

郑钢 ● 著

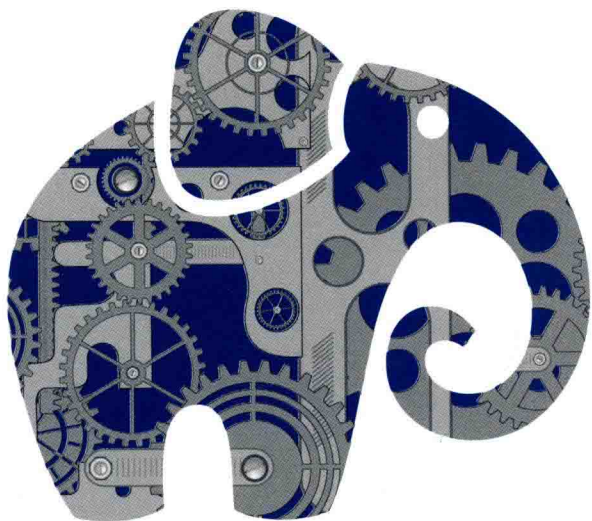
历时 26 个月，行文 30 余万字，  
用 6000 多行代码实现了一个完整的操作系统。

彻底剖析操作系统的原理，实现内核线程、特权级变换、  
用户进程、任务调度、文件系统等操作系统最基本的组成单元。

用实际代码解释了锁、信号量、生产者、消费者问题。

实现了一个简单的 shell，帮助大家理解内部命令、  
外部命令、管道等操作。

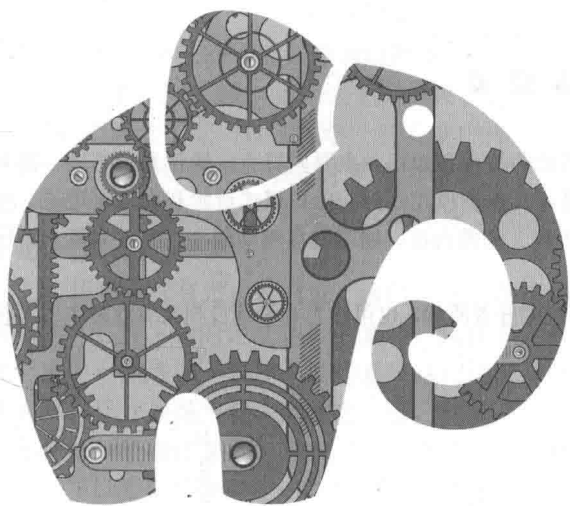
操作系统并不深奥，本书给予权威解读。



# 操作系统

# 真象 还原

郑钢  
◎ 著



操作系统

真象  
还原

藏书

人民邮电出版社  
北京

## 图书在版编目 (C I P) 数据

操作系统真象还原 / 郑钢著. — 北京 : 人民邮电出版社, 2016.3  
ISBN 978-7-115-41434-2

I. ①操… II. ①郑… III. ①操作系统—基本知识  
IV. ①TP316

中国版本图书馆CIP数据核字(2016)第016739号

## 内 容 提 要

本书共分 16 章, 讲解了开发一个操作系统需要的技术和知识, 主要内容有: 操作系统基础、部署工作环境、编写 MBR 主引导记录、完善 MBR 错误、保护模式入门、保护模式进阶和向内核迈进、中断、内存管理系统、线程、输入输出系统、用户进程、完善内核、编写硬盘驱动程序、文件系统、系统交互等核心技术。

本书适合程序员、系统底层开发人员、操作系统爱好者阅读, 也可作为大专院校相关专业师生用书和培训学校的教材。

- 
- ◆ 著 郑 钢
  - 责任编辑 张 涛
  - 责任印制 张佳莹 焦志炜
  - ◆ 人民邮电出版社出版发行 北京市丰台区成寿寺路 11 号  
邮编 100164 电子邮件 315@ptpress.com.cn  
网址 <http://www.ptpress.com.cn>  
三河市中晟雅豪印务有限公司印刷
  - ◆ 开本: 787×1092 1/16  
印张: 48.25  
字数: 1 281 千字 2016 年 3 月第 1 版  
印数: 1—2 500 册 2016 年 3 月河北第 1 次印刷
- 

定价: 108.00 元

读者服务热线: (010)81055410 印装质量热线: (010)81055316  
反盗版热线: (010)81055315

# 前 言

本书面向操作系统基础知识薄弱，但又想把操作系统搞清楚、喜欢刨根问底的技术人，在此向你们致敬，本书用诙谐幽默的语言，把深奥的操作系统尽量讲解清楚，读者在轻松阅读中就学通了深奥的知识，是一本难得的好书。

多数学习操作系统的读者都会有这样的感受：

(1) “太难了，对于操作系统这个庞然大物我简直无从下手”；

(2) “很后悔选了这门课（大学一些专业中操作系统是选修课），甚至不想学习计算机了”；

(3) “上课完全听不懂，我都不想继续听下去了”；

(4) “即使实验做出来了，由于只是完成了局部功能，我依然不明白操作系统是怎样运行起来的，甚至不知道自己在做什么”。

以上的感受我都有过，坦白说，这门课并不是很难，但想把这门课完全搞明白真不容易。我是个喜欢刨根问底的人，为了弄清楚这背后的真相，我花了大量时间学习课程之外的内容，甚至付出了惨痛的代价——大学中第一次考试不及格，操作系统这门课我是第二次才考过的。这确实很“讽刺”——操作系统不及格的人在写操作系统书籍！但转念一想，考试过了的同学并不代表能够写出操作系统，因为试卷上并不是在考如何写一个操作系统。和技术能力相比，卷面成绩并不重要。

想象一下，如果是爱因斯坦那样的天才给我们讲物理知识，我们会觉得物理更容易理解吗？肯定是不会的，因为在爱因斯坦眼中比较容易的内容也许对我们来说非常深奥，他用 B 解释 A 的时候也许会让我们更迷惑，因为 B 我们也不懂，这就是基础的问题了。幸运的是阅读本书时读者只要有 C 语言和部分汇编语言的基础就行了，涉及的其他方面的知识我都会详细介绍，并以更易懂的方式去解释技术难点，读者不必担心看不懂本书。

回忆一下学车的经历：教练让学员先踩离合器再挂档，然后再踩油门，车子就开动啦。如果学员总是学不会这些，有可能是学员根本不知道什么是离合器，或者不知道离合器的作用是什么。即使把这些操作背下来，也会对驾车感到心有余而力不足，可见，只有了解了背后的原理，才会知道自己在做什么，驾车才变得游刃有余。

以上情况对我们学习操作系统来说也同样存在，比如当老师介绍中断发生时的上下文保护时，我们更多的疑问不是如何保存 CPU 的上下文数据，而是想知道为什么在不同的特权级下会使用不同的栈，这背后的原理是什么，并且这是如何做到的。

诸如此类的疑问需要了解硬件原生支持的运行机制，因为很多操作都是硬件自动完成的，比如处理器进入 0 特权级时，会自动在任务状态段 TSS 中获得 0 特权级的栈地址，这不需要人工干涉，完全由处理器维护。我们想知道的是，硬件在背后自动完成了哪些工作，这样才便于我们理解操作系统的全貌。

