


计 算 机 科 学 丛 书

现代 操作系统

Modern
Operating
Systems

Andrew S. Tanenbaum 著
陈向群 等译

 机械工业出版社
China Machine Press

Prentice Hall

计算机科学丛书

现代操作系统

Andrew S. Tanenbaum 著

陈向群 等译



机械工业出版社
China Machine Press

本书共分两部分，第一部分详尽讲述了传统操作系统知识，包括进程、存储器管理、文件系统、I/O设备管理、死锁等内容；第二部分主要介绍了分布式操作系统，包括层次协议、远程过程调用、互斥操作、分布式文件系统等专题。

为加深概念的理解，本书还详细介绍了四个操作系统，包括两个传统的系统UNIX和MS-DOS；两个分布式系统Amoeba和Mach。此外还简要介绍了NFS、AFS、ISIS等其他几个系统。

本书体系完整、内容丰富、叙述清晰，是大学计算机及相关专业学生不可多得的教科书，对于从事计算机管理、开发、系统分析等职业的专业人员也是优秀的参考书。

Authorized translation from the English language edition published by Prentice Hall.

Copyright © 1999 by Prentice Hall. All rights reserved.

Chinese simplified language edition published by China Machine Press.

Copyright © 2000 by China Machine Press.

本书中文简体字版由美国Prentice Hall公司授权机械工业出版社独家出版。未经出版者书面许可，不得以任何方式复制或抄袭本书内容。

版权所有，侵权必究。

本书版权登记号：图字：01-98-0720

图书在版编目(CIP)数据

现代操作系统/坦尼鲍姆(Tanenbaum, A. S.)著；陈向群等译。—北京：机械工业出版社，1999.11

(计算机科学丛书)

书名原文：Modern Operating Systems

ISBN 7-111-07117-4

I. 现… II. ①坦… ②陈… III. 操作系统 IV. TN316

中国版本图书馆CIP数据核字(1999)第44617号

机械工业出版社(北京市西城区百万庄大街22号 邮政编码 100037)

责任编辑：陈 谊

北京第二外国语学院印刷厂印刷·新华书店北京发行所发行

1999年11月第1版·2000年4月第2次印刷

787mm × 1092mm 1/16 · 32印张

印数：6 001-9 000册

定价：40.00元

凡购本书，如有缺页、倒页、脱页，由本社发行部调换

译者序

当今，在知识经济蓬勃发展的浪潮中，软件产业占据着核心地位，而操作系统则是所有软件的基础，是软件的根本。正是由于这个原因，长久以来世界软件业巨头们从不放弃对操作系统的控制争霸。

操作系统是计算机科学与技术专业的一门专业基础课，是大学学生的必修课程。在学习过程中，教材起着相当大的作用。因此，要学好这门课，首先应该选择一本好的教材。

本书的作者Andrew S. Tanenbaum是一位经验丰富的教授，从事计算机教学多年，在操作系统、计算机网络方面颇有研究，出版了多本教材。本书是作者关于操作系统系列教材中的一本，书中既讨论了传统操作系统的原理，又系统介绍了分布式操作系统的方方面面。本书体系完整、结构合理、内容丰富、叙述清晰、适用面广。美国不少主要大学的计算机软件及相关专业都采用本书作为教材或主要教学参考书。

通过本书的学习，相信学生一定会对操作系统的功能、实现技术有全面的了解。

本书在翻译的过程中，得到很多人的支持帮助。参加翻译工作的还有张杰、张乃琳、叶松、李茹、翁念龙、王梦秋、何新明、高卫、胡英琨、陈宇、朱克耀，北京大学计算机科学技术系的部分学生也参加了翻译工作，其中徐颖、高磊、刘西川、潘颖、章苏、董方鹏等做了很多工作。

译者

前 言

过去，大多数的计算机都是独立工作的，大多数的操作系统也设计为在单处理机上运行。然而，这种状况正在发生迅速的变化：许多计算机都联网工作，因而分布式操作系统也变得越来越重要。作为操作系统课程的教材，本书的特别之处在于认识到了这个变化，并将分布式操作系统提高到与传统的单CPU操作系统同等地位来讲述。

在过去的15年中，我参与设计和实现了三个不同的操作系统：TSS-11(PDP-11)、MINIX(IBM PC、Atari、Amiga、Macintosh和SPARC)，还有Amoeba(80386、Sun-3、SPARC和VAX)。经过这些实践，我认识到应该在教材中强调那些在实际的操作系统中真正有用的东西。本书包含了所有在本科生的操作系统课程中应该包含的内容，包括进程、进程间通信、信号量、管程、消息传递、经典IPC问题、调度、交换技术、虚拟存储技术、分页算法、段式管理、文件系统、安全性、保护机制、I/O硬件与软件、死锁等。不过，本书分配给各专题的篇幅与一般的教材不尽相同，因为我认为学生应该学习那些在实践中真正有用的概念，而不仅仅是寻求理论上的完善。例如，CPU调度只占了一节，而不是整个一章。许多著作中都提出并分析了很多复杂的调度算法，而实际的操作系统中采用的一般都是最简单的优先级法或轮转法。

本书的第二部分主要介绍分布式系统。我们先用一章的篇幅简单介绍了分布式系统的软硬件配置，然后依次介绍各个专题：层次协议、客户机-服务器(Client-Server)模型、远程过程调用、组内通信、时钟同步、互斥操作、选举算法、原子事务、线程、分布式文件系统等等。在一门本科生课程中包含这么多分布式系统的内容可能有些特别，不过计算机行业发展得如此迅速，分布式系统会很快成为操作系统的主流，因此每个学生都应该熟悉它。

