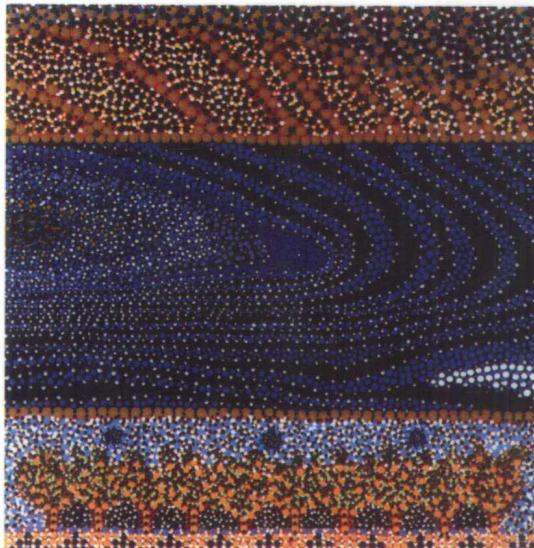


计 算 机 科 学 从 书

分布式操作系统 原理与实践

(美) Doreen L. Galli 著 徐良贤 唐英 毛家菊 金恩华 等译

DISTRIBUTED
OPERATING SYSTEMS
CONCEPTS & PRACTICE



DOREEN L. GALLI

Distributed Operating Systems
Concepts and Practice



机械工业出版社
China Machine Press



本书从概念和实践的角度详细论述了分布式操作系统的各个主要方面，还包含了实际操作系统的相关范例以及广泛算法的具体实例，其中提及的核心 Web 站点、ftp 站点和文献提供了大量参考资源。此外，Windows 2000 案例研究提供了一个实际商业解决方案的例子，附录中还有一个简单的 C/S 实际应用来演示关键的分布式计算程序设计概念，以加深概念并展示分布式操作系统设计人员所需进行的设计与操作。

本书内容丰富，结构合理，适于作为计算机及相关专业的本科生和研究生的教材，也是计算机从业人员掌握分布式操作系统原理的理想读物。

Doreen L. Galli; *Distributed Operating Systems: Concepts & Practice* (ISBN 0-13-0 79843-6).

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

Copyright © 2000 by Prentice Hall PTR, Inc.

All rights reserved.

Chinese simplified language edition published by China Machine Press.

Copyright © 2002 by China Machine Press.

本书中文简体字版由美国 Prentice Hall PTR 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：01-2001-2205

图书在版编目 (CIP) 数据

分布式操作系统：原理与实践 / (美) 加利 (Galli, D.L.) 著；徐良贤等译 . - 北京：
机械工业出版社，2003.1

(计算机科学丛书)

书名原文：Distributed Operating Systems: Concepts and Practice

ISBN 7-111-10952-X

I . 分… II . ①加… ②徐… III . 分布式操作系统 - 原理 IV . TP316

中国版本图书馆 CIP 数据核字 (2002) 第 070382 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：杨文

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

2003 年 1 月第 1 版第 1 次印刷

787mm × 1092mm 1/16 · 21.75 印张

印数：0 001 ~ 5 000 册

定价：38.00 元

凡购本书，如有倒页、脱页、缺页，由本社发行部调换

出版者的话

文艺复兴以降，源远流长的科学精神和逐步形成的学术规范，使西方国家在自然科学的各个领域取得了垄断性的优势；也正是这样的传统，使美国在信息技术发展的六十多年间名家辈出、独领风骚。在商业化的进程中，美国的产业界与教育界越来越紧密地结合，计算机学科中的许多泰山北斗同时身处科研和教学的最前线，由此而产生的经典科学著作，不仅擘划了研究的范畴，还揭橥了学术的源变，既遵循学术规范，又自有学者个性，其价值并不会因年月的流逝而减退。

近年，在全球信息化大潮的推动下，我国的计算机产业发展迅猛，对专业人才的需求日益迫切。这对计算机教育界和出版界都既是机遇，也是挑战；而专业教材的建设在教育战略上显得举足轻重。在我国信息技术发展时间较短、从业人员较少的现状下，美国等发达国家在其计算机科学发展的几十年间积淀的经典教材仍有许多值得借鉴之处。因此，引进一批国外优秀计算机教材将对我国计算机教育事业的发展起积极的推动作用，也是与世界接轨、建设真正的世界一流大学的必由之路。

机械工业出版社华章图文信息有限公司较早意识到“出版要为教育服务”。自1998年始，华章公司就将工作重点放在了遴选、移译国外优秀教材上。经过几年的不懈努力，我们与Prentice Hall, Addison-Wesley, McGraw-Hill, Morgan Kaufmann等世界著名出版公司建立了良好的合作关系，从它们现有的数百种教材中甄选出Tanenbaum, Stroustrup, Kernighan, Jim Gray等大师名家的一批经典作品，以“计算机科学丛书”为总称出版，供读者学习、研究及庋藏。大理石纹理的封面，也正体现了这套丛书的品位和格调。

“计算机科学丛书”的出版工作得到了国内外学者的鼎力襄助，国内的专家不仅提供了中肯的选题指导，还不辞劳苦地担任了翻译和审校的工作；而原书的作者也相当关注其作品在中国的传播，有的还专诚为其书的中译本作序。迄今，“计算机科学丛书”已经出版了近百个品种，这些书籍在读者中树立了良好的口碑，并被许多高校采用为正式教材和参考书籍，为进一步推广与发展打下了坚实的基础。

