

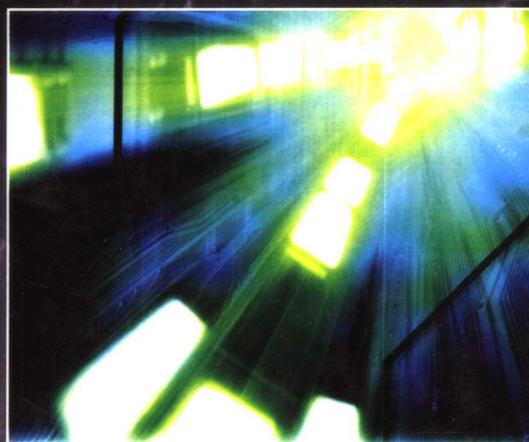


计 算 机 科 学 丛 书

原书第4版

分布式系统 概念与设计

(英) George Coulouris Jean Dollimore Tim Kindberg 著 金蓓弘 曹冬磊 等译



fourth edition

DISTRIBUTED SYSTEMS CONCEPTS AND DESIGN

George Coulouris
Jean Dollimore
Tim Kindberg



Distributed Systems Concepts and Design

Fourth Edition



机械工业出版社
China Machine Press

计 算 机 科 学 丛 书

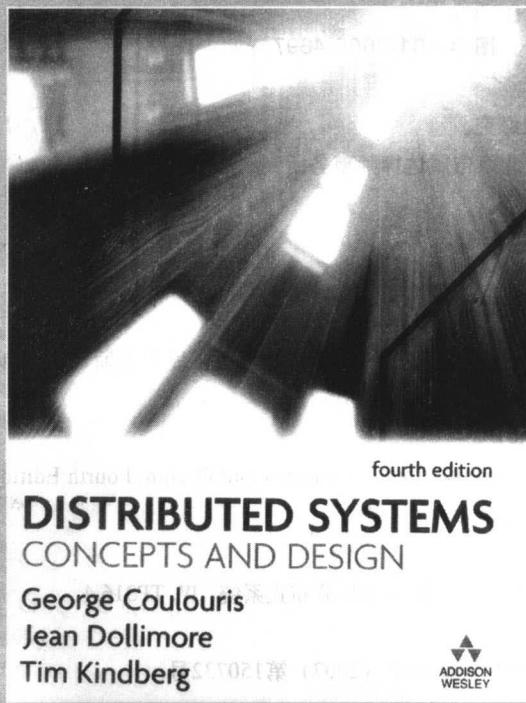
原书第4版

TP316.4/5=2

2008

分布式系统 概念与设计

(英) George Coulouris Jean Dollimore Tim Kindberg 著 金蓓弘 曹冬磊 等译



Distributed Systems
Concepts and Design
Fourth Edition

机械工业出版社
China Machine Press

本书旨在全面介绍因特网及其他常用分布式系统的原理、体系结构、算法和设计，内容涵盖分布式系统的相关概念、安全、数据复制、组通信、分布式文件系统、分布式事务等，以及相关的前沿主题，包括Web服务、网格、移动系统和无处不在系统等。

本书素材丰富、内容充实、深入浅出，每章后都有相关的习题，并有配套网站提供本书的学习和教学资源。本书可作为相关专业本科生及研究生的分布式系统课程的教材，也可供广大技术人员参考。

George Coulouris, Jean Dollimore, Tim Kindberg: *Distributed Systems: Concepts and Design*, Fourth Edition (ISBN: 0-321-26354-5).

Copyright © Addison Wesley Publishers Limited 1988, 1994, © Pearson Education Limited 2001, 2005.

This translation of *Distributed Systems: Concepts and Design*, Fourth Edition (ISBN: 0-321-26354-5) is published by arrangement with Pearson Education Limited.

All rights reserved.

本书中文简体字版由英国Pearson Education培生教育出版集团授权出版。

本书版权登记号：图字：01-2005-4697

版权所有，侵权必究。

本书法律顾问 北京市展达律师事务所

图书在版编目 (CIP) 数据

分布式系统：概念与设计（原书第4版）/（英）库劳里斯（Coulouris, G.）等著；金蓓弘等译。—北京：机械工业出版社，2008.1

（计算机科学丛书）

书名原文：Distributed Systems: Concepts and Design, Fourth Edition

ISBN 978-7-111-22438-9

I. 分… II. ①库… ②金… III. 分布式系统 IV. TP316.4

中国版本图书馆CIP数据核字（2007）第150732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：朱劼

北京京北制版厂印刷 新华书店北京发行所发行

2008年1月第1版第1次印刷

184mm×260mm · 36.25印张

定价：69.00元

凡购本书，如有倒页、脱页、缺页，由本社发行部调换

本社购书热线：(010) 68326294

译 者 序

随着网络技术的发展和计算机应用的深入，分布式系统已成为目前主流的软件系统。

本书介绍了分布式系统的概念、基本原理和核心技术，内容涉及通信、中间件、系统基础设施、分布式数据处理以及分布式算法等。通过阅读本书，读者既可以从系统层面上了解分布式系统构造的基本原理，又可以从算法层面上获知分布式系统构造的核心技术。

本书素材广泛、内容充实、叙述深入浅出、条理清楚，每章后都配有练习题，并有配套网站提供大量学习和教学资料。因此，本书可以作为高等院校高年级本科生和研究生“分布式计算”及相关课程的教材或参考书，也可供分布式计算领域的科研人员阅读、参考。

本书第3版和第4版都是由中国科学院软件研究所金蓓弘研究员组织和主持翻译的。

本书第3版由若干同仁通力合作、共同翻译而成。前言、第1、2、4、10、11章由金蓓弘翻译，第6、8、15、16、18章由李剑博士翻译，第12、13、14章由丁柯博士翻译，第5、17章由刘绍华博士翻译，第3、9章由阮彤博士翻译，第3章由王仲玉翻译，第7章由刘志军翻译，由金蓓弘通校了第3版全部译文。