为了加深对这些基本概念的理解，本书还详细介绍了四个操作系统，两个传统的，两个分布式的。其中两个传统的，即非分布式的系统是UNIX和MS-DOS；分布式操作系统是Amoeba和Mach。这四个例子每个各占一章，除此之外，还简要介绍了NFS、AFS、ISIS等其他几个系统。学完本书之后，读者会对操作系统有一个较完整的了解。

在此，我想指出，开始时我只想修订我的一本旧作“Operating System: Design and Implementation”，将其中有关MINIX的内容去掉，使之适合传统的“理论”教学，而不是“实验”课程。在修订的过程中，我意识到分布式系统的重要性，因而又加了七章这方面的内容。我还想在近期内修订那本MINIX的书，这样就会有两本更新后的书，一本适用于理论学习，一本适用于实践。

在本书的写作过程中，很多朋友审阅了初稿并提出了宝贵的意见，然而，本书还是存在着一些错误和不妥之处。错误似乎是不可避免的，无论有多少人审阅过书稿。欢迎发现错误的读者通过E-mail与我联系，地址是 ast@CS.VU.NL。另外，对Amoeba系统(见14章)感兴趣的公司或大学也请通过E-mail与我联系。

Andrew S. Tanenbaum (安德鲁·S. 坦尼鲍姆)

目 录

译者序
前言

第一部分 传统操作系统

第1章 引言	1
1.1 什么是操作系统	2
1.1.1 作为扩展机器的操作系统	2
1.1.2 作为资源管理器的操作系统	3
1.2 操作系统历史	3
1.2.1 第一代计算机(1945~1955): 真空管和插件板	3
1.2.2 第二代计算机(1955~1965): 晶体管和批处理系统	4
1.2.3 第三代计算机(1965~1980): 集成电路芯片和多道程序	5
1.2.4 第四代计算机(1980~1990): 个人计算机	7
1.3 操作系统基本概念	8
1.3.1 进程	8
1.3.2 文件	9
1.3.3 系统调用	11
1.3.4 外壳	12
1.4 操作系统结构	12
1.4.1 整体式系统	12
1.4.2 层次式系统	14
1.4.3 虚拟机	14
1.4.4 客户机/服务器系统	15
1.5 本书其他部分的概要	16
1.6 小结	17
习题	18
第2章 进程	19
2.1 进程介绍	19
2.1.1 进程模型	19
2.1.2 进程的实现	22
2.2 进程间通信	23
2.2.1 竞争条件	23

2.2.2 临界区	24
2.2.3 忙等待的互斥	24
2.2.4 睡眠与唤醒	27
2.2.5 信号量	29
2.2.6 事件计数器	31
2.2.7 管程	32
2.2.8 消息传递	34
2.2.9 原语等价	36
2.3 经典的IPC问题	39
2.3.1 哲学家进餐问题	39
2.3.2 读者-写者问题	41
2.3.3 睡眠的理发师问题	42
2.4 进程调度	44
2.4.1 时间片轮转调度	45
2.4.2 优先级调度	45
2.4.3 多级队列	46
2.4.4 最短作业优先	47
2.4.5 保证调度算法	48
2.4.6 策略与机制	48
2.4.7 两级调度法	48
2.5 小结	49
习题	49
第3章 存储管理	52
3.1 无交换或分页的存储管理	52
3.1.1 无交换或分页的单道程序	52
3.1.2 多道程序设计和内存利用率	53
3.1.3 固定分区的多道程序	55
3.2 交换	56
3.2.1 可变分区的多道程序设计	57
3.2.2 使用位图的内存管理	58
3.2.3 使用链表的内存管理	59
3.2.4 伙伴式的内存管理	60
3.2.5 交换空间的分配	61
3.2.6 交换系统分析	62
3.3 虚拟存储器	62
3.3.1 分页	63

3.3.2 页表	65	4.3.1 实现文件	113
3.3.3 分页硬件示例	67	4.3.2 实现目录	115
3.3.4 相联存储器	72	4.3.3 共享目录	117
3.4 页面置换算法	75	4.3.4 磁盘空间管理	118
3.4.1 最优页面置换算法	75	4.3.5 文件系统的可靠性	121
3.4.2 最近未使用页面置换算法	76	4.3.6 文件系统性能	124
3.4.3 先进先出页面置换算法	76	4.4 安全性	126
3.4.4 第二次机会页面置换算法	77	4.4.1 安全环境	126
3.4.5 时钟页面置换算法	77	4.4.2 著名的安全缺陷	127
3.4.6 最久未使用页面置换算法	78	4.4.3 Internet蠕虫	129
3.4.7 用软件模拟LRU	79	4.4.4 一般的安全性攻击	130
3.5 分页算法模型	80	4.4.5 安全性的设计原则	131
3.5.1 Belady异常现象	80	4.4.6 用户验证	132
3.5.2 栈式算法	81	4.5 保护机制	134
3.5.3 距离字符串	82	4.5.1 保护域	134
3.5.4 缺页率预测	83	4.5.2 存取控制表	136
3.6 分页系统中的设计问题	83	4.5.3 权限	137
3.6.1 工作集模型	83	4.5.4 保护模型	138
3.6.2 局部与全局分配策略	84	4.5.5 隐藏通道	139
3.6.3 页面尺寸	86	4.6 小结	140
3.6.4 实现时涉及的问题	86	习题	140
3.7 分段	89	第5章 I/O设备管理	143
3.7.1 纯分段系统的实现	91	5.1 I/O硬件组成原理	143
3.7.2 分段和分页结合: MULTICS	92	5.1.1 I/O设备	143
3.7.3 分段和分页结合: Intel的386	95	5.1.2 设备控制器	144
3.8 小结	98	5.1.3 直接存储器存取	145
习题	98	5.2 I/O软件原理	146
第4章 文件系统	101	5.2.1 I/O软件的目标	146
4.1 文件	101	5.2.2 中断处理程序	147
4.1.1 文件命名	101	5.2.3 设备驱动程序	148
4.1.2 文件结构	102	5.2.4 与设备无关的I/O软件	148
4.1.3 文件类型	103	5.2.5 用户空间的I/O软件	149
4.1.4 文件存取	105	5.3 磁盘	151
4.1.5 文件属性	105	5.3.1 磁盘硬件	151
4.1.6 文件操作	106	5.3.2 磁盘调度算法	151
4.1.7 存储映像文件	109	5.3.3 磁盘出错处理	153
4.2 目录	110	5.3.4 “每次一道”高速缓冲	154
4.2.1 层次目录系统	110	5.3.5 RAM 盘	154
4.2.2 路径名	111	5.4 时钟	155
4.2.3 目录操作	113	5.4.1 时钟硬件	155
4.3 文件系统的实现	113	5.4.2 时钟软件	156

