                       (-kernel/sys.c, 92)
58 extern int sys\_setpgid();          // 设置进程组 id。           (kernel/sys.c, 181)
59 extern int sys\_ulimit();           //                       (-kernel/sys.c, 97)
60 extern int sys\_uname();            // 显示系统信息。           (kernel/sys.c, 216)
61 extern int sys\_umask();            // 取默认文件创建属性码。   (kernel/sys.c, 230)
62 extern int sys\_chroot();           // 改变根系统。             (fs/open.c, 90)
63 extern int sys\_ustat();            // 取文件系统信息。         (fs/open.c, 19)
64 extern int sys\_dup2();             // 复制文件句柄。           (fs/fcntl.c, 36)
65 extern int sys\_getppid();          // 取父进程 id。            (kernel/sched.c, 353)
66 extern int sys\_getpgrp();          // 取进程组 id, 等于 getpgid(0)。 (kernel/sys.c, 201)
67 extern int sys\_setsid();          // 在新会话中运行程序。     (kernel/sys.c, 206)
68 extern int sys\_sigaction();       // 改变信号处理过程。       (kernel/signal.c, 63)
69 extern int sys\_sgetmask();         // 取信号屏蔽码。           (kernel/signal.c, 15)
70 extern int sys\_ssetmask();         // 设置信号屏蔽码。         (kernel/signal.c, 20)
71 extern int sys\_setreuid();         // 设置真实与/或有效用户 id。 (kernel/sys.c, 118)
72 extern int sys\_setregid();         // 设置真实与/或有效组 id。  (kernel/sys.c, 51)
73
// 系统调用函数指针表。用于系统调用中断处理程序(int 0x80), 作为跳转表。
74 fn_ptr sys\_call\_table[] = { sys\_setup, sys\_exit, sys\_fork, sys\_read,
75 sys\_write, sys\_open, sys\_close, sys\_waitpid, sys\_creat, sys\_link,
76 sys\_unlink, sys\_execve, sys\_chdir, sys\_time, sys\_mknod, sys\_chmod,
77 sys\_chown, sys\_break, sys\_stat, sys\_lseek, sys\_getpid, sys\_mount,
78 sys\_umount, sys\_setuid, sys\_getuid, sys\_stime, sys\_ptrace, sys\_alarm,
79 sys\_fstat, sys\_pause, sys\_utime, sys\_stty, sys\_gtty, sys\_access,
80 sys\_nice, sys\_ftime, sys\_sync, sys\_kill, sys\_rename, sys\_mkdir,
81 sys\_rmdir, sys\_dup, sys\_pipe, sys\_times, sys\_prof, sys\_brk, sys\_setgid,
82 sys\_getgid, sys\_signal, sys\_geteuid, sys\_getegid, sys\_acct, sys\_phys,
83 sys\_lock, sys\_ioctl, sys\_fcntl, sys\_mpx, sys\_setpgid, sys\_ulimit,
84 sys\_uname, sys\_umask, sys\_chroot, sys\_ustat, sys\_dup2, sys\_getppid,
85 sys\_getpgrp, sys\_setsid, sys\_sigaction, sys\_sgetmask, sys\_ssetmask,
86 sys\_setreuid, sys\_setregid };
87

```

## 10.31 tty.h 文件

### 10.31.1 功能描述

### 10.31.2 代码注释

列表 linux/include/linux/tty.h 文件

```

1 /*
2  * 'tty.h' defines some structures used by tty_io.c and some defines.
3  *
4  * NOTE! Don't touch this without checking that nothing in rs_io.s or
5  * con_io.s breaks. Some constants are hardwired into the system (mainly
6  * offsets into 'tty_queue'
7  */

```

```

8  /*
   * 'tty.h' 中定义了 tty_io.c 程序使用的某些结构和其它一些定义。
   *
   * 注意！在修改这里的定义时，一定要检查 rs_io.s 或 con_io.s 程序中不会出现问题。
   * 在系统中有些常量是直接写在程序中的（主要是一些 tty_queue 中的偏移值）。
   */
9  #ifndef TTY_H
10 #define TTY_H
11
12 #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。
13
14 #define TTY_BUF_SIZE 1024      // tty 缓冲区大小。
15
16 // tty 等待队列数据结构。
17 struct tty_queue {
18     unsigned long data;          // 等待队列缓冲区中当前数据指针（字符数[??]）。
19     unsigned long head;          // 缓冲区中数据头指针。
20     unsigned long tail;          // 缓冲区中数据尾指针。
21     struct task_struct * proc_list; // 等待进程列表。
22     char buf[TTY_BUF_SIZE];      // 队列的缓冲区。
23 };
24
25 // 以下定义了 tty 等待队列中缓冲区操作宏函数。（tail 在前，head 在后）。
26 // a 缓冲区指针前移 1 字节，并循环。
27 #define INC(a) ((a) = ((a)+1) & (TTY_BUF_SIZE-1))
28 // a 缓冲区指针后退 1 字节，并循环。
29 #define DEC(a) ((a) = ((a)-1) & (TTY_BUF_SIZE-1))
30 // 清空指定队列的缓冲区。
31 #define EMPTY(a) ((a).head == (a).tail)
32 // 缓冲区还可存放字符的长度（空闲区长度）。
33 #define LEFT(a) (((a).tail-(a).head-1)&(TTY_BUF_SIZE-1))
34 // 缓冲区中最后一个位置。
35 #define LAST(a) ((a).buf[(TTY_BUF_SIZE-1)&((a).head-1)])
36 // 缓冲区满（如果为 1 的话）。
37 #define FULL(a) (!LEFT(a))
38 // 缓冲区中已存放字符的长度。
39 #define CHARS(a) (((a).head-(a).tail)&(TTY_BUF_SIZE-1))
40 // 从 queue 队列项缓冲区中取一字符（从 tail 处，并且 tail+=1）。
41 #define GETCH(queue, c) \
42 (void) ({c=(queue).buf[(queue).tail];INC((queue).tail);})
43 // 往 queue 队列项缓冲区中放置一字符（在 head 处，并且 head+=1）。
44 #define PUTCH(c, queue) \
45 (void) ({(queue).buf[(queue).head]=(c);INC((queue).head);})
46
47 // 判断终端键盘字符类型。
48 #define INTR_CHAR(tty) ((tty)->termios.c_cc[VINTR]) // 中断符。
49 #define QUIT_CHAR(tty) ((tty)->termios.c_cc[VQUIT]) // 退出符。
50 #define ERASE_CHAR(tty) ((tty)->termios.c_cc[VERASE]) // 删除符。
51 #define KILL_CHAR(tty) ((tty)->termios.c_cc[VKILL]) // 终止符。
52 #define EOF_CHAR(tty) ((tty)->termios.c_cc[VEOF]) // 文件结束符。
53 #define START_CHAR(tty) ((tty)->termios.c_cc[VSTART]) // 开始符。

```








```

42 #define STOP_CHAR(tty) ((tty)->termios.c_cc[VSTOP]) // 结束符。
43 #define SUSPEND_CHAR(tty) ((tty)->termios.c_cc[VSUSP]) // 挂起符。
44
// tty 数据结构。
45 struct tty_struct {
46     struct termios termios; // 终端 io 属性和控制字符数据结构。
47     int pgrp; // 所属进程组。
48     int stopped; // 停止标志。
49     void (*write)(struct tty_struct * tty); // tty 写函数指针。
50     struct tty_queue read_q; // tty 读队列。
51     struct tty_queue write_q; // tty 写队列。
52     struct tty_queue secondary; // tty 辅助队列(存放规范模式字符序列),
53     }; // 可称为规范(熟)模式队列。
54
55 extern struct tty_struct tty_table[]; // tty 结构数组。
56
57 /*      intr=^C      quit=^/      erase=del      kill=^U
58      eof=^D      vtime=\0      vmin=\1      sxtc=\0
59      start=^Q      stop=^S      susp=^Z      eol=\0
60      reprint=^R      discard=^U      werase=^W      lnext=^V
61      eol2=\0
62 */
/* 中断 intr=^C      退出 quit=^|      删除 erase=del      终止 kill=^U
* 文件结束 eof=^D      vtime=\0      vmin=\1      sxtc=\0
* 开始 start=^Q      停止 stop=^S      挂起 susp=^Z      行结束 eol=\0
* 重显 reprint=^R      丢弃 discard=^U      werase=^W      lnext=^V
* 行结束 eol2=\0
*/
// 控制字符对应的 ASCII 码值。[8 进制]
63 #define INIT_C_CC "\003\034\177\025\004\0\1\0\021\023\032\0\022\017\027\026\0"
64
65 void rs_init(void); // 异步串行通信初始化。(kernel/chr_drv/serial.c, 37)
66 void con_init(void); // 控制终端初始化。(kernel/chr_drv/console.c, 617)
67 void tty_init(void); // tty 初始化。(kernel/chr_drv/tty_io.c, 105)
68
69 int tty_read(unsigned c, char * buf, int n); // (kernel/chr_drv/tty_io.c, 230)
70 int tty_write(unsigned c, char * buf, int n); // (kernel/chr_drv/tty_io.c, 290)
71
72 void rs_write(struct tty_struct * tty); // (kernel/chr_drv/serial.c, 53)
73 void con_write(struct tty_struct * tty); // (kernel/chr_drv/console.c, 445)
74
75 void copy_to_cooked(struct tty_struct * tty); // (kernel/chr_drv/tty_io.c, 145)
76
77 #endif
78

```

## 10.32 include/sys/目录中的文件

列表 linux/include/sys/目录下的文件

| Name  | Size       | Last modified (GMT) | Description |
|---|------------|---------------------|-------------|
|  <a href="#">stat.h</a>    | 1304 bytes | 1991-09-17 15:02:48 | m           |
|  <a href="#">times.h</a>   | 200 bytes  | 1991-09-17 15:03:06 | m           |
|  <a href="#">types.h</a>   | 805 bytes  | 1991-09-17 15:02:55 | m           |
|  <a href="#">utsname.h</a> | 234 bytes  | 1991-09-17 15:03:23 | m           |
|  <a href="#">wait.h</a>    | 560 bytes  | 1991-09-17 15:06:07 | m           |

## 10.33 stat.h 文件

### 10.33.1 功能描述

该头文件说明了函数 `stat()` 返回的数据及其结构类型，以及一些属性操作测试宏、函数原型。

### 10.33.2 代码注释

列表 linux/include/sys/stat.h 文件

```

1 #ifndef SYS\_STAT\_H
2 #define SYS\_STAT\_H
3
4 #include <sys/types.h>
5
6 struct stat {
7     dev\_t    st_dev;    // 含有文件的设备号。
8     ino\_t    st_ino;    // 文件 i 节点号。
9     umode\_t  st_mode;    // 文件属性（见下面）。
10    nlink\_t  st_nlink;   // 指定文件的连接数。
11    uid\_t    st_uid;    // 文件的用户（标识）号。
12    gid\_t    st_gid;    // 文件的组号。
13    dev\_t    st_rdev;    // 设备号（如果文件是特殊的字符文件或块文件）。
14    off\_t    st_size;    // 文件大小（字节数）（如果文件是常规文件）。
15    time\_t    st_atime;   // 上次（最后）访问时间。

```

---

```

16     time_t  st_mtime;    // 最后修改时间。
17     time_t  st_ctime;    // 最后节点修改时间。
18 };
19
    // 以下这些是 st_mode 值的符号名称。
    // 文件类型：
20 #define S_IFMT 00170000    // 文件类型。
21 #define S_IFREG 0100000    // 常规文件。
22 #define S_IFBLK 0060000    // 块特殊（设备）文件。
23 #define S_IFDIR 0040000    // 目录文件。
24 #define S_IFCHR 0020000    // 字符设备文件。
25 #define S_IFIFO 0010000    // FIFO 特殊文件。
    // 文件属性位：
26 #define S_ISUID 0004000    // 执行时设置用户 ID (set-user-ID) 。
27 #define S_ISGID 0002000    // 执行时设置组 ID。
28 #define S_ISVTX 0001000    // 对于目录，受限删除标志。
29
30 #define S_ISREG(m)        (((m) & S_IFMT) == S_IFREG)    // 测试是否常规文件。
31 #define S_ISDIR(m)        (((m) & S_IFMT) == S_IFDIR)    // 是否目录文件。
32 #define S_ISCHR(m)        (((m) & S_IFMT) == S_IFCHR)    // 是否字符设备文件。
33 #define S_ISBLK(m)        (((m) & S_IFMT) == S_IFBLK)    // 是否块设备文件。
34 #define S_ISFIFO(m)       (((m) & S_IFMT) == S_IFIFO)    // 是否 FIFO 特殊文件。
35
36 #define S_IRWXU 00700      // 宿主可以读、写、执行/搜索。
37 #define S_IRUSR 00400      // 宿主读许可。
38 #define S_IWUSR 00200      // 宿主写许可。
39 #define S_IXUSR 00100      // 宿主执行/搜索许可。
40
41 #define S_IRWXG 00070      // 组成员可以读、写、执行/搜索。
42 #define S_IRGRP 00040      // 组成员读许可。
43 #define S_IWGRP 00020      // 组成员写许可。
44 #define S_IXGRP 00010      // 组成员执行/搜索许可。
45
46 #define S_IRWXO 00007      // 其他人读、写、执行/搜索许可。
47 #define S_IROTH 00004      // 其他人读许可。
48 #define S_IWOTH 00002      // 其他人写许可。
49 #define S_IXOTH 00001      // 其他人执行/搜索许可。
50
51 extern int chmod(const char *_path, mode_t mode);    // 修改文件属性。
52 extern int fstat(int fildes, struct stat *stat_buf); // 取指定文件句柄的文件状态信息。
53 extern int mkdir(const char *_path, mode_t mode);    // 创建目录。
54 extern int mkfifo(const char *_path, mode_t mode);    // 创建管道文件。
55 extern int stat(const char *filename, struct stat *stat_buf); // 取指定文件名的文件状态信息。
56 extern mode_t umask(mode_t mask);    // 设置属性屏蔽码。
57
58 #endif
59

```

---

## 10.34 times.h 文件

### 10.34.1 功能描述

该头文件中主要定义了文件访问与修改时间结构 `tms`。它将由 `times()` 函数返回。其中 `time_t` 是在 `sys/types.h` 中定义的。还定义了一个函数原型 `times()`。

### 10.34.2 代码注释

列表 linux/include/sys/times.h 文件

---