随着学科建设的初步完善和教材改革的逐渐深化，教育界对国外计算机教材的需求和应用都步入一个新的阶段。为此，华章公司将加大引进教材的力度，在“华章教育”的总规划之下出版三个系列的计算机教材：针对本科生的核心课程，剔抉外版菁华而成“国外经典教材”系列；对影印版的教材，则单独开辟出“经典原版书库”；定位在高级教程和专业参考的“计算机科学丛书”还将保持原来的风格，继续出版新的品种。为了保证这三套丛书的权威性，同时也为了更好地为学校和老师们服务，华章公司聘请了中国科学院、北京大学、清华大学、国防科技大学、复旦大学、上海交通大学、南京大学、浙江大学、中国科技大学、哈尔滨工业大学、西安交通大学、中国人民大学、北京航空航天大学、北京邮电大学、中山大学、解放军理工大学、郑州大学、湖北工学院、中国国家信息安全测评认证中心等国内重点大学和科研机构在计算机的各个领域的著名学者组成“专家指导委员会”，为我们提供选题意见和出版监督。

“国外经典教材”是响应教育部提出的使用外版教材的号召，为国内高校的计算机本科教学度身订造的。在广泛地征求并听取丛书的“专家指导委员会”的意见后，我们最终选定了这20多种篇幅内容适度、讲解鞭辟入里的教材，其中的大部分已经被M.I.T.、Stanford、U.C. Berkley、C.M.U.等世界名牌大学采用。丛书不仅涵盖了程序设计、数据结构、操作系统、计算机体系结构、数据库、编译原理、软件工程、图形学、通信与网络、离散数学等国内大学计算机专业普遍开设的核心课程，而且各具特色——有的出自语言设计者之手、有的历三十年而不衰、有的已被全世界的几百所高校采用。在这些圆熟通博的名师大作的指引之下，读者必将在计算机科学的宫殿中由登堂而入室。

权威的作者、经典的教材、一流的译者、严格的审校、精细的编辑，这些因素使我们的图书有了质量的保证，但我们的目标是尽善尽美，而反馈的意见正是我们达到这一终极目标的重要帮助。教材的出版只是我们的后续服务的起点。华章公司欢迎老师和读者对我们的工作提出建议或给予指正，我们的联系方法如下：

电子邮件：hzedu@hzbook.com

联系电话：(010) 68995265

联系地址：北京市西城区百万庄南街1号

邮政编码：100037

译 者 序

近年来，分布式操作系统在国内外都获得了迅猛的发展，成为计算机科学与技术研究领域中的一个热点，同时在商业应用方面也取得了丰硕的成果。学习和研究世界上分布式系统最新的理论和技术，以便在该领域迅速跟上并超越国际水平，是我国计算机界义不容辞的任务。因此我们翻译了 Doreen L. Galli 博士的《分布式操作系统：原理与实践》一书。

本书涵盖了分布式操作系统的所有内容，全面地介绍了分布式系统环境中的内核、通信、内存管理、并发控制等基本概念和算法，并对进程管理、文件系统、事务管理和同步等方面高级概念和算法进行了详细的介绍和研究。

本书紧密结合当前的最新技术，并对这些技术进行了详细的阐述，比如基于对象的分布式操作系统、分布式操作系统的安全等，同时还给出了具有实际意义的研究实例——Windows 2000。它不仅是计算机专业本科生和研究生的一本很好的教材，对业界相关人士的研究开发工作也不无裨益。

本书的第1、2章及前言等由唐英翻译，第3、4、5章由毛家菊翻译，第6、7、8章由肖正光翻译，第9、10、11章由金恩华翻译，第12章由徐良贤翻译，附录程序中的注释部分由邱敏翻译。本书前半部分由金恩华初审，后半部分由唐英初审，并由唐英统排全书。最后，全书由徐良贤教授审校并统一给出了词汇表。

在翻译过程中我们尽量做到尊重原意、翻译准确，但由于水平有限，不当和疏漏之处在所难免，敬请广大读者不吝指正。

译 者
上海交通大学电子信息学院
计算机科学与技术系

前　　言

本书阐述了分布式计算的原理和实践，不仅适合学生学习，也可供从业人员及公司培训之用。在过去的十年中，计算机系统日益进步，大多数计算机都是连接在具有一致基础的某种网络之中的。小型企业中安装局域网越来越普遍，家庭安装的局域网也以越来越快的速度增加。软件技术必须跟上这一趋势，当前和未来的从业人员也是如此。按目前速度，所有计算机科研人员掌握分布式系统的工作原理只是个时间问题，因为大部分计算机及其应用都要用到这个技术。

本书读者

学习标准操作系统概念对计算机专业的大学生特别重要，同时，扩展分布式操作系统知识也是研究生和大学四年级学生以及业界工作人员迫切的和不断增长的要求。因此，很需要对分布式操作系统的原理、实际解决方案和方法进行研究。本书就能够满足学生和从业人员的这方面需要。