5.5 终端	158	7.2.1 设计目标	189
5.5.1 终端硬件	158	7.2.2 接口	190
5.5.2 存储映像终端	159	7.2.3 登录	190
5.5.3 输入软件	160	7.2.4 shell	191
5.5.4 输出软件	163	7.2.5 文件和目录	193
5.6 小结	164	7.2.6 应用程序	194
习题	165	7.3 UNIX的一些基础概念	195
第6章 死锁	167	7.3.1 进程	196
6.1 资源	167	7.3.2 内存管理模式	199
6.2 死锁定义	168	7.3.3 文件系统	200
6.2.1 死锁的条件	168	7.3.4 I/O设备	203
6.2.2 死锁模型	169	7.4 UNIX的系统调用	205
6.3 鸵鸟算法	170	7.4.1 进程管理系统调用	205
6.4 死锁检测和恢复	171	7.4.2 内存管理系统调用	208
6.4.1 单种资源类型下的死锁检测	171	7.4.3 有关文件和目录的系统调用	208
6.4.2 多种资源类型下的死锁检测	173	7.4.4 I/O系统调用	209
6.4.3 从死锁恢复	174	7.5 UNIX的实现	209
6.5 死锁避免	175	7.5.1 进程的实现	210
6.5.1 资源轨迹图	176	7.5.2 内存管理的实现	212
6.5.2 安全和不安全状态	177	7.5.3 文件系统的实现	215
6.5.3 单种资源的银行家算法	178	7.5.4 I/O实现	217
6.5.4 多种资源的银行家算法	178	7.6 小结	218
6.6 死锁预防	179	习题	219
6.6.1 破坏互斥条件	180	第8章 实例研究2: MS-DOS	221
6.6.2 破坏占有和等待条件	180	8.1 MS-DOS的历史	221
6.6.3 破坏不可剥夺条件	180	8.1.1 IBM PC	221
6.6.4 破坏循环等待条件	180	8.1.2 MS-DOS 1.0	222
6.7 其他问题	181	8.1.3 MS-DOS 2.0	223
6.7.1 两阶段加锁	181	8.1.4 MS-DOS 3.0	223
6.7.2 非资源死锁	182	8.1.5 MS-DOS 4.0	224
6.7.3 饥饿	182	8.1.6 MS-DOS 5.0	224
6.8 小结	182	8.2 MS-DOS 概述	225
习题	183	8.2.1 使用	226
第7章 实例研究1: UNIX	185	8.2.2 shell	228
7.1 UNIX的历史	185	8.2.3 配置	229
7.1.1 UNICS	185	8.3 MS-DOS 的基本概念	230
7.1.2 PDP-11 UNIX	186	8.3.1 进程	230
7.1.3 可移植的UNIX	186	8.3.2 内存模式	233
7.1.4 伯克利UNIX	187	8.3.3 文件系统	240
7.1.5 UNIX的标准化	187	8.3.4 I/O设备	241
7.2 UNIX概述	189	8.4 MS-DOS的系统调用	242