本书第4版新增了三章。其中，第10章由张发恩翻译，第16章由张英翻译，第19章由臧志翻译。由金蓓弘、曹冬磊博士通校了第4版全部译文。

由于时间和水平所限，翻译中不当之处在所难免，欢迎广大读者提出批评和指正。

感谢机械工业出版社华章分社在引进、编辑、出版本书中所做的努力。

译 者

2007年8月于北京

前　　言

在因特网和Web走向成熟、能够支持多种分布式系统之际，本书的第4版问世了。如今，分布式系统的规模已远远超过了本书第3版出版时的预期。

本书旨在介绍因特网和其他分布式系统所蕴涵的原理、体系结构、算法和设计。前两章是概念上的简介，概括分布式系统的特征和必须在设计中解决的挑战：可伸缩性、异构性、安全性和故障处理。这两章也给出了理解进程交互、故障和安全的抽象模型。后续几章关注连网、进程间通信、远程调用和中间件、操作系统和命名。

接着，我们论及一些比较成熟的主题，包括安全、数据复制、组通信、分布式文件系统、分布式事务、CORBA、分布式共享内存和多媒体系统。此外，还会讨论一些新的主题：Web服务、XML、网格、对等、移动和无处不在系统。与这些主题相关的算法将在相关主题中讨论。我们还将另辟几章讨论时序、协调和协定。

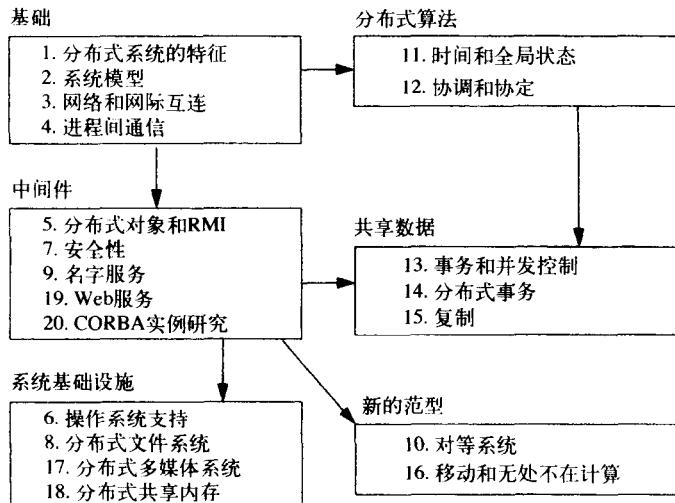
目的和读者群

本书可作为本科生教材和研究生的入门教材，也可作为自学教材。本书采用自顶向下的方法，首先叙述在分布式系统设计中要解决的问题，然后，通过抽象模型、算法和对广泛使用的系统实例进行详细研究的方式，描述成功开发系统的方法。本书覆盖的领域有足够的宽度和广度，使得读者在读完本书后能继续研究分布式系统文献中的大多数研究论文。

本书希望读者具有面向对象编程、操作系统以及基本的计算机体系结构知识。本书涵盖与分布式系统有关的计算机网络的知识，包括因特网、广域网、局域网和无线网的基本技术。本书中的大部分算法和接口用Java描述，有一小部分用ANSI C描述。为了表述上的简洁明了，还将使用一种从Java/C中派生出来的伪代码。

本书的组织

下图表明本书的各章可划分为六个主题。它说明了本书的结构，为教师、读者提供了一个推荐的导航路径，以便于他们理解分布式系统设计中的不同子领域。



目 录

译者序	
前言	
第1章 分布式系统的特征	1
1.1 简介	1
1.2 分布式系统的实例	2
1.2.1 因特网	2
1.2.2 企业内部网	3
1.2.3 移动计算和无处不在计算	3
1.3 资源共享和Web	5
1.4 挑战	10
1.4.1 异构性	11
1.4.2 开放性	11
1.4.3 安全性	12
1.4.4 可伸缩性	13
1.4.5 故障处理	14
1.4.6 并发性	15
1.4.7 透明性	15
1.5 小结	16
练习	17
第2章 系统模型	19
2.1 简介	19
2.2 体系结构模型	20
2.2.1 软件层	20
2.2.2 系统体系结构	22
2.2.3 变体	23
2.2.4 接口和对象	27
2.2.5 分布式体系结构的设计需求	27
2.3 基础模型	30
2.3.1 交互模型	31
2.3.2 故障模型	34
2.3.3 安全模型	36
2.4 小结	39
练习	40
第3章 网络和网际互连	42
3.1 简介	42
3.2 网络类型	44
3.3 网络原理	46
3.3.1 数据包的传输	47
3.3.2 数据流	47
3.3.3 交换模式	47
3.3.4 协议	48
3.3.5 路由	52
3.3.6 拥塞控制	54
3.3.7 网际互连	55
3.4 因特网协议	57
3.4.1 IP寻址	59
3.4.2 IP协议	60
3.4.3 IP路由	62
3.4.4 IPv6	65
3.4.5 移动IP	67
3.4.6 TCP和UDP	68
3.4.7 域名	69
3.4.8 防火墙	70
3.5 实例研究：以太网、WiFi、蓝牙和ATM	72
3.5.1 以太网	73
3.5.2 IEEE 802.11无线LAN	76
3.5.3 IEEE 802.15.1蓝牙无线PAN	78
3.5.4 异步传输模式网络	80
3.6 小结	82
练习	82
第4章 进程间通信	84
4.1 简介	84
4.2 因特网协议的API	85
4.2.1 进程间通信的特征	85
4.2.2 套接字	86
4.2.3 UDP数据报通信	87
4.2.4 TCP流通信	90
4.3 外部数据表示和编码	93
4.3.1 CORBA的公共数据表示	94
4.3.2 Java对象序列化	95
4.3.3 可扩展标记语言	97
4.3.4 远程对象引用	99
4.4 客户－服务器通信	100
4.5 组通信	105