操作系统受制于硬件的支持，很大程度上它的能力取决于硬件的能力，因此，要想全面理解操作系统，不仅需要了解上层软件的算法、原理、实现，还要了解很多硬件底层的内容。和硬件相关的知识是在微机接口电路中讲解的，而绝大多数读者在学习这门课时，根本不知道它有何用，只有学习操作系统课程时才用到它，因此，本书内容兼顾相关的硬件知识。

除硬件外，本书还把操作系统中的理论付诸于实践，让读者真正学到包含在操作系统中的实实在在的技术，比如在代码中实现了著名的生产者消费者问题，还有进程、线性、阻塞、信号量、锁、文件系统、目录、shell、管道等。各个章节的代码都可独立运行，方便调试，本书更让读者有成就感的是，我们最终完成的一个操作系统总共代码量只有几千行左右，极大地减少了操作系统源码阅读的工作量。

操作系统还是比较庞大的，因此，大部分介绍操作系统原理的书中，对各个部分都是分拆出来介绍的，这导致我们学习操作系统时犹如盲人摸象、管中窥豹。本书的封面是一个完整的大象的拼图，就像封面展示的那样，本书内容我们不再局部学习，而是把所有局部还原成一个整体，做出一个真正的操作系统。

为了让读者不再惧怕操作系统，同时也为了完成我自己的心愿，我辞职专心进行本书的编写，在此期间也曾拒绝了多份回报丰厚的工作，现在想想真是疯狂……苦了我的父母和女朋友，在这里跟你们说声抱歉，你们“纵容”我的偏执，真心不容易，辛苦啦，我爱你们！

感谢我在北京大学就读期间的 Linux 内核课程老师（同时也是我的研究生导师）荆琦教授和操作系统课程老师陈向群教授，很荣幸能够成为你们的学生，时至今日我常常回想起课堂上你们言传身教并为我解答问题的身影，你们渊博的知识和教学上严谨的态度深深影响了我，仅以此书向我这两位恩师致谢。

感谢父母给予我的理解和宽容，以后我一定加倍努力回报你们的养育之恩！

最后，感谢女朋友给予我的陪伴和照顾，在写此书的过程中我深深体会到：爱并不仅仅体现在相信对方一定能成功，更多是体现在支持对方去做想做的事，即使失败了也不会嫌弃。尽管在这漫长枯燥的 19 个月当中，如果没有你的“唠唠叨叨”本书早就写完了，但恰恰是这种“唠唠叨叨”下的不离不弃让我相信这世上还有真爱。

我爱你王小兔（对我女朋友的昵称），本书是我送给你的礼物。

本书中出现的“兄弟”“大伙儿”“同学”和“咱们”的称谓，是作者为了活泼写作风格故意为之，别无他意，在此说明一下。本书读者交流 QQ 群为：148177180，编辑联系邮箱：zhangtao@ptpress.com.cn。

作者

于北京大学图书馆

# 目 录

<b>第 0 章 一些你可能正感到迷惑的问题</b> .....1	<b>第 1 章 部署工作环境</b> .....42
0.1 操作系统是什么.....1	1.1 工欲善其事，必先利其器.....42
0.2 你想研究到什么程度.....2	1.2 我们需要哪些编译器.....42
0.3 写操作系统，哪些需要我来做.....2	1.2.1 世界顶级编译器 GCC.....42
0.4 软件是如何访问硬件的.....2	1.2.2 汇编语言编译器新贵 NASM.....43
0.5 应用程序是什么，和操作系统是如何配合到一起的.....3	1.3 操作系统的宿主环境.....43
0.6 为什么称为“陷入”内核.....4	1.3.1 什么是虚拟机.....44
0.7 内存访问为什么要分段.....4	1.3.2 盗梦空间般的开发环境，虚拟机中再装一个虚拟机.....45
0.8 代码中为什么分为代码段、数据段？这和内存访问机制中的段是一回事吗.....6	1.3.3 virtualBox 下载，安装.....46
0.9 物理地址、逻辑地址、有效地址、线性地址、虚拟地址的区别.....11	1.3.4 Linux 发行版下载.....46
0.10 什么是段重叠.....12	1.3.5 Bochs 下载安装.....46
0.11 什么是平坦模型.....12	1.4 配置 bochs.....48
0.12 cs、ds 这类 sreg 段寄存器，位宽是多少.....12	1.5 运行 bochs.....49
0.13 什么是工程，什么是协议.....13	<b>第 2 章 编写 MBR 主引导记录，让我们开始掌权</b> .....52
0.14 为什么 Linux 系统下的应用程序不能在 Windows 系统下运行.....14	2.1 计算机的启动过程.....52
0.15 局部变量和函数参数为什么要放在栈中.....14	2.2 软件接力第一棒，BIOS.....52
0.16 为什么说汇编语言比 C 语言快.....15	2.2.1 实模式下的 1MB 内存布局.....52
0.17 先有的语言，还是先有的编译器，第 1 个编译器是怎么产生的.....16	2.2.2 BIOS 是如何苏醒的.....54
0.18 编译型程序与解释型程序的区别.....19	2.2.3 为什么是 0x7c00.....56
0.19 什么是大端字节序、小端字节序.....19	2.3 让 MBR 先飞一会儿.....58
0.20 BIOS 中断、DOS 中断、Linux 中断的区别.....21	2.3.1 神奇好用的 \$ 和 \$\$，令人迷惑的 section.....58
0.21 Section 和 Segment 的区别.....25	2.3.2 NASM 简单用法.....60
0.22 什么是魔数.....29	2.3.3 请下一位选手 MBR 同学做准备.....60
0.23 操作系统是如何识别文件系统的.....30	<b>第 3 章 完善 MBR</b> .....65
0.24 如何控制 CPU 的下一条指令.....30	3.1 地址、section、vstart 浅尝辄止.....65
0.25 指令集、体系结构、微架构、编程语言.....30	3.1.1 什么是地址.....65
0.26 库函数是用户进程与内核的桥梁.....33	3.1.2 什么是 section.....67
0.27 转义字符与 ASCII 码.....37	3.1.3 什么是 vstart.....68
0.28 MBR、EBR、DBR 和 OBR 各是什么.....39	3.2 CPU 的实模式.....70
	3.2.1 CPU 的工作原理.....71
	3.2.2 实模式下的寄存器.....72
	3.2.3 实模式下内存分段由来.....76
	3.2.4 实模式下 CPU 内存寻址方式.....78
	3.2.5 栈到底是什么玩意儿.....81

