

Linux

内核完全剖析

赵 炯 编著



 机械工业出版社
CHINA MACHINE PRESS

Linux 内核完全剖析

赵 炯 编著



机械工业出版社

本书对早期 Linux 操作系统内核(v0.11)全部代码文件进行了详细的剖析,旨在让读者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。虽然选择的版本较低,但该内核已能够正常编译运行,并且其中已包括了 Linux 工作原理的精髓。书中首先以 Linux 源代码版本的变迁为主线,简要介绍了 Linux 系统的发展历史,同时着重说明了各个内核版本之间的主要区别和改进方面,给出了选择 0.11 版内核源代码作为研究对象的原因。然后概要介绍了 PC 机的硬件组成结构、编制内核使用的汇编语言和 C 语言扩展部分,并且重点说明了 80x86 处理器在保护模式下运行的编程方法。接着详细介绍了 Linux 内核源代码目录树组织结构,并依据该组织结构对所有内核程序和文件进行了注释和详细说明。为了加深读者对内核工作原理的理解,书中最后一章给出了围绕 Linux 0.11 系统的多个试验。试验中使用的相关程序均可从本书配套网站(www.oldlinux.org)上下载。

本书适合作为高校计算机专业学生学习操作系统课程的辅助和实践教材,也适合 Linux 爱好者作为学习内核工作原理的自学参考书籍,还可供一般技术人员作为开发嵌入式系统的参考书使用。

图书在版编目(CIP)数据

Linux 内核完全剖析/赵炯编著. —北京:机械工业出版社,2006.1

ISBN 7-111-18032-1

I. L... II. 赵... III. Linux 操作系统
IV. TP316.89

中国版本图书馆 CIP 数据核字(2005)第 145021 号

机械工业出版社(北京市百万庄大街 22 号 邮政编码 100037)

策 划:胡毓坚

责任编辑:车 忱

责任印制:石 冉

北京中兴印刷有限公司印刷

2006 年 1 月第 1 版·第 1 次印刷

787mm×1092mm $\frac{1}{16}$ ·56 印张·1466 千字

0 001—3 000 册

定价:79.00 元

凡购本图书,如有缺页、倒页、脱页,由本社发行部调换

本社购书热线电话:(010) 68326294

封面无防伪标均为盗版

序

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

本书的主要目标是用尽量少的篇幅,对完整的 Linux 内核源代码进行解剖,使读者对操作系统的基本功能和实际实现方式获得全面的理解。

本书读者应是一些知晓 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。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 Linux 7.0 使用的 2.2.16 稳定版等)进行描述,但由于目前 Linux 内核整个源代码已经非常大(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或原理性的说明,许多系统实现细节被忽略。所以并不能使读者对实际 Linux 内核有清晰而完整的理解。

Scott Maxwell 的《Linux 内核源代码分析》基本上是对 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。

John Lions 著的《莱昂氏 UNIX 源代码分析》一书虽然是学习 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废弃的 PDP-11 系列机的汇编语言编制的,因此在阅读和理解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。

国内出版的部分 Linux 源代码分析的图书也有类似缺点。

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,尤其是那些 Linux 系统初学者或刚学会如何使用 Linux 系统的人在使用这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。本人在多年的 Linux 内核学习过程中对此也深有体会。而现今流行的 Linux 系统非常庞大和复杂,因此不适合作为操作系统初学者的入门学习起点。这也是作者基于 Linux 早期内核版本写作本书的动机之一。

阅读早期内核的其他好处

目前,已经出现了不少基于 Linux 早期内核而开发的专门用于嵌入式系统的内核版本,如 μ Clinux 等(在 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出版类似本书的书籍。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方法,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. Beck 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,是理想的选择之一。

对于那些已经比较熟悉内核工作原理的人,为了能让自己在实际工作中对系统的实际运转机制不产生一种空中楼阁的感觉,因此也有必要阅读内核源代码。

当然,使用早期内核作为学习的对象也有不足之处,例如早期内核不支持虚拟文件系统 VFS,也不支持网络系统。但由于本书是 Linux 内核的入门教材,简单也正是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。