```

1 #ifndef TIMES\_H
2 #define TIMES\_H
3
4 #include <sys/types.h>    // 类型头文件。定义了基本的系统数据类型。
5
6 struct tms {
7     time\_t tms_ftime; // 用户使用的 CPU 时间。
8     time\_t tms_stime; // 系统（内核）CPU 时间。
9     time\_t tms_cutime; // 已终止的子进程使用的用户 CPU 时间。
10    time\_t tms_cstime; // 已终止的子进程使用的系统 CPU 时间。
11 };
12
13 extern time\_t times(struct tms * tp);
14
15 #endif
16

```

---

## 10.35 types.h 文件

### 10.35.1 功能描述

`types.h` 头文件中定义了基本的数据类型。所有的类型定义为适当的数学类型长度。另外，`size_t` 是无符号整数类型，`off_t` 是扩展的符号整数类型，`pid_t` 是符号整数类型。

### 10.35.2 代码注释

列表 linux/include/sys/types.h 文件

---

```

1 #ifndef SYS\_TYPES\_H
2 #define SYS\_TYPES\_H
3
4 #ifndef SIZE\_T
5 #define SIZE\_T
6 typedef unsigned int size\_t;    // 用于对象的大小（长度）。
7 #endif
8
9 #ifndef TIME\_T

```

---

```

10 #define TIME_T
11 typedef long time_t;           // 用于时间（以秒计）。
12 #endif
13
14 #ifndef PTRDIFF_T
15 #define PTRDIFF_T
16 typedef long ptrdiff_t;
17 #endif
18
19 #ifndef NULL
20 #define NULL ((void *) 0)
21 #endif
22
23 typedef int pid_t;             // 用于进程号和进程组号。
24 typedef unsigned short uid_t;  // 用于用户号（用户标识号）。
25 typedef unsigned char gid_t;  // 用于组号。
26 typedef unsigned short dev_t; // 用于设备号。
27 typedef unsigned short ino_t; // 用于文件序列号。
28 typedef unsigned short mode_t; // 用于某些文件属性。
29 typedef unsigned short umode_t; //
30 typedef unsigned char nlink_t; // 用于连接计数。
31 typedef int daddr_t;
32 typedef long off_t;           // 用于文件长度（大小）。
33 typedef unsigned char u_char; // 无符号字符类型。
34 typedef unsigned short ushort; // 无符号短整数类型。
35
36 typedef struct { int quot,rem; } div_t; // 用于 DIV 操作。
37 typedef struct { long quot,rem; } ldiv_t; // 用于长 DIV 操作。
38
39 struct ustat {
40     daddr_t f_tfree;
41     ino_t f_tinode;
42     char f_fname[6];
43     char f_fpack[6];
44 };
45
46 #endif
47

```

## 10.36 utsname.h 文件

### 10.36.1 功能描述

utsname.h 是系统名称结构头文件。其中定义了结构 utsname 以及函数原型 uname()。POSIX 要求字符数组长度应该是不指定的，但是其中存储的数据需以 null 终止。因此该版内核的 utsname 结构定义不符合要求（数组长度都被定义为 9）。

### 10.36.2 代码注释

列表 linux/include/sys/utsname.h 文件

```

1 #ifndef \_SYS\_UTSNAME\_H
2 #define \_SYS\_UTSNAME\_H
3
4 #include <sys/types.h>    // 类型头文件。定义了基本的系统数据类型。
5
6 struct utsname {
7     char sysname[9];    // 本版本操作系统的名称。
8     char nodename[9];    // 与实现相关的网络中节点名称。
9     char release[9];    // 本实现的当前发行级别。
10    char version[9];    // 本次发行的版本级别。
11    char machine[9];    // 系统运行的硬件类型名称。
12 };
13
14 extern int uname(struct utsname * utsbuf);
15
16 #endif
17

```

## 10.37 wait.h 文件

### 10.37.1 功能描述

该头文件描述了进程等待时信息。包括一些符号常数和 wait()、waitpid() 函数原型申明。

### 10.37.2 代码注释

列表 linux/include/sys/wait.h 文件

```

1 #ifndef \_SYS\_WAIT\_H
2 #define \_SYS\_WAIT\_H
3
4 #include <sys/types.h>
5
6 #define LOW(v)          ( (v) & 0377)          // 取低字节（8 进制表示）。
7 #define HIGH(v)        ( ((v) >> 8) & 0377)    // 取高字节。
8
9 /* options for waitpid, WUNTRACED not supported */
10 /* waitpid 的选项，其中 WUNTRACED 未被支持 */
11 #define WNOHANG         1          // 如果没有状态也不要挂起，并立刻返回。
12 #define WUNTRACED       2          // 报告停止执行的子进程状态。
13
14 #define WIFEXITED(s)     (!((s)&0xFF)          // 如果子进程正常退出，则为真。
15 #define WIFSTOPPED(s)   (((s)&0xFF)==0x7F)    // 如果子进程正停止着，则为 true。
16 #define WEXITSTATUS(s)  (((s)>>8)&0xFF)        // 返回退出状态。
17 #define WTERMSIG(s)      ((s)&0x7F)           // 返回导致进程终止的信号值（信号量）。
18 #define WSTOPSIG(s)      (((s)>>8)&0xFF)        // 返回导致进程停止的信号值。
19 #define WIFSIGNALED(s)   (((unsigned int)(s)-1 & 0xFFFF) < 0xFF) // 如果由于未捕捉到信号
                                                                    // 而导致子进程退出则为真。
20
21 // wait() 和 waitpid() 函数允许进程获取与其子进程之一的状态信息。各种选项允许获取已经终止或

```

---

```
// 停止的子进程状态信息。如果存在两个或两个以上子进程的状态信息，则报告的顺序是不指定的。  
// wait() 将挂起当前进程，直到其子进程之一退出（终止），或者收到要求终止该进程的信号，  
// 或者是需要调用一个信号句柄（信号处理程序）。  
// waitpid() 挂起当前进程，直到 pid 指定的子进程退出（终止）或者收到要求终止该进程的信号，  
// 或者是需要调用一个信号句柄（信号处理程序）。  
// 如果 pid= -1, options=0, 则 waitpid() 的作用与 wait() 函数一样。否则其行为将随 pid 和 options  
// 参数的不同而不同。（参见 kernel/exit.c, 142）  
20 pid_t wait(int *stat_loc);  
21 pid_t waitpid(pid_t pid, int *stat_loc, int options);  
22  
23 #endif  
24
```

---


















## 第 11 章 库文件(lib)

### 11.1 概述

---

| Name   | Size       | Last modified (GMT) | Description |
|--|------------|---------------------|-------------|
|  <a href="#">Makefile</a>   | 2602 bytes | 1991-12-02 03:16:05 |             |
|  <a href="#">_exit.c</a>    | 198 bytes  | 1991-10-02 14:16:29 |             |
|  <a href="#">close.c</a>    | 131 bytes  | 1991-10-02 14:16:29 |             |
|  <a href="#">ctype.c</a>    | 1202 bytes | 1991-10-02 14:16:29 |             |
|  <a href="#">dup.c</a>      | 127 bytes  | 1991-10-02 14:16:29 |             |
|  <a href="#">errno.c</a>  | 73 bytes   | 1991-10-02 14:16:29 |             |
|  <a href="#">execve.c</a> | 170 bytes  | 1991-10-02 14:16:29 |             |
|  <a href="#">malloc.c</a> | 7469 bytes | 1991-12-02 03:15:20 |             |
|  <a href="#">open.c</a>   | 389 bytes  | 1991-10-02 14:16:29 |             |
|  <a href="#">setsid.c</a> | 128 bytes  | 1991-10-02 14:16:29 |             |
|  <a href="#">string.c</a> | 177 bytes  | 1991-10-02 14:16:29 |             |
|  <a href="#">wait.c</a>   | 253 bytes  | 1991-10-02 14:16:29 |             |
|  <a href="#">write.c</a>  | 160 bytes  | 1991-10-02 14:16:29 |             |

---

## 11.2 Makefile 文件

### 11.2.1 功能描述

### 11.2.2 代码注释

列表 linux/lib/Makefile 文件

```

1 #
2 # Makefile for some libs needed in the kernel.
3 #
4 # Note! Dependencies are done automagically by 'make dep', which also
5 # removes any old dependencies. DON'T put your own dependencies here
6 # unless it's something special (ie not a .c file).
7 #
8 # 内核需要用到的 libs 库文件程序的 Makefile。
9 #
10 # 注意！依赖关系是由'make dep'自动进行的，它也会自动去除原来的依赖信息。不要把你自己的
11 # 依赖关系信息放在这里，除非是特别文件的（也即不是一个.c文件的信息）。
12
13 AR      =gar      # GNU 的二进制文件处理程序，用于创建、修改以及从归档文件中抽取文件。
14 AS      =gas      # GNU 的汇编程序。
15 LD      =gld      # GNU 的连接程序。
16 LDFLAGS =-s -x    # 连接程序所有的参数，-s 输出文件中省略所有符号信息。-x 删除所有局部符号。
17 CC      =gcc      # GNU C 语言编译器。
18 CFLAGS  =-Wall -O -fstrength-reduce -fomit-frame-pointer -fcombine-regs \
19         -finline-functions -mstring-insns -nostdinc -I../include
20 # C 编译程序选项。-Wall 显示所有的警告信息；-O 优化选项，优化代码长度和执行时间；
21 # -fstrength-reduce 优化循环执行代码，排除重复变量；-fomit-frame-pointer 省略保存不必要
22 # 的框架指针；-fcombine-regs 合并寄存器，减少寄存器类的使用；-finline-functions 将所有简
23 # 单短小的函数代码嵌入调用程序中；-mstring-insns Linux 自己添加的优化选项，以后不再使用；
24 # -nostdinc -I../include 不使用默认路径中的包含文件，而使用这里指定目录中的(../include)。
25
26 CPP      =gcc -E -nostdinc -I../include
27 # C 前处理选项。-E 只运行 C 前处理，对所有指定的 C 程序进行预处理并将处理结果输出到标准输
28 # 出设备或指定的输出文件中；-nostdinc -I../include 同前。
29
30 # 下面的规则指示 make 利用下面的命令将所有的.c 文件编译生成.s 汇编程序。该规则的命令
31 # 指使 gcc 采用 CFLAGS 所指定的选项对 C 代码编译后不进行汇编就停止（-S），从而产生与
32 # 输入的各个 C 文件对应的汇编代码文件。默认情况下所产生的汇编程序文件名是原 C 文件名
33 # 去掉.c 而加上.s 后缀。-o 表示其后是输出文件的名称。其中$.s（或$@）是自动目标变量，
34 # $<代表第一个先决条件，这里即是符合条件*.c 的文件。
35
36 .c.s:
37     $(CC) $(CFLAGS) \
38     -S -o $.s $<
39 # 下面规则表示将所有.s 汇编程序文件编译成.o 目标文件。22 行是实现该操作的具体命令。
40
41 .s.o:
42     $(AS) -c -o $.o $<
43
44 .c.o:          # 类似上面，*.c 文件→*.o 目标文件。不进行连接。
45     $(CC) $(CFLAGS) \
46     -c -o $.o $<

```

```

26 # 下面定义目标文件变量 OBJs。
27 OBJs = ctype.o _exit.o open.o close.o errno.o write.o dup.o setsid.o \
28       execve.o wait.o string.o malloc.o
29
30 # 在有了先决条件 OBJs 后使用下面的命令连接成目标 lib.a 库文件。
31 lib.a: $(OBJs)
32     $(AR) rcs lib.a $(OBJs)
33     sync
34
35 # 下面的规则用于清理工作。当执行'make clean'时，就会执行下面的命令，去除所有编译
36 # 连接生成的文件。'rm'是文件删除命令，选项-f 含义是忽略不存在的文件，并且不显示删除信息。
37 clean:
38     rm -f core *.o *.a tmp_make
39     for i in *.c;do rm -f `basename $$i .c`.s;done
40
41 # 下面得目标或规则用于检查各文件之间的依赖关系。方法如下：
42 # 使用字符串编辑程序 sed 对 Makefile 文件（即是本文件）进行处理，输出为删除 Makefile
43 # 文件中'### Dependencies'行后面的所有行（下面从 45 开始的行），并生成 tmp_make
44 # 临时文件（39 行的作用）。然后对 kernel/blk_drv/目录下的每个 C 文件执行 gcc 预处理操作。
45 # -M 标志告诉预处理程序输出描述每个目标文件相关性的规则，并且这些规则符合 make 语法。
46 # 对于每一个源文件，预处理程序输出一个 make 规则，其结果形式是相应源程序文件的目标
47 # 文件名加上其依赖关系--该源文件中包含的所有头文件列表。把预处理结果都添加到临时
48 # 文件 tmp_make 中，然后将该临时文件复制成新的 Makefile 文件。
49 dep:
50     sed '/\#\#\# Dependencies/q' < Makefile > tmp_make
51     (for i in *.c;do echo -n `echo $$i | sed 's,\.c,\.s,`" "; \
52         $(CPP) -M $$i;done) >> tmp_make
53     cp tmp_make Makefile
54
55 ### Dependencies:
56 _exit.s _exit.o : _exit.c ../include/unistd.h ../include/sys/stat.h \
57     ../include/sys/types.h ../include/sys/times.h ../include/sys/utsname.h \
58     ../include/utime.h
59 close.s close.o : close.c ../include/unistd.h ../include/sys/stat.h \
60     ../include/sys/types.h ../include/sys/times.h ../include/sys/utsname.h \
61     ../include/utime.h
62 ctype.s ctype.o : ctype.c ../include/ctype.h
63 dup.s dup.o : dup.c ../include/unistd.h ../include/sys/stat.h \
64     ../include/sys/types.h ../include/sys/times.h ../include/sys/utsname.h \
65     ../include/utime.h
66 errno.s errno.o : errno.c
67 execve.s execve.o : execve.c ../include/unistd.h ../include/sys/stat.h \
68     ../include/sys/types.h ../include/sys/times.h ../include/sys/utsname.h \
69     ../include/utime.h
70 malloc.s malloc.o : malloc.c ../include/linux/kernel.h ../include/linux/mm.h \
71     ../include/asm/system.h
72 open.s open.o : open.c ../include/unistd.h ../include/sys/stat.h \
73     ../include/sys/types.h ../include/sys/times.h ../include/sys/utsname.h \
74     ../include/utime.h ../include/stdarg.h
75 setsid.s setsid.o : setsid.c ../include/unistd.h ../include/sys/stat.h \
76     ../include/sys/types.h ../include/sys/times.h ../include/sys/utsname.h \
77     ../include/utime.h

```

---

```

67 string.s string.o : string.c ../include/string.h
68 wait.s wait.o : wait.c ../include/unistd.h ../include/sys/stat.h \
69 ../include/sys/types.h ../include/sys/times.h ../include/sys/utsname.h \
70 ../include/utime.h ../include/sys/wait.h
71 write.s write.o : write.c ../include/unistd.h ../include/sys/stat.h \
72 ../include/sys/types.h ../include/sys/times.h ../include/sys/utsname.h \
73 ../include/utime.h

```

---

## 11.3 \_exit.c 程序

### 11.3.1 功能描述

### 11.3.2 代码注释

列表 linux/lib/\_exit.c 程序

---