4.5.1 IP组播——组通信的实现	106	7.2 安全技术概述	174
4.5.2 组播的可靠性和排序	108	7.2.1 密码学	175
4.6 实例研究：UNIX中的进程间通信	108	7.2.2 密码学的应用	175
4.6.1 数据报通信	109	7.2.3 证书	177
4.6.2 流通信	110	7.2.4 访问控制	178
4.7 小结	110	7.2.5 凭证	180
练习	111	7.2.6 防火墙	181
第5章 分布式对象和远程调用	114	7.3 密码算法	181
5.1 简介	114	7.3.1 密钥（对称）算法	184
5.2 分布式对象间的通信	116	7.3.2 公钥（不对称）算法	186
5.2.1 对象模型	117	7.3.3 混合密码协议	188
5.2.2 分布式对象	117	7.4 数字签名	188
5.2.3 分布式对象模型	118	7.4.1 公钥数字签名	189
5.2.4 RMI的设计问题	120	7.4.2 密钥数字签名——MAC	189
5.2.5 RMI的实现	122	7.4.3 安全摘要函数	190
5.2.6 分布式无用单元收集	125	7.4.4 证书标准和证书权威机构	191
5.3 远程过程调用	126	7.5 密码实用学	192
5.4 事件和通知	129	7.5.1 密码算法的性能	192
5.4.1 分布式事件通知的参与者	131	7.5.2 密码学的应用和政治障碍	193
5.4.2 实例研究：Jini分布式事件规约	132	7.6 案例研究：Needham-Schroeder、 Kerberos、TLS和802.11 WiFi	194
5.5 实例研究：Java RMI	133	7.6.1 Needham-Schroeder认证协议	194
5.5.1 创建客户和服务器程序	136	7.6.2 Kerberos	195
5.5.2 Java RMI的设计和实现	138	7.6.3 使用安全套接字确保 电子交易安全	199
5.6 小结	139	7.6.4 IEEE 802.11 WiFi 安全 设计中的缺陷	201
练习	139	7.7 小结	203
第6章 操作系统支持	142	练习	204
6.1 简介	142	第8章 分布式文件系统	205
6.2 操作系统层	143	8.1 简介	205
6.3 保护	144	8.1.1 文件系统的特点	207
6.4 进程和线程	145	8.1.2 分布式文件系统的需求	208
6.4.1 地址空间	146	8.1.3 实例研究	209
6.4.2 新进程的生成	147	8.2 文件服务体系结构	210
6.4.3 线程	149	8.3 实例研究：SUN网络文件系统	214
6.5 通信和调用	157	8.4 实例研究：Andrew文件系统	222
6.5.1 调用性能	158	8.4.1 实现	223
6.5.2 异步操作	162	8.4.2 缓存的一致性	225
6.6 操作系统的体系结构	164	8.4.3 其他方面	227
6.7 小结	167	8.5 最新进展	228
练习	167	8.6 小结	232
第7章 安全性	169	练习	232
7.1 简介	169	第9章 名字服务	234
7.1.1 威胁和攻击	170		
7.1.2 保护电子事务	172		
7.1.3 设计安全系统	173		

9.1 简介	234	11.6.4 在同步系统中判定可能的 ϕ 和明确的 ϕ	295
9.2 名字服务和域名系统	236	11.7 小结	296
9.2.1 名字空间	237	练习	296
9.2.2 名字解析	239	第12章 协调和协定	298
9.2.3 域名系统	241	12.1 简介	298
9.3 目录服务	246	12.2 分布式互斥	300
9.4 实例研究：全局名字服务	246	12.3 选举	305
9.5 实例研究：X.500目录服务	248	12.4 组播通信	308
9.6 小结	251	12.4.1 基本组播	309
练习	252	12.4.2 可靠组播	310
第10章 对等系统	253	12.4.3 有序组播	312
10.1 简介	253	12.5 共识和相关问题	317
10.2 Napster及其遗留系统	256	12.5.1 系统模型和问题定义	317
10.3 对等中间件	257	12.5.2 同步系统中的共识问题	320
10.4 路由覆盖	259	12.5.3 同步系统中的拜占庭将军问题	320
10.5 路由覆盖实例研究： Pastry和Tapestry	261	12.5.4 异步系统的不可能性	323
10.5.1 Pastry	261	12.6 小结	324
10.5.2 Tapestry	266	练习	325
10.6 应用实例研究：Squirrel、 OceanStore和Ivy	267	第13章 事务和并发控制	327
10.6.1 Squirrel Web缓存	267	13.1 简介	327
10.6.2 OceanStore文件存储	269	13.1.1 简单的同步机制（无事务）	328
10.6.3 Ivy文件系统	272	13.1.2 事务的故障模型	329
10.7 小结	274	13.2 事务	329
练习	275	13.2.1 并发控制	332
第11章 时间和全局状态	277	13.2.2 事务放弃时的恢复	334
11.1 简介	277	13.3 嵌套事务	336
11.2 时钟、事件和进程状态	278	13.4 锁	337
11.3 同步物理时钟	279	13.4.1 死锁	342
11.3.1 同步系统中的同步	280	13.4.2 在加锁机制中增加并发度	345
11.3.2 同步时钟的Cristian方法	281	13.5 乐观并发控制	346
11.3.3 Berkeley算法	281	13.6 时间戳排序	349
11.3.4 网络时间协议	282	13.7 并发控制方法的比较	353
11.4 逻辑时间和逻辑时钟	284	13.8 小结	354
11.5 全局状态	286	练习	355
11.5.1 全局状态和一致割集	287	第14章 分布式事务	359
11.5.2 全局状态谓词、稳定性、 安全性和活性	288	14.1 简介	359
11.5.3 Chandy和Lamport的“快照”算法	289	14.2 平面分布式事务和嵌套分布式事务	359
11.6 分布式调试	291	14.3 原子提交协议	361
11.6.1 观察一致的全局状态	293	14.3.1 两阶段提交协议	362
11.6.2 判定可能的 ϕ	294	14.3.2 嵌套事务的两阶段提交协议	364
11.6.3 判定明确的 ϕ	294	14.4 分布式事务的并发控制	367