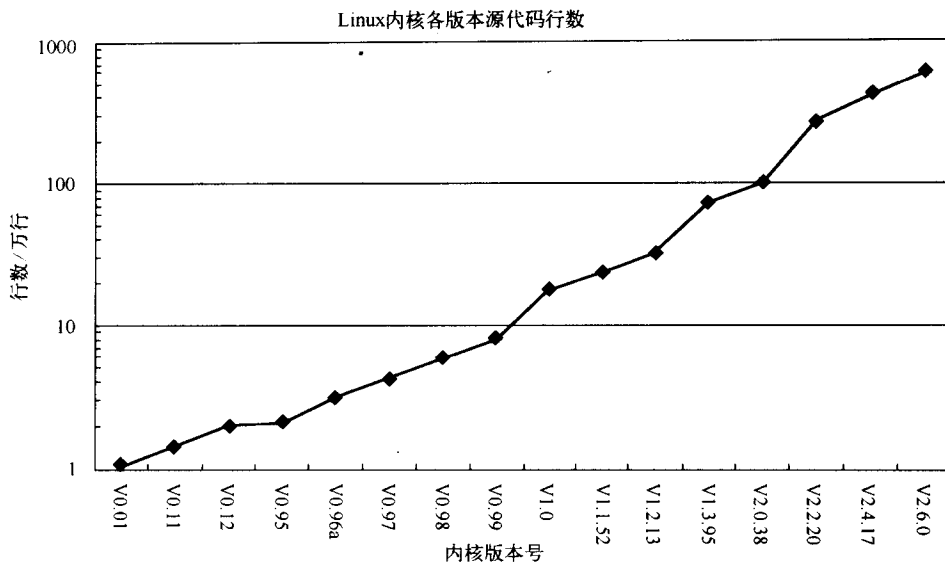
阅读完整源代码的重要性和必要性

正如 Linus 在一篇新闻组投稿上所说的,要理解一个操作系统的真正运行机制,一定要阅读其源代码。操作系统本身是一个整体,具有很多看似不重要的细节,但是若忽略这些细节,就不能真正了解一个实际操作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了解,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 Tanenbaum 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,才会对系统有一种豁然开朗的感觉,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,选择一个合适的 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 内核源代码量都有几百万行,2.6.0 版内核代码行数约为 592 万行,极其庞大,对这些版本进行完全剖析几乎是不可能的。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剖析清楚。麻雀虽小,五脏俱全。为了对所研究的系统有感性的了解,并能利用实验来加深对原理的理解,作者还专门重建了基于该内核的可运行的 Linux 0.11 系统。由于其中含有 GNU gcc 编译环境,因此使用该系统也能做一些简单的开发工作。

另外,使用该版本可以避免对现有较新内核版本中已经变得越来越复杂的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,读者必须具备一些基本的 C 语言知识和 Intel CPU 汇编语言知识。有关 C 语言最佳的参考资料仍然是 Brian W. Kernighan 和 Dennis M. Ritchie 编写的《The C Programming Language》一书。而汇编语言的资料则可以参考任意一本讲解与 Intel CPU 相关的汇编语言教材。另外还需要一些嵌入式汇编语言的知识。有关嵌入式汇编的权威信息都包含在 GNU gcc 编译器手册中。我们也可以从 Internet 上搜索到一些有关嵌入式汇编比较有价值的短文。本书中也包含了一些嵌入汇编基本的语法说明。

除此之外,还希望读者具备有关 80x86 处理器结构和编程的知识以及有关 80x86 硬件体系结构和接口编程的知识;读者还应具备使用 Linux 系统的简单技能。

另外,由于 Linux 系统内核最早是根据 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,会更易于理解内核源代码。

Linus 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,Tanenbaum 的书《操作系统:设计与实现》也具有较大的参考价值。

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。其

功能详尽并且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。Linux 内核 0.11 版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 Linux 内核中各部分的工作原理和过程了。

