

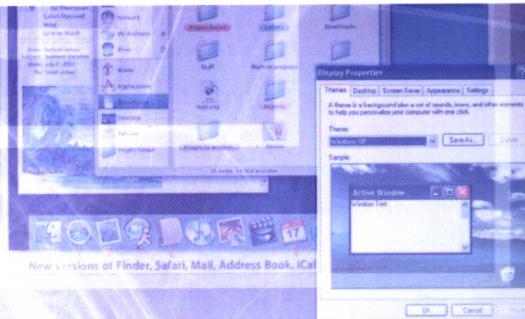
高等学校教材



计算机操作系统教程

(第三版)

左万历 周长林



高等教育出版社
HIGHER EDUCATION PRESS

高等学校教材

计算机操作系统教程

(第二版)

左万历 周长林

高等教育出版社

内容提要

本书主要讲述操作系统的基本概念、基本方法与实现技术。在经典内容的基础上,突出介绍了近年来操作系统的最新进展,如多线程、实时调度与多处理机调度、多处理机互斥、多级页表与倒置页表、RAID技术、快速文件系统、分布协同、微内核与嵌入式系统、操作系统安全等。主要章节后附有流行系统方法案例,并对UNIX系统做了全面分析。最后给出一个基于自动机的操作系统理论模型。

本书在选材和内容组织上进行了认真推敲,力求做到概念准确、层次清晰、系统性强、联系实际、富有启发性。本书第一版曾获得国家级教学成果二等奖、国家教委第三届优秀教材一等奖、国家教委科技进步三等奖等多种奖项,可作为高等学校计算机及相关专业操作系统课程教材,也可供相关技术人员阅读参考。

图书在版编目(CIP)数据

计算机操作系统教程/左万历, 周长林. —2 版.
—北京:高等教育出版社,2004.7

ISBN 7-04-012309-6

I. 计... II. ①左... ②周... III. 操作系统 - 高等学校 - 教材 IV. TP316

中国版本图书馆 CIP 数据核字(2004)第 056799 号

策划编辑 康兆华 责任编辑 康兆华 封面设计 张志 责任绘图 朱静
版式设计 王莹 责任校对 王效珍 责任印制 孔源

出版发行 高等教育出版社
社址 北京市西城区德外大街 4 号
邮政编码 100011
总机 010-82028899

购书热线 010-64054588
免费咨询 800-810-0598
网址 <http://www.hep.edu.cn>
<http://www.hep.com.cn>

经 销 新华书店北京发行所
印 刷 北京星月印刷厂
开 本 787×1092 1/16
印 张 21.5
字 数 480 000

版 次 1994 年 9 月第 1 版
2004 年 7 月第 2 版
印 次 2004 年 7 月第 1 次印刷
定 价 26.90 元

本书如有缺页、倒页、脱页等质量问题,请到所购图书销售部门联系调换。

版权所有 侵权必究

第一版前言

本书是根据国家教育委员会 1992 年颁发的计算机软件专业操作系统课程教学基本要求,结合作者多年教学经验和科研成果,参考国内外操作系统最新发展,联系目前我国操作系统教学实际编写而成的。在编写过程中,我们力求做到体系完整、结构合理、内容丰富、取舍得当、概念准确、叙述清晰、观点明确、深入浅出、适用面广、便于自学,适应 20 世纪 90 年代高等学校计算机专业操作系统教学的需要。

全书共分十四章:第一章是对操作系统的概括性描述;第二章给出操作系统的一个理论模型;第三章至第六章讲述资源管理,包括处理机管理、存储管理、文件管理、设备管理等;第七章和第八章讲述任务管理,包括作业管理和进程管理;第九章讨论操作系统的体系结构;第十章介绍网络与分布式操作系统;第十一章讲述并发程序设计;第十二章讲述操作系统管理;第十三章介绍 UNIX 操作系统;第十四章简述几个常见的操作系统。

考虑近年来操作系统在网络、分布式等方面有了较大的发展,为提高学生的适应能力,除基本内容外,本书还增加了一些新的内容,具体包括:(1) 网络与分布式操作系统,介绍网络和分布式操作系统的基本原理以及与传统操作系统之间的联系和差别;(2) 并发程序设计,在着重阐明操作系统对并发程序设计所应提供的支持的基础上,概述几种并发程序设计语言,并给出若干并发程序设计的例子;(3) 操作系统管理,讲述系统管理员和操作员所应掌握的知识;(4) UNIX 操作系统,描述 UNIX 操作系统所采用的主要技术和特点,既使理论联系实际,又使读者对这一著名的系统有一定程度的了解。这样一来,全书篇幅较长,如因学时限制不能全部讲授,对于增加的内容可以根据具体情况选择其中部分章节。