3.2.6	实模式下的 ret	84	4.4.4	分支预测	169
3.2.7	实模式下的 call	85	4.5	使用远跳转指令清空流水线, 更新段描述符缓冲寄存器	172
3.2.8	实模式下的 jmp	92	4.6	保护模式之内存段的保护	173
3.2.9	标志寄存器 flags	97	4.6.1	向段寄存器加载选择子时的保护	173
3.2.10	有条件转移	99	4.6.2	代码段和数据段的保护	174
3.2.11	实模式小结	101	4.6.3	栈段的保护	175
3.3	让我们直接对显示器说点什么吧	101	<b>第 5 章</b>	<b>保护模式进阶, 向内核迈进</b>	<b>177</b>
3.3.1	CPU 如何与外设通信——IO 接口	101	5.1	获取物理内存容量	177
3.3.2	显卡概述	105	5.1.1	学习 Linux 获取内存的方法	177
3.3.3	显存、显卡、显示器	106	5.1.2	利用 BIOS 中断 0x15 子功能 0xe820 获取内存	177
3.3.4	改进 MBR, 直接操作显卡	110	5.1.3	利用 BIOS 中断 0x15 子功能 0xe801 获取内存	179
3.4	bochs 调试方法	112	5.1.4	利用 BIOS 中断 0x15 子功能 0x88 获取内存	180
3.4.1	bochs 一般用法	113	5.1.5	实战内存容量检测	181
3.4.2	bochs 调试实例	118	5.2	启用内存分页机制, 畅游虚拟空间	186
3.5	硬盘介绍	122	5.2.1	内存为什么要分页	186
3.5.1	硬盘发展简史	122	5.2.2	一级页表	188
3.5.2	硬盘工作原理	123	5.2.3	二级页表	192
3.5.3	硬盘控制器端口	126	5.2.4	规划页表之操作系统与用户进程的关系	197
3.5.4	常用的硬盘操作方法	128	5.2.5	启用分页机制	198
3.6	让 MBR 使用硬盘	129	5.2.6	用虚拟地址访问页表	204
3.6.1	改造 MBR	130	5.2.7	快表 TLB (Translation Lookaside Buffer) 简介	206
3.6.2	实现内核加载器	134	5.3	加载内核	207
<b>第 4 章</b>	<b>保护模式入门</b>	<b>136</b>	5.3.1	用 C 语言写内核	207
4.1	保护模式概述	136	5.3.2	二进制程序的运行方法	211
4.1.1	为什么要有保护模式	136	5.3.3	elf 格式的二进制文件	213
4.1.2	实模式不是 32 位 CPU, 变成了 16 位	137	5.3.4	elf 文件实例分析	218
4.2	初见保护模式	137	5.3.5	将内核载入内存	222
4.2.1	保护模式之寄存器扩展	137	5.4	特权级深入浅出	229
4.2.2	保护模式之寻址扩展	140	5.4.1	特权级那点事	229
4.2.3	保护模式之运行模式反转	141	5.4.2	TSS 简介	230
4.2.4	保护模式之指令扩展	145	5.4.3	CPL 和 DPL 入门	232
4.3	全局描述符表	150	5.4.4	门、调用门与 RPL 序	235
4.3.1	段描述符	150	5.4.5	调用门的过程保护	240
4.3.2	全局描述符表 GDT、局部描述符表 LDT 及选择子	155	5.4.6	RPL 的前世今生	243
4.3.3	打开 A20 地址线	157	5.4.7	IO 特权级	248
4.3.4	保护模式的开关, CR0 寄存器的 PE 位	158	<b>第 6 章</b>	<b>完善内核</b>	<b>252</b>
4.3.5	让我们进入保护模式	158	6.1	函数调用约定简介	252
4.4	处理器微架构简介	165			
4.4.1	流水线	166			
4.4.2	乱序执行	168			
4.4.3	缓存	168			

6.2 汇编语言和 C 语言混合编程	256	8.1.1 makefile 是什么	357
6.2.1 浅析 C 库函数与系统调用	256	8.1.2 makefile 基本语法	358
6.2.2 汇编语言和 C 语言共同协作	259	8.1.3 跳到目标处执行	360
6.3 实现自己的打印函数	261	8.1.4 伪目标	361
6.3.1 显卡的端口控制	261	8.1.5 make: 递归式推导目标	362
6.3.2 实现单个字符打印	265	8.1.6 自定义变量与系统变量	363
6.3.3 实现字符串打印	275	8.1.7 隐含规则	365
6.3.4 实现整数打印	277	8.1.8 自动化变量	366
6.4 内联汇编	281	8.1.9 模式规则	367
6.4.1 什么是内联汇编	281	8.2 实现 assert 断言	367
6.4.2 汇编语言 AT&T 语法简介	281	8.2.1 实现开、关中断的函数	367
6.4.3 基本内联汇编	283	8.2.2 实现 ASSERT	370
6.4.4 扩展内联汇编	284	8.2.3 通过 makefile 来编译	372
6.4.5 扩展内联汇编之机器模式简介	294	8.3 实现字符串操作函数	374
<b>第 7 章 中断</b>	<b>298</b>	8.4 位图 bitmap 及其函数的实现	377
7.1 中断是什么, 为什么要有中断	298	8.4.1 位图简介	377
7.2 操作系统是中断驱动的	299	8.4.2 位图的定义与实现	378
7.3 中断分类	299	8.5 内存管理系统	381
7.3.1 外部中断	299	8.5.1 内存池规划	381
7.3.2 内部中断	301	8.5.2 内存管理系统第一步, 分配页 内存	388
7.4 中断描述符表	304	<b>第 9 章 线程</b>	<b>398</b>
7.4.1 中断处理过程及保护	306	9.1 实现内核线程	398
7.4.2 中断发生时的压栈	308	9.1.1 执行流	398
7.4.3 中断错误码	310	9.1.2 线程到底是什么	399
7.5 可编程中断控制器 8259A	311	9.1.3 进程与线程的关系、区别简述	402
7.5.1 8259A 介绍	311	9.1.4 进程、线程的状态	405
7.5.2 8259A 的编程	314	9.1.5 进程的身份证——PCB	405
7.6 编写中断处理程序	319	9.1.6 实现线程的两种方式——内核或 用户进程	406
7.6.1 从最简单的中断处理程序 开始	319	9.2 在内核空间实现线程	409
7.6.2 改进中断处理程序	335	9.2.1 简单的 PCB 及线程栈的实现	409
7.6.3 调试实战: 处理器进入中断时 压栈出栈完整过程	339	9.2.2 线程的实现	413
7.7 可编程计数器/定时器 8253 简介	346	9.3 核心数据结构, 双向链表	417
7.7.1 时钟——给设备打拍子	346	9.4 多线程调度	421
7.7.2 8253 入门	348	9.4.1 简单优先级调度的基础	421
7.7.3 8253 控制字	349	9.4.2 任务调度器和任务切换	425
7.7.4 8253 工作方式	350	<b>第 10 章 输入输出系统</b>	<b>439</b>
7.7.5 8253 初始化步骤	353	10.1 同步机制——锁	439
7.8 提高时钟中断的频率, 让中断来得更 猛烈一些	354	10.1.1 排查 GP 异常, 理解原子操作	439
<b>第 8 章 内存管理系统</b>	<b>357</b>	10.1.2 找出代码中的临界区、互斥、 竞争条件	444
8.1 makefile 简介	357	10.1.3 信号量	445