同济大学
赵炯 博士

目 录

序	
第 1 章 概述	1
1.1 Linux 的诞生和发展	1
1.1.1 UNIX 操作系统的诞生	1
1.1.2 MINIX 操作系统	1
1.1.3 GNU 计划	2
1.1.4 POSIX 标准	2
1.1.5 Linux 操作系统的诞生	3
1.1.6 Linux 操作系统版本的变迁	5
1.1.7 Linux 名称的由来	7
1.1.8 早期 Linux 系统开发的主要贡献者	7
1.2 内容综述	8
第 2 章 微型计算机组成结构	13
2.1 微型计算机组成原理	13
2.2 I/O 端口寻址和访问控制方式	15
2.2.1 I/O 端口和寻址	15
2.2.2 接口访问控制	17
2.3 主存储器、BIOS 和 CMOS 存储器	17
2.3.1 主存储器	17
2.3.2 基本输入/输出程序 BIOS	18
2.3.3 CMOS 存储器	19
2.4 控制器和控制卡	19
2.4.1 中断控制器	19
2.4.2 DMA 控制器	21
2.4.3 定时/计数器	21
2.4.4 键盘控制器	21
2.4.5 串行控制卡	23
2.4.6 显示控制	24
2.4.7 软盘和硬盘控制器	26
第 3 章 内核编程语言和环境	28
3.1 as86 汇编器	28
3.1.1 as86 汇编语言语法	29
3.1.2 as86 汇编语言程序	30
3.1.3 as86 汇编语言程序的编译和链接	32
3.1.4 as86 和 ld86 使用方法和选项	34
3.2 GNU as 汇编	34
3.2.1 编译 as 汇编语言程序	35
3.2.2 as 汇编语法	36
3.2.3 指令语句、操作数和寻址	38
3.2.4 区与重定位	40
3.2.5 符号	42
3.2.6 as 汇编命令	43
3.2.7 编写 16 位代码	45
3.2.8 as 汇编器命令行选项	45
3.3 C 语言程序	45
3.3.1 C 程序编译和链接	45
3.3.2 嵌入汇编	46
3.3.3 圆括号中的组合语句	50
3.3.4 寄存器变量	50
3.3.5 内联函数	51
3.4 C 与汇编程序的相互调用	53
3.4.1 C 函数调用机制	53
3.4.2 在汇编程序中调用 C 函数	58
3.4.3 在 C 程序中调用汇编函数	60
3.5 Linux 0.11 目标文件格式	62
3.5.1 目标文件格式	62
3.5.2 Linux 0.11 中的目标文件格式	65
3.5.3 链接程序输出	67
3.5.4 链接程序预定义变量	68
3.5.5 System.map 文件	70
3.6 Make 程序和 Makefile 文件	71
第 4 章 80x86 保护模式及其编程	74
4.1 80x86 系统寄存器和系统指令	74
4.1.1 标志寄存器	74
4.1.2 内存管理寄存器	75
4.1.3 控制寄存器	76
4.1.4 系统指令	79
4.2 保护模式内存管理	80
4.2.1 内存寻址	80
4.2.2 地址变换	80
4.2.3 保护	83
4.3 分段机制	84
4.3.1 段的定义	84
4.3.2 段描述符表	86
4.3.3 段选择符	88