```

1  /*
2   * linux/lib/_exit.c
3   *
4   * (C) 1991 Linus Torvalds
5   */
6
7  #define LIBRARY          // 定义一个符号常量，见下行说明。
8  #include <unistd.h>       // Linux 标准头文件。定义了各种符号常数和类型，并声明了各种函数。
                          // 如定义了__LIBRARY__，则还包括系统调用号和内嵌汇编_syscall0()等。
9
10  // 内核使用的程序(退出)终止函数。
11  // 直接调用系统中断 int 0x80，功能号__NR_exit。
12  // 参数: exit_code - 退出码。
13  volatile void exit(int exit_code)
14  {
15      // %0 - eax(系统调用号__NR_exit); %1 - ebx(退出码 exit_code)。
16      __asm__ ("int $0x80"::"a" (__NR_exit), "b" (exit_code));
17  }

```

---

### 11.3.3 相关信息

参见 include/unistd.h 中的说明。

## 11.4 close.c 程序

### 11.4.1 功能描述

## 11.4.2 代码注释

列表 linux/lib/close.c 程序

---

```

1  /*
2  *  linux/lib/close.c
3  *
4  *  (C) 1991  Linus Torvalds
5  */
6
7  #define  LIBRARY
8  #include <unistd.h>          // Linux 标准头文件。定义了各种符号常数和类型，并声明了各种函数。
                                // 如定义了__LIBRARY__，则还包括系统调用号和内嵌汇编_syscall0()等。
9
10 // 关闭文件函数。
11 // 下面该调用宏函数对应：int close(int fd)。直接调用了系统中断 int 0x80，参数是__NR_close。
12 // 其中 fd 是文件描述符。
13 syscall1(int, close, int, fd)
14

```

---

## 11.5 ctype.c 程序

### 11.5.1 功能描述

### 11.5.2 代码注释

列表 linux/lib/ctype.c 程序

---

```

1  /*
2  *  linux/lib/ctype.c
3  *
4  *  (C) 1991  Linus Torvalds
5  */
6
7  #include <ctype.h>          // 字符类型头文件。定义了一些有关字符类型判断和转换的宏。
8
9  char ctmp;                // 一个临时字符变量，供 ctype.h 文件中转换字符宏函数使用。
10 // 字符特性数组(表)，定义了各个字符对应的属性，这些属性类型(如_C等)在 ctype.h 中定义。
11 // 用于判断字符是控制字符(_C)、大写字母(_U)、小写字母(_L)等所属类型。
12 unsigned char ctype[] = {0x00,          /* EOF */
13 C, C, C, C, C, C, C, C,          /* 0-7 */
14 C, C, S, C, S, C, S, C, S, C, C,          /* 8-15 */
15 C, C, C, C, C, C, C, C,          /* 16-23 */
16 C, C, C, C, C, C, C, C,          /* 24-31 */
17 S, SP, P, P, P, P, P, P, P,          /* 32-39 */
18 P, P, P, P, P, P, P, P,          /* 40-47 */
19 D, D, D, D, D, D, D, D,          /* 48-55 */
20 D, D, P, P, P, P, P, P,          /* 56-63 */
21 P, U, X, U, X, U, X, U, X, U, X, U,          /* 64-71 */
22 U, U, U, U, U, U, U, U,          /* 72-79 */

```

---

|           |   |                      |
|-----------|---|----------------------|
| <u>21</u> | <u>U</u> , <u>U</u> , <u>U</u> , <u>U</u> , <u>U</u> , <u>U</u> , <u>U</u> , <u>U</u> ,   | <i>/* 80-87 */</i>   |
| <u>22</u> | <u>U</u> , <u>U</u> , <u>U</u> , <u>P</u> , <u>P</u> , <u>P</u> , <u>P</u> , <u>P</u> ,   | <i>/* 88-95 */</i>   |
| <u>23</u> | <u>P</u> , <u>L</u>   <u>X</u> , <u>L</u>   <u>X</u> , <u>L</u>   <u>X</u> , <u>L</u>   <u>X</u> , <u>L</u>   <u>X</u> , <u>L</u>   <u>X</u> , <u>L</u> , | <i>/* 96-103 */</i>  |
| <u>24</u> | <u>L</u> , <u>L</u> , <u>L</u> , <u>L</u> , <u>L</u> , <u>L</u> , <u>L</u> , <u>L</u> , <u>L</u> ,  | <i>/* 104-111 */</i> |
| <u>25</u> | <u>L</u> , <u>L</u> , <u>L</u> , <u>L</u> , <u>L</u> , <u>L</u> , <u>L</u> , <u>L</u> , <u>L</u> ,  | <i>/* 112-119 */</i> |
| <u>26</u> | <u>L</u> , <u>L</u> , <u>L</u> , <u>P</u> , <u>P</u> , <u>P</u> , <u>P</u> , <u>C</u> ,   | <i>/* 120-127 */</i> |
| <u>27</u> | 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,   | <i>/* 128-143 */</i> |
| <u>28</u> | 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,   | <i>/* 144-159 */</i> |
| <u>29</u> | 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,   | <i>/* 160-175 */</i> |
| <u>30</u> | 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,   | <i>/* 176-191 */</i> |
| <u>31</u> | 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,   | <i>/* 192-207 */</i> |
| <u>32</u> | 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,   | <i>/* 208-223 */</i> |
| <u>33</u> | 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,   | <i>/* 224-239 */</i> |
| <u>34</u> | 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0 } ;  | <i>/* 240-255 */</i> |
| <u>35</u> |   |                      |
| <u>36</u> |   |                      |

## 11.6 dup.c 程序

### 11.6.1 功能描述

该程序包括一个创建文件描述符拷贝的函数 `dup()`。在成功返回之后，新的和原来的描述符可以交替使用。它们共享锁定、文件读写指针以及文件标志。例如，如果文件读写位置指针被其中一个描述符使用 `lseek()` 修改过之后，则对于另一个描述符来讲，文件读写指针也被改变。该函数使用数值最小的未使用描述符来建立新描述符。但是这两个描述符并不共享执行时关闭标志 (`close-on-exec`)。

### 11.6.2 代码注释

列表 linux/lib/dup.c 程序

```

1  /*
2   *  linux/lib/dup. c
3   *
4   *  (C) 1991  Linus Torvalds
5   */
6
7  #define  LIBRARY
8  #include <unistd.h>          // Linux 标准头文件。定义了各种符号常数和类型，并声明了各种函数。
                                // 如定义了__LIBRARY__，则还包括系统调用号和内嵌汇编_syscall0()等。
9
10  ///// 复制文件描述符函数。
11  // 下面该调用宏函数对应：int dup(int fd)。直接调用了系统中断 int 0x80，参数是__NR_dup。
12  // 其中 fd 是文件描述符。
13  syscall11(int, dup, int, fd)

```

## 11.7 errno.c 程序

### 11.7.1 功能描述

该程序仅定义了一个出错号变量 `errno`。用于在函数调用失败时存放出错号。

### 11.7.2 代码注释

列表 linux/lib/errno.c 程序

---

```

1 /*
2  * linux/lib/errno.c
3  *
4  * (C) 1991 Linus Torvalds
5  */
6
7 int errno;
8

```

---

## 11.8 execve.c 程序

### 11.8.1 功能描述

### 11.8.2 代码注释

列表 linux/lib/execve.c 程序

---

```

1 /*
2  * linux/lib/execve.c
3  *
4  * (C) 1991 Linus Torvalds
5  */
6
7 #define LIBRARY
8 #include <unistd.h>          // Linux 标准头文件。定义了各种符号常数和类型，并声明了各种函数。
                                // 如定义了 __LIBRARY__，则还包括系统调用号和内嵌汇编 _syscall0() 等。
9
10 // 加载并执行子进程 (其它程序) 函数。
11 // 下面该调用宏函数对应: int execve(const char * file, char ** argv, char ** envp)。
12 // 参数: file - 被执行程序文件名; argv - 命令行参数指针数组; envp - 环境变量指针数组。
13 // 直接调用了系统中断 int 0x80，参数是 __NR_execve。参见 include/unistd.h 和 fs/exec.c 程序。
14 syscall13(int, execve, const char *, file, char **, argv, char **, envp)
15

```

---

## 11.9 malloc.c 程序

### 11.9.1 功能描述

该程序中主要包括内存分配函数 `malloc()`。该函数使用了存储桶(bucket)的原理对分配的内存进行管理。基本思想是对不同请求的内存块大小(长度),使用存储桶目录(下面简称目录)分别进行处理。比如对于请求内存块的长度在 32 字节或 32 字节以下但大于 16 字节时,就使用存储桶目录第二项对应的存储桶描述符链表分配内存块。其基本结构示意图见下图所示。该函数目前一次所能分配的最大内存长度是一个内存页面,即 4096 字节。

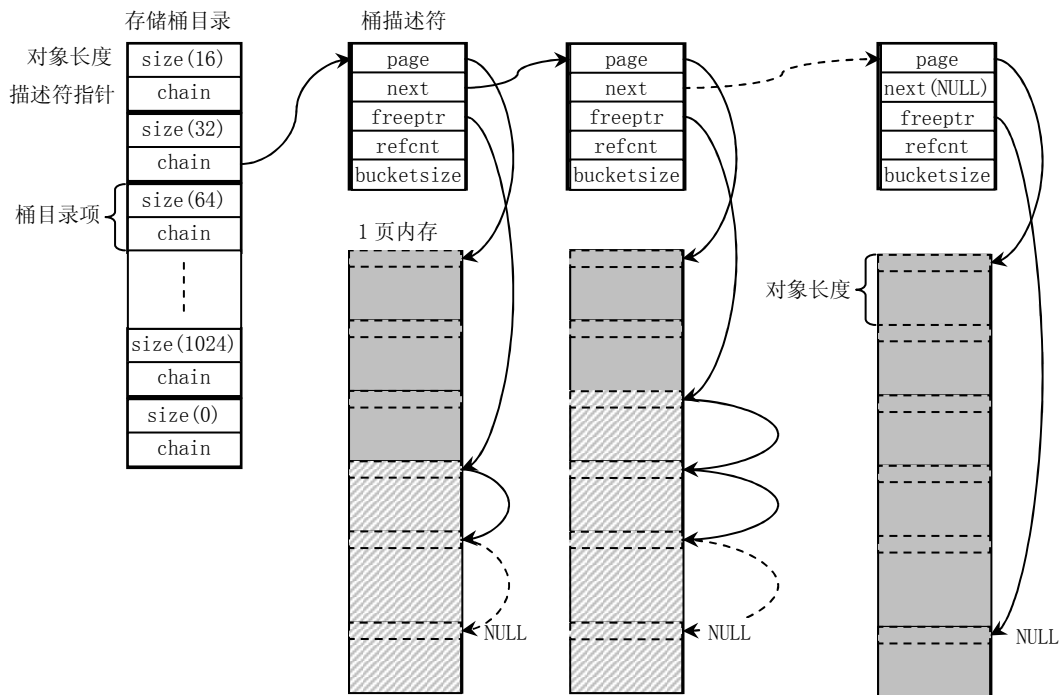


图 使用存储桶原理进行内存分配管理的结构示意图

在第一次调用 `malloc()` 函数时,首先要建立一个页面的空闲存储桶描述符(下面简称描述符)链表,其中存放着还未使用或已经使用完毕而收回的描述符。该链表结构示意图见下图所示。其中 `free_bucket_desc` 是链表头指针。从链表中取出或放入一个描述符都是从链表头开始操作。当取出一个描述符时,就将链表头指针所指向的头一个描述符取出;当释放一个空闲描述符时也是将其放在链表头处。

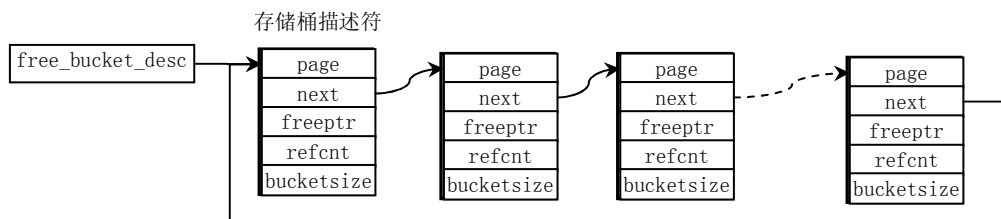


图 空闲存储桶描述符链表结构示意图



malloc() 函数执行的基本步骤如下:

1. 首先搜索目录, 寻找适合请求内存块大小的目录项对应的描述符链表。当目录项的对象字节长度大于请求的字节长度, 就算找到了相应的目录项。如果搜索完整个目录都没有找到合适的目录项, 则说明用户请求的内存块太大。
2. 在目录项对应的描述符链表中查找具有空闲空间的描述符。如果某个描述符的空闲内存指针 freeptr 不为 NULL, 则表示找到了相应的描述符。如果没有找到具有空闲空间的描述符, 那么我们就需要新建一个描述符。新建描述符的过程如下:
  - a. 如果空闲描述符链表头指针还是 NULL 的话, 说明是第一次调用 malloc() 函数, 此时需要 init\_bucket\_desc() 来创建空闲描述符链表。
  - b. 然后从空闲描述符链表头处取一个描述符, 初始化该描述符, 令其对象引用计数为 0, 对象大小等于对应目录项指定对象的长度值, 并申请一内存页面, 让描述符的页面指针 page 指向该内存页, 描述符的空闲内存指针 freeptr 也指向页开始位置。
  - c. 对该内存页面根据本目录项所用对象长度进行页面初始化, 建立所有对象的一个链表。也即每个对象的头部都存放一个指向下一个对象的指针, 最后一个对象的开始处存放一个 NULL 指针值。
  - d. 然后将该描述符插入到对应目录项的描述符链表开始处。
3. 将该描述符的空闲内存指针 freeptr 复制为返回给用户的内存指针, 然后调整该 freeptr 指向描述符对应内存页面中下一个空闲对象位置, 并使该描述符引用计数值增 1。

free\_s() 函数用于回收用户释放的内存块。基本原理是首先根据该内存块的地址换算出该内存块对应页面的地址(用页面长度进行模运算), 然后搜索目录中的所有描述符, 找到对应该页面的描述符。将该释放的内存块链入 freeptr 所指向的空闲对象链表中, 并将描述符的对象引用计数值减 1。如果引用计数值此时等于零, 则表示该描述符对应的页面已经完全空出, 可以释放该内存页面并将该描述符收回到空闲描述符链表中。

## 11.9.2 代码注释

列表 linux/lib/malloc.c 程序