14.4.3 乐观并发控制	368	16.4.2 感知体系结构	434
14.5 分布式死锁	369	16.4.3 位置感知	438
14.6 事务恢复	374	16.4.4 小结和前景	441
14.6.1 日志	375	16.5 安全和私密性	442
14.6.2 影子版本	377	16.5.1 背景	442
14.6.3 为何恢复文件需要事务 状态和意图列表	378	16.5.2 一些解决办法	443
14.6.4 两阶段提交协议的恢复	378	16.5.3 小结和前景	447
14.7 小结	380	16.6 自适应	447
练习	381	16.6.1 内容的上下文敏感自适应	448
第15章 复制	383	16.6.2 适应变化的系统资源	449
15.1 简介	383	16.6.3 小结和前景	450
15.2 系统模型和组通信	385	16.7 Cooltown实例研究	450
15.2.1 系统模型	385	16.7.1 Web存在	451
15.2.2 组通信	386	16.7.2 物理超链接	452
15.3 容错服务	390	16.7.3 互操作和eSquirt协议	454
15.3.1 被动（主备份）复制	392	16.7.4 小结和前景	455
15.3.2 主动复制	393	16.8 小结	455
15.4 高可用服务的实例研究：gossip 体系结构、Bayou和Coda	394	练习	456
15.4.1 gossip体系结构	395	第17章 分布式多媒体系统	458
15.4.2 Bayou系统和操作变换方法	401	17.1 简介	458
15.4.3 Coda文件系统	402	17.2 多媒体数据的特征	461
15.5 复制数据上的事务	407	17.3 服务质量管理	462
15.5.1 复制事务的体系结构	407	17.3.1 服务质量协商	464
15.5.2 可用拷贝复制	409	17.3.2 许可控制	467
15.5.3 网络分区	410	17.4 资源管理	468
15.5.4 带验证的可用拷贝	411	17.5 流适应	469
15.5.5 法定数共识方法	411	17.5.1 调整	470
15.5.6 虚拟分区算法	413	17.5.2 过滤	471
15.6 小结	415	17.6 实例研究：Tiger视频文件服务器	471
练习	415	17.7 小结	474
第16章 移动计算和无处不在计算	417	练习	474
16.1 简介	417	第18章 分布式共享内存	476
16.2 关联	423	18.1 简介	476
16.2.1 发现服务	424	18.1.1 消息传递机制和DSM	477
16.2.2 物理关联	427	18.1.2 DSM的实现方法	478
16.2.3 小结和前景	428	18.2 设计和实现问题	479
16.3 互操作	428	18.2.1 结构	479
16.3.1 易变系统的面向数据编程	429	18.2.2 同步模型	480
16.3.2 间接关联和软状态	432	18.2.3 一致性模型	481
16.3.3 小结和前景	433	18.2.4 更新选项	483
16.4 感知和上下文敏感	433	18.2.5 粒度	485
16.4.1 传感器	434	18.2.6 系统颤簸	485
		18.3 顺序一致性和lvy实例研究	485
		18.3.1 系统模型	486

18.3.2 写失效	487	一种网格应用	518
18.3.3 失效协议	488	19.7.2 数据密集型科学应用的特征	518
18.3.4 一个动态分布式管理器算法	489	19.7.3 开放的网格服务体系结构	519
18.3.5 系统颤簸	490	19.7.4 一些网格应用的例子	521
18.4 释放一致性和Munin实例研究	491	19.7.5 Globus工具包	522
18.4.1 内存访问	491	19.8 小结	523
18.4.2 释放一致性	492	练习	524
18.4.3 Munin	493	第20章 CORBA实例研究	526
18.5 其他一致性模型	494	20.1 简介	526
18.6 小结	495	20.2 CORBA RMI	527
练习	496	20.2.1 CORBA客户和服务器实例	529
第19章 Web服务	498	20.2.2 CORBA体系结构	532
19.1 简介	498	20.2.3 CORBA接口定义语言	534
19.2 Web服务	499	20.2.4 CORBA远程对象引用	537
19.2.1 SOAP	501	20.2.5 CORBA语言映射	538
19.2.2 Web服务与分布式对象 模型的比较	504	20.2.6 CORBA与Web的集成	538
19.2.3 在Java中使用SOAP	505	20.3 CORBA服务	539
19.2.4 Web服务和CORBA的比较	508	20.3.1 CORBA名字服务	540
19.3 服务描述和Web服务接口定义语言	509	20.3.2 CORBA事件服务	542
19.4 Web服务使用的目录服务	512	20.3.3 CORBA通知服务	543
19.5 XML安全性	513	20.3.4 CORBA安全服务	544
19.6 Web服务的协作	516	20.4 小结	544
19.7 实例研究：网格	517	练习	545
19.7.1 World-Wide Telescope——		索引	548
		参考文献	[⊕]

第1章 分布式系统的特征

分布式系统是其组件分布在连网的计算机上，组件之间通过传递消息进行通信和动作协调的系统。该定义导出了分布式的下列特征：组件的并发性、缺乏全局时钟、组件故障的独立性。

我们给出分布式的三个例子：

- 因特网。
- 企业内部网，它是因特网的一部分，一般由一个机构负责管理。
- 移动计算和无处不在计算。

资源共享是构造分布式的动力。资源可以由服务器管理，由客户访问，或它们被封装成对象，由其他客户对象访问。作为一个资源共享的例子，我们将讨论Web并介绍它的主要特征。