10.1.4	线程的阻塞与唤醒	447	12.2.1	系统调用实现框架	527
10.1.5	锁的实现	449	12.2.2	增加 0x80 号中断描述符	527
10.2	用锁实现终端输出	452	12.2.3	实现系统调用接口	528
10.3	从键盘获取输入	456	12.2.4	增加 0x80 号中断处理例程	528
10.3.1	键盘输入原理简介	456	12.2.5	初始化系统调用和实现 sys_getpid	530
10.3.2	键盘扫描码	457	12.2.6	添加系统调用 getpid	531
10.3.3	8042 简介	463	12.2.7	在用户进程中的系统调用	532
10.3.4	测试键盘中断处理程序	465	12.2.8	系统调用之栈传递参数	534
10.4	编写键盘驱动	468	12.3	让用户进程“说话”	536
10.4.1	转义字符介绍	468	12.3.1	可变参数的原理	536
10.4.2	处理扫描码	469	12.3.2	实现系统调用 write	538
10.5	环形输入缓冲区	476	12.3.3	实现 printf	539
10.5.1	生产者与消费者问题简述	476	12.3.4	完善 printf	542
10.5.2	环形缓冲区的实现	478	12.4	完善堆内存管理	545
10.5.3	添加键盘输入缓冲区	481	12.4.1	malloc 底层原理	545
10.5.4	生产者与消费者实例测试	482	12.4.2	底层初始化	548
<b>第 11 章</b>	<b>用户进程</b>	<b>485</b>	12.4.3	实现 sys_malloc	550
11.1	为什么要有任务状态段 TSS	485	12.4.4	内存的释放	555
11.1.1	多任务的起源, 很久很久 以前	485	12.4.5	实现 sys_free	558
11.1.2	LDT 简介	486	12.4.6	实现系统调用 malloc 和 free	562
11.1.3	TSS 的作用	488	<b>第 13 章</b>	<b>编写硬盘驱动程序</b>	<b>566</b>
11.1.4	CPU 原生支持的任务切换 方式	492	13.1	硬盘及分区表	566
11.1.5	现代操作系统采用的任务 切换方式	495	13.1.1	创建从盘及获取安装的 磁盘数	566
11.2	定义并初始化 TSS	497	13.1.2	创建磁盘分区表	567
11.3	实现用户进程	501	13.1.3	磁盘分区表浅析	571
11.3.1	实现用户进程的原理	501	13.2	编写硬盘驱动程序	578
11.3.2	用户进程的虚拟地址空间	501	13.2.1	硬盘初始化	578
11.3.3	为进程创建页表和 3 特权 级栈	502	13.2.2	实现 thread_yield 和 idle 线程	582
11.3.4	进入特权级 3	505	13.2.3	实现简单的休眠函数	584
11.3.5	用户进程创建的流程	506	13.2.4	完善硬盘驱动程序	585
11.3.6	实现用户进程——上	507	13.2.5	获取硬盘信息, 扫描分区表	590
11.3.7	bss 简介	513	<b>第 14 章</b>	<b>文件系统</b>	<b>595</b>
11.3.8	实现用户进程——下	515	14.1	文件系统概念简介	595
11.3.9	让进程跑起来——用户进程的 调度	519	14.1.1	inode、间接块索引表、文件 控制块 FCB 简介	595
11.3.10	测试用户进程	520	14.1.2	目录项与目录简介	597
<b>第 12 章</b>	<b>进一步完善内核</b>	<b>523</b>	14.1.3	超级块与文件系统布局	599
12.1	Linux 系统调用浅析	523	14.2	创建文件系统	601
12.2	系统调用的实现	527	14.2.1	创建超级块、i 结点、目录项	601
			14.2.2	创建文件系统	603
			14.2.3	挂载分区	609

14.3 文件描述符简介 .....	612	14.14.1 显示当前工作目录的原理及 基础代码 .....	679
14.3.1 文件描述符原理 .....	612	14.14.2 实现 sys_getcwd .....	681
14.3.2 文件描述符的实现 .....	614	14.14.3 实现 sys_chdir 改变工作目录 .....	683
14.4 文件操作相关的基础函数 .....	615	14.15 获得文件属性 .....	684
14.4.1 inode 操作有关的函数 .....	616	14.15.1 ls 命令的幕后功臣 .....	684
14.4.2 文件相关的函数 .....	620	14.15.2 实现 sys_stat .....	685
14.4.3 目录相关的函数 .....	623	<b>第 15 章 系统交互</b> .....	<b>687</b>
14.4.4 路径解析相关的函数 .....	628	15.1 fork 的原理与实现 .....	687
14.4.5 实现文件检索功能 .....	630	15.1.1 什么是 fork .....	687
14.5 创建文件 .....	633	15.1.2 fork 的实现 .....	689
14.5.1 实现 file_create .....	633	15.1.3 添加 fork 系统调用与实现 init 进程 .....	695
14.5.2 实现 sys_open .....	636	15.2 添加 read 系统调用, 获取键盘输入 .....	696
14.5.3 在文件系统中创建第 1 个 文件 .....	639	15.3 添加 putchar、clear 系统调用 .....	697
14.6 文件的打开与关闭 .....	640	15.4 实现一个简单的 shell .....	699
14.6.1 文件的打开 .....	640	15.4.1 shell 雏形 .....	699
14.6.2 文件的关闭 .....	642	15.4.2 添加 Ctrl+u 和 Ctrl+l 快捷键 .....	701
14.7 实现文件写入 .....	643	15.4.3 解析键入的字符 .....	703
14.7.1 实现 file_write .....	643	15.4.4 添加系统调用 .....	705
14.7.2 改进 sys_write 及 write 系统 调用 .....	648	15.4.5 路径解析转换 .....	708
14.7.3 把数据写入文件 .....	650	15.4.6 实现 ls、cd、mkdir、ps、rm 等 命令 .....	712
14.8 读取文件 .....	651	15.5 加载用户进程 .....	717
14.8.1 实现 file_read .....	651	15.5.1 实现 exec .....	717
14.8.2 实现 sys_read 与功能验证 .....	653	15.5.2 让 shell 支持外部命令 .....	723
14.9 实现文件读写指针定位功能 .....	655	15.5.3 加载硬盘上的用户程序执行 .....	724
14.10 实现文件删除功能 .....	657	15.5.4 使用户进程支持参数 .....	727
14.10.1 回收 inode .....	657	15.6 实现系统调用 wait 和 exit .....	731
14.10.2 删除目录项 .....	660	15.6.1 wait 和 exit 的作用 .....	731
14.10.3 实现 sys_unlink 与功能验证 .....	663	15.6.2 孤儿进程和僵尸进程 .....	732
14.11 创建目录 .....	665	15.6.3 一些基础代码 .....	733
14.11.1 实现 sys_mkdir 创建目录 .....	666	15.6.4 实现 wait 和 exit .....	737
14.11.2 创建目录功能验证 .....	669	15.6.5 实现 cat 命令 .....	741
14.12 遍历目录 .....	671	15.7 管道 .....	745
14.12.1 打开目录和关闭目录 .....	671	15.7.1 管道的原理 .....	745
14.12.2 读取 1 个目录项 .....	673	15.7.2 管道的设计 .....	747
14.12.3 实现 sys_readdir 及 sys_ rewinddir .....	674	15.7.3 管道的实现 .....	748
14.13 删除目录 .....	676	15.7.4 利用管道实现进程间通信 .....	752
14.13.1 删除目录与判断空目录 .....	676	15.7.5 在 shell 中支持管道 .....	754
14.13.2 实现 sys_rmdir 及功能验证 .....	677	<b>参考文献</b> .....	<b>760</b>
14.14 任务的工作目录 .....	679		

## 第0章 一些你可能正感到迷惑的问题

正如计算机中数组下标是从0开始的，我们的内容也从0开始，尽量做到低基础学习（负责地说，不是0基础，而且还只是尽量），解释一些学习过程中经常被问到的问题。

### 0.1 操作系统是什么

我并没有给你提供教科书上对操作系统的定义，因为解释得太抽象了，看了之后似乎只是获得一些感性认识，好奇心强的读者反而会产生更多迷惑。为了说清楚问题，让我给您举个例子。