```

1 /*
2  * malloc.c — a general purpose kernel memory allocator for Linux.
3  *
4  * Written by Theodore Ts'o (tytso@mit.edu), 11/29/91
5  *
6  * This routine is written to be as fast as possible, so that it
7  * can be called from the interrupt level.
8  *
9  * Limitations: maximum size of memory we can allocate using this routine
10 *   is 4k, the size of a page in Linux.
11 *
12 * The general game plan is that each page (called a bucket) will only hold
13 * objects of a given size. When all of the object on a page are released,
14 * the page can be returned to the general free pool. When malloc() is
15 * called, it looks for the smallest bucket size which will fulfill its
16 * request, and allocate a piece of memory from that bucket pool.
17 *
18 * Each bucket has as its control block a bucket descriptor which keeps
19 * track of how many objects are in use on that page, and the free list
20 * for that page. Like the buckets themselves, bucket descriptors are

```

```

21 * stored on pages requested from get_free_page(). However, unlike buckets,
22 * pages devoted to bucket descriptor pages are never released back to the
23 * system. Fortunately, a system should probably only need 1 or 2 bucket
24 * descriptor pages, since a page can hold 256 bucket descriptors (which
25 * corresponds to 1 megabyte worth of bucket pages.) If the kernel is using
26 * that much allocated memory, it's probably doing something wrong. :-)
27 *
28 * Note: malloc() and free() both call get_free_page() and free_page()
29 * in sections of code where interrupts are turned off, to allow
30 * malloc() and free() to be safely called from an interrupt routine.
31 * (We will probably need this functionality when networking code,
32 * particularly things like NFS, is added to Linux.) However, this
33 * presumes that get_free_page() and free_page() are interrupt-level
34 * safe, which they may not be once paging is added. If this is the
35 * case, we will need to modify malloc() to keep a few unused pages
36 * "pre-allocated" so that it can safely draw upon those pages if
37 * it is called from an interrupt routine.
38 *
39 * Another concern is that get_free_page() should not sleep; if it
40 * does, the code is carefully ordered so as to avoid any race
41 * conditions. The catch is that if malloc() is called re-entrantly,
42 * there is a chance that unnecessary pages will be grabbed from the
43 * system. Except for the pages for the bucket descriptor page, the
44 * extra pages will eventually get released back to the system, though,
45 * so it isn't all that bad.
46 */
47
/*
 * malloc.c - Linux 的通用内核内存分配函数。
 *
 * 由 Theodore Ts'o 编制 (tytso@mit.edu), 11/29/91
 *
 * 该函数被编写成尽可能地快, 从而可以从中断层调用此函数。
 *
 * 限制: 使用该函数一次所能分配的最大内存是 4k, 也即 Linux 中内存页面的大小。
 *
 * 编写该函数所遵循的一般规则是每页(被称为一个存储桶)仅分配所要容纳对象的大小。
 * 当一页上的所有对象都释放后, 该页就可以返回通用空闲内存池。当 malloc() 被调用
 * 时, 它会寻找满足要求的最小的存储桶, 并从该存储桶中分配一块内存。
 *
 * 每个存储桶都有一个作为其控制用的存储桶描述符, 其中记录了页面上有多少对象正被
 * 使用以及该页上空闲内存的列表。就象存储桶自身一样, 存储桶描述符也是存储在使用
 * get_free_page() 申请到的页面上的, 但是与存储桶不同的是, 桶描述符所占用的页面
 * 将不再会释放给系统。幸运的是一个系统大约只需要 1 到 2 页的桶描述符页面, 因为一
 * 个页面可以存放 256 个桶描述符(对应 1MB 内存的存储桶页面)。如果系统为桶描述符分
 * 配了许多内存, 那么肯定系统什么地方出了问题☺。
 *
 * 注意! malloc() 和 free() 两者关闭了中断的代码部分都调用了 get_free_page() 和
 * free_page() 函数, 以使 malloc() 和 free() 可以安全地被从中断程序中调用
 * (当网络代码, 尤其是 NFS 等被加入到 Linux 中时就可能需要这种功能)。但前
 * 提是假设 get_free_page() 和 free_page() 是可以安全地在中断级程序中使用的,
 * 这在一旦加入了分页处理之后就很可能不是安全的。如果真是这种情况, 那么我
 * 们就需要修改 malloc() 来“预先分配”几页不用的内存, 如果 malloc() 和 free()

```

```

*      被从中断程序中调用时就可以安全地使用这些页面。
*
*      另外需要考虑到的是 get_free_page() 不应该睡眠；如果会睡眠的话，则为了防止
*      任何竞争条件，代码需要仔细地安排顺序。 关键在于如果 malloc() 是可以重入地
*      被调用的话，那么就会存在不必要的页面被从系统中取走的机会。除了用于桶描述
*      符的页面，这些额外的页面最终会释放给系统，所以并不是象想象的那样不好。
*/

```

```

48 #include <linux/kernel.h> // 内核头文件。含有一些内核常用函数的原形定义。
49 #include <linux/mm.h>     // 内存管理头文件。含有页面大小定义和一些页面释放函数原型。
50 #include <asm/system.h>    // 系统头文件。定义了设置或修改描述符/中断门等的嵌入式汇编宏。

```

```
51
```

```
// 存储桶描述符结构。
```

```

52 struct bucket_desc { /* 16 bytes */
53     void *page; // 该桶描述符对应的内存页面指针。
54     struct bucket_desc *next; // 下一个描述符指针。
55     void *freeptr; // 指向本桶中空闲内存位置的指针。
56     unsigned short refcnt; // 引用计数。
57     unsigned short bucket_size; // 本描述符对应存储桶的大小。
58 };

```

```
59
```

```
// 存储桶描述符目录结构。
```

```

60 struct bucket_dir { /* 8 bytes */
61     int size; // 该存储桶的大小(字节数)。
62     struct bucket_desc *chain; // 该存储桶目录项的桶描述符链表指针。
63 };

```

```
64
```

```
65 /*
```

```

66  * The following is the where we store a pointer to the first bucket
67  * descriptor for a given size.
68  *
69  * If it turns out that the Linux kernel allocates a lot of objects of a
70  * specific size, then we may want to add that specific size to this list,
71  * since that will allow the memory to be allocated more efficiently.
72  * However, since an entire page must be dedicated to each specific size
73  * on this list, some amount of temperance must be exercised here.
74  *
75  * Note that this list must be kept in order.
76  */

```

```
/*
```

```
* 下面是我们存放第一个给定大小存储桶描述符指针的地方。
```

```
*
```

```

* 如果 Linux 内核分配了许多指定大小的对象，那么我们就希望将该指定的大小加到
* 该列表(链表)中，因为这样可以使内存的分配更有效。但是，因为一页完整内存页面
* 必须用于列表中指定大小的所有对象，所以需要总数方面的测试操作。

```

```
*/
```

```
// 存储桶目录列表(数组)。
```

```

77 struct bucket_dir bucket_dir[] = {
78     { 16, (struct bucket_desc *) 0}, // 16 字节长度的内存块。
79     { 32, (struct bucket_desc *) 0}, // 32 字节长度的内存块。
80     { 64, (struct bucket_desc *) 0}, // 64 字节长度的内存块。
81     { 128, (struct bucket_desc *) 0}, // 128 字节长度的内存块。
82     { 256, (struct bucket_desc *) 0}, // 256 字节长度的内存块。

```

```

83     { 512, (struct bucket\_desc *) 0},    // 512 字节长度的内存块。
84     { 1024, (struct bucket\_desc *) 0},    // 1024 字节长度的内存块。
85     { 2048, (struct bucket\_desc *) 0},    // 2048 字节长度的内存块。
86     { 4096, (struct bucket\_desc *) 0},    // 4096 字节(1 页)内存。
87     { 0, (struct bucket\_desc *) 0}};    /* End of list marker */
88
89 /*
90  * This contains a linked list of free bucket descriptor blocks
91  */
92 /*
93  * 下面是含有空闲桶描述符内存块的链表。
94  */
95 struct bucket\_desc *free\_bucket\_desc = (struct bucket\_desc *) 0;
96
97 /*
98  * This routine initializes a bucket description page.
99  */
100 /*
101  * 下面的子程序用于初始化一页桶描述符页面。
102  */
103 // 初始化桶描述符。
104 // 建立空闲桶描述符链表，并让 free_bucket_desc 指向第一个空闲桶描述符。
105 static inline void init\_bucket\_desc()
106 {
107     struct bucket\_desc *bdesc, *first;
108     int i;
109
110     // 申请一页内存，用于存放桶描述符。如果失败，则显示初始化桶描述符时内存不够出错信息，死机。
111     first = bdesc = (struct bucket\_desc *) get\_free\_page();
112     if (!bdesc)
113         panic("Out of memory in init_bucket_desc()");
114     // 首先计算一页内存中可存放的桶描述符数量，然后对其建立单向连接指针。
115     for (i = PAGE\_SIZE/sizeof(struct bucket\_desc); i > 1; i--) {
116         bdesc->next = bdesc+1;
117         bdesc++;
118     }
119     /*
120      * This is done last, to avoid race conditions in case
121      * get_free_page() sleeps and this routine gets called again...
122      */
123     /*
124      * 这是在最后处理的，目的是为了在 get_free_page() 睡眠时该子程序又被
125      * 调用而引起的竞争条件。
126      */
127     // 将空闲桶描述符指针 free_bucket_desc 加入链表中。
128     bdesc->next = free\_bucket\_desc;
129     free\_bucket\_desc = first;
130 }
131
132 // 分配动态内存函数。
133 // 参数: len - 请求的内存块长度。
134 // 返回: 指向被分配内存的指针。如果失败则返回 NULL。
135 void *malloc(unsigned int len)

```

```

118 {
119     struct bucket\_dir      *bdir;
120     struct bucket\_desc    *bdesc;
121     void                  *retval;
122
123     /*
124      * First we search the bucket_dir to find the right bucket change
125      * for this request.
126      */
127     /*
128      * 首先我们搜索存储桶目录 bucket_dir 来寻找适合请求的桶大小。
129      */
130     // 搜索存储桶目录，寻找适合申请内存块大小的桶描述符链表。如果目录项的桶字节数大于请求的字节
131     // 数，就找到了对应的桶目录项。
132     for (bdir = bucket\_dir; bdir->size; bdir++)
133         if (bdir->size >= len)
134             break;
135     // 如果搜索完整个目录都没有找到合适大小的目录项，则表明所请求的内存块大小太大，超出了该
136     // 程序的分配限制(最长为 1 个页面)。于是显示出错信息，死机。
137     if (!bdir->size) {
138         printk("malloc called with impossibly large argument (%d)\n",
139             len);
140         panic("malloc: bad arg");
141     }
142     /*
143      * Now we search for a bucket descriptor which has free space
144      */
145     /*
146      * 现在我们来搜索具有空闲空间的桶描述符。
147      */
148     cli(); /* Avoid race conditions */ /* 为了避免出现竞争条件，首先关中断 */
149     // 搜索对应桶目录项中描述符链表，查找具有空闲空间的桶描述符。如果桶描述符的空闲内存指针
150     // freeptr 不为空，则表示找到了相应的桶描述符。
151     for (bdesc = bdir->chain; bdesc; bdesc = bdesc->next)
152         if (bdesc->freeptr)
153             break;
154     /*
155      * If we didn't find a bucket with free space, then we'll
156      * allocate a new one.
157      */
158     /*
159      * 如果没有找到具有空闲空间的桶描述符，那么我们就需要新建一个该目录项的描述符。
160      */
161     if (!bdesc) {
162         char      *cp;
163         int       i;
164
165         // 若 free_bucket_desc 还为空时，表示第一次调用该程序，则对描述符链表进行初始化。
166         // free_bucket_desc 指向第一个空闲桶描述符。
167         if (!free\_bucket\_desc)
168             init\_bucket\_desc();
169         // 取 free_bucket_desc 指向的空闲桶描述符，并让 free_bucket_desc 指向下一个空闲桶描述符。
170         bdesc = free\_bucket\_desc;

```

```

153         free_bucket_desc = bdesc->next;
// 初始化该新的桶描述符。令其引用数量等于 0；桶的大小等于对应桶目录的大小；申请一内存页面，
// 让描述符的页面指针 page 指向该页面；空闲内存指针也指向该页开头，因为此时全为空闲。
154         bdesc->refcnt = 0;
155         bdesc->bucket_size = bdir->size;
156         bdesc->page = bdesc->freeptr = (void *) cp = get_free_page();
// 如果申请内存页面操作失败，则显示出错信息，死机。
157         if (!cp)
158             panic("Out of memory in kernel malloc()");
159         /* Set up the chain of free objects */
        /* 在该页空闲内存中建立空闲对象链表 */
// 以该桶目录项指定的桶大小为对象长度，对该页内存进行划分，并使每个对象的开始 4 字节设置
// 成指向下一对象的指针。
160         for (i=PAGE_SIZE/bdir->size; i > 1; i--) {
161             *((char **) cp) = cp + bdir->size;
162             cp += bdir->size;
163         }
// 最后一个对象开始处的指针设置为 0(NULL)。
// 然后让该桶描述符的下一描述符指针字段指向对应桶目录项指针 chain 所指的描述符，而桶目录的
// chain 指向该桶描述符，也即将该描述符插入到描述符链链头处。
164         *((char **) cp) = 0;
165         bdesc->next = bdir->chain; /* OK, link it in! */ /* OK, 将其链入! */
166         bdir->chain = bdesc;
167     }
// 返回指针即等于该描述符对应页面的当前空闲指针。然后调整该空闲空间指针指向下一个空闲对象，
// 并使描述符中对应页面中对象引用计数增 1。
168     retval = (void *) bdesc->freeptr;
169     bdesc->freeptr = *((void **) retval);
170     bdesc->refcnt++;
// 最后开放中断，并返回指向空闲内存对象的指针。
171     sti(); /* OK, we're safe again */ /* OK, 现在我们又安全了*/
172     return(retval);
173 }
174
175 /*
176  * Here is the free routine. If you know the size of the object that you
177  * are freeing, then free_s() will use that information to speed up the
178  * search for the bucket descriptor.
179  *
180  * We will #define a macro so that "free(x)" is becomes "free_s(x, 0)"
181  */
// 下面是释放子程序。如果你知道释放对象的大小，则 free_s() 将使用该信息加速
// 搜寻对应桶描述符的速度。
// 我们将定义一个宏，使得 "free(x)" 成为 "free_s(x, 0)"。
//// 释放存储桶对象。
// 参数: obj - 对应对象指针; size - 大小。
182 void free_s(void *obj, int size)
183 {
184     void *page;
185     struct bucket_dir *bdir;

```