构造分布式的挑战是处理其组件的异构性、开放性（允许增加或替换组件）、安全性、可伸缩性（用户数量增加时能正常运行的能力）、故障处理、组件的并发性和透明性问题。

1.1 简介

计算机网络无处不在。因特网也是其中的一个，因为它是许多种网络组成的。移动电话网、协作网、企业网、校园网、家庭网、车内网，所有这些，既可单独使用，又可相互结合，它们具有相同的本质特征，这些特征使得它们可以放在分布式的标题下来研究。本书旨在解释影响系统设计者和实现者的连网的计算机的特征，给出已有的可帮助完成设计和实现分布式的任务的主要概念和技术。

我们把分布式的定义成一个其硬件或软件组件分布在连网的计算机上，组件之间通过传递消息进行通信和动作协调的系统。这个简单的定义覆盖了所有可部署连网计算机的系统。

由一个网络连接的计算机可能在空间上的距离不等。它们可能分布在地球上不同的洲，也可能在同一栋楼或同一个房间里。分布式的有如下显著特征：

并发：在一个计算机网络中，执行并发程序是常见的行为。用户可以在各自的计算机上工作，在必要时共享诸如Web页面或文件之类的资源。系统处理共享资源的能力会随着网络资源（例如，计算机）的增加而提高。在本书的许多地方将描述有效实施这种额外能力的方法。对共享资源的并发执行的程序的协调也是一个重要和重复提及的主题。

缺乏全局时钟：在程序需要协作时，通过交换消息来协调它们的动作。密切的协作通常取决于对程序动作发生的时间的共识。但是，事实证明，网络上的计算机与时钟同步所达到的准确性是有限的，即没有一个正确时间的全局概念。这是由于通信仅仅是通过网络发送消息这个事实带来的直接结果。定时问题和它们的解决方案将在第11章描述。

故障独立性：所有的计算机系统都可能出故障，一般由系统设计者负责为可能的故障设计结果。分布式的可能以新的方式出现故障。网络故障导致网上互连的计算机的隔离，但这并不意味着它们停止运行，事实上，计算机上的程序不能检测到网络是出现故障还是网络运行得比通常慢。类似的，计算机的故障或系统中程序的异常终止（崩溃），并不能马上被与它通信的其他组件感知。系统的每个组件会单独地出现故障，而其他组件还在运行。分布式的这个特征所带来的后果将是本书的一个反复提及的主题。

构造和使用分布式的动力来源于对共享资源的期望。“资源”一词是相当抽象的，但它很

好地描述了能在连网的计算机系统中共享的事物的范围。它涉及的范围从硬件组件（如硬盘、打印机）到软件定义的实体（如文件、数据库和所有的数据对象）。它包括来自数字摄像机的视频流和移动电话呼叫所表示的音频连接。

本章主要论述分布式系统的本质，以及成功部署分布式系统所面临的挑战，1.2节展示了分布式系统的一些重要的例子，以及构造系统所需的组件和这些组件的作用。1.3节描述了在万维网环境中资源共享系统的设计。1.4节则阐述了分布式系统设计者所要面对的重要挑战：异构性、开放性、安全性、可伸缩性、故障处理、并发性和对透明性的要求。

1.2 分布式系统的实例

本节选用的实例基于大家熟悉和广泛使用的计算机网络，包括因特网、企业内部网和新兴的基于移动设备的网络技术。我们举这些例子主要是说明由计算机网络支持的服务和应用具有很广的范围，然后从这些系统开始讨论支持系统实现的技术问题。

1.2.1 因特网

因特网是一个巨大的由多种类型计算机网络互连的集合。图1-1摘取了因特网的部分典型组成。因特网上的计算机程序通过传递消息进行交互，采用了一种公共的通信手段。因特网通信机制（因特网协议）的设计和构造是一项重大的技术成果，它使得一个在某处运行的程序能给另一个地方的程序发送消息。



图1-1 因特网的典型部分

因特网也是一个非常大的分布式系统。它使得世界各地的用户能利用诸如万维网、电子邮件和文件传送等服务。（有时，不确切地说，Web等同于因特网。）服务集是开放的，它能够通过服务器计算机和新的服务的增加而被扩展。图1-1还展示了许多企业内部网——由公司和其他组织操作的子网。因特网服务提供商（ISP）是给个体用户和小型组织提供调制解调器链接和其他类型连接的公司，使他们能获得因特网上的服务；同时提供诸如电子邮件和Web主机等本地服务。企业内部网通过主干网实现互相链接。主干网是具有高传送能力的网络链接，通常采用卫星连接、光缆和其他高宽带线路。

因特网能够提供多媒体服务，使用户能获得包括音乐、广播和TV频道在内的音频和视频数据，并召开电话和视频会议。由于目前因特网没有提供为单个数据流预留网络容量所必需的设施，因此它处理多媒体数据这类特殊通信需求的能力还很有限。第17章将讨论分布式多媒体系统的需求。

因特网的实现和它支持的服务已经解决了分布式系统的许多问题（包括在1.4节中定义的大多数问题）。本书将着重阐述这些解决方案，并在适当的时候说明它们的适用范围和局限性。

1.2.2 企业内部网

企业内部网是因特网的一部分，它是独立管理的，具有一个可被配置来执行本地安全策略的边界。图1-2给出了一个典型的企业内部网。它由几个通过主干网连接的局域网（LAN）组成。每个企业内部网的网络配置都由管理企业内部网的组织负责，这种管理的范围差异很广，可以从单个场地的LAN到（可能分布在不同的国家）属于同一个公司或组织的若干部门的若干LAN。