4.3.4 段描述符	90	5.3 Linux 内核对内存的管理和使用 ...	143
4.3.5 代码和数据段描述符类型	93	5.3.1 物理内存	143
4.3.6 系统描述符类型	94	5.3.2 内存地址空间概念	144
4.4 分页机制	95	5.3.3 内存分段机制	145
4.4.1 页表结构	96	5.3.4 内存分页管理	147
4.4.2 页表项格式	98	5.3.5 CPU 多任务和保护方式	150
4.4.3 虚拟存储	98	5.3.6 虚拟地址、线性地址和物理地址 之间的关系	150
4.5 保护	99	5.3.7 用户申请内存的动态分配	154
4.5.1 段级保护	99	5.4 Linux 系统的中断机制	155
4.5.2 访问数据段时的特权级检查	102	5.4.1 中断操作原理	155
4.5.3 代码段之间转移控制时的特权级 检查	103	5.4.2 80x86 微机的中断子系统	156
4.5.4 页级保护	108	5.4.3 中断向量表	157
4.5.5 组合页级和段级保护	109	5.4.4 Linux 内核的中断处理	157
4.6 中断和异常处理	110	5.4.5 标志寄存器的中断标志	159
4.6.1 异常和中断向量	110	5.5 Linux 的系统调用	159
4.6.2 中断源和异常源	110	5.5.1 系统调用接口	159
4.6.3 异常分类	111	5.5.2 系统调用处理过程	160
4.6.4 程序或任务的重新执行	112	5.5.3 Linux 系统调用的参数传递方式 ...	161
4.6.5 开启和禁止中断	112	5.6 系统时间和定时	161
4.6.6 异常和中断的优先级	112	5.6.1 系统时间	161
4.6.7 中断描述符表	113	5.6.2 系统定时	162
4.6.8 IDT 描述符	114	5.7 Linux 进程控制	162
4.6.9 异常与中断处理	115	5.7.1 任务数据结构	163
4.6.10 中断处理任务	117	5.7.2 进程运行状态	167
4.6.11 错误码	118	5.7.3 进程初始化	168
4.7 任务管理	119	5.7.4 创建新进程	169
4.7.1 任务的结构和状态	119	5.7.5 进程调度	170
4.7.2 任务的执行	120	5.7.6 终止进程	171
4.7.3 任务管理数据结构	120	5.8 Linux 系统中堆栈的使用方法 ...	172
4.7.4 任务切换	123	5.8.1 初始化阶段	172
4.7.5 任务链	125	5.8.2 任务的堆栈	173
4.7.6 任务地址空间	126	5.8.3 任务内核态堆栈与用户态堆栈 之间的切换	175
4.8 保护模式编程初始化	127	5.9 Linux 0.11 采用的文件系统	176
4.8.1 进入保护模式时的初始化操作 ...	127	5.10 Linux 内核源代码的目录结构 ...	176
4.8.2 模式切换	129	5.10.1 内核主目录 linux	177
4.9 一个简单的多任务内核实例	130	5.10.2 引导启动程序目录 boot	177
4.9.1 多任务程序结构和工作原理	131	5.10.3 文件系统目录 fs	177
4.9.2 引导启动程序 boot.s	134	5.10.4 头文件主目录 include	179
4.9.3 多任务内核程序 head.s	135	5.10.5 内核初始化程序目录 init	180
第 5 章 Linux 内核体系结构	141	5.10.6 内核程序主目录 kernel	180
5.1 Linux 内核模式	141	5.10.7 内核库函数目录 lib	183
5.2 Linux 内核系统体系结构	142		

5.10.8 内存管理程序目录 mm	183	8.5.2 代码注释	264
5.10.9 编译内核工具程序目录 tools ...	183	8.5.3 其他信息	272
5.11 内核系统与应用程序的关系 ...	183	8.6 mktime.c 程序	274
5.12 linux/Makefile 文件	184	8.6.1 功能描述	274
5.12.1 功能描述	184	8.6.2 代码注释	274
5.12.2 代码注释	185	8.6.3 闰年的计算方法	276
第 6 章 引导启动程序	190	8.7 sched.c 程序	276
6.1 总体功能	190	8.7.1 功能描述	276
6.2 bootsect.s 程序	192	8.7.2 代码注释	278
6.2.1 功能描述	192	8.7.3 其他信息	290
6.2.2 代码注释	192	8.8 signal.c 程序	295
6.2.3 其他信息	200	8.8.1 功能描述	295
6.3 setup.s 程序	201	8.8.2 代码注释	301
6.3.1 功能描述	201	8.8.3 进程信号说明	305
6.3.2 代码注释	203	8.9 exit.c 程序	306
6.3.3 其他信息	210	8.9.1 功能描述	306
6.4 head.s 程序	219	8.9.2 代码注释	307
6.4.1 功能描述	219	8.10 fork.c 程序	312
6.4.2 代码注释	221	8.10.1 功能描述	312
6.4.3 其他信息	228	8.10.2 代码注释	314
第 7 章 初始化程序	232	8.10.3 任务状态段(TSS)信息	318
7.1 main.c 程序	232	8.11 sys.c 程序	320
7.1.1 功能描述	232	8.11.1 功能描述	320
7.1.2 代码注释	235	8.11.2 代码注释	321
7.1.3 其他信息	242	8.12 vsprintf.c 程序	326
7.2 环境初始化工作	244	8.12.1 功能描述	326
第 8 章 内核代码	246	8.12.2 代码注释	327
8.1 总体功能	246	8.12.3 vsprintf()的格式字符串	332
8.1.1 中断处理程序	247	8.12.4 与当前版本的区别	334
8.1.2 系统调用处理相关程序	248	8.13 printk.c 程序	334
8.1.3 其他通用类程序	248	8.13.1 功能描述	334
8.2 Makefile 文件	249	8.13.2 代码注释	335
8.2.1 功能简介	249	8.14 panic.c 程序	336
8.2.2 文件注释	249	8.14.1 功能描述	336
8.3 asm.s 程序	251	8.14.2 代码注释	336
8.3.1 功能描述	251	第 9 章 块设备驱动程序	338
8.3.2 代码注释	252	9.1 总体功能	339
8.3.3 Intel 保留中断向量的定义	256	9.1.1 块设备请求项和请求队列	339
8.4 traps.c 程序	257	9.1.2 块设备操作方式	341
8.4.1 功能描述	257	9.2 Makefile 文件	342
8.4.2 代码注释	257	9.2.1 功能描述	342
8.5 system_call.s 程序	262	9.2.2 代码注释	342
8.5.1 功能描述	262	9.3 blk.h 文件	344