教材的第二章处于比较特殊的地位,它给出了操作系统的一个理论模型。考虑到易于讲解和接受,此模型不是严格形式化的。本章的目的是使学生在学习操作系统具体内容之前,先对操作系统的基本知识——资源、进程、资源共享、进程并发执行等有一个理性的认识,以指导其后各章具体内容的学习。由于本章具有较强的独立性,可以安排在第九章之后讲授,作为对前面具体内容的理论总结。本章也可以不讲,这不会影响其他章节的学习。

对于具体教学安排提出如下建议,谨供参考:

(1) 用 48~54 个学时讲授基本内容,包括第一章及第三章至第九章。其中某些节(标记※号的)可以选讲。有些内容可以由学生自学。

(2) 如有更多的学时,可以根据教学目标选讲其余各章,其中第十三章应当优先介绍。

本书由周长林负责总体策划以及第二章的编写,其余各章均由左万历执笔。

限于作者水平,错误与不妥之处在所难免,恳请读者给予批评指正。

南京大学计算机科学系谭耀铭副教授在百忙之中审阅了本书，并提出许多宝贵意见和有益的建议，在此深表谢意。

作 者

1993年8月于吉林大学

第二版前言

《计算机操作系统教程》(高等教育出版社,1994年9月第一版)出版至今已近10年,这期间操作系统从理论到实践上都经历了一个较大的发展过程。本书第二版的主要目标是根据新的教学大纲,在保持原教程经典风格的同时,对操作系统的教学内容进行系统而全面的更新,以适应新世纪国内操作系统课程教学需要。

新版《计算机操作系统教程》增加了线程、实时调度与多处理器调度、自旋锁与多处理器互斥、条件临界区、饥饿与饿死、两阶段锁、同步方法、多级页表与倒置页表、伙伴堆内存分配、内存映射文件、日志结构文件、快速文件系统、Ext2fs、DMA、RAID技术、稳定存储器、进程迁移、RMI、分布式文件系统、保护与安全、嵌入式与微内核、面向对象操作系统设计方法等新内容,同时删除了作业控制语言等过时内容和并发程序设计等外延性成分,保证新版教程在篇幅上增加适当,适合国内操作系统课程教学学时要求。

为使理论与实际更好地结合,新版教程在每章后增加了流行系统(如Windows 2000 /XP 和Linux 系统)中采用的具体算法和实现技术。全书最后比较全面地介绍了UNIX 系统,给出了完整的数据结构以及核心算法。

新版教程还对部分章节的次序进行了重新安排和组织,将作业、进程、线程合并,作为主体在第二章中集中讲述;将中断与处理机调度这两个相互承接、密不可分的内容独立出来构成第三章;将互斥、同步与通讯、死锁与饥饿提前,分别构成第四章和第五章,使章节篇幅适中,逻辑关系连贯;将保护与安全纳入第十章“操作系统管理”;另外,原版第二章“操作系统理论”经过精炼后移到最后,更加方便教学取舍。

新版教程对重要术语进行了认真推敲,并以定义形式给出,明确了“系统开销”、“忙式等待”等概念,给出了地址映射的数学意义。新版教程对经典内容进行了认真的审定,使其更加严谨。除涉及管程语言和Ada 语言的成分外,所有算法均由Pascal 语言改为C 语言描述。另外,新版教程特别注重能力和素质培养,各章中都安排有启发性问题,留给学生充分的思考空间。

作为高等教育出版社立体化精品课程教材,作者开发了与之配套的远程教学系统,并开通了相关教学网站(<http://info.jlu.edu.cn/~cs/los>)。为了方便教学,作者还制作了与新版教材配套的多媒体教学课件,可以通过下载方式获取。我们将对网站进行动态维护并对课件进行周期性版本更新。与教材配套的操作系统习题解答与实验指导也将在近期完成。

新版教程力求更加准确地讲述操作系统的基本概念、基本原理、设计方法和实现技术,全面反映操作系统近年来的最新发展,深入介绍流行操作系统的核心算法。教材中部分章节引用了参考文献中列出的国外著作中的一些内容,谨在此向作者致以衷心的感谢和深深的敬意。