8.4.1 进程管理系统调用	243	10.1.3 网络层	282
8.4.2 内存管理系统调用	243	10.1.4 传输层	283
8.4.3 文件与目录系统调用	244	10.1.5 会话层	283
8.4.4 I/O系统调用	244	10.1.6 表示层	283
8.5 MS-DOS的实现	244	10.1.7 应用层	284
8.5.1 进程的实现	245	10.2 客户-服务器模型	284
8.5.2 内存管理的实现	246	10.2.1 客户和服务端	284
8.5.3 文件系统的实现	247	10.2.2 一个客户和服务器的实例	285
8.5.4 I/O的实现	250	10.2.3 寻址	288
8.6 小结	252	10.2.4 阻塞与非阻塞原语	289
习题	253	10.2.5 缓冲和非缓冲原语	291
		10.2.6 可靠和不可靠原语	293
		10.2.7 实现客户-服务器模型	294
		10.3 远程过程调用	295
		10.3.1 基本的RPC操作	296
		10.3.2 参数传递	298
		10.3.3 动态联编	301
		10.3.4 出现差错时的RPC语义	303
		10.3.5 与实现相关的论题	306
		10.3.6 问题域	313
		10.4 组通信	315
		10.4.1 组通信入门	315
		10.4.2 与设计有关的问题	316
		10.4.3 ISIS中的组通信	322
		10.5 小结	324
		习题	325
		第11章 分布式系统中的同步	327
		11.1 时钟同步	327
		11.1.1 逻辑时钟	328
		11.1.2 物理时钟	331
		11.1.3 时钟同步算法	333
		11.2 互斥	336
		11.2.1 集中式算法	336
		11.2.2 分布式算法	337
		11.2.3 令牌环算法	339
		11.2.4 三种算法的比较	340
		11.3 选举算法	341
		11.3.1 Bully算法	341
		11.3.2 环算法	342
		11.4 原子事务	343
		11.4.1 原子事务简介	343
第9章 分布式操作系统概述	255		
9.1 目标	255		
9.1.1 分布式系统相对于集中式系统的优点	255		
9.1.2 分布式系统相对于独立PC机的优点	256		
9.1.3 分布式系统的缺点	257		
9.2 硬件上的概念	258		
9.2.1 总线型多处理机	259		
9.2.2 交换型多处理机	260		
9.2.3 总线型多计算机	261		
9.2.4 交换型多计算机	262		
9.3 软件上的概念	262		
9.3.1 网络操作系统和网络文件系统	263		
9.3.2 真正的分布式系统	269		
9.3.3 多处理机分时系统	269		
9.4 设计上的考虑	271		
9.4.1 透明性	271		
9.4.2 灵活性	273		
9.4.3 可靠性	274		
9.4.4 性能	275		
9.4.5 可扩展性	276		
9.5 小结	277		
习题	277		
第10章 分布式系统中的通信问题	279		
10.1 分层协议	279		
10.1.1 物理层	281		
10.1.2 数据链路层	282		

11.4.2 事务模型	344	13.2.6 经验教训	411
11.4.3 实现	347	13.3 分布式文件系统的发展趋势	412
11.4.4 并发控制	350	13.3.1 新的硬件	412
11.5 分布式系统中的死锁	352	13.3.2 可伸缩性	414
11.5.1 分布式死锁的检测	353	13.3.3 广域网	414
11.5.2 分布式死锁的预防	356	13.3.4 移动用户	415
11.6 小结	357	13.3.5 容错	415
习题	357	13.4 小结	415
第12章 分布式系统中的进程及 处理器	359	习题	416
12.1 线程	359	第14章 实例研究3: AMOEBa	418
12.1.1 线程的引入	359	14.1 Amoeba简介	418
12.1.2 线程使用	360	14.1.1 Amoeba的历史	418
12.1.3 线程包的设计问题	362	14.1.2 研究目的	418
12.1.4 线程包的实现	365	14.1.3 Amoeba的体系结构	419
12.1.5 线程和RPC	367	14.1.4 Amoeba微内核	420
12.1.6 一个线程包的例子	368	14.1.5 Amoeba服务程序	422
12.2 系统模型	371	14.2 Amoeba中的对象与权限字	422
12.2.1 工作站模型	371	14.2.1 权限字	423
12.2.2 使用空闲工作站	373	14.2.2 对象的保护	424
12.2.3 处理器池模型	376	14.2.3 标准操作	425
12.2.4 混合模型	378	14.3 Amoeba中的进程管理	425
12.3 处理器分配	378	14.3.1 进程	426
12.3.1 分配模型	378	14.3.2 线程	427
12.3.2 处理器分配算法的设计原则	380	14.4 Amoeba的内存管理	428
12.3.3 处理器分配算法的实现问题	381	14.4.1 段的管理	428
12.3.4 处理器分配算法实例	382	14.4.2 段的映像	429
12.4 分布式系统中的调度	386	14.5 Amoeba的通信	429
12.5 小结	387	14.5.1 远程过程调用	430
习题	387	14.5.2 Amoeba的组通信	432
第13章 分布式文件系统	389	14.5.3 快速局域网协议	438
13.1 分布式文件系统设计	389	14.6 Amoeba 服务器	443
13.1.1 文件服务的接口	389	14.6.1 文件服务器	443
13.1.2 目录服务器接口	391	14.6.2 目录服务器	446
13.1.3 文件共享的语义	394	14.6.3 复制服务器	449
13.2 分布式文件系统的实现	396	14.6.4 运行服务器	449
13.2.1 文件的使用	396	14.6.5 引导服务器	451
13.2.2 系统结构	397	14.6.6 TCP/IP服务器	451
13.2.3 缓存区处理	400	14.6.7 其他服务器	451
13.2.4 复制性	404	14.7 小结	451
13.2.5 例子: Andrew文件系统	407	习题	452
		第15章 实例研究4: Mach操作系统	453

15.1 概述	453	15.4.1 端口	470
15.1.1 Mach的历史	453	15.4.2 发送和接收消息	474
15.1.2 Mach的设计目标	454	15.4.3 网络信息服务器	477
15.1.3 Mach的微内核	454	15.5 Mach 中的BSD UNIX 仿真	479
15.1.4 Mach 的BSD UNIX服务器	455	15.6 Amoeba与Mach的比较	480
15.2 Mach 中的进程管理	456	15.6.1 基本原理	480
15.2.1 进程	456	15.6.2 对象	481
15.2.2 线程	458	15.6.3 进程	481
15.2.3 调度	460	15.6.4 存储模式	482
15.3 Mach的存储管理	462	15.6.5 通信	482
15.3.1 虚拟存储	462	15.6.6 服务程序	483
15.3.2 内存共享	464	15.7 小结	484
15.3.3 外部存储管理器	466	习题	484
15.3.4 Mach的分布式共享内存	469	附录A 阅读材料及参考文献	486
15.4 MACH中的通信	469	附录B C语言简介	494