```

186     struct bucket\_desc      *bdesc, *prev;
187
188     /* Calculate what page this object lives in */
189     /* 计算该对象所在的页面 */
190     page = (void *) ((unsigned long) obj & 0xfffff000);
191     /* Now search the buckets looking for that page */
192     /* 现在搜索存储桶目录项所链接的桶描述符，寻找该页面 */
193
194     //
195     for (bdir = bucket\_dir; bdir->size; bdir++) {
196         prev = 0;
197         /* If size is zero then this conditional is always false */
198         /* 如果参数 size 是 0，则下面条件肯定是 false */
199         if (bdir->size < size)
200             continue;
201         // 搜索对应目录项中链接的所有描述符，查找对应页面。如果某描述符页面指针等于 page 则表示找到
202         // 了相应的描述符，跳转到 found。如果描述符不含有对应 page，则让描述符指针 prev 指向该描述符。
203         for (bdesc = bdir->chain; bdesc; bdesc = bdesc->next) {
204             if (bdesc->page == page)
205                 goto found;
206             prev = bdesc;
207         }
208         // 若搜索了对应目录项的所有描述符都没有找到指定的页面，则显示出错信息，死机。
209         panic("Bad address passed to kernel free_s()");
210 found:
211         // 找到对应的桶描述符后，首先关中断。然后将该对象内存块链入空闲块对象链表中，并使该描述符
212         // 的对象引用计数减 1。
213         cli(); /* To avoid race conditions */ /* 为了避免竞争条件 */
214         *((void **)obj) = bdesc->freeptr;
215         bdesc->freeptr = obj;
216         bdesc->refcnt--;
217         // 如果引用计数已等于 0，则我们就可以释放对应的内存页面和该桶描述符。
218         if (bdesc->refcnt == 0) {
219             /*
220             * We need to make sure that prev is still accurate. It
221             * may not be, if someone rudely interrupted us....
222             */
223             /*
224             * 我们需要确信 prev 仍然是正确的，若某程序粗鲁地中断了我们
225             * 就有可能不是了。
226             */
227             // 如果 prev 已经不是搜索到的描述符的前一个描述符，则重新搜索当前描述符的前一个描述符。
228             if ((prev && (prev->next != bdesc)) ||
229                 (!prev && (bdir->chain != bdesc)))
230                 for (prev = bdir->chain; prev; prev = prev->next)
231                     if (prev->next == bdesc)
232                         break;
233             // 如果找到该前一个描述符，则从描述符链中删除当前描述符。
234             if (prev)
235                 prev->next = bdesc->next;
236             // 如果 prev==NULL，则说明当前一个描述符是该目录项首个描述符，也即目录项中 chain 应该直接
237             // 指向当前描述符 bdesc，否则表示链表有问题，则显示出错信息，死机。因此，为了将当前描述符
238             // 从链表中删除，应该让 chain 指向下一个描述符。

```

```

220         else {
221             if (bdir->chain != bdesc)
222                 panic("malloc bucket chains corrupted");
223             bdir->chain = bdesc->next;
224         }
225     // 释放当前描述符所操作的内存页面，并将该描述符插入空闲描述符链表开始处。
226     free\_page((unsigned long) bdesc->page);
227     bdesc->next = free\_bucket\_desc;
228     free\_bucket\_desc = bdesc;
229 }
230 // 开中断，返回。
231 sti();
232 return;
233 }

```

## 11.10 open.c 程序

### 11.10.1 功能描述

`open()` 系统调用用于将一个文件名转换成一个文件描述符。当调用成功时，返回的文件描述符将是进程没有打开的最小数值的描述符。该调用创建一个新的打开文件，并不与任何其它进程共享。在执行 `exec` 函数时，该新的文件描述符将始终保持着打开状态。文件的读写指针被设置在文件开始位置。

参数 `flag` 是 `O_RDONLY`、`O_WRONLY`、`O_RDWR` 之一，分别代表文件只读打开、只写打开和读写打开方式，可以与其它一些标志一起使用。（参见 `fs/open.c`，138 行）

### 11.10.2 代码注释

列表 linux/lib/open.c 程序

```

1  /*
2   *  linux/lib/open.c
3   *
4   *  (C) 1991  Linus Torvalds
5   */
6
7  #define  LIBRARY
8  #include <unistd.h>    // Linux 标准头文件。定义了各种符号常数和类型，并声明了各种函数。
                        // 如定义了 __LIBRARY__，则还包括系统调用号和内嵌汇编 _syscall0() 等。
9  #include <stdarg.h>    // 标准参数头文件。以宏的形式定义变量参数列表。主要说明了一个
                        // 类型 (va_list) 和三个宏 (va_start, va_arg 和 va_end)，用于
                        // vsprintf、vprintf、vfprintf 函数。
10
    ///// 打开文件函数。
    // 打开并有可能创建一个文件。
    // 参数: filename - 文件名; flag - 文件打开标志; ...
    // 返回: 文件描述符，若出错则置出错码，并返回-1。

```



```

11 int open(const char * filename, int flag, ...)
12 {
13     register int res;
14     va_list arg;
15
16     // 利用 va_start() 宏函数, 取得 flag 后面参数的指针, 然后调用系统中断 int 0x80, 功能 open 进行
17     // 文件打开操作。
18     // %0 - eax(返回的描述符或出错码); %1 - eax(系统中断调用功能号 __NR_open);
19     // %2 - ebx(文件名 filename); %3 - ecx(打开文件标志 flag); %4 - edx(后随参数文件属性 mode)。
20     va_start(arg, flag);
21     __asm__( "int $0x80"
22             : "=a" (res)
23             : "" ( __NR_open), "b" (filename), "c" (flag),
24             "d" (va_arg(arg, int)));
25     // 系统中断调用返回值大于或等于 0, 表示是一个文件描述符, 则直接返回之。
26     if (res >= 0)
27         return res;
28     // 否则说明返回值小于 0, 则代表一个出错码。设置该出错码并返回-1。
29     errno = -res;
30     return -1;
31 }

```

## 11.11 setsid.c 程序

### 11.11.1 功能描述

该程序包括一个 setsid() 系统调用函数。如果调用的进程不是一个组的领导时, 该函数用于创建一个新会话。则调用进程将成为该新会话的领导、新进程组的组领导, 并且没有控制终端。调用进程的组 id 和会话 id 被设置成进程的 PID(进程标识符)。调用进程将成为新进程组和新会话中的唯一进程。

### 11.11.2 代码注释

列表 linux/lib/setsid.c 程序

```

1  /*
2   * linux/lib/setsid.c
3   *
4   * (C) 1991 Linus Torvalds
5   */
6
7  #define LIBRARY
8  #include <unistd.h> // Linux 标准头文件。定义了各种符号常数和类型, 并声明了各种函数。
9                      // 如定义了 __LIBRARY__, 则还包括系统调用号和内嵌汇编 _syscall10() 等。
10
11  // 创建一个会话并设置进程组号。
12  // 下面系统调用宏对应于函数: pid_t setsid()。
13  // 返回: 调用进程的会话标识符(session ID)。
14  _syscall10(pid_t, setsid)
15

```

## 11.12 string.c 程序

### 11.12.1 功能描述

所有字符串操作函数已经在 string.h 中实现，因此 string.c 程序仅包含 string.h 头文件。

### 11.12.2 代码注释

列标 linux/lib/string.c 程序

---

```

1  /*
2   *  linux/lib/string.c
3   *
4   *  (C) 1991  Linus Torvalds
5   */
6
7  #ifndef __GNUC__           // 需要 GNU 的 C 编译器编译。
8  #error I want gcc!
9  #endif
10
11 #define extern
12 #define inline
13 #define  LIBRARY
14 #include <string.h>
15

```

---

## 11.13 wait.c 程序

### 11.13.1 功能描述

该程序包括函数 waitpid() 和 wait()。这两个函数允许进程获取与其子进程之一的状态信息。各种选项允许获取已经终止或停止的子进程状态信息。如果存在两个或两个以上子进程的状态信息，则报告的顺序是不指定的。

wait() 将挂起当前进程，直到其子进程之一退出（终止），或者收到要求终止该进程的信号，或者是需要调用一个信号句柄（信号处理程序）。

waitpid() 挂起当前进程，直到 pid 指定的子进程退出（终止）或者收到要求终止该进程的信号，或者是需要调用一个信号句柄（信号处理程序）。

如果 pid = -1, options = 0, 则 waitpid() 的作用与 wait() 函数一样。否则其行为将随 pid 和 options 参数的不同而不同。（参见 kernel/exit.c, 142）

### 11.13.2 代码注释

列表 linux/lib/wait.c 程序

---

```

1  /*
2  *  linux/lib/wait.c
3  *
4  *  (C) 1991  Linus Torvalds
5  */
6
7  #define  LIBRARY
8  #include <unistd.h>          // Linux 标准头文件。定义了各种符号常数和类型，并声明了各种函数。
                                // 如定义了__LIBRARY__，则还包括系统调用号和内嵌汇编_syscall0()等。
9  #include <sys/wait.h>       // 等待调用头文件。定义系统调用 wait() 和 waitpid() 及相关常数符号。
10
    // 等待进程终止系统调用函数。
    // 该下面宏结构对应于函数：pid_t waitpid(pid_t pid, int * wait_stat, int options)
    //
    // 参数：pid - 等待被终止进程的进程 id，或者是用于指定特殊情况的其它特定数值；
    //        wait_stat - 用于存放状态信息；options - WNOHANG 或 WUNTRACED 或是 0。
11  syscall3(pid_t, waitpid, pid_t, pid, int *, wait_stat, int, options)
12
    // 等待进程终止系统调用。直接调用 waitpid() 函数。
13  pid_t wait(int * wait_stat)
14  {
15      return waitpid(-1, wait_stat, 0);
16  }
17

```

---

## 11.14 write.c 程序

### 11.14.1 功能描述

该程序中包括一个向文件描述符写操作函数 write()。该函数向文件描述符指定的文件写入 count 字节的数据到缓冲区 buf 中。

### 11.14.2 代码注释

列表 linux/lib/write.c 程序

---

```

1  /*
2  *  linux/lib/write.c
3  *
4  *  (C) 1991  Linus Torvalds
5  */
6
7  #define  LIBRARY
8  #include <unistd.h>          // Linux 标准头文件。定义了各种符号常数和类型，并声明了各种函数。
                                // 如定义了__LIBRARY__，则还包括系统调用号和内嵌汇编_syscall0()等。
9

```

---

---

//// 写文件系统调用函数。

// 该宏结构对应于函数: int write(int fd, const char \* buf, off\_t count)

// 参数: fd - 文件描述符; buf - 写缓冲区指针; count - 写字节数。

// 返回: 成功时返回写入的字节数(0 表示写入 0 字节); 出错时将返回-1, 并且设置了出错号。

10 syscall3(int, write, int, fd, const char \*, buf, off\_t, count)

11

---

## 第 12 章 建造工具(tools)

### 12.1 概述

Linux 内核源代码中的 tools 目录中包含一个生成内核磁盘映象文件的工具程序 build.c，该程序将单独编译成可执行文件，在 linux/ 目录下的 Makefile 文件中被调用运行，用于将所有内核编译代码连接和合并成一个可运行的内核映象文件 image。具体方法是对 boot/ 中的 bootsect.s、setup.s 使用 8086 汇编器进行编译，分别生成各自的执行模块。再对源代码中的其它所有程序使用 GNU 的编译器 gcc/gas 进行编译，并连接成模块 system。然后使用 build 工具将这三块组合成一个内核映象文件 image。基本编译连接/组合结构如下图 12.1 所示。

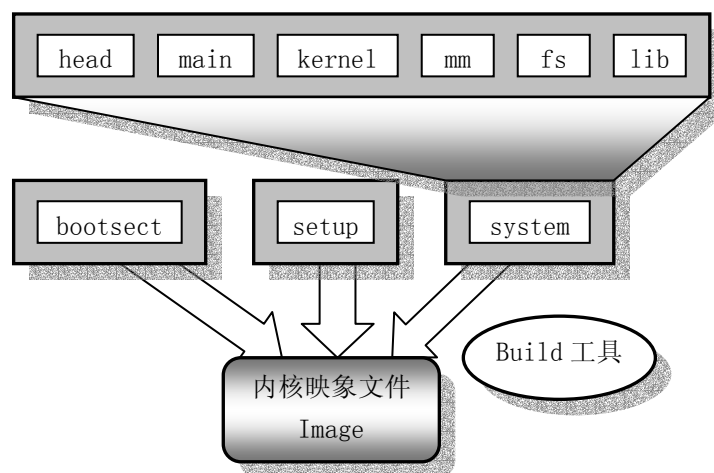


图 12.1 内核编译连接/组合结构

### 12.2 build.c 程序

#### 12.2.1 功能概述

build 程序使用 4 个参数，分别是 bootsect 文件名、setup 文件名、system 文件名和可选的根文件系统设备文件名。

程序首先检查命令行上最后一个根设备文件名可选参数，若其存在，则读取该设备文件的状态信息结构 (stat)，取出设备号。若命令行上不带该参数，则使用默认值。

然后对 bootsect 文件进行处理，读取该文件的 minix 执行头部信息，判断其有效性，然后读取随后 512 字节的引导代码数据，判断其是否具有可引导标志 0xAA55，并将前面获取的根设备号写入到 508, 509 位移处，最后将该 512 字节代码数据写到 stdout 标准输出，由 Make 文件重定向到 Image 文件。

接下来以类似的方法处理 setup 文件。若该文件长度小于 4 个扇区，则用 0 将其填满为 4 个扇区的长度，并写到标准输出 stdout 中。

最后处理 system 文件。该文件是使用 GCC 编译器产生，所以其执行头部格式是 GCC 类型的，与 linux 定义的 a.out 格式一样。在判断执行入口点是 0 后，就将数据写到标准输出 stdout 中。若其代码数据长度超过 128KB，则显示出错信息。

## 12.2.2 代码注释

列表 linux/tools/build.c 程序

---

```

1  /*
2  *  linux/tools/build.c
3  *
4  *  (C) 1991  Linus Torvalds
5  */
6
7  /*
8  *  This file builds a disk-image from three different files:
9  *
10 *  - bootsect: max 510 bytes of 8086 machine code, loads the rest
11 *  - setup: max 4 sectors of 8086 machine code, sets up system parm
12 *  - system: 80386 code for actual system
13 *
14 *  It does some checking that all files are of the correct type, and
15 *  just writes the result to stdout, removing headers and padding to
16 *  the right amount. It also writes some system data to stderr.
17 */
18 /*
19  *  该程序从三个不同的程序中创建磁盘映象文件：
20  *
21  *  - bootsect: 该文件的 8086 机器码最长为 510 字节，用于加载其它程序。
22  *  - setup: 该文件的 8086 机器码最长为 4 个磁盘扇区，用于设置系统参数。
23  *  - system: 实际系统的 80386 代码。
24  *
25  *  该程序首先检查所有程序模块的类型是否正确，并将检查结果在终端上显示出来，
26  *  然后删除模块头部并扩充大正确的长度。该程序也会将一些系统数据写到 stderr。
27 */
28
29 /*
30 *  Changes by tytso to allow root device specification
31 */
32 /*
33 *  tytso 对该程序作了修改，以允许指定根文件设备。
34 */
35
36 #include <stdio.h>          /* fprintf */          /* 使用其中的 fprintf() */
37 #include <string.h>         /* 字符串操作 */
38 #include <stdlib.h>         /* contains exit */        /* 含有 exit() */
39 #include <sys/types.h>      /* unistd.h needs this */    /* 供 unistd.h 使用 */
40 #include <sys/stat.h>       /* 文件状态信息结构 */

```