本书第一版曾获得国家级教学成果二等奖、国家教委第三届优秀教材一等奖、国家教委科技

进步三等奖等多项奖励。新版教材得到高等教育出版社立体化精品课程教材建设项目的支
持，在此深表谢意。

本书第二版由左万历负责整体规划并执笔前十二章，第十三章由周长林编写。由于作者水
平所限，加之时间仓促，若有错误与不当之处，敬请读者不吝赐教。作者的电子信箱为：wanli@jlu.edu.cn。

作 者

2004年3月于长春

目 录

第一章 操作系统概述	1	1.6 操作系统的界面形式	15
1.1 操作系统的概念	1	1.6.1 交互终端命令	15
1.1.1 操作系统的地位	1	1.6.2 图形用户界面	15
1.1.2 操作系统的作用	1	1.6.3 作业控制语言	15
1.1.3 操作系统的定义	2	1.6.4 系统调用命令	16
1.2 操作系统的历史	2	1.7 操作系统的运行机理	16
1.2.1 操作系统的产生	2	1.8 系统举例	17
1.2.2 操作系统的完善	4	1.8.1 Linux 系统	17
1.2.3 操作系统的发展	5	1.8.2 Windows 2000/XP 系统	17
1.3 操作系统的特性	6	习题一	17
1.3.1 程序并发性	6	第二章 进程、线程与作业	19
1.3.2 资源共享性	6	2.1 多道程序设计	19
1.4 操作系统的分类	7	2.1.1 单道程序设计的缺点	19
1.4.1 多道批处理操作系统	7	2.1.2 多道程序设计的提出	20
1.4.2 分时操作系统	8	2.1.3 多道程序设计的问题	21
1.4.3 实时操作系统	8	2.2 进程的引入	22
1.4.4 通用操作系统	9	2.2.1 进程的概念	22
1.4.5 单用户操作系统	9	2.2.2 进程状态及状态转换	23
1.4.6 网络操作系统	10	2.2.3 进程控制块	23
1.4.7 分布式操作系统	10	2.2.4 进程的组成与上下文	24
1.4.8 多处理器操作系统	11	2.2.5 进程的队列	25
1.4.9 嵌入式操作系统	11	2.2.6 进程的类型和特性	25
1.4.10 智能卡操作系统	12	2.2.7 进程间的相互联系与相互作用	26
1.5 操作系统的硬件环境	12	2.2.8 进程的创建与撤销	26
1.5.1 定时装置	12	2.2.9 进程与程序的联系和差别	27
1.5.2 系统栈	13	2.3 线程与轻进程	27
1.5.3 特权指令与非特权指令	13	2.3.1 线程的引入	27
1.5.4 处理机状态及状态转换	13	2.3.2 线程的概念	28
1.5.5 地址映射机构	14	2.3.3 线程的结构	28
1.5.6 存储保护设施	14	2.3.4 线程控制块	28
1.5.7 中断装置	14	2.3.5 线程的实现	29
1.5.8 通道与 DMA 控制器	14	2.3.6 线程的应用	31

2.4 作业.....	32	4.2.1 共享变量与临界区	67
2.4.1 批处理作业	32	4.2.2 临界区与进程互斥	68
2.4.2 交互式作业	32	4.2.3 进程互斥的实现	69
2.5 系统举例.....	34	4.2.4 多处理器环境下的互斥	75
2.5.1 Java 线程	34	4.3 进程同步.....	76
2.5.2 Linux 进程与线程	35	4.3.1 进程同步的概念	76
2.5.3 Windows 2000/XP 进程、线程 与线程	35	4.3.2 进程同步机制	77
习题二	37	4.3.3 信号灯与 PV 操作	78
第三章 中断与处理器调度	39	4.3.4 条件临界区	82
3.1 中断与中断系统.....	39	4.3.5 管程	83
3.1.1 中断概念	39	4.3.6 会合	92
3.1.2 中断装置	39	4.4 进程高级通讯.....	99
3.1.3 中断处理程序	43	4.4.1 进程通讯的概念	99
3.2 处理机调度.....	49	4.4.2 进程通讯的模式	99
3.2.1 处理机调度算法	49	4.4.3 直接方式	100
3.2.2 处理机调度时机	53	4.4.4 间接方式	103
3.2.3 处理机调度过程	54	4.5 系统举例	105
3.3 调度级别与多级调度	55	4.5.1 Java 中的管程	105
3.3.1 交换与中级调度	55	4.5.2 Linux 进程通讯	105
3.3.2 作业与高级调度	56	4.5.3 Windows 2000/XP 并发控制	107
3.4 实时调度.....	57	习题四	108
3.4.1 最早截止期优先调度	58	第五章 死锁与饥饿	111
3.4.2 速率单调调度	58	5.1 死锁的概念	111
3.5 多处理器调度	59	5.2 死锁的类型	112
3.5.1 自调度	59	5.2.1 竞争资源引起的死锁	112
3.5.2 组调度	60	5.2.2 进程通讯引起的死锁	112
3.6 系统举例	60	5.2.3 其他原因引起的死锁	112
3.6.1 Linux 进程调度	60	5.3 死锁的条件	112
3.6.2 Windows 2000/XP 线程调度	61	5.4 死锁的处理	113
习题三	63	5.5 资源分配图	113
第四章 互斥、同步与通讯	65	5.5.1 资源分配图的定义	113
4.1 并发进程	65	5.5.2 资源分配图的约简	115
4.1.1 顺序程序及其特性	65	5.6 死锁的预防	115
4.1.2 并发程序及其特性	65	5.6.1 预先分配策略	115
4.1.3 与时间有关的错误	66	5.6.2 有序分配策略	116
4.2 进程互斥	67	5.7 死锁的避免	117
		5.7.1 安全状态与安全进程序列	117