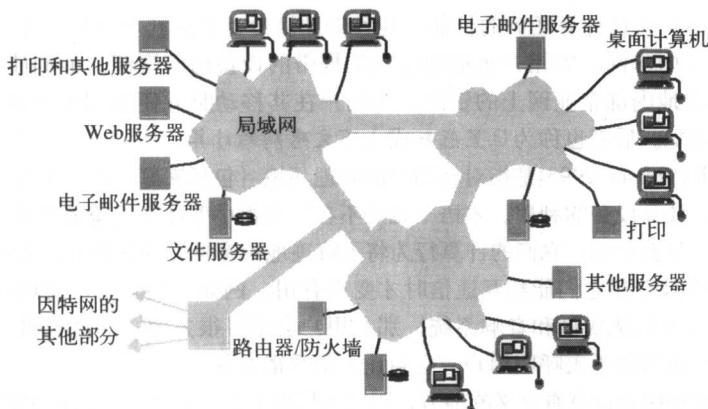


图1-2 典型企业内部网

企业内部网通过路由器连接到因特网，因此企业内部网内的用户能使用因特网上的服务（如Web或电子邮件）。企业内部网也允许其他企业内部网的用户访问它提供的服务。许多组织需要保护他们自己的服务以免其他地方可能有恶意的用户未经授权便使用。例如，公司不希望保密的信息被竞争对手获取；医院不希望病人病历被爆光。公司也希望自己免受病毒入侵这样的有害程序影响，以免他人对企业内部网内的计算机的攻击并摧毁有用的数据。

防火墙的作用是通过防止未授权消息进出网络来保护企业内部网。防火墙是通过过滤进出的消息来实现的，例如根据消息的源地址或目的地址进行过滤。一个防火墙可能仅允许与电子邮件和Web访问有关的消息进出它所保护的企业内部网。

一些组织根本不希望将他们的内部网连接到因特网上。例如，警察机关和其他安全法律执行机构可能至少有几个内部网和外部世界隔离；一些军事组织在战争时期会将它们的内部网与因特网断连。但即使是这样的组织，也希望从大量的采用因特网通信协议的应用和系统软件中受益。这些组织通常采用的解决方案是像上面描述的那样操作企业内部网，但不与因特网相连。这样一个企业内部网可以没有防火墙，或者，从另一个角度来看，这有可能是最有效的防火墙——与因特网没有任何物理连接。

在设计用于企业内部网的组件时，会出现下列主要问题：

- 需要文件服务以便用户能共享数据，文件服务的设计将在第8章讨论。
- 防火墙试图阻止对服务的合法访问——当需要在企业内部网和外部用户之间共享资源时，防火墙必须增加细粒度的安全机制。这部分内容将在第7章讨论。
- 用于软件安装和支持的花销是一个重要的问题。通过使用诸如网络计算机和瘦客户这样的系统体系结构，能减少这些开销。这方面的内容参见第2章。

1.2.3 移动计算和无处不在计算

设备小型化和无线网络方面的技术进步已经逐步使得小型和便携式计算设备集成到分布式系统中。这些设备包括：

- 笔记本电脑。
- 手持设备，包括个人数字助理（PDA）、移动电话、传呼机、摄像机和数码相机。
- 可穿戴设备，如具有类似PDA功能的智能手表。
- 嵌入在家电（如洗衣机、高保真音响系统、汽车和冰箱）中的设备。

这些设备大多数具有可携带性，再加上它们可以在不同地方方便地连接到网络的能力，使得移动计算成为可能。移动计算，也叫游牧计算[Kleinrock 1997]，是指用户在移动或参观某处（而不是在通常环境下）执行计算任务的性能。在移动计算中，远离其本地的企业内部网（指工作环境或其住处的企业内部网）的用户也能通过他们携带的设备访问资源。他们能继续访问因特网，继续访问在他们本地内部企业网上的资源。为用户在其移动时方便地利用周围资源（如打印机）的设备也在不断增加。后者也称为位置感知或上下文感知的计算。

无处不在计算[Weiser 1993]是指对在用户的物理环境（包括家庭、办公室和其他地方）中存在的多个小型便宜的计算设备的利用。术语“无处不在”意指小型计算设备最终将在日常不会引人注意的物品中普及。也就是说，它们的计算行为将无痕迹地紧密捆绑到这些日常物品的物理功能上。

各处的计算机只有在它们能相互通信时才变得有用。例如，如果用户能通过一个“通用远程控制”设备控制家里的洗衣机和音响系统，那么用户会觉得很方便。而洗衣机在完成洗衣后能通过一个智能报警器或智能手表呼叫用户，也会让人觉得很方便。

无处不在计算和移动计算有交叉的地方，因为从原理上说，移动用户能利用遍布各处的计算机。但一般而言，它们是不同的。无处不在计算能让呆在家里或医院这样单一的环境中的用户受益。类似地，即使移动计算只涉及常见的分立的计算机和设备（如笔记本电脑和打印机），它还是有优势的。

图1-3显示了一个正在访问一个组织的用户。该图显示出用户本地的内部网和用户正在访问的内部网。两个企业内部网通过因特网相连。

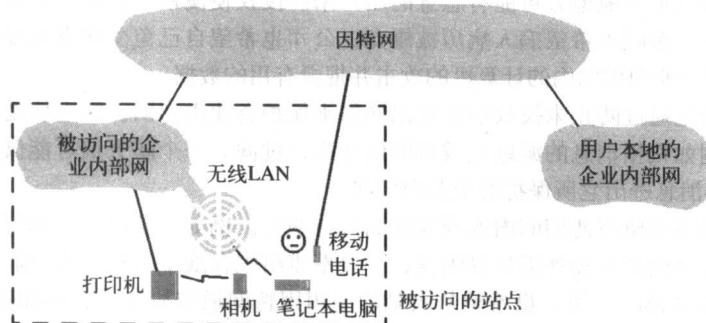


图1-3 分布式系统中的便携式设备和手持设备

用户可以使用三种无线连接。笔记本电脑可以连接到被访问组织的无线LAN。无线LAN覆盖方圆几百米的范围（即建筑物的一层）。它通过网关连接到被访问组织企业内部网。用户还有一部连到因特网的移动电话，在电话中这些信息按页显示在电话的显示屏上。最后，用户携带一台数码相机，它能通过一个个人局域无线网络（其覆盖范围大约为10m）与打印机这样的设备通信。