第一部分 传统操作系统

第1章 引言

离开了软件的计算机就成为一堆废铜烂铁。有了软件，计算机可以对信息进行存储、处理和检索，检查文档拼写错误、玩探险游戏以及处理许多其他有意义的事务。计算机软件大致分为两类：系统软件和应用软件。系统软件用于管理计算机本身及应用程序；应用软件实现用户所需要的功能。操作系统(operating system)是最基本的系统软件，它控制计算机的所有资源并提供应用程序开发的基础。

现代计算机系统包含一个或多个处理器、若干主存(常称为“磁芯存储器”，尽管磁芯存储器早在十多年前就不再使用了)、时钟、终端、磁盘、网络接口及其他输入/输出设备，总而言之是一个非常复杂的系统。为了正确地使用它们，需要编写一个程序来管理所有这些部件，编写这样的程序，即使不考虑优化也是一件很困难的事情。如果每个程序员都必须处理诸如磁盘如何工作，以及每读一个磁盘块时，可能导致操作出错的有几十种因素等硬件控制问题，那么很多程序几乎没法写成。

在许多年以前，人们就认识到必须要找到某种方法，将程序员从复杂的硬件控制中脱离出来。解决的方法是在硬件的基础之上加载一层软件来管理整个系统。并为用户提供一个接口或虚拟机(Virtual machine)，从而更容易理解和进行程序设计。这一层软件就是操作系统，也就是本书所论述的主题。

这一种处理方式如图1-1所示。其底层是硬件，它本身可能由两层或多层构成。最底一层是物理设备，包括集成电路芯片、连线、电源、阴极射线管和相关的物理装置。至于它们是如何构造的和如何工作的则属于电子工程师的范围。

其次是直接控制设备并向上一层提供更清晰的一个接口的很原始的软件。该软件被称为微程序(microprogram)，通常存放在只读存储器中，它实际上是一个解释器，先取得诸如ADD，MOVE和JUMP等机器语言指令，然后通过一序列的细小动作执行这些指令。例如，为了执行一条ADD指令，微程序必须先确定要做加法的数据的位置，取数，然后相加，最后在某处存放结果。由微程序解释执行的这一套指令集称为机器语言(machine language)。机器语言并不真正是硬件的组成部分，但硬件制造商通常在手册中给出有关机器语言的完整描述，所以不少人将它看作真正的“计算机”。在有些机器里，微程序由硬件实现，它并不是明晰的一层。

典型的机器语言有50到300条指令，大多数在机器里从事数据传送、算术运算和数值比较

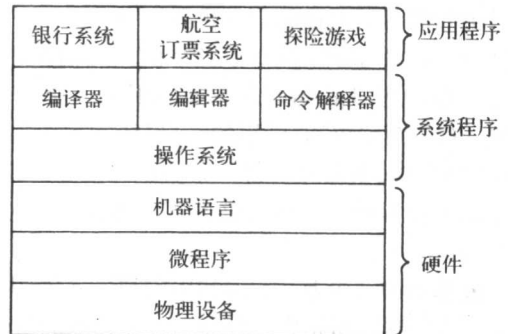


图1-1 计算机系统由硬件、系统程序和应用程序组成

等操作。在这个层次上，通过向特定的设备寄存器(device register)写入数值来控制输入/输出设备。例如，硬盘读操作可通过把磁盘地址、内存地址、字节数和操作类型(读/写)等值装入特定的寄存器来完成。实际上，往往需要更多的操作参数，而操作完成后设备返回的状态也非常复杂。进一步说，对许多I/O设备而言，在程序设计中时序的作用非常重要。

操作系统的主要功能之一就是将这些所有的复杂性隐藏起来，并为程序员的工作提供更加方便的一套指令。例如，在概念上“从文件中读数据块”比考虑“移动磁头，等待放置”之类的细节来得简单、方便。

在操作系统的顶层是其他系统软件。其中有命令解释器(shell)、编译器、编辑器及类似的独立于应用的程序。重要的是这些程序本身并不是操作系统的组成部分，尽管通常由计算机厂商提供它们。这一点很重要，也很微妙。操作系统是专指在核心态(kernel mode)，或称管态(supervisor mode)下运行的软件，它受硬件保护而免遭用户的篡改(不考虑某些老式微处理器，它们根本没有硬件的保护)。编译器和编辑器运行在用户态(user mode)。如果某一个编译器用户不喜欢，他尽可以自己重写一个，但是用户却不可以自己换一个磁盘中断处理程序，因为这是操作系统的一部分，通常由硬件保护，以防止用户试图对它进行修改。

最后，在系统软件的上面是应用软件。这些软件是用户编写的，用来解决诸如商业数据处理、工程计算或者电子游戏等特定的问题。

1.1 什么是操作系统

多数计算机用户都有一些使用操作系统的体验，但要给出操作系统的准确定义却很困难。部分原因是操作系统执行两个相对独立的任务，并取决于从什么角度看待，读者多半听说过其中一个或另一个的功能。下面我们逐项进行讨论。

1.1.1 作为扩展机器的操作系统

正如前面所叙述的，在机器语言一级上，多数计算机的体系结构(architecture)(指令集、存储组织、I/O和总线结构)是原始的而且编程是很困难的，尤其是输入输出操作。要更细致地考察这一点，可以简要地考虑如何用NEC PD765控制器芯片来进行软盘I/O操作，这个芯片是IBM PC和多数PC机采用的。(尽管在本书中，我们交叉使用“软盘”与“磁盘”这两个名词)PD765有16条命令，每一条命令向一个设备寄存器装入长度从1到9字节的特定数据。这些命令用于读写数据、移动磁头臂、格式化磁道，以及初始化、检测状态、复位、校准控制器及设备。