```

28 #include <linux/fs.h>                                /* 文件系统 */
29 #include <unistd.h>      /* contains read/write */ /* 含有 read()/write() */
30 #include <fcntl.h>        /* 文件操作模式符号常数 */
31
32 #define MINIX_HEADER 32      // minix 二进制模块头部长度为 32 字节。
33 #define GCC_HEADER 1024     // GCC 头部信息长度为 1024 字节。
34
35 #define SYS_SIZE 0x2000     // system 文件最长节数(字节数为 SYS_SIZE*16=128KB)。
36
37 #define DEFAULT_MAJOR_ROOT 3 // 默认根设备主设备号 - 3 (硬盘)。
38 #define DEFAULT_MINOR_ROOT 6 // 默认根设备次设备号 - 6 (第 2 个硬盘的第 1 分区)。
39
40 /* max nr of sectors of setup: don't change unless you also change
41  * bootsect etc */
42 /* 下面指定 setup 模块占的最大扇区数：不要改变该值，除非也改变 bootsect 等相应文件。
43 #define SETUP_SECTS 4      // setup 最大长度为 4 个扇区 (4*512 字节)。
44
45 #define STRINGIFY(x) #x    // 用于出错时显示语句中表示扇区数。
46
47 //// 显示出错信息，并终止程序。
48 void die(char * str)
49 {
50     fprintf(stderr, "%s\n", str);
51     exit(1);
52 }
53
54 // 显示程序使用方法，并退出。
55 void usage(void)
56 {
57     die("Usage: build bootsect setup system [rootdev] [> image]");
58 }
59
60 int main(int argc, char ** argv)
61 {
62     int i, c, id;
63     char buf[1024];
64     char major_root, minor_root;
65     struct stat sb;
66
67     // 如果程序命令行参数不是 4 或 5 个，则显示程序用法并退出。
68     if ((argc != 4) && (argc != 5))
69         usage();
70     // 如果参数是 5 个，则说明带有根设备名。
71     if (argc == 5) {
72         // 如果根设备名是软盘 ("FLOPPY")，则取该设备文件的状态信息，若出错则显示信息，退出。
73         if (strcmp(argv[4], "FLOPPY")) {
74             if (stat(argv[4], &sb)) {
75                 perror(argv[4]);
76                 die("Couldn't stat root device.");
77             }
78         }
79         // 若成功则取该设备名状态结构中的主设备号和次设备号。
80         major_root = MAJOR(sb.st_rdev);
81         minor_root = MINOR(sb.st_rdev);

```

```

74         } else {
// 否则让主设备号和次设备号取 0。
75             major_root = 0;
76             minor_root = 0;
77         }
// 若参数只有 4 个，则让主设备号和次设备号等于系统默认的根设备。
78     } else {
79         major_root = DEFAULT\_MAJOR\_ROOT;
80         minor_root = DEFAULT\_MINOR\_ROOT;
81     }
// 在标准错误终端上显示所选择的根设备主、次设备号。
82     fprintf(stderr, "Root device is (%d, %d)\n", major_root, minor_root);
// 如果主设备号不等于 2(软盘)或 3(硬盘)，也不等于 0(取系统默认根设备)，则显示出错信息，退出。
83     if ((major_root != 2) && (major_root != 3) &&
84         (major_root != 0)) {
85         fprintf(stderr, "Illegal root device (major = %d)\n",
86             major_root);
87         die("Bad root device --- major #");
88     }
// 初始化 buf 缓冲区，全置 0。
89     for (i=0; i<sizeof buf; i++) buf[i]=0;
// 以只读方式打开参数 1 指定的文件(bootsect)，若出错则显示出错信息，退出。
90     if ((id=open(argv[1], O\_RDONLY, 0))<0)
91         die("Unable to open 'boot'");
// 读取文件中的 minix 执行头部信息(参见列表后说明)，若出错则显示出错信息，退出。
92     if (read(id, buf, MINIX\_HEADER) != MINIX\_HEADER)
93         die("Unable to read header of 'boot'");
// 0x0301 - minix 头部 a_magic 魔数; 0x10 - a_flag 可执行; 0x04 - a_cpu, Intel 8086 机器码。
94     if (((long *) buf)[0]!=0x04100301)
95         die("Non-Minix header of 'boot'");
// 判断头部长字段 a_hdrlen (字节) 是否正确。(后三字节正好没有用，是 0)
96     if (((long *) buf)[1]!=MINIX\_HEADER)
97         die("Non-Minix header of 'boot'");
// 判断数据段长 a_data 字段(long)内容是否为 0。
98     if (((long *) buf)[3]!=0)
99         die("Illegal data segment in 'boot'");
// 判断堆 a_bss 字段(long)内容是否为 0。
100    if (((long *) buf)[4]!=0)
101        die("Illegal bss in 'boot'");
// 判断执行点 a_entry 字段(long)内容是否为 0。
102    if (((long *) buf)[5] != 0)
103        die("Non-Minix header of 'boot'");
// 判断符号表长字段 a_sym 的内容是否为 0。
104    if (((long *) buf)[7] != 0)
105        die("Illegal symbol table in 'boot'");
// 读取实际代码数据，应该返回读取字节数为 512 字节。
106    i=read(id, buf, sizeof buf);
107    fprintf(stderr, "Boot sector %d bytes. \n", i);
108    if (i != 512)
109        die("Boot block must be exactly 512 bytes");
// 判断 boot 块 0x510 处是否有可引导标志 0xAA55。
110    if ((*(unsigned short *) (buf+510)) != 0xAA55)
111        die("Boot block hasn't got boot flag (0xAA55)");

```



```

// 引导块的 508, 509 偏移处存放的是根设备号。
112     buf[508] = (char) minor_root;
113     buf[509] = (char) major_root;
// 将该 boot 块 512 字节的数据写到标准输出 stdout, 若写出字节数不对, 则显示出错信息, 退出。
114     i=write(1, buf, 512);
115     if (i!=512)
116         die("Write call failed");
// 最后关闭 bootsect 模块文件。
117     close (id);
118
// 现在开始处理 setup 模块。首先以只读方式打开该模块, 若出错则显示出错信息, 退出。
119     if ((id=open(argv[2], O_RDONLY, 0))<0)
120         die("Unable to open 'setup'");
// 读取该文件中的 minix 执行头部信息(32 字节), 若出错则显示出错信息, 退出。
121     if (read(id, buf, MINIX_HEADER) != MINIX_HEADER)
122         die("Unable to read header of 'setup'");
// 0x0301 - minix 头部 a_magic 魔数; 0x10 - a_flag 可执行; 0x04 - a_cpu, Intel 8086 机器码。
123     if (((long *) buf)[0]!=0x04100301)
124         die("Non-Minix header of 'setup'");
// 判断头部长度字段 a_hdrlen (字节) 是否正确。(后三字节正好没有用, 是 0)
125     if (((long *) buf)[1]!=MINIX_HEADER)
126         die("Non-Minix header of 'setup'");
// 判断数据段长 a_data 字段(long)内容是否为 0。
127     if (((long *) buf)[3]!=0)
128         die("Illegal data segment in 'setup'");
// 判断堆 a_bss 字段(long)内容是否为 0。
129     if (((long *) buf)[4]!=0)
130         die("Illegal bss in 'setup'");
// 判断执行点 a_entry 字段(long)内容是否为 0。
131     if (((long *) buf)[5] != 0)
132         die("Non-Minix header of 'setup'");
// 判断符号表长字段 a_sym 的内容是否为 0。
133     if (((long *) buf)[7] != 0)
134         die("Illegal symbol table in 'setup'");
// 读取随后的执行代码数据, 并写到标准输出 stdout。
135     for (i=0 ; (c=read(id, buf, sizeof buf))>0 ; i+=c )
136         if (write(1, buf, c)!=c)
137             die("Write call failed");
//关闭 setup 模块文件。
138     close (id);
// 若 setup 模块长度大于 4 个扇区, 则算出错, 显示出错信息, 退出。
139     if (i > SETUP_SECTS*512)
140         die("Setup exceeds " STRINGIFY(SETUP_SECTS)
141             " sectors - rewrite build/boot/setup");
// 在标准错误 stderr 显示 setup 文件的长度值。
142     fprintf(stderr, "Setup is %d bytes. \n", i);
// 将缓冲区 buf 清零。
143     for (c=0 ; c<sizeof(buf) ; c++)
144         buf[c] = '\0';
// 若 setup 长度小于 4*512 字节, 则用\0 将 setup 补足为 4*512 字节。
145     while (i<SETUP_SECTS*512) {
146         c = SETUP_SECTS*512-i;
147         if (c > sizeof(buf))

```

```

148         c = sizeof(buf);
149         if (write(1, buf, c) != c)
150             die("Write call failed");
151         i += c;
152     }
153
154     // 下面处理 system 模块。首先以只读方式打开该文件。
155     if ((id=open(argv[3], O_RDONLY, 0))<0)
156         die("Unable to open 'system'");
157     // system 模块是 GCC 格式的文件，先读取 GCC 格式的头部结构信息(linux 的执行文件也采用该格式)。
158     if (read(id, buf, GCC_HEADER) != GCC_HEADER)
159         die("Unable to read header of 'system'");
160     // 该结构中的执行代码入口点字段 a_entry 值应为 0。
161     if (((long *) buf)[5] != 0)
162         die("Non-GCC header of 'system'");
163     // 读取随后的执行代码数据，并写到标准输出 stdout。
164     for (i=0; (c=read(id, buf, sizeof buf))>0; i+=c )
165         if (write(1, buf, c)!=c)
166             die("Write call failed");
167     // 关闭 system 文件，并向 stderr 上打印 system 的字节数。
168     close(id);
169     fprintf(stderr, "System is %d bytes. \n", i);
170     // 若 system 代码数据长度超过 SYS_SIZE 节（或 128KB 字节），则显示出错信息，退出。
171     if (i > SYS_SIZE*16)
172         die("System is too big");
173     return(0);
174 }
175

```

### 12.2.3 相关信息

#### 可执行文件头部数据结构

Minix 可执行文件 a.out 的头部结构如下所示（参见 minix 2.0 源代码 01400 行开始）：