9.3.1 功能描述	344	10.7.1 功能描述	478
9.3.2 代码注释	345	10.7.2 代码注释	479
9.4 hd.c 程序	348	10.8 tty_ioctl.c 程序	489
9.4.1 功能描述	348	10.8.1 功能描述	489
9.4.2 代码注释	351	10.8.2 代码注释	490
9.4.3 其他信息	362	10.8.3 波特率与波特率因子	495
9.5 ll_rw_blk.c 程序	371	第 11 章 数学协处理器	497
9.5.1 功能描述	371	11.1 Makefile 文件	497
9.5.2 代码注释	371	11.1.1 功能描述	497
9.6 ramdisk.c 程序	376	11.1.2 代码注释	497
9.6.1 功能描述	376	11.2 math_emulate.c 程序	499
9.6.2 代码注释	378	11.2.1 功能描述	499
9.7 floppy.c 程序	382	11.2.2 代码注释	499
9.7.1 功能描述	382	第 12 章 文件系统	501
9.7.2 代码注释	382	12.1 总体功能	502
9.7.3 其他信息	396	12.1.1 MINIX 文件系统	502
第 10 章 字符设备驱动程序	409	12.1.2 文件类型、属性和目录项	506
10.1 总体功能	409	12.1.3 高速缓冲区	510
10.1.1 终端驱动程序基本原理	409	12.1.4 文体系统底层函数	511
10.1.2 Linux 支持的终端设备类型	410	12.1.5 文件中数据的访问操作	511
10.1.3 终端基本数据结构	412	12.1.6 文件和目录管理系统调用	513
10.1.4 规范模式和非规范模式	415	12.1.7 360 KB 软盘中文件系统实例 分析	513
10.1.5 控制台终端和串行终端设备	416	12.2 Makefile 文件	517
10.1.6 终端驱动程序接口	419	12.2.1 功能描述	517
10.2 Makefile 文件	419	12.2.2 代码注释	517
10.2.1 功能描述	419	12.3 buffer.c 程序	519
10.2.2 代码注释	419	12.3.1 功能描述	520
10.3 keyboard.S 程序	422	12.3.2 代码注释	525
10.3.1 功能描述	422	12.4 bitmap.c 程序	536
10.3.2 代码注释	422	12.4.1 功能描述	536
10.3.3 其他信息	436	12.4.2 代码注释	537
10.4 console.c 程序	440	12.5 truncate.c 程序	542
10.4.1 功能描述	440	12.5.1 功能描述	542
10.4.2 代码注释	440	12.5.2 代码注释	543
10.4.3 其他信息	459	12.6 inode.c 程序	544
10.5 serial.c 程序	466	12.6.1 功能描述	544
10.5.1 功能描述	466	12.6.2 代码注释	546
10.5.2 代码注释	466	12.7 super.c 程序	556
10.5.3 异步串行通信控制器 UART	468	12.7.1 功能描述	556
10.6 rs_io.s 程序	474	12.7.2 代码注释	557
10.6.1 功能描述	474	12.8 namei.c 程序	565
10.6.2 代码注释	474		
10.7 tty_io.c 程序	478		