让我们扯点远的……在盘古开天之际，除动物以外，世界上只有土地、荒草、树木、石头等资源。人们为了躲避天灾、野兽攻击等危险，开始住进了山洞，为了获取食物，用石头和树木等材料打造一些武器。当时所有人都在做这些相同的事。这就是没有组织的人类社会，所有人都在重复“造轮子”。

后来各个地区有了自己权威性的部落，部落都专门找人打造武器，谁需要武器就直接申请领取便可，大部分人不需要自己打造武器了。后来嫌打猎太麻烦了，干脆养一些家畜好了，直接供给人们，谁需要可以过来交换。这就是把大家的重复性劳动集中到了一起，让人们可以专注于自己的事情。

再后来，部落之间为了通信，开始有信使了，这是最原始的通信方式。到后来发展到有社会组织，通信越来越频繁了，干脆搞个驿站吧，谁需要通信，直接写信，由驿站代为送达。

随着人口越来越多，社会组织需要了解到到底有多少人，为了方便人口管理，于是就在各地建了“户籍办事”处，人们的生老病死都要到那里登记申报。

说到这我估计您已经猜出我所说的了，上面提到的部落其实就是最原始的操作系统雏形，它将大家都需要的工作找专人负责，大家不用重复劳动。而以上的社会组织其实就是代表现代操作系统，除了把重复性工作集中化，还有了自己的管理策略。

把上面的例子再具体一下，人们想狩猎时，可不可以自己先打造武器，然后拿着自己的武器去狩猎？当然可以，自己制造武器完全没有问题，但部落既然有现成的武器可用，何必自己再费事呢。另外，部落担不担心你随意制造武器会对他人造成伤害？当然会，所以部落不允许你自己制造武器了，人们只有申请的资格，给不给还是要看人家部落的意愿。这就是操作系统提供给用户进程一些系统调用，当用户进程需要某个资源时，直接调用便可，不用自己再费尽心思考虑硬件的事情了，由操作系统把资源获取到后交给用户进程，用户进程可以专注于自己的工作。但操作系统为了保护计算机系统不被损害，不允许用户进程直接访问硬件资源，比如用户进程将操作系统所占据的内存恶意覆盖了，操作系统也就不复存在了，没有操作系统的话，计算机将会瘫痪无法运作。

当人们想和远方的朋友说话时，虽然可以徒步走到亲朋好友身边再对其表达想说的话，但社会组织已经给提供了邮局和电话，何必自己再大老远跑一趟呢。这就是操作系统（社会组织）提供的资源。两个人想在一起生活，要不要一定先结婚呢？完全不用，领不领证都不会阻碍人们在一起生活，但是社会组织为了方便人口管理做了额外约束。不领证的话，至少社会组织无法预测未来人口数量趋势，无法做出宏观调控，甚至这是找到你家人的一种方法。这就如Linux系统中的内存管理，分别要记录哪些页是Active，哪些是“脏页”。不记录会不会影响程序执行，当然不会，记录这些状态还不是为了更好地管理内存吗。

以上说的社会组织和人们之间的关系，正是操作系统和用户进程的关系，希望大家能对操作系统有个初步印象，后面的实践中我们将实例化各个部分。

## 0.2 你想研究到什么程度

学无止境，学习没有说到头的那天。学习到任何程度都是存有疑惑的，就像中学和大学都讲物理，但学的深度不一样，各个阶段都会产生疑问。我们只是基于一些公认的知识，使其作为学习的起点，并以此展开上层的研究。

比如我对太空很感兴趣，大伙儿都知道地球围绕太阳做周期性公转，后来又知道电子围绕原子核来做周期性公转运动，这和地球绕太阳公转的行为如出一辙，甚至我在想太阳是不是相当于原子核，地球相当于一个电子，我们只是生活在一个电子上……而我们身体里有那么多的原子和电子，对那些我们身体中更为细微的生物来说，我们的身体是不是一个宇宙，无尽的猜想，无尽的疑惑。想法虽然有些荒诞，但基于现有科技目前谁也无法证明这是错的，而且近期已有科学文献证明人的大脑就像个宇宙。如果无止境地刨根问底下去，虽然会对底层科学更加清晰，但这对上层知识的学习非常不利，从而我们需要一个公设，我们认为原子是不可再分的，没有更微小的对象了，一切理论研究以此为基础展开。比如乘法是基于加法的，我们研究 $3 \times 4$ 等于多少，必须要承认 $1+1$ 等于 $2$ ，并认为其为真理，不用再去质疑 $1+1$ 为什么等于 $2$ 了，这就是我们的公设，至于为什么 $1+1$ 等于 $2$ ，还是由专门研究基础科学的学者们去探究吧。

学习操作系统也一样，不必纠结于硬件内部是如何工作的，我们只要认为给硬件一个输入，硬件就会给我一个输出就行了，因为即使你学到了硬件内部电子电路，随着你不断进步，钻研不断深入，也许有一天你的求知欲到了物理领域，并产生了物理科学方面的质疑……这让我想到一个笑话，某人准备去买自行车，结果被销售人员不断劝说，加点钱就能买摩托啦，等决定买摩托时，销售人员又说既然都决定买摩托车了，不如再加点买汽车吧，给出了各种汽车方面的优势，欲望需求不断升级，不断被销售劝说，最后居然花了几百万元买车，最后才想起自己是来买自行车的，甚至他还没有驾照……于是，咱们赶紧就此打住，我们是来学操作系统的。

你想学到哪个程度呢，你的公设是什么，要不咱们还是走一步说一步吧。

## 0.3 写操作系统，哪些需要我来做

首先应该明确，在计算机中有分层的概念，也就是说，计算机是一个大的组合物，由各个部分组成一个系统。每个部分就是一层功能模块，各司其职，它只完成一定的工作，并将自己的工作结果（也就是输出）交给下一层的模块，这里的模块指的是各种外设、硬件。

这样，各种工作成果不断累加，通过这种流水线式的上下游协作，便实现了所谓的系统。可见，系统就是各种功能组合到一起后，产生最终输出的组合物。就像人的身体，胃负责搅拌食物，将这些食物变食糜后交给小肠，因为小肠只能处理流食，所以上游的输出一定要适合作为下游的输入，是不是有点类似管道操作了，哈哈，分工协作是大自然的安排，并不是只有计算机世界才有。我们人类的思想是大自然安排好的，所以人类创造的事物也是符合大自然规律的。

好，赶紧回到正题，操作系统是管理资源的软件，操作系统能做什么，取决于主机上硬件的功能。就像用 Maya 造一个人体模型出来，首先我得知道 Maya 这个软件提供曲线曲面各种建模方法才行，换句话说，对于人体建模，你不可能想到用 QQ，因为它不是干这个的。我想说的是硬件不支持的话，操作系统也没招……操作系统一直是所谓的底层，拥有至高无上的控制权，一副牛气轰轰的样子，原来也要依仗他人啊。是啊，操作系统毕竟是软件，而软件的逻辑是需要作用在硬件上才能体现出来的。

所以说，写操作系统需要了解硬件，这些硬件提供了软件方面的接口，这样我们的操作系统通过软件（计算机指令）就能够控制硬件。我们需要做的就是知道如何通过计算机指令来控制硬件，参考硬件手册这下少不了啦。