```

struct exec {
    unsigned char a_magic[2];    // 执行文件魔数。
    unsigned char a_flags;      // 标志（参见后面说明）。
    unsigned char a_cpu;        // cpu 标识号。
    unsigned char a_hdrlen;     // 头部长度。
    unsigned char a_unused;     // 保留给将来使用。
    unsigned short a_version;    // 版本信息（目前未用）。
    long          a_text;        // 代码段长度，字节数。
    long          a_data;        // 数据段长度，字节数。
    long          a_bss;         // 堆长度，字节数。
    long          a_entry;       // 执行入口点地址。
    long          a_total;       // 分配的内存总量。
    long          a_syms;        // 符号表大小。
    ...
};

```

其中标志字段定义为：

---

|        |      |                  |
|--------|------|------------------|
| A_UZP  | 0x01 | // 未映射的 0 页（页数）。 |
| A_PAL  | 0x02 | // 以页边界调整的可执行文件。 |
| A_NSYM | 0x04 | // 新型符号表。        |
| A_EXEC | 0x10 | // 可执行文件。        |
| A_SEP  | 0x20 | // 代码和数据是分开的。    |

CPU 标识号为：

|          |      |                      |
|----------|------|----------------------|
| A_NONE   | 0x00 | // 未知。               |
| A_I8086  | 0x04 | // Intel i8086/8088。 |
| A_M68K   | 0x0B | // Motorola m68000。  |
| A_NS16K  | 0x0C | // 国家半导体公司 16032。    |
| A_I80386 | 0x10 | // Intel i80386。     |
| A_SPARC  | 0x17 | // Sun 公司 SPARC。     |

GCC 执行文件头部结构信息参见 linux/include/a.out.h 文件。



## 参考文献

- [1] Intel Co. INTEL 80386 Programmer's Reference Manual 1986, INTEL CORPORATION, 1987.
- [2] James L. Turley. Advanced 80386 Programming Techniques. Osborne McGraw-Hill, 1988.
- [3] Brian W. Kernighan, Dennis M. Ritchie. The C programming Language. Prentice-Hall 1988.
- [4] Leland L. Beck. System Software: An Introduction to Systems Programming, 3<sup>rd</sup>. Addison-Wesley, 1997.
- [5] Richard M. Stallman, Using and Porting the GNU Compiler Collection, For GCC Version 2.95, the Free Software Foundation, 1998.
- [6] The Open Group Base Specifications Issue 6 IEEE Std 1003.1-2001, The IEEE and The Open Group.
- [7] David A Rusling, The Linux Kernel, 1999. <http://www.tldp.org/>
- [8] Linux Kernel Source Code, <http://www.kernel.org/>
- [9] Digital co.ltd. VT100 User Guide, <http://www.vt100.net/>
- [10] FreeBSD Online Manual, <http://www.freebsd.org/cgi/man.cgi>
- [11] Maurice J. Bach 著, 陈葆珏, 王旭, 柳纯录, 冯雪山译, UNIX 操作系统设计. 机械工业出版社, 2000 年 4 月.
- [12] John Lions 著, 尤晋元译, 莱昂氏 UNIX 源代码分析, 机械工业出版社, 2000 年 7 月.
- [13] Andrew S. Tanenbaum 著 王鹏, 尤晋元等译, 操作系统: 设计与实现 (第 2 版), 电子工业出版社, 1998 年 8 月.
- [14] Alessandro Rubini, Jonathan 著, 魏永明, 骆刚, 姜君译, Linux 设备驱动程序, 中国电力出版社, 2002 年 11 月.
- [15] Daniel P. Bovet, Marco Cesati 著, 陈莉君, 冯锐, 牛欣源 译, 深入理解 LINUX 内核, 中国电力出版社 2001 年.
- [16] 张载鸿. 微型机(PC 系列)接口控制教程, 清华大学出版社, 1992 年.
- [17] 李风华, 周利华, 赵丽松. MS-DOS 5.0 内核剖析. 西安电子科技大学出版社, 1992 年.
- [18] RedHat 7.3 操作系统在线手册. <http://www.plinux.org/cgi-bin/man.cgi>

## 附录 1 80x86 保护运行模式

80386 是一个高级的 32 位微处理器，专门用于多任务的操作系统，并为需要高性能的应用所设计。32 位的寄存器和数据通道支持 32 位的寻址方式和数据类型，处理器可以寻址最高可达 4GB 的物理内存以及 64TB ( $2^{46}$  字节) 的虚拟内存。芯片上的内存管理包括地址转换寄存器、高级多任务硬件、保护机制以及分页虚拟内存机制。下面针对系统编程，概要说明使用 80386 的这些基本原理。

设计用于系统编程的系统寄存器主要包括以下几类:

内存管理寄存器:

调试寄存器:

测试寄存器。

|    |   |   |   |   |   |   |   |    |   |   |   |    |   |   |   |  |  |   |   |    |   |   |   |   |   |   |   |   |   |   |   |   |
|----|---|---|---|---|---|---|---|----|---|---|---|----|---|---|---|--|--|---|---|----|---|---|---|---|---|---|---|---|---|---|---|---|
| 31 |   |   |   |   |   |   |   | 23 |   |   |   | 15 |   |   |   |  |  |   |   | 7  |   |   |   |   |   |   |   | 0 |   |   |   |   |
| 0  | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0  | 0 | 0 | 0 | 0  | 0 | V | R |  |  | 0 | N | IO | 0 | D | I | T | S | Z | 0 | A | 0 | P | 1 | C |
|    |   |   |   |   |   |   |   |    |   |   |   |    |   | M | F |  |  |   | T | PL | F | F | F | F | F | F |   | F |   | F |   | F |

内存管理寄存器有 4 个，用于分段内存管理：

- GDTR - 全局描述符表寄存器(Global Descriptor Table Register);
- LDTR - 局部描述符表寄存器(Local Descriptor Table Register);
- IDTR - 中断描述符表寄存器(Interrupt Descriptor Table Register);
- TR - 任务寄存器。

80386 共有 4 个控制寄存器，分别是 CR0、CR1、CR2 和 CR3。格式见下图所示。

464

|                                      |                |  |  |  |  |  |  |   |   |     |
|--------------------------------------|----------------|--|--|--|--|--|--|---|---|-----|
| 页异常线性地址<br>Page Fault Linear Address |                |  |  |  |  |  |  |   |   | CR2 |
| 保留<br>Reserved                       |                |  |  |  |  |  |  |   |   | CR1 |
| P                                    | 保留<br>Reserved |  |  |  |  |  |  | E | T | CR0 |
| G                                    |                |  |  |  |  |  |  | T | S |     |
|                                      |                |  |  |  |  |  |  | E | M |     |
|                                      |                |  |  |  |  |  |  | P | E |     |

控制寄存器 CR0 含有系统整体的控制标志。其中：

- PE - 保护模式开启位 (Protection Enable, 比特位 0)。如果设置了该比特位, 就会使处理器开始在保护模式下运行。
- MP - 协处理器存在标志 (Math Present, 比特位 1)。用于控制 WAIT 指令的功能, 以配合协处理的运行。
- EM - 仿真控制 (Emulation, 比特位 2)。指示是否需要仿真协处理器的功能。
- TS - 任务切换 (Task Switch, 比特位 3)。每当任务切换时处理器就会设置该比特位, 并且在解释协处理器指令之前测试该位。
- ET - 扩展类型 (Extention Type, 比特位 4)。该位指出了系统中所含有的协处理器类型 (是 80287 还是 80387)。
- PG - 分页操作 (Paging, 比特位 31)。该位指示出是否使用页表将线性地址变换成物理地址。

## 内存管理

内存管理主要涉及处理器的内存寻址机制。80x86 使用两步将一个分段形式的逻辑地址转换为实际物理内存地址。

- 段变换, 将一个由段选择符和段内偏移构成的逻辑地址转换为一个线性地址;
- 页变换, 将线性地址转换为对应的物理地址。该步是可选的。

在分页机制开启时, 通过将前面所述的段转换和页转换组合在一起, 即实现了从逻辑地址到物理地址的两个转换阶段。

### 段变换

下图示出了处理器是如何将一个逻辑地址转换为线性地址的。在转换过程中 CPU 使用了以下一些数据结构:

- 段描述符 (Segment Descriptors);
- 描述符表 (Descriptor tables);
- 选择符 (Selectors);
- 段寄存器 (Segment Registers)。

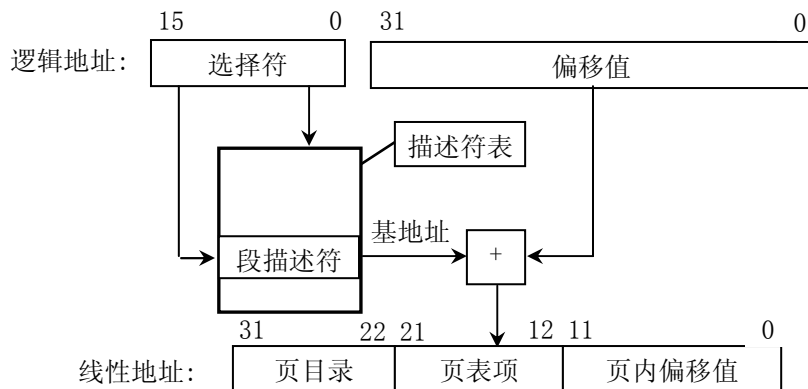


图 段变换示意图

## 段描述符

段描述符向 CPU 提供了将逻辑地址映射为线性地址所必要的信息。描述符是由程序编译器、链接器、加载器或操作系统创建的。下图示出了描述符的两种一般格式。所有种类的描述符都具有这两种格式之一。段描述符的各个字段的含义如下：

|                        |                |                        |                     |                        |
|------------------------|----------------|------------------------|---------------------|------------------------|
| 31                     | 23             | 15                     | 7                   | 0                      |
| 基地址 (BASE)<br>位 31..24 | G X 0 A<br>V L | 限长 (LIMIT)<br>位 19..16 | P DPL 1 TYPE A      | 基地址 (BASE)<br>位 23..16 |
| 段基地址 (BASE) 位 15..0    |                |                        | 段限长 (LIMIT) 位 15..0 |                        |
|                        |                |                        |                     |                        |

a. 用于程序代码段和数据段的描述符

|                        |                |                        |                     |                        |
|------------------------|----------------|------------------------|---------------------|------------------------|
| 31                     | 23             | 15                     | 7                   | 0                      |
| 基地址 (BASE)<br>位 31..24 | G X 0 A<br>V L | 限长 (LIMIT)<br>位 19..16 | P DPL 0 TYPE        | 基地址 (BASE)<br>位 23..16 |
| 段基地址 (BASE) 位 15..0    |                |                        | 段限长 (LIMIT) 位 15..0 |                        |
|                        |                |                        |                     |                        |

b. 用于特殊系统段的描述符

图 描述符的一般格式

**基地址 (BASE)：**定义段在 4GB 线性空间中的位置。处理器会将基地址的三个部分组合成一个 32 位的值。  
**段限长 (LIMIT)：**定义了段的最大长度。处理器将组合段限长的两个部分形成一个 20 位的值。处理器会依据颗粒度 (Granularity) 位字段的值来解释段限长域的实际含义：

1. 当以 1 字节为单元时，则定义了最高可为 1MB 字节的长度；
2. 当以 4KB 字节为单元时，则定义了最高可为 4GB 字节的长度。在加载时限长值将左移 12 位。

**颗粒度 (Granularity)：**指定了限长字段值代表的单元含义。当为 0 时，限长单元值为 1 字节；当该位为 1 时，限长的单元值为 4KB 字节。

**类型 (TYPE)：**用于区分各种不同类型的描述符。

**描述符特权级 (Descriptor Privilege Level - DPL)：**用于保护机制。共有 4 级：0-3。0 级是最高特权级，3 级是最低特权级。

**段存在位 (Segment-Present bit - P)：**如果该位为零，则该描述符无效，不能用于地址变换过程。当指向该描述符的选择符被加载到段寄存器中时，处理器就会发出一个异常信号。

**访问位 (Accessed bit - A)：**当处理器访问过该段时就会设置该比特位。

## 描述符表

段描述符是保存在描述符表中的，有两类描述符表：

- 全局描述符表 (Global descriptor table - GDT)；
- 局部描述符表 (Local descriptor table - LDT)。



描述符表是由 8 字节构成的描述符项的内存中的一个数组，见下图所示。描述符表的长度是可变的，最多可以含有 8192 ( $2^{13}$ ) 个描述符。但是对于 GDT 表，其第一个描述符（索引 0）是不用的。

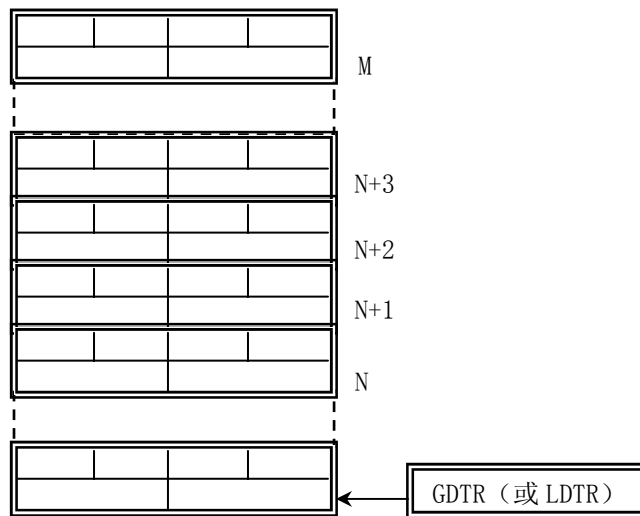
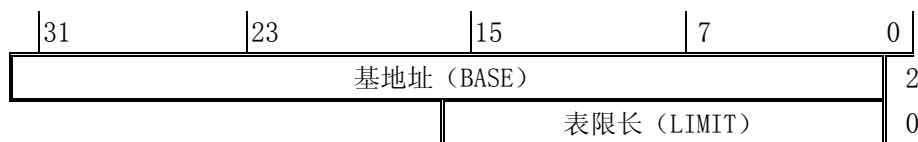


图 描述符表示意图

处理器是通过使用 GDTR 和 LDTR 寄存器来定位 GDT 表和当前的 LDT 表。这两个寄存器以线性地址的方式保存了描述符表的基地址和表的长度。指令 lgdt 和 sgdt 用于访问 GDTR 寄存器；指令 llidt 和 slidt 用于访问 LDTR 寄存器。lgdt 使用的是内存中一个 6 字节操作数来加载 GDTR 寄存器的。头两个字节代表描述符表的长度，后 4 个字节是描述符表的基地址。然而请注意，访问 LDTR 寄存器的指令 llidt 所使用的操作数却是一个 2 字节的操作数，表示全局描述符表 GDT 中一个描述符项的选择符。该选择符所对应的 GDT 表中的描述符项应该对应一个局部描述符表。选择符的含义见下面说明。



## 选择符(Selectors)

逻辑地址的选择符部分是用于指定一描述符的，它是通过指定一描述符表并且索引其中的一个描述符项完成的。下图示出了选择符的格式。各字段的含义为：

索引值(Index)：用于选择指定描述符表中 8192 个描述符中的一个。处理器将该索引值乘上 8(描述符的字节长度)，并加上描述符表的基地址即可访问表中指定的段描述符。

表指示器(Table Indicator - TI)：指定选择符所引用的描述符表。值为 0 表示指定 GDT 表，值为 1 表示指定当前的 LDT 表。

请求者的特权级(Requestor's Privilege Level - RPL)：用于保护机制。



由于 GDT 表的第一项(索引值为 0)没有被使用，因此一个具有索引值 0 和表指示器值也为 0 的选择符

(也即指向 GDT 的第一项的选择符) 可以用作为一个空(null)选择符。当一个段寄存器(不能是 CS 或 SS)加载了一个空选择符时, 处理器并不会产生一个异常。但是若使用这个段寄存器访问内存时就会产生一个异常。对于初始化还未使用的段寄存器以陷入意外的引用来说, 这个特性是很有用的。

## 段寄存器

处理器将描述符中的信息保存在段寄存器中, 因而可以避免在每次访问内存时查询描述符表。

每个段寄存器都有一个“可见”部分和一个“不可见”部分, 见下图所示。这些段地址寄存器的可见部分是由程序来操作的, 就好象它们只是简单的 16 位寄存器。不可见部分则是由处理器来处理的。

对这些寄存器的加载操作使用的是普通程序指令, 这些指令可以分为两类:

1. 直接加载指令; 例如, MOV, POP, LDS, LSS, LGS, LFS。这些指令显式地引用了指定的段寄存器。
2. 隐式加载指令; 例如, 远调用 CALL 和远跳转 JMP。这些指令隐式地引用了 CS 段寄存器, 并用新值加载到 CS 中。

程序使用这些指令会把 16 位的选择符加载到段寄存器的可见部分, 而处理器则会自动地从描述符表中将一个描述符的基地址、段限长、类型以及其它信息加载到段寄存器中的不可见部分中去。

|    | 16 位可见部分 | 隐藏部分 |
|----|----------|------|
| CS |          |      |
| SS |          |      |
| DS |          |      |
| ES |          |      |
| FS |          |      |
| GS |          |      |

## 页变换(翻译)

在地址变换的第二阶段, CPU 将线性地址转换为物理地址。地址变换的这个阶段实现了基于分页的虚拟内存系统和分页级保护的基本功能。页变化这一步是可选的, 页变换仅在设置了 CR0 的 PG 比特位后才起作用, 该比特位是在软件初始化时由操作系统设置的。如果操作系统需要实现多个虚拟 8086 任务、基于分页的保护机制或基于分页的虚拟内存, 那么就一定要设置该位。

### 页框(帧)(Page Frame)

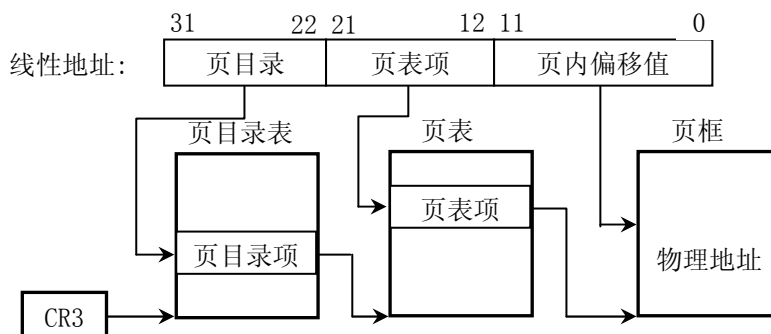
页框是一个物理内存地址连续的 4K 字节单元。它以字节为边界, 大小固定。