12.8.1 功能描述	565	使用分配	640
12.8.2 代码注释	566	13.1.4 页面出错异常处理	640
12.9 file_table.c 程序	588	13.1.5 写时复制(copy on write)机制	640
12.9.1 功能描述	588	13.1.6 需求加载(Load on demand)机制	641
12.9.2 代码注释	588	13.2 Makefile 文件	642
12.10 block_dev.c 程序	589	13.2.1 功能描述	642
12.10.1 功能描述	589	13.2.2 代码注释	642
12.10.2 代码注释	590	13.3 memory.c 程序	643
12.11 file_dev.c 程序	592	13.3.1 功能描述	643
12.11.1 功能描述	592	13.3.2 代码注释	645
12.11.2 代码注释	592	13.4 page.s 程序	661
12.12 pipe.c 程序	595	13.4.1 功能描述	661
12.12.1 功能描述	595	13.4.2 代码注释	662
12.12.2 代码注释	596	13.4.3 页出错异常处理	663
12.13 char_dev.c 程序	599	第 14 章 头文件	664
12.13.1 功能描述	599	14.1 include/目录下的文件	664
12.13.2 代码注释	599	14.2 a.out.h 文件	666
12.14 read_write.c 程序	602	14.2.1 功能描述	666
12.14.1 功能描述	602	14.2.2 代码注释	667
12.14.2 代码注释	602	14.2.3 a.out 执行文件格式	672
12.14.3 用户程序读写操作过程	605	14.3 const.h 文件	676
12.15 open.c 程序	608	14.3.1 功能描述	676
12.15.1 功能描述	608	14.3.2 代码注释	676
12.15.2 代码注释	608	14.4 ctype.h 文件	677
12.16 exec.c 程序	614	14.4.1 功能描述	677
12.16.1 功能描述	614	14.4.2 代码注释	677
12.16.2 代码注释	617	14.5 errno.h 文件	678
12.17 stat.c 程序	629	14.5.1 功能描述	678
12.17.1 功能描述	629	14.5.2 代码注释	678
12.17.2 代码注释	629	14.6 fcntl.h 文件	680
12.18 fcntl.c 程序	631	14.6.1 功能描述	680
12.18.1 功能描述	631	14.6.2 代码注释	680
12.18.2 代码注释	632	14.7 signal.h 文件	682
12.19 ioctl.c 程序	634	14.7.1 功能描述	682
12.19.1 功能描述	634	14.7.2 文件注释	682
12.19.2 代码注释	634	14.8 stdarg.h 文件	684
第 13 章 内存管理	636	14.8.1 功能描述	684
13.1 总体功能	636	14.8.2 代码注释	684
13.1.1 内存分页管理机制	636	14.9 stddef.h 文件	685
13.1.2 Linux 中物理内存的管理和 分配	639	14.9.1 功能描述	685
13.1.3 Linux 内核对线性地址空间的		14.9.2 代码注释	686
		14.10 string.h 文件	687

14.10.1 功能描述	687	14.25 head.h 文件	733
14.10.2 代码注释	687	14.25.1 功能描述	733
14.11 termios.h 文件	696	14.25.2 代码注释	733
14.11.1 功能描述	696	14.26 kernel.h 文件	733
14.11.2 代码注释	697	14.26.1 功能描述	733
14.11.3 控制字符 TIME、MIN	702	14.26.2 代码注释	733
14.12 time.h 文件	703	14.27 mm.h 文件	734
14.12.1 功能描述	703	14.27.1 功能描述	734
14.12.2 代码注释	704	14.27.2 代码注释	734
14.13 unistd.h 文件	705	14.28 sched.h 文件	735
14.13.1 功能描述	705	14.28.1 功能描述	735
14.13.2 代码注释	705	14.28.2 代码注释	736
14.14 utime.h 文件	711	14.29 sys.h 文件	742
14.14.1 功能描述	711	14.29.1 功能描述	742
14.14.2 代码注释	711	14.29.2 代码注释	742
14.15 include/asm/目录下的文件	711	14.30 tty.h 文件	744
14.16 io.h 文件	711	14.30.1 功能描述	744
14.16.1 功能描述	711	14.30.2 代码注释	744
14.16.2 代码注释	712	14.31 include/sys/目录中的文件	746
14.17 memory.h 文件	712	14.32 stat.h 文件	747
14.17.1 功能描述	712	14.32.1 功能描述	747
14.17.2 代码注释	713	14.32.2 代码注释	747
14.18 segment.h 文件	713	14.33 times.h 文件	748
14.18.1 功能描述	713	14.33.1 功能描述	748
14.18.2 代码注释	713	14.33.2 代码注释	749
14.19 system.h 文件	715	14.34 types.h 文件	749
14.19.1 功能描述	715	14.34.1 功能描述	749
14.19.2 代码注释	717	14.34.2 代码注释	749
14.20 include/linux/目录下的文件	719	14.35 utsname.h 文件	750
14.21 config.h 文件	720	14.35.1 功能描述	750
14.21.1 功能描述	720	14.35.2 代码注释	750
14.21.2 代码注释	720	14.36 wait.h 文件	751
14.22 fdreg.h 头文件	722	14.36.1 功能描述	751
14.22.1 功能描述	722	14.36.2 代码注释	751
14.22.2 文件注释	723	第 15 章 库文件	753
14.23 fs.h 文件	725	15.1 Makefile 文件	754
14.23.1 功能描述	725	15.1.1 功能描述	754
14.23.2 代码注释	725	15.1.2 代码注释	754
14.24 hdreg.h 文件	730	15.2 _exit.c 程序	756
14.24.1 功能描述	730	15.2.1 功能描述	756
14.24.2 代码注释	730	15.2.2 代码注释	756
14.24.3 硬盘分区表	732	15.2.3 相关信息	756