## 0.4 软件是如何访问硬件的

硬件是各种各样的，发展速度还是非常快的。各个硬件都有自己的个性，操作系统不可能及时更新各

种硬件的驱动方法吧。比如，刚出来某个新硬件，OS 开发者们便开始为其写驱动，这不太现实，会把人累死的。于是乎，便出现了各种硬件适配设备，这就是 IO 接口。接口其实就是标准，大家生产出来的硬件按照这个标准工作就实现了通用。

硬件在输入输出上大体分为串行和并行，相应的接口也就是串行接口和并行接口。串行硬件通过串行接口与 CPU 通信，反过来也是，CPU 通过串行接口与串行设备数据传输。并行设备的访问类似，只不过是通过并行接口进行的。

访问外部硬件有两个方式。

(1) 将某个外设的内存映射到一定范围的地址空间中，CPU 通过地址总线访问该内存区域时会落到外设的内存中，这种映射让 CPU 访问外设的内存就如同访问主板上的物理内存一样。有的设备是这样做的，比如显卡，显卡是显示器的适配器，CPU 不直接和显示器交互，它只和显卡通信。显卡上有片内存叫显存，它被映射到主机物理内存上的低端 1MB 的 0xB8000~0xBFFFF。CPU 访问这片内存就是访问显存，往这片内存上写字节便是往屏幕上打印内容。看上去这么高大上的做法是怎么实现的，这个我们就不关心了，前面说过，计算机中处处是分层，我们要充分相信上一层的工作。

(2) 外设是通过 IO 接口与 CPU 通信的，CPU 访问外设，就是访问 IO 接口，由 IO 接口将信息传递给另一端的外设，也就是说，CPU 从来不知道有这些设备的存在，它只知道自己操作的 IO 接口，你看，处处体现着分层。

于是问题来了，如何访问到 IO 接口呢，答案就是 IO 接口上面有一些寄存器，访问 IO 接口本质上就是访问这些寄存器，这些寄存器就是人们常说的端口。这些端口是人家 IO 接口给咱们提供的接口。人家接口电路也有自己的思维（系统），看到寄存器中写了什么就做出相应的反应。接口提供接口，哈哈，有意思。不过这是人家的约定，没有约定就乱了，各干各的，大家都累，咱们只要遵循人家的规定就能访问成功。

## 0.5 应用程序是什么，和操作系统是如何配合到一起的

应用程序是软件（似乎是废话，别急，往后看），操作系统也是软件。CPU 会将它们一视同仁，甚至，CPU 不知道自己在执行的程序是操作系统，还是一般应用软件，CPU 只知道去 cs: ip 寄存器中指向的内存取指令并执行，它不知道什么是操作系统，也无需知道。

操作系统是人想出来的，为了让自己管理计算机方便而创造出来的一套管理办法。

应用程序要用某种语言编写，而语言又是编译器来提供的。其实根本就没有什么语言，有的只是编译器。是编译器决定怎样解释某种关键字及某种语法。语言只是编译器和大家的约定，只要写入这样的代码，编译器便将其翻译成某种机器指令，翻译成什么样取决于编译器的行为，和语言无关，比如说 C 语言的 printf 函数，它的功能不是说一定要把字符打印到屏幕上，这要看编译器对这种关键字的处理。

编译器提供了一套库函数，库函数中又有封装的系统调用，这样的代码集合称之为运行库。C 语言的运行库称为 C 运行库，就是所谓的 CRT (C Runtime Library)。

应用程序加上操作系统提供功能才算是完整的程序。由于有了操作系统的支持，一些现成的东西已经摆在那了，但这些是属于操作系统的，不是应用程序的，所以咱们平时所写的应用程序只是半成品，需要调用操作系统提供好的函数才能完整地做成一件事，而这个函数便是系统调用。

用户态与内核态是对 CPU 来讲的，是指 CPU 运行在用户态（特权 3 级）还是内核态（特权 0 级），很多人误以为是对用户进程来讲的。

用户进程陷入内核态是指：由于内部或外部中断发生，当前进程被暂时终止执行，其上下文被内核的中断程序保存起来后，开始执行一段内核的代码。是内核的代码，不是用户程序在内核的代码，用户代码怎么可能在内核中存在，所以“用户态与内核态”是对 CPU 来说的。

当应用程序陷入内核后，它自己已经下 CPU 了，以后发生的事，应用程序完全不知道，它的上下文环境已经被保存到自己的 0 特权级栈中了，那时在 CPU 上运行的程序已经是内核程序了。所以要清楚，内核代码并不是成了应用程序的内核化身，操作系统是独立的部分，用户进程永远不会因为进入内核态而变身为操作系统了。

应用程序是通过系统调用来和操作系统配合完成某项功能的，有人可能会问：我写应用程序时从来没写什么系统调用的代码啊。这是因为你用到的标准库帮你完成了这些事，库中提供的函数其实都已经封装好了系统调用，你需要跟下代码才会看到。其实也可以跨过标准库直接执行系统调用，对于 Linux 系统来说，直接嵌入汇编代码“`int 0x80`”便可以直接执行系统调用，当然要提前设置好系统调用子功能号，该子功能号用寄存器 `eax` 存储。

会不会有人又问，编译器怎么知道系统调用接口是什么，哈哈，您想啊，下载编译器时，是不是要选择系统版本，编译器在设计时也要知道自己将来运行在哪个系统平台上，所以这都是和系统绑定好的，各个操作系统都有自己的系统调用号，编译器厂商在代码中已经把宿主系统的系统调用号写死了，没什么神奇的。

## 0.6 为什么称为“陷入”内核

前面提到了用户进程陷入内核，这个好解释，如果把软件分层的话，最外圈是应用程序，里面是操作系统，如图 0-1 所示。

应用程序处于特权级 3，操作系统内核处于特权级 0。当用户程序欲访问系统资源时（无论是硬件，还是内核数据结构），它需要进行系统调用。这样 CPU 便进入了内核态，也称管态。看图中凹下去的部分，是不是有陷进去的感觉，这就是“陷入内核”。



▲图 0-1 陷入内核

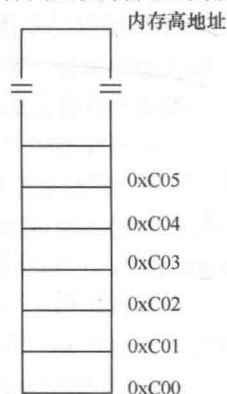
## 0.7 内存访问为什么要分段

按理说咱们应该先看看段是什么，不过了解段是什么之前，先看看内存是什么样子，如图 0-2 所示。

内存按访问方式来看，其结构就如同上面的长方形带子，地址依次升高。为了解释问题更明白，我们假设还在实模式下，如果读者不清楚什么是实模式也不要紧，这并不影响理解段是什么，故暂且先忽略。

内存是随机读写设备，即访问其内部任何一处，不需要从头开始找，只要直接给出其地址便可。如访问内存 `0xC00`，只要将此地址写入地址总线便可。问题来了，分段是内存访问机制，是给 CPU 用的访问内存的方式，只有 CPU 才关注段，那为什么 CPU 要用段呢，也就是为什么 CPU 非得将内存分成一段一段的才能访问呢？