5.7.2 银行家算法	117	6.6.1 Linux 存储管理	169
5.8 死锁的发现	120	6.6.2 Windows 2000/XP 存储管理	170
5.8.1 死锁检测算法	120	习题六	173
5.8.2 死锁检测时刻	122	第七章 文件系统	175
5.9 死锁的恢复	123	7.1 文件与文件系统	175
5.10 鸵鸟算法	123	7.1.1 文件	175
5.11 有关问题的讨论	124	7.1.2 文件系统	176
5.11.1 关于充要性算法	124	7.2 文件的访问方式	176
5.11.2 关于消耗型资源问题	124	7.2.1 顺序访问	176
5.11.3 关于两阶段封锁	124	7.2.2 随机访问	176
5.12 饥饿与活锁	125	7.3 文件的组织	177
5.13 死锁与饥饿的例子	126	7.3.1 文件的逻辑组织	177
习题五	129	7.3.2 文件的物理组织	178
第六章 存储管理	132	7.4 文件目录	184
6.1 存储管理的功能	132	7.4.1 文件控制块与目录项	184
6.1.1 存储分配	132	7.4.2 文件目录与目录文件	184
6.1.2 存储共享	132	7.4.3 单级目录与多级目录	185
6.1.3 存储保护	133	7.4.4 文件目录的改进	185
6.1.4 存储扩充	133	7.4.5 根目录与当前目录	186
6.1.5 地址映射	133	7.4.6 文件目录的查找	187
6.2 内存资源管理	134	7.5 文件的共享	187
6.2.1 内存分区	134	7.5.1 文件共享的目的	187
6.2.2 内存分配	134	7.5.2 文件共享的模式	187
6.2.3 碎片与紧凑	136	7.5.3 文件共享的实现	188
6.3 存储管理方式	137	7.6 文件的保护、保密与安全	188
6.3.1 单一连续区存储管理	137	7.6.1 文件的保护	188
6.3.2 分页式存储管理	139	7.6.2 文件的保密	189
6.3.3 分段式存储管理	145	7.6.3 文件的安全	190
6.3.4 段页式存储管理	150	7.7 文件系统的实现	191
6.4 外存管理技术	153	7.7.1 内存所需的表目	191
6.4.1 外存空间划分	153	7.7.2 外存空间的管理	192
6.4.2 外存空间分配	154	7.8 文件系统的界面	194
6.5 虚拟存储系统	154	7.9 日志结构文件系统	196
6.5.1 虚拟页式存储系统	155	7.10 内存映射文件	197
6.5.2 虚拟段式存储系统	162	7.11 系统举例	198
6.5.3 虚拟段页式存储系统	166	7.11.1 Linux 文件系统	198
6.6 系统举例	169	7.11.2 Windows 2000/XP 的 NTFS	199