利用适当的系统基础设施，用户能用他们携带的设备完成一些简单的任务。当用户连接到被访问的站点时，他能通过移动电话从Web服务器上取得最新的股票价格。在与访问企业开会时，通过把数码相机的照片直接发送到会议室的一台可用的打印机上，用户就能展示最近的照片。这仅仅要求相机和打印机之间具有无线连接。原则上，用户可以利用无线LAN或是有线的以太网链接从笔记本电脑上把文件发送到同一台打印机。

移动计算和无处不在计算是一个热门的研究领域，第16章将继续讨论这两个主题。

1.3 资源共享和Web

用户已经习惯了资源共享带来的好处，以至于很容易忽视它们的重要性。大家通常共享硬件资源（如打印机）、数据资源（如文件）和具有特定功能的资源（如搜索引擎）。

从硬件资源的观点看，大家共享设备（如打印机和磁盘）可以减少花费，但共享与用户应用、日常工作和社会活动有关的更高层的资源对用户的意义更大。例如，用户关心以共享数据库或Web页面方式出现的共享数据，而不是实现上述服务的硬盘和处理器。类似地，用户关心诸如搜索引擎或货币转换器的共享资源，而不关心提供这些服务的服务器。7

实际上，资源共享的模式随其工作范围和与用户工作的密切程度的不同而不同。一种极端是，Web上的搜索引擎给全世界的用户提供工具，而用户之间并不需要直接接触；另一种极端是，在计算机支持协同工作（CSCW）中，若干直接进行合作的用户在一个小型封闭的小组中共享诸如文档之类的资源。用户在地理上的分布以及用户之间进行共享的模式决定了系统必须提供协调用户动作的机制。

我们使用术语服务表示计算机系统中管理相关资源并提供功能给用户和应用的一个单独的部分。例如，我们通过文件服务访问共享文件；通过打印服务发送文件到打印机；通过电子支付服务购买商品。仅仅通过服务提供的操作可以实现对服务的访问。例如，一个文件服务提供对文件的read、write和delete操作。

服务将资源访问限制为一组定义良好的操作，这属于标准的软件工程实践，同时它也反映出分布式系统的物理组织。分布式的资源是物理地封装在计算机内的，其他计算机只能通过通信才能访问。为了实现有效的共享，每个资源必须由一个程序管理，这个程序提供通信接口使得对资源进行可靠和一致的访问和更新。

大多数读者很熟悉术语服务器，它指的是在连网的计算机上的一个运行程序（一个进程），这个程序接收来自其他计算机上正在运行的程序请求，执行一个服务并适当地响应。发出请求的进程称为客户。请求以消息的形式从客户发送到服务器，应答以消息的形式从服务器发送到客户。当客户发送一个要执行的操作请求，就称客户调用那个服务器上的操作。客户和服务器之间的完整交互，即从客户发送一个请求到它接收到服务器的应答，称为一个远程调用。

一个进程可能既是客户又是服务器，因为服务器有时调用其他服务器上的操作。术语“客户”和“服务器”仅仅是针对在一个请求中扮演的角色而言。就它们扮演的角色不同这点而言，客户是主动的，服务器是被动的；服务器是连续运行的，而客户所持续的时间只是客户所属的那部分应用程序持续的时间。

注意，默认情况下，术语“客户”和“服务器”指的是进程而不是运行客户或服务器的计算机，虽然在日常用法中这些术语也指计算机。另一个不同（见第5章）是在用面向对象语言实现的分布式系统中，资源被封装成对象，并由客户对象访问，这时，称一个客户对象调用了一个服务器对象上的方法。

许多（但不是所有的）分布式系统可以完全用客户和服务器交互的形式来构造，万维网、电子邮件和连网的打印机都满足这种模式。第2章将讨论除客户—服务器系统之外的其他系统类型。

一个正在执行的Web浏览器是一个客户的例子。Web浏览器与Web服务器通信，从服务器上请求Web页面。下面将详细讨论Web。8

万维网

万维网[www.w3.org I, Berners-Lee 1991]是一个不断发展的系统，用于发布和访问因特网上的资源和服务。通过常用的Web浏览器，用户可以检索和查看多种类型的文档、收听音频文件、观看视频文件、与无数服务进行交互。

Web是1989年在瑞士的欧洲原子能研究中心(CERN)诞生的，作为通过因特网连接的物理学家之间交换文档用的工具[Berners-Lee 1999]。Web的一个关键特征是它在所存储的文档中提供了超文本结构，超文本结构反映了用户对知识组织的要求。这意味着文档包含链接(或超链接)，链接指向其他存储在Web上的文档和资源。

对Web用户来说，当他遇到文档中的一幅图像或一段文字时，它很可能伴有到相关文档和其他资源的链接。链接的结构可以简单，也可以复杂，可加入的资源集是无限的，即链接的Web确实是世界范围的。Bush[1945]在五十年前就设想出了超文本结构，因特网的发展使得这个想法能在世界范围内得到证实。

Web是一个开放的系统，它可以被扩展，并且在不妨碍已有功能的前提下用新的方法实现扩展(见1.4.2节)。

首先，它的操作是基于被自由发布和广泛实现的通信标准和文档标准的。例如，浏览器的类型是多种多样的。在多数情况下，每种浏览器可以在多个平台上实现；有多种Web服务器实现。一种构造的浏览器能从不同构造的服务器中检索资源。所以，用户能访问大多数设备(从移动电话到桌面计算机)上的浏览器。

其次，相对于能在其上发布和共享的“资源”的类型而言，Web是开放的。在Web上，最简单的资源是一个Web页面或其他能保存在文件中并提交给用户的内容，如程序文件、介质文件和PostScript和PDF格式的文件。如果有人新发明了一种图像存储格式，那么这种格式的图像能马上在Web上发布。用户需要一种查看这种新格式图像的工具，而浏览器以“帮助者”应用和“插件程序”的形式来支持新的内容显示功能。