目标

本书从概念和实践的角度详细论述了分布式操作系统的各个主要方面，并且还包含了实际操作系统的相关范例，以加深概念和展示分布式系统设计人员所必须作出的决定。一些操作系统如 Amoeba、Clouds 和 Chorus（JavaOS 的技术基础）在全书中作为范例来讲解，此外，Windows 2000 案例研究还提供了一个实际商业解决方案的例子。针对分布式计算的各个方面，本书用 CORBA、DCOM、NFS、LDAP、X.500、Kerberos、RSA、DES、SSH 和 NTP 等技术说明实际的解决方案。附录中还有一个简单的 C/S 应用来演示关键的分布式计算程序设计概念，如 INET 套接字、pthread 和通过互斥操作实现同步。

总之，本书着重于分布式的原理、理论和实践。本书是为计算机从业人员、大学四年级学生和研究生编写的，并且假定读者已经学过基本的操作系统课程。我们希望这本书不仅对渴望充电的业界人员是有益的，而且对将本书作为教材来学习的学生来说也是有价值的。

内容组织和教学特点

本书分成两部分。从第 1 章到第 6 章为第一部分，提供分布式计算的基础知识，从第 7 章到第 11 章为第二部分，详述这些主题并更深入地研究高级分布式操作系统的主题。书中的教学特点如下：

1. 供进一步深入理解的详细说明。这些详细说明中包含诸如复杂算法和更深入的例子等内容。
2. 超过 150 幅的图表，以图解方式帮助阐明概念。
3. Windows 2000 案例研究，展示一个实际的商业解决方案。
4. 面向项目的习题（带有斜体数字），作为亲身体验。
5. 习题建立在前面几章的概念之上。
6. 参考资料来源包括：

- A. 供进一步学习的综述资料。
 - B. 研究论文。
 - C. 核心 web 站点和 ftp 站点。
7. 一个简化的分布式应用程序，以展示分布式程序设计的关键概念。
 8. 综合词汇表，集中给出主要的定义。
 9. 缩写字的完整列表，以帮助阅读并为快速查询提供一个集中存放地点。
 10. 各章小结。
 11. 完整的索引。
 12. 本书的 web 站点为 www.prenhall.com/galli。

对教师的建议

本书为教师提供了最大的灵活性，并具有一定的教学特点，以便教师能够根据班级需要和教学任务选择具体内容。在写本书时惟一的假定是读者已经学过了基础的操作系统导论课程。有些题材可能在操作系统导论课程中有所论述，但在本书中有时忽略或简单讲述，有的题材常常难于掌握或者可能已经淡忘，然而对分布式系统是很重要的，本书会在适当的地方插入这些内容。这些材料不必在课堂上讲解，但写进本书以保证学生在学习更高深的分布式题材时具有必备的基础知识。下面是一些建议，指导需要侧重实践的课程以及要求侧重研究的课程如何使用本书。需要侧重学习这两方面的研究生课程可以采用这两类建议。更多的信息可以在作者的 Prentice Hall 网站中获取，地址是 www.prenhall.com/galli。

侧重实践

下面是几个建议，适用于侧重实践的课程。

1. 给个别学生或一组学生一个或多个“项目练习”。这些练习在相关各章最后的练习题中用斜体数字指明。如果他们的设计和实现是在课堂上口述，则可以获得更多的实践经验。
2. 阅读有关实际实现的所有详细说明。
3. 在课堂上学习 Windows 2000 案例研究。
4. 建立一个个人或小组项目来研究 Windows 2000 的分布式特点。
5. 安排学生扩充或修改外科手术调度程序。可以简单地改变进程间通信的类型，也可以复杂地利用同样的分布式概念创建另一个程序。

侧重研究

下面是几个建议，适用于侧重研究的课程。

1. 让个别学生或一组学生就分布式操作系统的一个专题写一篇论文。每章后面的参考文献是较好的起点，这些练习可以包括一个口头讲解。
2. 展示一些从相关 RFC 中找到的讲授材料或者每章后面的研究论文。这些资料可以从 web 站点中获取，并将它们列在要求学生阅读的读物中。
3. 要求学生在网上查找相关的 RFC 文献和每章后面的论文，并作出摘要。
4. 在每章后面引用的参考论文中选择一部分，并创建一个活页簿以在整个课程中与此书进行同步学习。研究机构中的很多书店可为此提供所需的版权。

目 录

译者序	
前言	
第1章 分布式系统引论	1
1.1 什么是操作系统	1
1.2 什么是分布式系统	2
1.2.1 流行的网络拓扑和特点	2
1.2.2 ISO/OSI 参考模型	6
1.2.3 分布式计算模型	8
1.2.4 分布式与集中式解决方案	10
1.2.5 网络与分布式操作系统	10
1.3 什么是实时系统	11
1.3.1 实时事件的特点	11
1.3.2 影响分布式实时应用的网络特性	12
1.4 什么是并行系统	13
1.4.1 并行体系结构	13
1.4.2 并行软件范例	16
1.5 分布式应用实例	16
1.6 小结	18
1.7 参考文献	18
习题	19
第2章 内核	21
2.1 内核类型	21
2.2 进程和线程	22
2.2.1 多线程进程介绍	24
2.2.2 多线程进程范例	24
2.2.3 多线程支持	25
2.3 进程管理	26
2.3.1 进程类型	27
2.3.2 负荷分布和进程迁移	28
2.4 进程调度	30
2.4.1 识别用于调度的进程	30
2.4.2 调度器的组织	32
2.5 小结	33
2.6 参考文献	33
习题	34
第3章 进程间通信	37
3.1 选择因素	37
3.2 消息传递	37
3.2.1 阻塞原语	38
3.2.2 非阻塞原语	40
3.2.3 进程地址	40
3.3 管道	42
3.3.1 非命名管道	43
3.3.2 命名管道	43
3.4 套接字	44
3.4.1 UNIX 套接字	45
3.4.2 Java 对套接字的支持	48
3.5 远程过程调用	50
3.5.1 参数类型	50
3.5.2 数据类型支持	50
3.5.3 参数整理	50
3.5.4 RPC 绑定	51
3.5.5 RPC 认证	52
3.5.6 RPC 调用语义	52
3.5.7 SUN 的 ONC RPC	52
3.6 小结	53
3.7 参考文献	53
习题	54
第4章 内存管理	56
4.1 集中式内存管理回顾	56
4.1.1 虚拟内存	56
4.1.2 页面和段	56
4.1.3 页替换算法	58
4.2 简单内存模型	59
4.3 共享内存模型	59
4.3.1 共享内存性能	60
4.3.2 高速缓存一致性	61
4.4 分布式共享内存	61
4.4.1 分布式共享数据的方法	61
4.4.2 DSM 性能问题	66
4.5 内存迁移	66
4.6 小结	68
4.7 参考文献	69