最基本的命令是READ和WRITE。它们均需要13个参数，所有这13个参数被封装在9个字节中。这些参数所指定的信息有：欲读取的磁盘块地址、磁道的扇区数、物理介质的记录格式、扇区间隙，以及对已删除数据地址标识的处理方法等。如果不懂这些晦涩的指令也没关系，它本身相当深奥。当操作结束时，控制器芯片在7个字节中返回23个状态及出错字段。这样好象还不够，编写软盘控制的程序员还要保持注意步进电机的开关状态。如果电机关闭则在读写数据前要先启动它(有一段较长的启动延时)。电机不能长时间处于开启状态，否则将损坏软盘。程序员还必须在较长的启动延迟和可能对软盘造成损坏(和丢失数据)之间作出折衷。

很清楚，一般程序员并不想涉足软盘编程的这些具体细节(也包括硬盘，它也非常复杂且与软盘不同)。程序员需要的是一种简单的，高度抽象的可以与之打交道的设备。对磁盘而言，典型的抽象是包含了一组已命名文件的一张磁盘。每个文件可以打开、读写，最后关闭。至于其中的一些细节，如是否应使用修改的频率模块化数据记录格式、当前步进电机的状态等则在

对用户的抽象中隐藏。

将硬件细节与程序员隔离开来，并提供一个可以读写的命名文件的简洁的方式，当然就是操作系统。与操作系统屏蔽了磁盘硬件，提供了一个简明的、面向文件的接口类似，其他许多包括定时器、存储器管理等低层硬件的特性也被隐藏了。在每种情况中，操作系统提供的抽象都较底层硬件本身更简单、更易使用。

从这个角度看，操作系统的作用是为用户提供一台等价的扩展机器(extended machine)，或称虚拟机，它比底层硬件更容易编程。至于操作系统是如何做到这一点的，这正是本书所讨论的内容。

1.1.2 作为资源管理器的操作系统

把操作系统看作是提供给用户的基本的方便接口的概念是一种自顶向下的观点。按照另一种自底向上的观点，操作系统则用来管理一个复杂系统的各个部分。现代计算机包含处理器、存储器、时钟、磁盘、终端、磁带设备、网络接口、激光打印机以及许多其他设备。从这个角度看，操作系统的任务是在相互竞争的程序之间有序地控制对处理器、存储器以及其他I/O接口设备的分配。

假设在一台计算机上运行的三个程序试图同时在同一台打印机上输出计算结果。那么头几行可能是程序1的输出，下几行是程序2的输出，然后又是程序3的输出等等。最终结果将是一团糟。操作系统采用将打印输出送到磁盘上缓冲区的方法，可以避免这种混乱。在一个程序结束后，操作系统可以将暂存在磁盘上的文件送到打印机输出，同时其他程序可以继续产生新的输出结果，很明显，这些程序的输出并没有真正送到打印机。

当计算机有多个用户时，特别是用户还需要共享诸如磁带机和照片输出机等昂贵的资源。就更有必要管理和保护存储器、I/O设备以及其他设备。经济考虑是一方面，还要让用户共享信息。总之，这种观点认为操作系统的主要任务是跟踪谁在使用什么资源、满足资源请求、记录使用状况，以及协调各程序和用户对资源使用请求的冲突。

1.2 操作系统历史

操作系统经历了多年的发展过程，对此我们进行一下简要的回顾。在历史上由于操作系统与其运行的计算机系统结构联系非常密切，我们将按照计算机的换代历程介绍操作系统的发展。

第一台真正的数字计算机是英国数学家Charles Babbage(1792-1871)设计的。尽管Babbage为建造他的“分析机”投入了毕生精力，但却没能让它成功地运行起来，因为它是纯机械式的，当时的技术不可能使分析机的零部件达到需要的精度。当然，这个分析机没有操作系统。

1.2.1 第一代计算机(1945~1955): 真空管和插件板

从Babbage失败之后一直到第二次世界大战，数字计算机的建造几乎没有什么进展。40年代中期，哈佛大学的Howard Aiken、普林斯顿高等研究院的John Neumann(冯·诺依曼)、宾夕法尼亚大学的J.Presper Eckert和William Mauchley、德国的Konrad Zuse、以及其他一些人，都使用真空管成功地建造了运算机器。这些非常巨大的机器，使用了数万个真空管，占据了几个房间，然而其运算速度却比现在最便宜的家用计算机还要慢得多。

在早期，都有一个小组专门设计、制造、编程、操作和维护每台机器。程序设计全部采用

机器语言，通过在一些插板上的硬连线来控制其基本功能。没有程序设计语言(甚至没有汇编语言)，操作系统更是前所未闻。使用机器的方式是程序员提前在墙上机时表上预约一段时间，然后到机房将他的插件板插到计算机里，在接下来的几小时里计算自己的题目，期盼着在这段时间中，二万多个真空管没有烧断的。这时实际上所有的题目都是数值计算问题，例如计算正弦和余弦函数表。

到了50年代早期，有了改进，出现了穿孔卡片，这时就可以将程序写在卡片上然后读入计算机而不用插板，但其他过程则依然如旧。

1.2.2 第二代计算机(1955~1965): 晶体管和批处理系统

50年代晶体管的发明极大地改变了整个状况。计算机已经很可靠，厂商可以成批地生产计算机并卖给用户。用户可以指望计算机长时间运行，完成一些有用的工作。此时，设计人员、生产人员、操作人员、程序人员和维护人员之间第一次有了明确的分工。