说来话长，现实行业中有很多问题都是历史遗留问题，计算机行业也不能例外。分段是从 CPU 8086 开始的，限于技术和经济，那时候电脑还是非常昂贵的东西，所以 CPU 和寄存器等宽度都是 16 位的，并不是像今天这样寄存器已经扩展到 64 位，当然编译器用的最多的还是 32 位。16 位寄存器意味着其可存储的数字范围是 2 的 16 次方，即 65536 字节，64KB。那时的计算机没有虚拟地址之说，只有物理地址，访问任何存储单元都直接给出物理地址。



▲图 0-2 内存示例

编译器在编译程序时，肯定要根据 CPU 访问内存的规则将代码编译成机器指令，这样编译出来的程序才能在该 CPU 上运行无误，所以说，在直接以绝对物理地址访问内存的 CPU 上运行程序，该程序中指令的地址也必须得是绝对物理地址。总之，要想在该硬件上运行，就要遵从该硬件的规则，操作系统和编译器也无一例外。

若加载程序运行，不管其是内核程序，还是用户程序，程序中的地址若都是绝对物理地址，那该程序必须放在内存中固定的地方，于是，两个编译出来地址相同的用户程序还真没法同时运行，只能运行一个。于是伟大的计算机前辈们用分段的方式解决了这一问题，让 CPU 采用“段基址+段内偏移地址”的方式来访问任意内存。这样的好处是程序可以重定位了，尽管程序指令中给的是绝对物理地址，但终究可以同时运行多个程序了。

什么是重定位呢，简单来说就是将程序中指令的地址改写成另外一个地址，但该地址处的内容还是原

地址处的内容。

CPU 采用“段基址+段内偏移地址”的形式访问内存，就需要专门提供段基址寄存器，这些是 cs、ds、es 等。程序中需要用到哪块内存，只要先加载合适的段到段基址寄存器中，再给出相对于该段基址的偏移地址便可，CPU 中的地址单元会将这两个地址相加后的结果用于内存访问，送上地址总线。

注意，很多读者都觉得段基址一定得是 65536 的倍数（16 位段基址寄存器的容量），这个真的不用，段基址可以是任意的。这就是段可以重叠的原因。

举个例子，看图 0-2，假设段基址为 0xC00，要想访问物理内存 0xC01，就要将用 0xC00: 0x01 的方式来访问才行。若将段基址改为 0xc01，还是访问 0xC01，就要用 0xC01: 0x00 的方式来访问。同样，若想访问物理内存 0xC04，段基址和段内偏移的组合可以是：0xC01: 0x03、0xC02: 0x02、0xC00: 0xC04 等，总之要想访问某个物理地址，只要凑出合适的段基址和段内偏移地址，其和为该物理地址就行了。这时估计有人会问这样行不行，0xC05: -1，能这样提问的同学都是求知欲极强的，可以自己试一下。

说了这么多，我想告诉你的是只要程序分了段，把整个段平移到任何位置后，段内的地址相对于段基址是不变的，无论段基址是多少，只要给出段内偏移地址，CPU 就能访问到正确的指令。于是加载用户程序时，只要将整个段的内容复制到新的位置，再将段基址寄存器中的地址改成该地址，程序便可准确无误地运行，因为程序中用的是段内偏移地址，相对于新的段基址，该偏移地址处的内存内容还是一样的，如图 0-3 所示。

所以说，程序分段首先是为了重定位，我说的是首先，下面还有其他理由呢。

偏移地址也要存入寄存器，而那时的寄存器是 16 位的，也就是一个段最多可以访问到 64KB。而那时的内存再小也有 1MB，改变段基址，由一个段变为另一个段，就像一个段在内存中飘移，采用这种在内存中来回挪位置的方式可以访问到任意内存位置。

所以说，程序分段又是为了将大内存分成可以访问的小段，通过这样变通的方法便能够访问到所有内存了。

但想一想，1M 是 2 的 20 次方，1MB 内存需要 20 位的地址才能访问到，如何做到用 16 位寄存器访问 20 位地址空间呢？

在 8086 的寻址方式中，有基址寻址，这是用基址寄存器 bx 或 bp 来提供偏移地址的，如“mov [bx], 0x5;”指令便是将立即数 0x5 存入 ds: bx 指向的内存。

大家看，bx 寄存器是 16 位的，它最大只能表示 0~0xFFFF 的地址空间，即 64KB，也就是单一的一个寄存器无法表示 20 位的地址空间——1MB。也许有人会说，段基址和段内偏移地址都搞到最大，都为 0xFFFF，对不起，即使不溢出的话，其结果也只是由 16 位变成了 17 位，即两个 n 位的数字无论多大，其相加的结果也超过 n+1 位，因为即使是两个相同的数相加，其结果相当于乘以 2，也就是左移一位而已，依然无法访问 20 位的地址空间。也许读者又有好建议了：CPU 的寻址方式又不是仅仅这一种，上面的限制是因为寄存器是 16 位，只要不全部通过寄存器不就行了吗。既然段寄存器必须得用，那就在偏移地址上下功夫，不要把偏移地址写在寄存器里了，把它直接写成 20 位立即数不就行啦。例如 mov ax, [0x12345]，这样最终的地址是 ds+0x12345，肯定是 20 位，解决啦。不错，这种是直接寻址方式，至少道理上讲得通，这是通过编程技巧来突破这一瓶颈的，能想到这一点我觉得非常 nice。但是作为一个严谨的 CPU，既然宣称支持了通过寄存器来寻址，那就要能够自圆其说才行，不能靠程序员的软实力来克服 CPU 自身的缺陷。于是，一个大胆的想法出现了。

16 位的寄存器最多访问到 64KB 大小的内存。虽然 1MB 内存中可容纳 1MB/64KB=16 个最大段，但这只是可以容纳而已，并不是说可以访问到。16 位的寄存器超过 0xffff 后将会回卷到 0，又从 0 重新开始。20 位宽度的内存地址空间必然只能由 20 位宽度的地址来访问。问题又来了，在当时只有 16 位寄存器的情况下是如何做到访问 20 位地址空间的呢？

这是因为 CPU 设计者在地址处理单元中动了手脚，该地址部件接到“段基址+段内偏移地址”的地址后，自动将段基址乘以 16，即左移了 4 位，然后再和 16 位的段内偏移地址相加，这下地址变成了 20 位了吧，行

内存段重定位演示



▲图 0-3 段的重定位

啦，有了 20 位的地址便可以访问 20 位的空间，可以在 1MB 空间内自由翱翔了。

## 0.8 代码中为什么分为代码段、数据段？这和内存访问机制中的段是一回事吗

首先，程序不是一定要分段才能运行的，分段只是为了使程序更加优美。就像用饭盒装饭菜一样，完全可以将很多菜和米饭混合在一起，或者搅拌成一体，哈哈，但这样可能就没什么胃口啦。如果饭盒中有好多小格子，方便将不同的菜和饭区分存放，这样会让我们胃口大开增加食欲。

x86 平台的处理器是必须要用分段机制访问内存的，正因为如此，处理器才提供了段寄存器，用来指定待访问的内存段起始地址。我们这里讨论的程序代码中的段（用 section 或 segment 来定义的段，不同汇编编译器提供的关键字有所区别，功能是一样的）和内存访问机制中的段本质上是一回事。在硬件的内存访问机制中，处理器要用硬件——段寄存器，指向软件——程序代码中用 section 或 segment 以软件形式所定义的内存段。