Web的发展已超越这些简单的数据资源而开始包含服务，如电子化的商品购买。Web一直在发展，但其基本的体系结构没有改变。Web基于以下三个主要的标准技术组件：

- 超文本标记语言(HTML)是页面在Web浏览器上显示时指定其内容和布局的语言。
- 统一资源定位器(URL)用于识别保存成Web一部分的文档和其他资源。第9章将讨论有关的Web标识符的其他术语。
- 具有标准交互规则(超文本传送协议HTTP)的客户—服务器系统体系结构，浏览器和其他客户可利用标准交互规则从Web服务器上获取文档和其他资源。图1-4给出了一些Web服务器和向它们发送请求的浏览器。用户可以定位和管理位于因特网上任何地方的他们自己的Web服务器，这是一个很重要的特征。

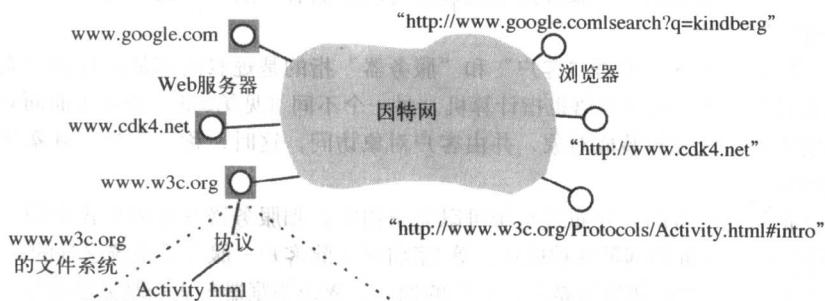


图1-4 Web服务器和Web浏览器举例

接下来我们依次讨论这些组件，并解释用户获取Web页面并单击页面上的链接时，浏览器和Web服务器的操作。

HTML 超文本标记语言将这些内容[www.w3.org II]用于指定组成Web页面内容的文本和图像，指定它们以何种布局方式和组织方式将这些内容显示给用户。Web页面包含结构化的成分，

如标题、段落、表格和图像。HTML也用于指定链接和与链接相关联的资源。

用户可使用标准的文本编辑器手写生成HTML，或用能识别HTML的“所见即所得型”编辑器，根据用户给出的一个图示布局生成HTML。下面是一段典型的HTML文本：

```
<IMG SRC = "http://www.cdk4.net/WebExample/Images/earth.jpg "> 1  
<P> 2  
Welcome to Earth! Visitors may also be interested in taking a look at the 3  
<A HREF = "http://www.cdk4.net/WebExample/moon.html "> Moon </A>. 4  
<P> 5
```

这段HTML文本保存在一个Web服务器可以访问的文件（例如earth.html文件）中。浏览器从Web服务器（本例中是一个位于名为www.cdk4.net的计算机上的服务器）中检索这个文件的内容，浏览器读取从服务器返回的内容后，把它变成格式化的文本和图像，以大家熟悉的方式放到Web页面上。只能由浏览器（不是服务器）解释HTML文本，但是服务器确实通知了浏览器它所返回的内容的类型，用于区分html文件和其他文件（如PostScript文件）。服务器能从文件的扩展名“.html”中推断出内容类型。

注意，HTML的指令（即标记）放在尖括号里，如〈P〉。例子中的第一行确定了一个包含图片显示的文件，图片的URL是http://www.cdk4.net/WebExample/Images/earth.jpg。第二行的指令表示开始新段落，第三行和第四行包含要在Web页面上以标准的段落格式显示的文本，第五行的指令表示该段落结束。10

其中，第四行指定了Web页面上的一个链接。它包含词“Moon”，该词位于两个匹配的HTML标记〈A HREF…〉和〈/A〉中间。这些标记之间的文本在Web页面上显示时是以链接的形式出现的。大多数浏览器在默认情况下给链接的文本加下划线，所以，用户看到的上面的段落将是：

Welcome to Earth! Visitors may also be interested in taking a look at the Moon.

浏览器记录了链接的显示文本和包含在<A HREF…>标记中的URL之间的关联，在这个例子中是：

<http://www.cdk4.net/WebExample/moon.html>

当用户单击文本时，浏览器获取由相应URL识别的资源，并将它显示给用户。在这个例子中，资源是一个HTML文件，它指定了关于月亮的一个Web页面。

URL 统一资源定位器[www.w3.org III]的作用是识别资源。在Web体系结构文档中使用的术语是统一资源标识符（URI），在不引起混淆的前提下，本书使用更为人们所熟悉的术语URL。浏览器检查URL以便从Web服务器上访问相应的资源。有时用户在浏览器中键入一个URL。更常见的方法是用户单击一个链接或选择一个书签，由浏览器查找相应的URL；或当浏览器去取一个Web页面里的内嵌资源（如一个图像）时，由浏览器查找相应的URL。

按绝对完整的格式，每一个URL有两个不可或缺的组成部分：

模式：模式特定的位置

第一个成分“模式”声明了URL的类型。要求URL能识别各种资源。例如，mailto:joe@anISP.net标识出一个用户的电子邮件地址；ftp://ftp.downloadIt.com/software/aProg.exe标识一个用文件传送协议（FTP）获取而不是用更常用的HTTP协议获取的文件。模式的其他例子有“nntp”（用于指定一个Usenet新闻组）和“mid”（用于标识一个邮件消息）。

从Web可访问（利用URL中的模式指示器）的资源类型的角度来说，它是开放的。如果有人发明了一种新的有用的“widget”资源（可能用它专有的寻址方案定位widget，用它专有的协议访问widget）那么大家就能使用widget: …格式的URL。当然，浏览器必须具备使用新的“widget”协议的能力，这一点可通过增加一个插件实现。