这些计算机安装在专门的空调房间里，由专人操作。只有少数大公司、主要的政府部门或大学才承受得了数百万美元的标价。要运行一个作业(job，即一个或一组程序)，程序员首先将程序写在纸上(用FORTRAN语言或汇编语言)，然后穿孔成卡片。再将卡片盒带到输入室，交给操作员。

计算机运行完当前任务后，其计算结果从打印机上输出，操作员到打印机上撕下运算结果并送到输出室，稍后程序员就可取到结果。然后，操作员从已送到输入室的卡片盒中读入另一个任务。如果需要FORTRAN编译器，操作员还要从文件柜把它取来读入计算机。许多机时被操作员在机房里走来走去而浪费掉了。

由于当时计算机非常昂贵，人们很自然地想办法减少机时的浪费。通常采用的方案就是批处理系统(batch system)。其思想是：在输入室收集全部的作业，然后用一台相对便宜的计算机，将它们读到磁带上，如IBM 1401计算机，它适用于读卡片、拷贝磁带和输出打印，但不适用于作数值运算。另外用较昂贵的计算机，如IBM 7094来完成真正的计算。情况如图1-2所示。

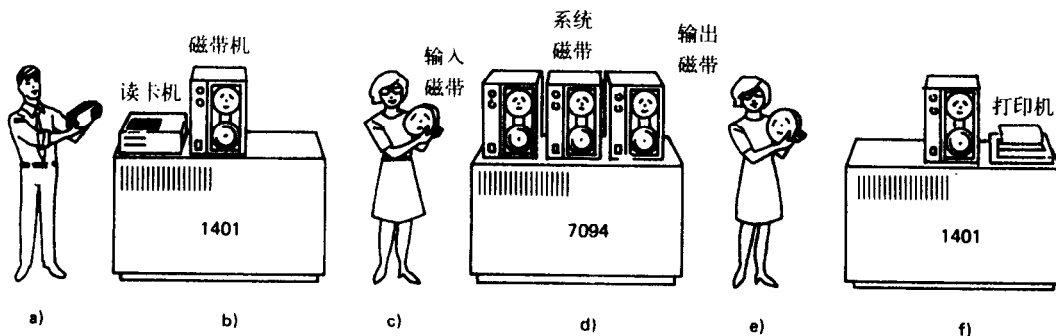


图1-2 一种早期的批处理系统

- a)程序员将卡片拿到1401机处 b)1401机将批处理作业读到磁带上 c)操作员将输入带送至7094机
d)7094机进行计算 e)操作员将输出带送到1401机 f)1401机打印输出

在一个小时的批量作业收集之后，磁带倒转后被送到机房里并装到磁带上。随后操作员装入一个特殊的程序(现代操作系统的前身)，从磁带上它读入第一个作业并运行。其输出写到第二盘磁带上，并不打印。每个作业结束后，操作系统自动地从磁带上读入下一个作业并运行。当一批作业完全结束后，操作员取下输入和输出磁带，将输入磁带换成下一批作业，并把输出

磁带拿到一台1401机器上进行脱机(off line, 即, 不与主计算机联机)打印。

典型的输入作业结构如图1-3所示。一开始是张\$JOB卡片, 它标识出所需的最大运行时间(分钟)、计费帐号, 以及程序员的名字。接着是\$FORTRAN卡片, 通知操作系统从系统磁带上装入FORTRAN语言编译器。之后是待编译的源程序, 然后是\$LOAD卡片, 通知操作系统装入编译好的目标程序(目标程序通常保存在临时磁带上以便于马上使用)。接着是\$RUN卡片, 告诉操作系统运行该程序并使用随后的数据。最后, \$END卡片标识作业结束。这些基本的控制卡片是现代作业控制语言和命令解释器的先驱。

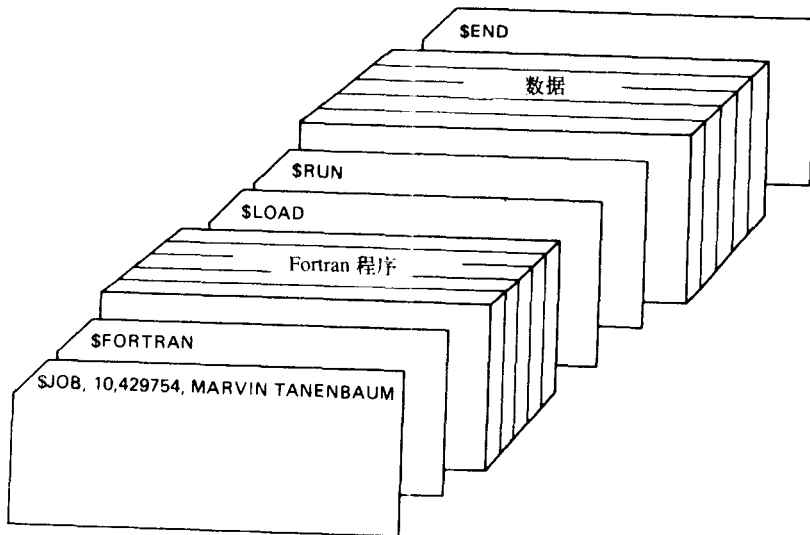


图1-3 典型的FMS作业结构

第二代计算机主要用于科学与工程计算, 例如解偏微分方程。这些题目大多用FORTRAN语言和汇编语言编写。典型的操作系统是FMS(FORTRAN Monitor System, FORTRAN监控系统)和IBSYS(IBM为7094机配备的操作系统)。