分段是必然的，只是在平坦模型下，硬件段寄存器中指向的内存段为最大的 4GB，而在多段模式下编程，硬件段寄存器中指向的内存段大小不一。

对于在代码中的分段，有的是操作系统做的，有的是程序员自己划分的。如果是在多段模型下编程，我们必然会在源码中定义多个段，然后需要不断地切换段寄存器所指向的段，这样才能访问到不同段中的数据，所以说，在多段模型下的程序分段是程序员人为划分的。如果是在平坦模型下编程，操作系统将整个 4GB 内存都放在同一个段中，我们就不需要来回切换段寄存器所指向的段。对于代码中是否要分段，这取决于操作系统是否在平坦模型下。

一般的高级语言不允许程序员自己将代码分成各种各样的段，这是因为其所用的编译器是针对某个操作系统编写的，该操作系统采用的是平坦模型，所以该编译器要编译出适合此操作系统加载运行的程序。由于处理器支持了具有分页机制的虚拟内存，操作系统也采用了分页模型，因此编译器会将程序按内容划分成代码段和数据段，如编译器 gcc 会把 C 语言写出的程序划分成代码段、数据段、栈段、.bss 段、堆等部分。这会由操作系统将编译器编译出来的用户程序中的各个段分配到不同的物理内存上。对于目前咱们用高级语言编码来说，我们之所以不用关心如何将程序分段，正是由于编译器按平坦模型编译，而程序所依赖的操作系统又采用了虚拟内存管理，即处理器的分页机制。像汇编这种低级语言允许程序员为自己的程序分段，能够灵活地编排布局，这就属于人为将程序分成段了，也就是采用多段模型编程。

这么说似乎不是很清楚，一会再用例子和大伙儿解释就明白了。在这之前，先和大家明确一件事。

CPU 是个自动化程度极高的芯片，就像心脏一样，给它一个初始的收缩，它将永远地跳下去。突然想到 Intel 的广告词：给你一颗奔腾的心。

只要给出 CPU 第一个指令的起始地址，CPU 在它执行本指令的同时，它会自动获取下一条的地址，然后重复上述过程，继续执行，继续取址。假如执行的每条指令都正确，没有异常发生的话，我想它可以运行到世界的尽头，能让它停下来的唯一条件就是断电。

它为什么能够取得下一条指令地址？也就是说为什么知道下一条指令在哪里。这是因为程序中的指令都是挨着的，彼此之间无空隙。有同学可能会问，程序中不是有对齐这回事吗？为了对齐，编译器在程序中塞了好多 0。是的，对齐确实是让程序中出现了好多空隙，但这些空隙是数据间的空隙，指令间不存在空隙，下一条指令的地址是按照前面指令的尺寸大小排下来的，这就是 Intel 处理器的程序计数器 cs: eip 能够自动获得下一条指令的原理，即将当前 eip 中的地址加上当前指令机器码的大小便是内存中下一条指令的起始地址。即使指令间有空隙或其他非指令的数据，这也仅仅是在物理上将其断开了，依然可以用 jmp 指令将非指令部分跳过以保持指令在逻辑上连续，我们在后面会通过实例验证这一原理。

为了让程序内指令接连不断地执行，要把指令全部排在一起，形成一片连续的指令区域，这就是代码段。这样 CPU 肯定能接连不断地执行下去。指令是由操作码和操作数组成的，这对于数据也一样，程序运行不仅要有操作码，也得有操作数，操作数就是指程序中的数据。把数据连续地并排在一起存储形成的



段落，就称为数据段。

指令大小是由实际指令的操作码决定的，也就是说 CPU 在译码阶段拿到了操作码后，就知道实际指令所占的大小。其实说来说去，本质上就是在解释地址是怎么来的。这部分在第 3 章中的“什么是地址”节中有详解。

给大家演示个小例子，代码没有实际意义，是我随便写的，只是为方便大家理解指令的地址，代码如下。

#### code\_seg.S

```
1    mov ds,ax
2    mov ax,[var]
3 label:
4    jmp label
5    var dw 0x99
```

本示例一共就 5 行，简单纯粹为演示。将其编译为二进制文件，程序内容是：

```
8E D8 A1 07 00 EB FE 99 00
```

就这 9 个字节的内容，有没有觉得一阵晕炫。如果没有，目测读者兄弟的技术水平远在我之上，请略过本书。

其实这 9 个字节的内容就是机器码。为了让大家理解得更清晰，给大家列个机器码和源码对照表，见表 0-1。

表 0-1 机器码和源码对照表

地 址	机 器 码	源 码
00000000	8ED8	mov ds, ax
00000002	A10700	mov ax, [0x7]
00000005	EBFE	jmp short 0x5
00000007		var dw 0x99
00000008		

表 0-1 第 1 行，地址 0 处的指令是“mov ds, ax”，其机器码是 8ED8，这是十六进制表示，可见其大小是 2 字节。前面说过，下一条指令的地址是按照前面指令的尺寸排下来的，那第 2 行指令的起始地址是 0+2=2。在第 2 行的地址列中，地址确实是 2。这不是我故意写上去的，编译器真的就是这样编排的。第 2 列的指令是“mov ax, [0x7]”（0x7 是变量 var 经过编译后的地址），其机器码是 A10700，这是 3 字节大小。所以第 3 条指令的地址是 2+3=5。后面的指令地址也是这样推算的。程序虽然很短，但麻雀虽小，五脏俱全，完美展示了程序中代码紧凑无缝隙的布局。

现在大伙儿明白为什么 CPU 能源源不断获取到指令了吧，如前所述，原因首先是指令是连续紧凑的，其次是通过指令机器码能够判断当前指令长度，当前指令地址+当前指令长度=下一条指令地址。

上面给出的例子，其指令在物理上是连续的，其实在 CPU 眼里，只要指令逻辑上是连续的就可以，没必要一定得是物理上连续。所以，明确一点，即使数据和代码在物理上混在一起，程序也是可以运行的，这并不意味着指令被数据“断开”了。只要程序中有指令能够跨过这些数据就行啦，最典型的就是用 jmp 跳过数据区。

比如这样的汇编代码：

```
1    jmp start          ;跳转到第三行的 start, 这是 CPU 直接执行的指令
2    var dd 1          ;定义变量 var 并赋值为 1。分配变量不是 CPU 的工作
                          ;汇编器负责分配空间并为变量编址
3    start:           ;标名为 start, 会被汇编器翻译为某个地址
4    mov ax,0         ;将 ax 赋值为 0
```

这几行代码没有实际意义，是为了解释清楚问题，咱们只要关注在第 2 行的定义变量 var 之前为什么要 jmp start。如果将上面的汇编代码按纯二进制编译，如果不加第 1 行的 jmp，CPU 也许会发出异常，显示无效指令，也许不知道执行到哪里去了。因为 CPU 只会执行 cs: ip 中的指令，这两个寄存器记录的是下一条待执行指令的地址，下一个地址 var 处的值为 1，显然我们从定义中看出这只是数据，但指令和数据都是二进制数字，CPU 可分不出这是指令，还是数据。保不准某些“数据”误打误撞恰恰是某种指令也说不定。既然 var 是我们定义的数据，那么必须加上 jmp start 跳过这个 var 所占的空间才可以。