习题七.....	201	习题八.....	223
第八章 设备与 I/O 管理	202	第九章 网络操作系统与分布式	
8.1 设备的分类	202	操作系统.....	225
8.1.1 输入/输出型设备与存储型设备	202	9.1 计算机网络	225
8.1.2 块型设备与字符型设备	202	9.1.1 网络的概念	225
8.1.3 独占型设备与共享型设备	202	9.1.2 网络的组成	225
8.2 设备的物理特性	203	9.1.3 网络的分类	226
8.2.1 输入/输出型设备的物理特性	203	9.1.4 网络的拓扑	226
8.2.2 存储型设备的物理特性	203	9.2 通信与协议	228
8.3 I/O 传输方式	206	9.3 网络服务	229
8.3.1 程序控制查询方式	206	9.3.1 远程登录	229
8.3.2 中断驱动方式	206	9.3.2 远程文件传输	229
8.3.3 DMA 方式	206	9.4 计算模型	230
8.3.4 通道方式	207	9.4.1 数据迁移	230
8.4 设备分配与去配	209	9.4.2 计算迁移	230
8.4.1 独占型设备的分配与去配	209	9.5 事件定序	232
8.4.2 共享型设备的分配与去配	210	9.5.1 前发生关系	232
8.5 设备驱动	211	9.5.2 全序关系	233
8.5.1 通道程序	211	9.6 进程互斥	233
8.5.2 设备启动	211	9.6.1 集中方式	234
8.5.3 中断处理	211	9.6.2 分布方式	234
8.6 设备调度	212	9.6.3 令牌传递方式	235
8.7 缓冲技术	214	9.7 进程同步与进程通讯	235
8.7.1 缓冲技术的引入	214	9.7.1 消息传递	235
8.7.2 硬缓冲与软缓冲	214	9.7.2 套接字	236
8.7.3 私用缓冲与公共缓冲	214	9.7.3 远程过程调用	237
8.7.4 缓冲池及其管理	214	9.7.4 远程方法启用	239
8.7.5 缓冲技术的实现	215	9.8 死锁处理	239
8.8 输入/输出进程.....	218	9.8.1 死锁预防	239
8.9 RAID 技术	218	9.8.2 死锁检测	240
8.9.1 RAID 级别.....	219	9.9 资源管理	240
8.9.2 硬件 RAID 与软件 RAID	220	9.9.1 集中方式	240
8.10 虚拟设备.....	220	9.9.2 分布方式	241
8.10.1 虚拟设备的引入	220	9.9.3 层次方式	241
8.10.2 虚拟设备的实现	221	9.10 分布式文件系统.....	241
8.11 稳定存储	222	9.10.1 一般结构	242
8.12 系统举例	222	9.10.2 命名与透明性	242

9.10.3 远程文件存取	243	11.4.4 面向对象设计方法	270
9.10.4 有状态服务与无状态服务	243	11.5 系统举例	271
9.10.5 缓存策略	243	习题十一	273
9.11 系统举例	244	第十二章 UNIX 实例分析	275
习题九	245	12.1 历史回顾	275
第十章 操作系统管理	246	12.2 系统结构	275
10.1 操作系统使用	246	12.2.1 内核部分	277
10.1.1 操作系统生成	246	12.2.2 外壳部分	277
10.1.2 操作系统装入	247	12.3 进程管理	277
10.1.3 操作系统初启	247	12.3.1 进程组成	277
10.1.4 操作系统运行	247	12.3.2 进程控制块	278
10.2 操作系统维护	248	12.3.3 进程状态与状态转换	280
10.2.1 改正性维护	249	12.3.4 进程调度	281
10.2.2 适应性维护	250	12.3.5 进程互斥	282
10.2.3 完善性维护	250	12.3.6 进程同步	282
10.3 操作系统保护	250	12.3.7 进程通讯	282
10.3.1 域结构	251	12.4 存储管理	285
10.3.2 访问矩阵	251	12.4.1 存储管理方式	285
10.4 操作系统安全	253	12.4.2 存储分配算法	286
10.4.1 闯入与身份认证	253	12.4.3 进程空间扩充	288
10.4.2 程序威胁	255	12.4.4 交换技术	288
10.4.3 安全策略	258	12.4.5 虚拟页式存储管理	288
10.4.4 可信系统	259	12.5 文件系统	289
习题十	260	12.5.1 文件类型	289
第十一章 操作系统设计	261	12.5.2 文件体系	290
11.1 操作系统设计目标	261	12.5.3 文件结构	290
11.2 操作系统基本内核	262	12.5.4 文件目录与连接	291
11.2.1 内核的基本组成	262	12.5.5 文件系统映射	292
11.2.2 内核各部分关系	262	12.5.6 文件卷的安装	293
11.3 操作系统体系结构	263	12.5.7 磁盘空间管理	294
11.3.1 基于共享变量结构	264	12.5.8 inode 区域管理	295
11.3.2 基于信件传递结构	264	12.5.9 快速文件系统	297
11.3.3 微内核结构	265	12.5.10 NFS 网络文件系统	298
11.4 操作系统设计方法	266	12.6 设备管理	300
11.4.1 模块接口法	266	12.6.1 设备分配	300
11.4.2 核扩充法	266	12.6.2 缓冲与缓存	300
11.4.3 层次结构法	266	12.6.3 预先读与延迟写	302