习题	69	6.5.4 DCOM 中支持的线程模型	95
第 5 章 并发控制	71	6.5.5 DCOM 的安全策略	96
5.1 互斥和临界区	71	6.6 CORBA 概述	97
5.2 信号量	72	6.6.1 CORBA 的 ORB	97
5.2.1 信号量的缺点	73	6.6.2 CORBA 的对象适配器	98
5.2.2 信号量评估	74	6.6.3 CORBA 的消息模型	100
5.3 管程	74	6.6.4 遵从 CORBA 标准	100
5.3.1 条件变量	74	6.6.5 CORBA 到 COM 的映射	100
5.3.2 管程评估	75	6.7 小结	100
5.4 锁	75	6.8 参考文献	101
5.4.1 轮转	76	习题	101
5.4.2 原子操作和硬件支持	77	第 7 章 分布式进程管理	102
5.5 软件锁控制	78	7.1 分布式调度算法选择	102
5.5.1 集中式锁管理器	78	7.1.1 调度层次	102
5.5.2 分布式锁管理器	79	7.1.2 负荷分布目标	103
5.6 令牌传递互斥	80	7.1.3 调度的有效目标	103
5.7 死锁	80	7.1.4 处理器绑定时间	104
5.7.1 防止死锁	81	7.2 调度算法的方法	106
5.7.2 避免死锁	82	7.2.1 使用点数方法	106
5.7.3 忽略死锁	82	7.2.2 图论方法	107
5.7.4 检测死锁	82	7.2.3 探查	109
5.8 小结	83	7.2.4 调度队列	110
5.9 参考文献	84	7.2.5 随机学习	111
习题	84	7.3 协调者选举	112
第 6 章 基于对象的操作系统	86	7.4 孤儿进程	114
6.1 对象介绍	86	7.4.1 孤儿进程清除	114
6.1.1 对象定义	86	7.4.2 子进程限额	116
6.1.2 对象的评价	87	7.4.3 进程版本号	116
6.2 Clouds 对象方法	88	7.5 小结	117
6.2.1 Clouds 的对象	88	7.6 参考文献	118
6.2.2 Clouds 的线程	89	习题	118
6.2.3 Clouds 内存存储	89	第 8 章 分布式文件系统	120
6.3 Chorus V3 和 COOL V2	90	8.1 分布式名字服务	120
6.3.1 基层:COOL 内存管理	90	8.1.1 文件类型	120
6.3.2 通用运行时系统层:COOL 对象	91	8.1.2 位置透明	121
6.3.3 特定语言运行时系统层	92	8.1.3 全局命名与名字透明	123
6.4 Amoeba	92	8.2 分布式文件服务	125
6.4.1 Amoeba 对象的标识和保护	92	8.2.1 文件多样性	126
6.4.2 Amoeba 的对象通信	92	8.2.2 文件修改通知	128
6.5 分布式组件对象模型	93	8.2.3 文件服务实现	128
6.5.1 标记	94	8.2.4 文件复制	129
6.5.2 远程方法调用	95	8.3 分布式目录服务	130
6.5.3 资源回收	95	8.3.1 目录结构	131