15.3 close.c 程序	757	17.1.2 配置文件*.bxrc	782
15.3.1 功能描述	757	17.2 在 Bochs 中运行 Linux 0.1x	
15.3.2 代码注释	757	系统	785
15.4 ctype.c 程序	757	17.2.1 软件包中文件说明	785
15.4.1 功能描述	757	17.2.2 运行 Linux 0.1x 系统	787
15.4.2 代码注释	757	17.3 访问磁盘映像文件中的信息	791
15.5 dup.c 程序	758	17.3.1 使用 WinImage 工具软件	791
15.5.1 功能描述	758	17.3.2 利用现有 Linux 系统	792
15.5.2 代码注释	758	17.4 编译运行简单内核示例程序	793
15.6 errno.c 程序	759	17.5 利用 Bochs 调试内核	795
15.6.1 功能描述	759	17.5.1 运行 Bochs 调试程序	796
15.6.2 代码注释	759	17.5.2 定位内核中的变量或数据结构	801
15.7 execve.c 程序	759	17.6 创建磁盘映像文件	802
15.7.1 功能描述	759	17.6.1 利用 Bochs 软件自带的 Image	
15.7.2 代码注释	759	生成工具	803
15.8 malloc.c 程序	760	17.6.2 在 Linux 系统下使用 dd 命令	
15.8.1 功能描述	760	创建 Image 文件	804
15.8.2 代码注释	762	17.6.3 利用 WinImage 创建 DOS 格式的	
15.9 open.c 程序	768	软盘 Image 文件	804
15.9.1 功能描述	768	17.7 制作根文件系统	805
15.9.2 代码注释	769	17.7.1 根文件系统和根文件设备	805
15.10 setuid.c 程序	769	17.7.2 创建文件系统	806
15.10.1 功能描述	769	17.7.3 Linux-0.11 的 Bochs 配置文件	808
15.10.2 代码注释	770	17.7.4 在 hdc.img 上建立根文件系统	810
15.11 string.c 程序	770	17.7.5 使用硬盘 Image 上的根文件	
15.11.1 功能描述	770	系统	811
15.11.2 代码注释	770	17.8 在 Linux 0.11 系统中编译 0.11	
15.12 wait.c 程序	771	内核	812
15.12.1 功能描述	771	17.9 在 Redhat Linux 9 系统中编译	
15.12.2 代码注释	771	Linux 0.11 内核	813
15.13 write.c 程序	772	17.9.1 修改 makefile 文件	814
15.13.1 功能描述	772	17.9.2 修改汇编程序中的注释	814
15.13.2 代码注释	772	17.9.3 内存位置对齐语句 align 值的	
第 16 章 建造工具	773	修改	815
16.1 build.c 程序	773	17.9.4 修改嵌入宏汇编程序	815
16.1.1 功能概述	773	17.9.5 C 程序变量在汇编语句中的	
16.1.2 代码注释	774	引用表示	815
16.2 MINIX 可执行文件头部数据		17.9.6 保护模式下调试显示函数	815
结构	778	17.10 内核引导启动 + 根文件系统	
第 17 章 实验环境设置与使用方法	780	组成的集成盘	816
17.1 Bochs 仿真软件系统	781	17.10.1 集成盘制作原理	817
17.1.1 设置 Bochs 系统	781	17.10.2 集成盘的制作过程	818