12.7 系统调用.....	303	13.4.2 进程的执行	316
12.7.1 有关进程的系统调用命令	303	13.4.3 进程与资源的关系	316
12.7.2 有关文件的系统调用命令	306	13.4.4 进程的互斥	317
12.8 外壳语言.....	309	13.5 资源管理.....	318
习题十二.....	310	13.5.1 主要资源管理思想概述	318
第十三章 操作系统理论.....	312	13.5.2 互斥机制与资源管理	319
13.1 前言.....	312	13.6 进程管理.....	324
13.1.1 操作系统理论所处的地位	312	13.6.1 进程同步	324
13.1.2 操作系统理论的描述形式	312	13.6.2 进程通讯	325
13.1.3 操作系统理论的主要内容	312	13.6.3 进程死锁	326
13.2 并发程序.....	313	13.7 虚拟资源.....	326
13.2.1 并发程序的概念	313	13.8 操作系统理论的形式化.....	326
13.2.2 并发程序的不确定性	313	13.8.1 资源	326
13.2.3 不确定性带来的问题	313	13.8.2 进程	327
13.3 资源.....	314	13.8.3 指针选择	327
13.3.1 资源的概念	314	13.8.4 有关理论问题	327
13.3.2 资源的分类	315	13.9 本章小结.....	328
13.4 进程.....	316	习题十三.....	328
13.4.1 进程的定义	316	参考文献.....	329

第一章 操作系统概述

1.1 操作系统的概念

关于什么是操作系统,目前尚无统一的定义。这里从操作系统在整个计算机系统中所处的地位以及所起的作用来给出关于操作系统的一个非形式化的描述。

1.1.1 操作系统的地位

计算机系统是由硬件和软件两部分构成的。软件又分成系统软件与应用软件两类,操作系统是一个最基本也是最重要的系统软件。从虚拟机的观点来看,软件是划分为层次的。系统软件位于低层,应用软件位于高层。当然,系统软件和应用软件都可以进一步分层。如果将系统软件进一步分层的话,操作系统位于系统层次中的最底层,如图 1.1 所示。

据此可以看出,操作系统是与计算机硬件关系最为密切的一个基本软件,是对硬件机器的第一次扩充。

注意图 1.1 所表示的层次关系具有穿透性:高层软件可以调用所有低于其所在层次的软件,并可以直接与硬件打交道,每个软件层都在原有层次的基础上另外增加一层新的界面。例如,应用程序以目标代码形式运行时可以与操作系统和硬件直接打交道(调用操作系统或执行硬件指令),操作系统之上的系统库可以被应用程序调用,库函数又可以调用操作系统,如图 1.2 所示。

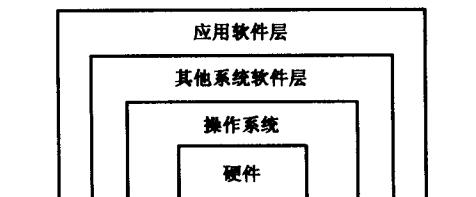


图 1.1 虚拟机层次

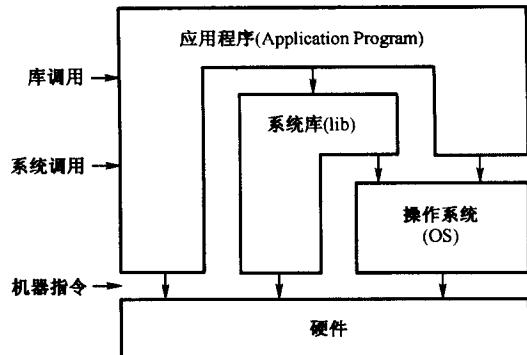


图 1.2 操作系统地位

1.1.2 操作系统的作用

操作系统有以下两个重要的作用。

1. 管理系统中的各种资源

一个多道计算机系统可同时为多个用户服务,也就是说在系统中同时有多个程序在执行。

这些程序在执行的过程中可能会要求使用系统中的各种资源,例如,当程序运行时需要处理机资源,输出时需要打印机资源。多个程序的资源需求经常会发生冲突,如程序 P₁ 和程序 P₂ 可能会同时想要使用打印机输出。如果对这些程序的资源需求不加以管理,则会造成混乱甚至可能损坏设备。也就是说,在系统中需要有一个资源仲裁者,由它负责资源在各个程序之间的调度,保证系统中的各种资源得以有效的利用。这个资源仲裁者就是操作系统。