8.3.2 目录管理	131	10.2.2 物理时间的同步	160
8.3.3 目录操作	131	10.2.3 集中式物理时间服务	161
8.4 网络文件系统	132	10.2.4 分布式物理时间服务	163
8.4.1 NFS 文件服务	132	10.3 网络时间协议	164
8.4.2 NFS 目录服务	133	10.3.1 NTP 体系结构	164
8.4.3 NFS 名字服务	134	10.3.2 NTP 设计目标	165
8.5 X.500	134	10.3.3 NTP 同步模式	166
8.5.1 X.500 文件和名字服务:信息模 型	135	10.3.4 简单网络时间协议	169
8.5.2 X.500 的目录服务:目录模型	135	10.4 逻辑时钟	169
8.6 小结	135	10.4.1 超前关系	169
8.7 参考文献	136	10.4.2 逻辑顺序	170
习题	137	10.4.3 带有逻辑时钟的总体排序	172
第 9 章 事务管理和一致性模型	139	10.5 小结	172
9.1 事务管理的动机	139	10.6 参考文献	173
9.1.1 更新遗失	139	习题	173
9.1.2 检索的不一致	140		
9.2 事务的 ACID 特性	143	第 11 章 分布式安全	175
9.3 一致性模型	145	11.1 加密和数字签名	175
9.3.1 严格一致性模型	145	11.1.1 对称加密	176
9.3.2 顺序一致性模型	145	11.1.2 非对称加密	179
9.3.3 偶然一致性模型	146	11.2 身份认证	183
9.3.4 PRAM 一致性模型	147	11.2.1 证书表	183
9.3.5 处理器一致性模型	147	11.2.2 集中式证书分送中心	186
9.3.6 弱一致性模型	148	11.3 访问控制(防火墙)	189
9.3.7 释放一致性模型	150	11.3.1 包过滤网关	189
9.3.8 懒释放一致性	151	11.3.2 代理服务	190
9.3.9 入口一致性模型	151	11.3.3 防火墙体系结构	191
9.4 两阶段提交协议	152	11.4 小结	192
9.4.1 准备提交阶段	153	11.5 参考文献	192
9.4.2 提交阶段	153	习题	193
9.5 嵌套事务	154		
9.6 事务实现中的问题	156	第 12 章 实例研究:Windows 2000	195
9.6.1 预读写	156	12.1 概述:Windows 2000 设计	196
9.6.2 中途退出的多米诺效应	156	12.2 内核模式综述	197
9.6.3 保证恢复能力	156	12.2.1 内核对象	199
9.7 小结	156	12.2.2 硬件抽象层	200
9.8 参考文献	157	12.2.3 设备驱动程序	200
习题	157	12.2.4 执行程序	200
第 10 章 分布式同步	159	12.3 即插即用	201
10.1 全局时间介绍	159	12.4 Windows 2000 中的 NT 文件系统	204
10.2 物理时钟	159	12.4.1 访问控制表	204
10.2.1 获得准确的物理时间	159	12.4.2 再解析点	205
		12.4.3 存储管理	206
		12.5 活动目录	206
		12.5.1 名字空间	208

12.5.2 通过修改日志实现复制和 可扩展性	209
12.5.3 微软的索引服务器和 HTTP 支 持	210
12.6 微软管理控制台	212
12.7 集群服务	213
12.7.1 集群服务概况	213
12.7.2 集群抽象	213
12.7.3 集群服务体系结构	214
12.7.4 为应用程序配置的集群服务	215
12.8 Windows 2000 安全性	215
12.8.1 安全配置编辑器	215
12.8.2 加密文件系统	217
12.8.3 微软安全支持提供者接口	218
12.9 HYDRA——一个瘦客户	219
12.10 小结	220
12.11 参考文献	220
习题	220
附录 A 外科手术调度程序	222
缩写词表	288
术语表	292
参考文献目录	303
索引	318

第1章 分布式系统引论

计算机系统的进展从来没有像开发与创建分布式计算这样迅猛。因此对所有计算机研究人员深入理解系统的要求迅速增涨，这种一度限于对计算机专家的要求，已经在分布式计算技术中起着前所未有的作用 [Nev95]。在基本的集中式计算环境中，简单了解一下操作系统的概念就够了。但在分布式环境中的开发人员往往不仅要懂得分布式和实时概念，甚至并行系统，还要在更大的范围中实现它们 [SHFECB93]。个人经验显示，在高级环境中开发一个应用程序，开发团队常常要花费 50% 以上的时间来实现这些高级操作系统概念。

1.1 什么是操作系统

操作系统是计算机的任务管理者，它控制、管理、协调和调度所有的处理过程和相应资源。资源主要包括：

- ◆ CPU。
- ◆ 存储模块。
- ◆ 媒体存储器（磁盘驱动器、DVD[⊖]、CD – ROM 等）。
- ◆ 声卡和显卡。
- ◆ 总线和互联网络。
- ◆ 调制解调器和网卡。
- ◆ 显示器、终端、打印机和其他输出设备。
- ◆ 键盘、鼠标、扫描仪和其他输入设备。

操作系统的管理任务主要包括：

- ◆ 过程（或对象）管理。
- ◆ 通信。
- ◆ 存储访问/管理。
- ◆ 资源调度。
- ◆ 信息和资源安全。
- ◆ 数据完整性。
- ◆ 定时。

对于这些任务和部件，我们只希望系统能创建一个统一完整的视图或“图像”，传统上，操作系统的这些管理任务都由软件来负责，但在分布式系统中并不完全如此。我们讲解主要的基本操作系统概念，即使它们在分布式环境中是由硬件实现的。

为了形象地说明这个问题对于计算机科学家的重要性和复杂性，我们可以把操作系统看作音乐指挥。当然乐章越复杂，指挥的工作量和复杂性会迅速增加。如果同时指挥多个管弦

[⊖] DVD 曾经是指数字影碟（Digital Video Disk）或者数字多功能影碟（Digital Versatile Disk）。其格式已结合到现在的 DVD 中，现在的 DVD 已经不是一个字母缩写的组合了，虽然人们经常认为它是某种缩写。

乐队演奏不同的乐章，指挥的工作会变得更复杂。同时演出地点可能不同，并且还可能有时间限制！很显然这是一个需全面解决的任务。但不管多复杂的任务我们都可以每次集中解决一个问题，一次克服一个障碍，并最终将揭示出全局的图像。在我们知道最终结果以前，整个图像会是清晰的，展现的是音乐般的美丽与和谐，而不是混乱、噪声或者杂乱无章的现象。作为第一步，让我们首先在 1.2 节到 1.4 节中分别定义分布式、实时以及并行计算的概念。