1.2.3 第三代计算机(1965~1980): 集成电路芯片和多道程序

60年代初期, 大多数计算机厂商都有两条完全不同并且不兼容的生产线。一条是基于字的大型科学用计算机, 诸如IBM 7094, 主要用于科学和工程计算。另一条是基于字符的商用计算机, 诸如IBM 1401, 银行和保险公司主要用它从事磁带归档和打印服务。

对厂商来说, 开发和维护两种完全不同的产品是昂贵的。另外, 许多新的计算机用户一开始时只需要一台小计算机, 后来可能需要一台较大的计算机, 而且希望能够更快地执行原有程序。

IBM公司试图通过引入System/360来一次性地解决这两个问题。360是一个软件兼容的计算机系列, 其低档机与1401相当, 高档机则比7094功能强很多。这些计算机只在价格和性能(最大存储器容量、处理器速度、允许的I/O设备数量等)上有差异。由于所有的计算机都有相同的体系结构和指令集, 因此在理论上, 为一种型号机器编写的程序可以在其他所有型号的机器上运行。而且360被设计成既可用于科学计算, 又可用于商业计算, 这样一个系列的计算机可以满足所有用户的要求。在随后的几年里, IBM使用更现代的技术陆续推出了360的后续机型, 如著名的370、4300、3080和3090系列。

360是第一个采用小规模集成电路的主流机型, 与采用分立晶体管制造的第二代计算机相

比，其性能价格比有很大提高。360很快就获得了成功，其他主要厂商也很快采纳了系列兼容机的思想。这些计算机的后代目前仍在大型的计算中心里使用。

“单一家族”思想的最大优点同时也是其最大的缺点。原因所有的软件，包括操作系统，都要能够在所有机器上运行，从小的代替1401把卡片拷贝到磁带上的机器，到用于替代7094进行气象预报及其他繁重计算的大型机；从只能带很少外部设备的机器到有很多外设的机器；从商业领域到科学计算领域等。总之，它要有效地适用于所有的用途。

IBM(别人也不行)无法写出同时满足这些需求相互冲突的软件。其结果是一个庞大的又极其复杂的操作系统诞生了，它比FMS大了约2~3数量级规模。其中包含有数千名程序员写的数百万行汇编语言代码，不可避免地也有成千上万处错误，这就不断地导致发行新的版本试图更正这些错误。每个新版本在订正老错误的同时又引入新错误，所以随着时间的流逝，错误的数量可能大致保持不变。

OS/360的设计者之一Fred Brooks 后来写过一本书(Brooks, 1975)描述他在开发OS/360过程中的经验。我们不可能在这里复述该书的全部内容，不过其封面是一群史前动物陷入泥潭而不能自拔。1991年出版的Silberschatz等人的著作的封面也表达了类似的观点。

抛开OS/360的庞大和存在的问题，OS/360和其他公司的类似的第三代操作系统的确合理地满足了大多数用户的要求。同时，它们也使第二代操作系统所缺乏的几项关键技术得到广泛应用。其中最重要的可能是多道程序设计(multiprogramming)。在7094机上，若当前作业因等待磁带或其他I/O而暂停时，CPU就只能简单地踏步直至该I/O完成。对于CPU操作密集的科学计算问题，I/O操作较少，因此浪费的时间很少。然而对于商业数据处理，I/O操作等待的时间通常占到80%~90%，所以必须采取某种措施减少CPU时间的浪费。

解决办法是将内存分几个部分，每一部分存放不同的作业，如图1-4所示。当一个作业等待I/O操作完成时，另一个作业可以使用CPU。如果内存中可以存放足够多的作业，则CPU利用率可以接近100%。在主存中同时驻留多个作业需要特殊的硬件来对其进行保护，以避免作业的信息被窃取或受到攻击，360及其他第三代计算机都配有此类硬件。

第三代计算机的另一个特性是，卡片被拿到机房后能够很快地将一个作业从卡片读入磁盘。于是任何时刻当一个作业运行结束，操作系统就能将一个新作业从磁盘读出，装入空出来的内存区域运行，这种技术叫做spooling (Simultaneous Peripheral Operation On Line, 同时的外部设备联机操作，亦称假脱机技术)，该技术同时也用于输出。当采用了spooling技术后，就不再需要IBM 1401机，也不必再将磁带搬来搬去。

第三代操作系统很适于大型科学计算和繁忙的商务数据处理，但其实质上仍然是批处理系统。许多程序员很怀念第一代计算机的使用方式，他们那时可以几个小时独占一台机器，可以随时调试他们的程序。而对第三代计算机，从一个作业提交到运算结果取回往往长达数小时，更有甚者，一个逗号的误用会导致编译失败，可能浪费程序员半天时间。

程序员们希望很快得到响应，这种需求导致了分时系统(timesharing)的出现。它实际上是多道程序的一个变种，每个用户都有一个联机终端。在分时系统中，假设有20个用户登录，其中17个在思考或谈论或喝咖啡，则CPU可给那三个需要的作业轮流分配服务。由于调试程序的用户常常只发出简短的命令(如编译一个五页的源文件)，而很少有长的费时命令(如上百万条记录的文件排序)，所以计算机能够为许多用户提供交互式快速的服务，同时在CPU空闲时还能

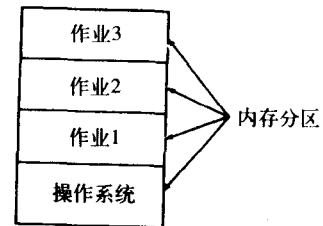


图1-4 内存中有三个作业的一个多道程序系统