2. 为用户提供良好的界面

早期的计算机是没有操作系统的,那时使用计算机需要大量的手工操作,既繁琐又费时。读者可以想象,如果没有操作系统,要运行一个用 C 语言编写的源程序将是多么困难。有了操作系统之后,原来需要由人来做的许多繁琐而费时的工作就由操作系统代替完成,这使得用户能够非常方便地使用计算机系统。例如,若要运行一个用 C 语言编写的源程序,用户只需在终端上键入几个命令或者点击几次鼠标便可完成。可以说,操作系统的产生是计算机发展过程中具有历史性意义的一步。

随着硬件成本的下降,计算机已经走入家庭和办公自动化领域,多数计算机的使用者不是计算机专业人员,这时界面的友好性比资源的利用效率更具实际的意义。目前商业化操作系统提供的 GUI(Graphic User Interface,图形用户界面)就是在此背景下的产物。

1.1.3 操作系统的定义

根据前面关于操作系统地位和作用的描述,可以给出关于操作系统的一个非形式化定义如下。

定义 操作系统是位于硬件层之上、所有其他软件层之下的一一个系统软件,是管理系统中各种软件和硬件资源、使其得以充分利用并方便用户使用计算机系统的程序集合。

应当指出,关于操作系统的概念是难以用几句话来概括的,读者需要在下面的学习过程中认真地加以体会。

1.2 操作系统的歷史

为使读者了解操作系统的形成、完善和发展,现在简略地回顾一下操作系统的歷史。

应当指出,由于操作系统是直接建造于硬件之上的,它的演变必然与计算机系统结构的演变有着密切的联系。可以说,操作系统的发展与硬件系统结构的发展是相互促进、相互影响的。一方面,为了方便而有效地使用硬件,导致了操作系统的产生。另一方面,为了利于操作系统的构造,硬件也经历了不断改进的过程。此外,由于操作系统为上层软件及用户提供界面,它的演变必然反映出上层软件及用户对于操作系统的使用要求。

简而言之,操作系统经历了从无到有、由功能单纯到功能完整的演变过程,并且还处于进一步的发展之中,下面分别加以叙述。

1.2.1 操作系统的产生

计算机操作系统在从无到有的产生过程中经历了如下几个主要阶段。

1. 手工操作阶段(20世纪40年代)

计算机诞生的初期并没有操作系统,人们采用手工操作方式使用计算机,典型的作业(job)处理步骤如下:首先将程序和数据通过手工操作记录在穿孔纸带上;然后将程序纸带放到光电输入机上并通过控制台开关启动光电机将程序输入内存;继而再通过控制台开关启动程序由第一条指令开始执行;程序在运行的过程中通常需要人工干预,如将数据纸带放到光电输入机上、出错时显示错误地址并修改指令等;最后运行结果在电传打印机上输出。显然,这种操作方式有如下两个缺点:(1)用户在其作业处理的整个过程中独享系统中的全部资源;(2)手工操作所需的时间很长。

这种操作方式在计算机速度较慢的情况下是可以容忍的,但是当计算机速度大幅度提高之后,就暴露出了严重的缺点。例如,假设一个作业在速度为每秒1 000次的机器上运行需要1 hr,手工操作所需要的时间为4 min,则手工操作时间与程序运行时间之比为1:15;若计算机的速度提高到每秒600 000次,同一程序的运行时间只需6 s,而手工操作时间不变,仍为4 min,则手工操作时间与程序运行时间之比为40:1,就是说,手工操作时间远远大于程序运行时间。因而,缩短手工操作时间在以晶体管为代表的第二代计算机出现后便成了亟待解决的问题。

其他软件在此阶段所取得的成就是汇编系统和汇编语言的出现,它在一定程度上减轻了用户使用计算机的负担。

2. 批处理阶段(20世纪50年代)

为了缩短手工操作的时间,人们自然想到的是使作业到作业之间的过渡摆脱人的干预,实现自动化,如此便出现了批处理。批处理经历了两个阶段:即联机批处理阶段和脱机批处理阶段。