1.2 什么是分布式系统

分布式系统是一组异构计算机和处理器通过网络联接在一起，如图 1-1 所示。整组机器紧密配合工作，完成一个共同的目标。分布式操作系统的目标是为文件系统、名字空间、时间、安全和资源访问提供一个共同的、一致的全局视图。为了提供这个共同视图，所有成员系统都有许多限制和要求，因此分布式系统通常也称为紧耦合系统（每个规则总有例外！）。如果异构计算机系统通过网络互联起来，并且不是紧耦合的，那么通常称为**网络系统**。网络系统不提供共同的全局视图，也不对成员系统提出明显的要求。分布式系统或网络系统中的部件可以是简单的集中式系统部件，或者如 1.3 节所示的有实时限制的部件，甚至如 1.4 节所介绍的更复杂的并行系统。而且，分布式系统有时也会同时结合集中式的部件、实时性的部件和并行性的部件，在 1.5 节中可看到这样的实例。

由于 PC 机的每美元计算能力在迅速提高，网络和分布式系统都在不断普及。1.2.1 节介绍流行的网络；1.2.2 节提供一个由国际标准化组织（International Standards Organization, ISO）定义的开放式

3 系统互联（Open System Interconnection, OSI）参考模型；1.2.3 节介绍分布式计算模型；1.2.4 节讨论在一个分布式系统中决定用分布式还是集中式解决方案的有关问题；最后，1.2.5 节描述在分布式环境中可用的各种计算模型。

1.2.1 流行的网络拓扑和特点

虽然整本书和整个课程都是关于网络的，但本节只对网络简要介绍一下。希望这点最基本的介绍足以用来理解本书要学习的操作系统概念。

网络有两大基本类：局域网（LAN）和广域网（WAN）。典型的局域网由一个单位所拥有，范围只有几公里。尽管许多现代广域网（比如 ATM 网）已经大大降低了错误率，但在 20 世纪 90 年代中期之前，局域网的错误率往往比广域网低一千倍。一个节点可能有几个子网或更小的局域网，大的局域网可以由多个小局域网组成。小局域网可能以如下方式连接，

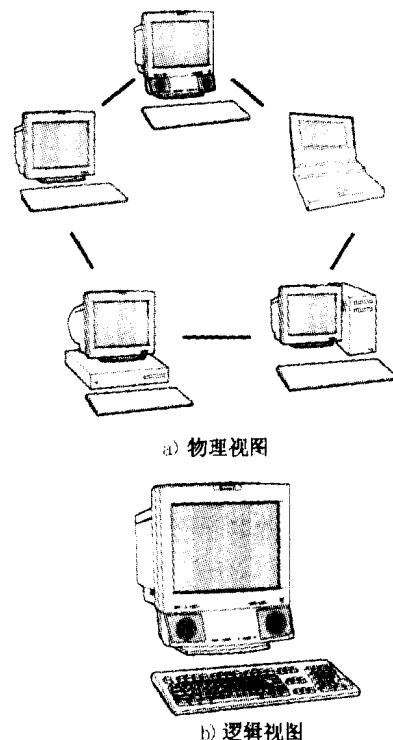


图 1-1 网络环境中的计算机

其算法可参见详细说明 1.1。

详细说明 1.1

连接局域网子网的路由器算法

令 p1 和 p2 为设备的两个端口。

令 OUTPUT (a, b) 将端口 b 接到的消息内容输出到端口 a 所连的网络。

令 DESTINATION (a, b) 为真，当且仅当端口 a 接到的消息必须使用端口 b 来送达目的地。

令 DIF_PROTOCOL (a, b) 为真，当且仅当端口 a 和端口 b 使用不同的网络协议并且 DESTINATION (a, b) 为真。

令 CONVERT (a, b) 将端口 a 接到的消息转换为使用端口 b 所用协议的消息，并发送到端口 b。

转发器算法：

```
While ()
{
    OUTPUT (P1, P2);           // P2 接到的所有消息都输出到 P1
    OUTPUT (P2, P1);
}
```

网桥算法：

```
While ()
{
    if DESTINATION (P1, P2);
        Then OUTPUT (P1, P2);           // 仅在到达目的地需要经过的情况下转发
    if DESTINATION (P2, P1);
        Then OUTPUT (P2, P1);
}
```

路由器算法：

```
While ()
{
    If DIF_PROTOCOL (P1, P2);           // 仅在到达目的地需要经过的情况下转发
        Then CONVERT (P1, P2);           // 只在必须的情况下转换
    Else If DESTINATION (P1, P2);
        Then OUTPUT (P1, P2);           // 仅在到达目的地需要经过的情况下转发
    If DIF_PROTOCOL (P2, P1);
        Then CONVERT (P2, P1);
    Else If DESTINATION (P1, P2);
        Then OUTPUT (P1, P2);
}
```

1. **转发器**：一个非智能型的设备，只是简单地将一个网络里的东西全部转发到另一个网络中去。这两个网络必须使用相同的协议，也就是说它们是同类型的网络。
2. **网桥**：一个智能设备，只将一个子网里的数据转发到目标子网，或者到达目标所需要经过的子网。这些网络必须实现相同的网络协议。
3. **路由器**：比网桥更先进之处在于它能连接使用不同协议的局域网网段。路由器能同时连接两个以上的网络。
4. **主干**：不包含用户的局域网，它只连接其他网络而不连接用户。如图 1-2 所示。