### 线性地址(Linear Address)

线性地址通过指定一个页表、页表中的某一页以及该页中的偏移值, 从而间接地指向对应的物理地址。下图示出了线性地址的格式。

|              |    |             |    |                 |   |
|--------------|----|-------------|----|-----------------|---|
| 31           | 22 | 21          | 12 | 11              | 0 |
| 页目录<br>(DIR) |    | 页<br>(PAGE) |    | 偏移值<br>(OFFSET) |   |

下图示出了处理器将一个线性地址转换成物理地址的方法。通过使用两级页表, 处理器将一个线性地址的页目录字段(DIR)、页字段(PAGE)和偏移字段(OFFSET)翻译成对应的物理地址。寻址机制使用线性地址的页目录字段作为页目录中的索引值、使用页表字段作为页目录所指定页表中的索引值、使用偏移字段作为页表所确定的内存页中的字节偏移值。



## 页表(Page Table)

页表只是一个简单的 32 位页指示器的数组。页表本身也是一页内存，因此它含有 4K 字节的内存，可容纳 1K 个 32 位的项。

这里使用了两级页表来定位一页内存页。最高层是页目录，页目录可定位最多 1K 个第二级页表，而每个二级页表可以定位最多 1K 内存页。因此，一个页目录定位的所有页表可以寻址 1M 内存页 ( $2^{20}$ )。由于每一页内存含有 4K 字节 ( $2^{12}$ )，最终一个页目录所指定的页表可以寻址 80386 的整个物理地址空间 ( $2^{20} * 2^{12} = 2^{32}$ )。

当前页目录的物理地址是存储在 CPU 控制寄存器 CR3 中的，因此该寄存器也被称为页目录基地址寄存器 (page directory base register - PDBR)。内存管理软件可以选择对所有的任务只使用一个页目录，或每个任务使用一个页目录，也可以组合两个任务使用一个页目录。

## 页表项(Page-Table Entries)

各级页表所使用的页表项是相同的，其格式见下图所示。

|                                       |    |               |   |   |   |   |   |   |          |          |   |
|---------------------------------------|----|---------------|---|---|---|---|---|---|----------|----------|---|
| 31                                    | 12 | 11            | 0 |   |   |   |   |   |          |          |   |
| 页框地址 位 31..12<br>(PAGE FRAME ADDRESS) |    | 可用<br>(AVAIL) | 0 | 0 | D | A | 0 | 0 | U<br>/ S | R<br>/ W | P |

其中，页框地址 (PAGE FRAME ADDRESS) 指定了一页内存的物理起始地址。因为内存页是位于 4K 边界上的，所以其低 12 比特总是 0。在一个页目录中，页表项的页框地址是一个页表的起始地址；在第二级页表中，页表项的页框地址是包含期望内存操作的页框的地址。

存在位 (PRESENT - P) 确定了一个页表项是否可以用于地址转换过程。P=1 表示该项可用。当目录表项或第二级表项的 P=0 时，则该表项时无效的，不能用于地址转换过程。此时该表项的其它所有比特位都可供程序使用；处理器不对这些位进行测试。

当 CPU 试图使用一个页表项进行地址转换时，如果此时任意一级页表项的 P=0，则处理器就会发出页异常信号。对于支持分页虚拟内存的软件系统中，页不存在 (page-not-present) 异常处理程序就可以把所请求的页加入到物理内存中。此时导致异常的指令就可以被重新执行。

已访问 (Accessed - A) 和已修改 (Dirty - D) 比特位提供了有关页使用的信息。除了页目录项中的已修改位，这些比特位将由硬件置位，但不复位。

在对一页内存进行读或写操作之前，处理器将设置相关的目录和二级页表项的已访问位。在向一个二级页表项所涵盖的地址进行写操作之前，处理器将设置该二级页表项的已修改位，而页目录项中的已修改位是不用的。当需求的内存超出实际物理内存容量时，支持分页虚拟内存的操作系统可以使用这些位来确定

那些页可以从内存中取走。操作系统必须负责检测和复位这些比特位。

读/写位 (Read/Write - R/W) 和用户/超级用户位 (User/Supervisor - U/S) 并不用于地址转换，但用于分页级的保护机制，是由处理器在地址转换过程中同时操作的。

### 页转换高速缓冲

为了最大地提高地址转换的效率，处理器将最近所使用的页表数据存放在芯片上的高速缓冲中。操作系统设计人员必须在当前页表改变时刷新高速缓冲，可使用以下两种方式之一：

1. 通过使用 MOV 指令重新加载 CR3 页目录基址寄存器；
2. 通过执行一个任务切换。

## 多任务 (Multitasking)

为了提供有效的、受保护的多任务机制，80x86 使用了一些特殊的数据结构。支持多任务运行的寄存器和数据结构主要有任务状态段 (Task State Segment) 和任务寄存器 (Task register)。使用这些数据结构，CPU 可以快速地从一个任务的执行切换到另一个任务，并保存原有任务的内容。

### 任务状态段 (Task State Segment - TSS)

处理器管理一个任务的所有信息存储在一个特殊类型的段中，即任务状态段 TSS。下图给出了 TSS 的格式。其中的字段可分为两类：

- 处理器只读其中信息的静态字段集 (图中灰色部分)；
- 每次任务切换时处理器将会更新的动态字段集。

| 31                              | 23 | 15                              | 7 | 0 |  |  |  |  |  |  |  |  |  |    |
|---------------------------------|----|---------------------------------|---|---|--|--|--|--|--|--|--|--|--|----|
| I/O 映射图基地址 (MAP BASE)           |    | 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 |   |   |  |  |  |  |  |  |  |  |  | 64 |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 |    | 局部描述符表 (LDT) 的选择符               |   |   |  |  |  |  |  |  |  |  |  | 60 |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 |    | GS                              |   |   |  |  |  |  |  |  |  |  |  | 5C |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 |    | FS                              |   |   |  |  |  |  |  |  |  |  |  | 58 |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 |    | DS                              |   |   |  |  |  |  |  |  |  |  |  | 54 |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 |    | SS                              |   |   |  |  |  |  |  |  |  |  |  | 50 |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 |    | CS                              |   |   |  |  |  |  |  |  |  |  |  | 4C |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 |    | ES                              |   |   |  |  |  |  |  |  |  |  |  | 48 |
| EDI                             |    |                                 |   |   |  |  |  |  |  |  |  |  |  | 44 |
| ESI                             |    |                                 |   |   |  |  |  |  |  |  |  |  |  | 40 |
| EBP                             |    |                                 |   |   |  |  |  |  |  |  |  |  |  | 3C |
| ESP                             |    |                                 |   |   |  |  |  |  |  |  |  |  |  | 38 |
| EBX                             |    |                                 |   |   |  |  |  |  |  |  |  |  |  | 34 |
| EDX                             |    |                                 |   |   |  |  |  |  |  |  |  |  |  | 30 |
| ECX                             |    |                                 |   |   |  |  |  |  |  |  |  |  |  | 2C |
| EAX                             |    |                                 |   |   |  |  |  |  |  |  |  |  |  | 28 |
| EFLAGS                          |    |                                 |   |   |  |  |  |  |  |  |  |  |  | 24 |
| 指令指针 (EIP)                      |    |                                 |   |   |  |  |  |  |  |  |  |  |  | 20 |
| 页目录基地址寄存器 CR3 (PDBR)            |    |                                 |   |   |  |  |  |  |  |  |  |  |  | 1C |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 |    | SS2                             |   |   |  |  |  |  |  |  |  |  |  | 18 |
| ESP2                            |    |                                 |   |   |  |  |  |  |  |  |  |  |  | 14 |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 |    | SS1                             |   |   |  |  |  |  |  |  |  |  |  | 10 |

|                                 |  |  |  |  |  |  |  |  |  |                 |  |  |  |  |  |    |
|---------------------------------|--|--|--|--|--|--|--|--|--|-----------------|--|--|--|--|--|----|
| ESP1                            |  |  |  |  |  |  |  |  |  |                 |  |  |  |  |  | 0C |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 |  |  |  |  |  |  |  |  |  | SS0             |  |  |  |  |  | 08 |
| ESP0                            |  |  |  |  |  |  |  |  |  |                 |  |  |  |  |  | 04 |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 |  |  |  |  |  |  |  |  |  | 前一执行任务 TSS 的描述符 |  |  |  |  |  | 00 |

任务状态段 TSS 可以处于线性空间的任何位置。TSS 与其它段一样，也是使用段描述符来定义的。访问 TSS 的描述符会导致任务切换。因此，在大多数系统中都将描述符的 DPL（描述符特权级）字段设置为最高特权级 0，这样就可以只允许可信的软件执行任务的切换。TSS 的描述符只能放在全局描述符表 GDT 中。

## 任务寄存器

任务寄存器（Task Register - TR）的作用与一般段寄存器的类似，它通过指向 TSS 来确定当前执行的任务。它也有 16 位的可见部分和不可见部分。可见部分中的选择符用于在 GDT 表中选择一个 TSS 描述符，处理器使用不可见部分来存放描述符中的基地址和段限长值。指令 LTR 和 STR 用于修改和读取任务寄存器中的可见部分，指令所使用的操作数是一 16 位的选择符。

另外，还有一种提供对 TSS 间接、受保护引用的任务门描述符（Task Gate Descriptor）。这种描述符是在一般段描述符格式的基地址位 15..0 字段（第 3、4 字节）中存放的是一个 TSS 描述符的选择符，并利用其中的特权级字段（DPL）来控制使用描述符执行任务切换的权限。见下面有关中断描述符表 IDT 描述符中的说明。

在以下 4 种情况下，CPU 会切换执行的任务：

1. 当前任务执行了一条引用 TSS 描述符的 JMP 或 CALL 指令；
2. 当前任务执行了一条引用任务门的 JMP 或 CALL 指令；
3. 引用了中断描述符表（IDT）中任务门的中断或异常；
4. 当嵌套任务标志 NT 置位时，当前任务执行了一个 IRET 指令。

## 中断和异常

中断和异常是一种特殊类型的控制转换。它们改变了正常程序流而去处理其它的事件（例如外部事件、出错报告或异常条件）。中断与异常的主要区别在于中断常用于处理 CPU 外部的异步事件，而异常则是处理 CPU 在执行过程中本身检测到的问题。

外部中断源有两种：由 CPU 的 INTR 引脚输入的可屏蔽中断和 NMI 引脚输入的不可屏蔽中断。同样，异常也有两类：由 CPU 检测到的出错、陷阱或放弃事件以及编程设置的“软中断”（如 INT 3 指令等）。

处理器使用标识号（中断号）来识别每种类型的中断或异常。处理器所能识别的不可屏蔽中断 NMI 和异常的标识号是预先确定的，范围是 0 到 31（0x00-0x1f）。目前这些号码并没有全都使用，未确定的号码由 Intel 公司留作今后使用。

可屏蔽中断的标识号由外部中断控制器（如 8259A 可编程中断控制器）确定，并在 CPU 的中断识别阶段通知 CPU。8259A 所分配的中断号可以通过编程指定，可使用的标识号范围是 32 到 255（0x20-0xff）。Linux 系统将 32-47 分配给了可屏蔽中断，余下的 48-255 用来标识其它软中断。当 Linux 只使用了号码 128（0x80）作为系统调用的中断向量号。

## 中断描述符表

中断描述符表（Interrupt Descriptor Table - IDT）将每个中断或异常标识号与处理相应事件程序指令的一个描述符相关联。与 GDT 和 LDT 相似，IDT 是一个 8 字节描述符数组，但其第 1 项可以含有一个描述符。处理器通过将中断号异常号乘上 8 即可索引 IDT 中对应的描述符。IDT 可以位于物理内存的任何地方。处理器是使用 IDT 寄存器（IDTR）来定位 IDT 的。修改和复制 IDT 的指令是 LIDT 和 SIDT。与 GDT

表的操作一样，IDT 也是使用 6 字节数据的内存地址作为操作数的。前两个字节表示表的限长，后 4 个字节是表的线性基地址。

## IDT 描述符

在中断描述符表 IDT 中可以含有三类描述符中的任意一种：

- 任务门 (Task gates);
- 中断门 (Interrupt gates);
- 陷阱门 (Trap gates);

下图给出了任务门、中断门和陷阱门描述符的格式。

|                     |    |       |     |           |
|---------------------|----|-------|-----|-----------|
| 31                  | 23 | 15    | 7   | 0         |
| (未使用)               |    | P     | DPL | 0 0 1 0 1 |
| TSS 段选择符 (SELECTOR) |    | (未使用) |     |           |

80X86 任务门描述符

|                      |    |                      |     |           |
|----------------------|----|----------------------|-----|-----------|
| 31                   | 23 | 15                   | 7   | 0         |
| 偏移值(OFFSET) 位 31..16 |    | P                    | DPL | 0 1 1 1 0 |
| 段选择符 (SELECTOR)      |    | 偏移值 (OFFSET) 位 15..0 |     |           |

80X86 中断门描述符

|                      |    |                      |     |           |
|----------------------|----|----------------------|-----|-----------|
| 31                   | 23 | 15                   | 7   | 0         |
| 偏移值(OFFSET) 位 31..16 |    | P                    | DPL | 0 1 1 1 1 |
| 段选择符 (SELECTOR)      |    | 偏移值 (OFFSET) 位 15..0 |     |           |

80X86 陷阱门描述符

## 中断任务和中断过程

正如 CALL 指令能调用一个过程或任务一样，一个中断或异常也能“调用”中断处理程序，该程序是一个过程或一个任务。当响应一个中断或异常时，CPU 使用中断或异常的标识号来索引 IDT 表中的描述符。如果 CPU 索引到一个中断门或陷阱门时，它就调用处理过程；如果是一个任务门，它就引起任务切换。

中断门或陷阱门间接地指向一个过程，该过程将在当前执行任务上下文中执行。门描述符中的段选择符指向 GDT 或当前 LDT 中的一个可执行段的描述符。门描述符中的偏移字段值指向中断或异常处理过程的开始处。

80X86 执行一个中断或异常处理过程的方式与 CALL 指令调用一个过程的方式非常相似，只是两者在使用堆栈上略有不同。中断会在把原指令指针压入堆栈之前，把原标志寄存器 EFLAGS 的内容也推入堆栈中。对于与段有关的异常，CPU 还会将一个错误码压入异常处理程序的堆栈上。

对于中断过程处理结束的返回操作，中断返回指令 IRET 与 RET 相似，但是 IRET 为了去除压入堆栈的 EFLAGS 值，ESP 会多递增 4 个字节。

中断门与陷阱门的区别在于对中断允许标志 IF 的影响。由中断门向量引起的中断会复位 IF，因为可以避免其它中断干扰当前中断的处理。随后的 IRET 指令会从堆栈上恢复 IF 的原值；而通过陷阱门产生的中断不会改变 IF。

IDT 表中的任务门描述符间接地指向一个任务状态段 TSS。任务门描述符中的段选择符指向 GDT 表中的一个 TSS 描述符。当产生的中断或异常指向 IDT 中的一个任务门描述符，就会导致任务切换，从而会在独立的任务中处理中断。Linux 系统中并没有使用任务门描述符。

# 索引