(1) 联机批处理:早期的批处理是联机的。其工作原理如下:操作员将若干个作业合成为一批,并将其卡片依次放到读卡机上,监督程序通过内存将这一批作业传送到磁带机上,输入完毕后监督程序开始处理这一批作业。它自动地将第一个作业读入内存,并对该作业的程序进行汇编或编译,然后将产生的目标程序与所需要的例行子程序连接装配在一起,继而执行该程序,计算完成之后输出其结果。第一个作业处理完毕后立即开始处理第二个作业,如此重复直到所有作业处理完毕。此时,监督程序将第二批作业由读卡机传送到磁带机上,并按上述步骤处理。这样,监督程序不间断地处理各个作业,实现了作业之间转换的自动化,大大地缩短了手工操作时间。不过,联机批处理也有一个缺点,即作业由读卡机到磁带机的传输需要处理机完成,由于设备的传输速度远远低于处理机的速度,在此传输过程中处理机仍会浪费较多的时间。

(2) 脱机批处理:为了克服联机批处理的缺点,引入了脱机批处理。它的思想是把输入/输出操作交给一个功能较为单纯的卫星机去做,使主机从繁琐耗时的输入/输出操作中解脱出来,其基本原理如图1.3所示。待处理的作业由卫星机负责经读卡机传送到输入磁带上,主机由输入磁带读入作业并加以处理,其结果送到输出磁带上,最后由卫星机负责将输出磁带上的结果信息在打印机上输出。

批处理系统是操作系统的雏形。在此阶段,其他软件也有了相应的发展,如输入/输出标准程序、高级语言编译程序、连接装配程序等。

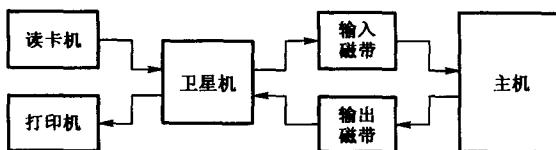


图 1.3 脱机批处理

3. 执行系统阶段(20世纪60年代初期)

批处理较之手工操作前进了一大步,但它仍有一些缺点:如需要额外的卫星机、磁带机的装卸需要手工操作等。

在20世纪60年代初期,硬件在两方面取得了重要的进展,一是通道的引入,二是通道中断主机功能的出现。这是操作系统发展史上的重要事件,它推进操作系统进入执行系统阶段。

通道也称I/O处理器,它具有自己的指令系统和运控部件,与处理机共享内存资源。通道可以受处理机的委托执行通道程序,完成输入/输出操作。通道的I/O操作可以同处理机的计算工作完全并行,并在I/O操作完成时向处理机发出中断请求。这样,作业由读卡机到磁带机的传输以及运行结果由磁带机到打印机的传输可由通道完成,这既非联机,也非脱机,称做“假脱机”或“伪脱机”。通道取代了卫星机,也免去了手工装卸磁带的麻烦。

执行系统是操作系统的初级阶段,它为操作系统的最终形成奠定了基础。

1.2.2 操作系统的完善

操作系统由形成到完善经历了如下几个主要的发展阶段。

1. 多道批处理系统(20世纪60年代初期)

执行系统出现不久,人们就发现在内存中同时存放多道作业是有利的。当一道作业因等待I/O传输完成等原因暂时不能运行时,系统可将处理机资源分配给另外一个可以运行的程序。如此便产生了多道批处理操作系统。

多道批处理的出现是操作系统发展史上一个革命性的变革,它将多道程序设计的概念引入操作系统中。本书后面将会讲到,多道程序设计与传统的单道程序设计相比具有本质上的差别。它的引入给操作系统的理论及实践等各个方面都带来了许多新的研究课题。

2. 分时系统(20世纪60年代初/中期)

手工操作是一种联机操作方式,其效率很低。批处理系统否定并代替了手工操作,是一种脱机操作方式。执行系统及多道批处理系统是批处理系统的进一步发展,属于更高级的脱机处理方式。但是,多道批处理系统出现不久,人们便发现仍有联机操作的必要,这个要求首先是由程序员提出的。对于脱机操作来说,程序员无法了解其作业的执行情况并对其进行动态控制,如果作业在处理过程中出现错误,程序员不能对其进行及时的修改,必须等待批处理结果输出后才能从输出报告中得知错误所在,并对其进行修改,然后再次提交批作业,如此可能需要重复多次,使得作业的处理周期较长。也就是说,脱机方式非常不利于程序的动态调试。

为达到联机操作的目标,出现了分时系统。分时系统由一个主机和若干个与其相连的终端