17.10.3 运行集成盘系统	821	内核	825
17.11 从硬盘启动:利用 shoelace 引导软件	821	17.12.3 调试方法和步骤	826
17.11.1 shoelace 程序设置路径	822	附录	830
17.11.2 设置过程	822	附录 A 内核数据结构	830
17.11.3 问题和解决方法	822	附录 B ASCII 码表	838
17.12 利用 GDB 和 Bochs 调试内核 源代码	824	附录 C 常用 C0、C1 控制字符表	840
17.12.1 编译带 gdbstub 的 Bochs 系统 ..	824	附录 D 常用转义序列和控制序列 ..	841
17.12.2 编译带调试信息的 Linux 0.11		附录 E 第 1 套键盘扫描码集	843
		索引	845
		参考文献	881

第 1 章 概 述

本章首先回顾 Linux 操作系统的诞生和发展,读者由此可以理解本书选择 Linux 系统早期版本作为学习对象的一些原因。然后具体说明选择早期 Linux 内核版本的优点和不足之处,以及如何进一步学习。最后对各章的内容进行简要介绍。

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 10 月 5 日(这是第一次正式向外公布的时间)。此后借助于 Internet,经过全世界各地计算机爱好者的共同努力,现已成为当今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、发展和成长过程依赖于以下五个重要支柱:UNIX 操作系统、MINIX 操作系统、GNU 计划、POSIX 标准和 Internet。下面根据这五个基本线索来追寻一下 Linux 的酝酿、开发以及最初的发展。首先分别介绍其中的四个基本要素,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 的帮助下推出比较完善的 1.0 版本这段时间的发展经过。

当然,目前 Linux 内核版本已经开发到了 2.6.x 版。而大多数 Linux 系统中所用到的内核是稳定的 2.6.12 版内核(其中第 2 个数字若是奇数则表示是正在开发的版本,不能保证系统的稳定性)。

1.1.1 UNIX 操作系统的诞生

UNIX 操作系统是美国贝尔实验室的 Ken Thompson 和 Dennis Ritchie 于 1969 年夏天在 DEC PDP-7 小型计算机上开发的一个分时操作系统。

1969 年夏天, 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,趁他夫人回家乡加利福尼亚度假期间,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 开发的。Tanenbaum 在荷兰阿姆斯特丹的 Vrije 大学数学与计算机科学系工作,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 多篇文章和 5 部计算机著作。

Tanenbaum 虽出生在美国纽约,却是荷兰侨民(1914 年他的祖辈来到美国)。他在纽约上中学、麻省理工学院上大学、加州大学 Berkeley 分校念博士学位。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。荷兰首都阿姆斯特丹是个常年阴雨绵绵的城市,但对于 Tanenbaum 来说,这最好不过了,因为在这样的环境下

他就可以经常待在家中摆弄他的计算机了。

MINIX 是他 1987 年编制的,主要用于学生学习操作系统原理。到 1991 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。目前 MINIX 系统可以从许多 FTP 上下载。

对于 Linux 系统,他后来曾表示对其开发者 Linus 的称赞。但他认为 Linux 的发展很大原因是由于他为了保持 MINIX 的小型化,能让学生在一个学期内就能学完,因而没有接纳全世界许多人对 MINIX 的扩展要求。因此在这样的前提下 Linus 编写了 Linux 系统。当然 Linus 也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 语言和汇编语言编写的系统源代码。这是第一次使得有抱负的程序员或黑客能够阅读操作系统的源代码。在当时,这种源代码是软件商们一直小心守护着的秘密。

1.1.3 GNU 计划

GNU 计划和自由软件基金会 FSF(the Free Software Foundation)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使用。虽然这些系统通常被称作“Linux”,但是 Stallman 认为,严格地说,它们应该被称为 GNU/Linux 系统。

到 20 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 能够诞生的重要因素。

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和运行。它是在 20 世纪 80 年代一个 UNIX 用户组(usr/group)的工作基础上取得的。该 UNIX 用户组原来试图将 AT&T 的 System V 操作系统和 Berkeley CSRG 的 BSD 操作系统的调用接口之间的区别重新调和集成。并于 1984 年定制出了 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 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),即以后经常提到的 POSIX.1 标准。

到 1989 年,POSIX 的工作被转移至 ISO/IEC 社团,并由 15 个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。因此工作组期望对系统的其他功能也制定出标准。这样 IEEE POSIX 的工作就开始了。刚开始有十个批准的计划在进行,有近 300 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