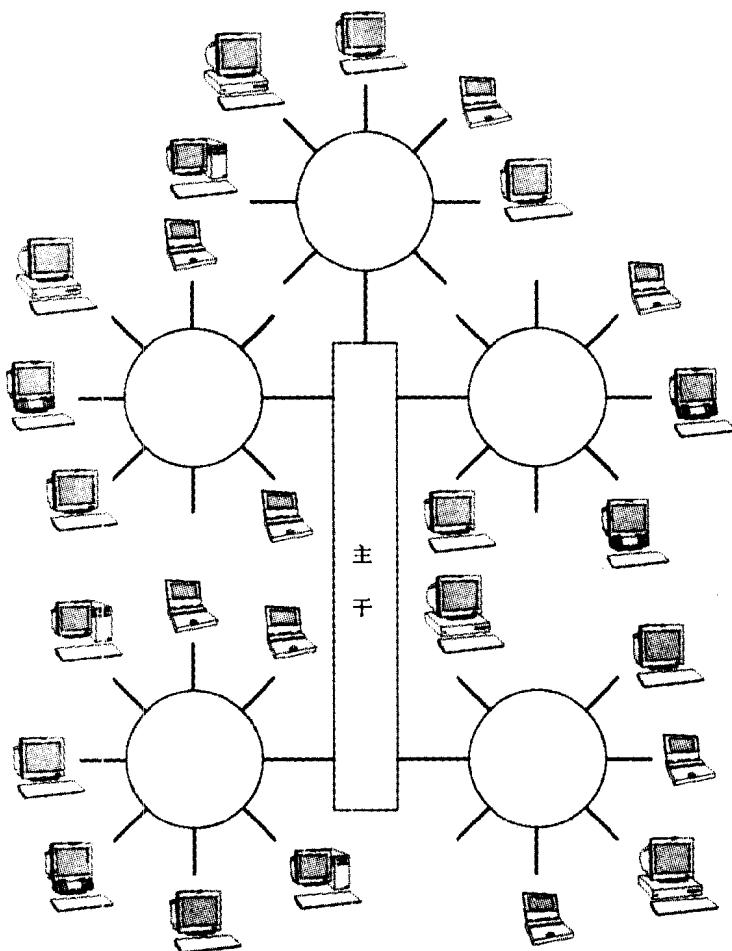


图 1-2 通过主干来连接多个局域网子网

局域网可以是有线的或无线的。有线局域网通过物理线缆连接，无线局域网通过无形的通信通道如红外线或无线电连接。有线局域网有三种常用拓扑，如图 1-3 所示。其中最常用的两种拓扑是以太网（在 IEEE802.3 标准中的正式命名是 CSMA/CD [IEEE85a]）和令牌环网（在 IEEE802.5 标准中指明 [IEEE85b]）。由于这些拓扑提供的操作服务类型根本不同，它们之间没有直接竞争。以太网类似于将邮票贴在信封上，而令牌环网类似于次日发送并要求回执。基本以太网运行在 10Mbps 上，运行在 100Mbps 上的快速以太网是当前流行的版本，而千兆以太网则运行在 1000Mbps 上。以太网以尽可能快的速度发送信息，但不确认信息是否

被正确地接收。而且标准以太网协议不支持信息优先权，也不保证传输时间。对大多数应用来说，它在达到容量的 50% 之前运行得很好，之后它开始迅速变慢，并且几乎所有的 CSMA/CD 网络在超过约 60% 的容量时将崩溃而无法传送任何信息。

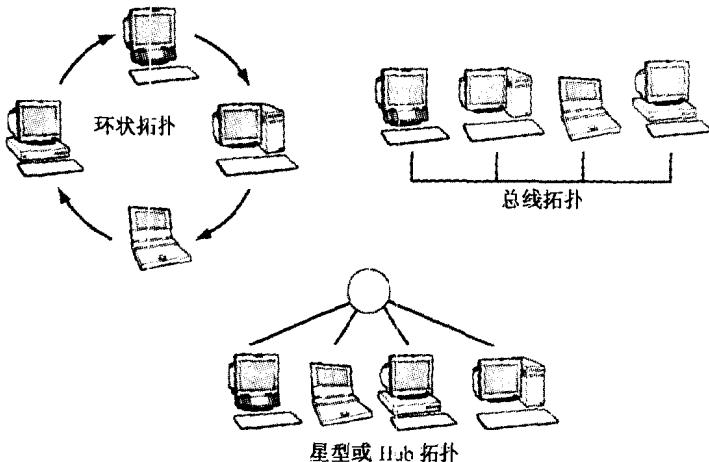


图 1-3 常用有线局域网拓扑结构

令牌环模拟带回执的次日传送，能确认所有信息被正确送达。它允许优先级并保证传输时间，这使它对实时应用特别有吸引力。与以太网不同，令牌环网在高负荷时也不会崩溃，而是继续顺利地发送信息。这个协议是考虑优先权的公平访问。令牌环网依所用种类不同，运行在 4 到 16Mbps 之上，跟 10 Base[⊖] 以太网同一个档次，但实现成本是以太网的两倍。

光纤分布式数据接口（Fiber Distributed Data Interface, FDDI）是另一种常用网络协议，经常用于主干或高速局域网。FDDI 操作在 100Mbps，其实现是建立在两个反相令牌环网上。这个协议设计支持同步信息流，因此对实时应用很有吸引力。

无线局域网由于不断增长的移动需求和有线局域网固定的重配置费用而变得越来越普遍。无线局域网通常有一个较高的初装成本，但维护成本较低。无线局域网可以是基于无线电，也可以是基于红外线。基于无线电的局域网当前可达 4Mbps，并且容易安装，因为无线电信号能穿过不透明的物体传播，但这也意味着网络信号传到建筑物之外。因此基于无线电的局域网有时不如线局域网或红外线局域网安全。不过，由于传播频谱技术的利用，信号可以以噪音的方式传播或者令信号不断跳频传播，这样无线电局域网也可以更加安全。此外，有线局域网所采用的各种安全手段如第 11 章的加密技术，也可以用于无线电局域网。红外线局域网如今可达 10Mbps，但红外线信号不能穿过不透明物体。这一点使红外线局域网很难安装，因为发送器和接收器必须在一条直线上并且相互可见。而且红外线信号可能受到光噪音的干扰，影响了可靠性。

许多分布式系统并不是在局域网，而是在广域网。要将一个局域网连接到一个广域网，需要用到网关。信息发送到网络上在传输之前被分成更小的片，称做帧或包。有两个常用的方法可以用来在广域网上发送这些包。第一个方法是电路交换，用于公用交换电话网（public-switched telephone network, PSTN）。这类网络在进行任何数据传输之前提供一个固定的数据

[⊖] 10Base 以太网一般指运行在 10Mbps 下的以太网络。

速率通道，这个保留的通道完全用于当前的数据传输。这个方法需要大量的设置时间并降低了网络吞吐量，但能保证一个信息在其传输过程中只有传输延迟。第二个方法是包交换。每个信息被分成称为包的小单元，包的大小有固定的上限。它们可以使用一个虚拟电路以面向连接的方式发送，也可以使用数据报以无连接的方式发送。虚拟电路要求同一个信息的每个包都沿着同样的路径按正常顺序传输，这通常要求先发送一个设置包。相反，数据报对相关包的传输路径和次序不作限制，在目的地把收到的数据报依正常次序重新排列，并再组装成原始信息。这个方法可靠性较弱，但更健壮和高效。

广域网常用的网络协议包括帧中继和异步传输模式（asynchronous transfer mode, ATM）。帧中继允许网络消息包含拥塞控制和速率信息，以便在网络过于繁忙时丢掉该消息。ATM 网络以速度和支持多媒体文件传输而著称，尽管这些特点并非 ATM 网络所独有。ATM 网中的每个包尺寸相同，都是 53 字节。ATM 网把这些 53 字节的包称作单元（cell）。通过使用不同的帧格式，ATM 适应层（ATM adaptation layer, AAL）支持固定速率和可变位速率传输，以及面向连接和无连接传输。它在设计时考虑了支持实时多媒体应用。当然，像 Internet 这样的广域网使用了所有网络类型的组合。

1.2.2 ISO/OSI 参考模型

ISO/OSI 参考模型 [ISO84] 描述了在各种网络模型中所执行的所有计算任务的一种可能划分。这个参考模型如图 1-4 所示划分成七层，层数越低越接近网络。第一到第三层依赖网

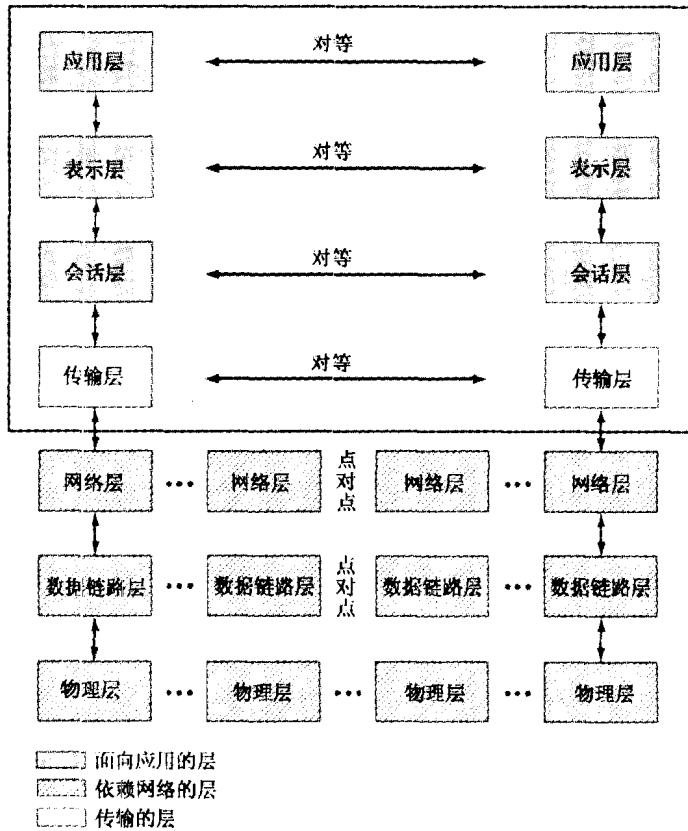


图 1-4 ISO/OSI 参考模型