

计算机科学本科核心课程教材

# Practical Distributed Processing



## 分布式处理实践

Phillip J. Brooke Richard F. Paige 著

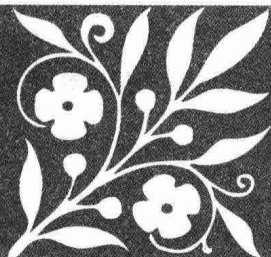
孙燕 等译



清华大学出版社

计算机科学本科核心课程教材

# Practical Distributed Processing



## 分布式处理实践

清华大学出版社  
北京

English reprint edition Copyright © 2009 by **Springer-Verlag and TSINGHUA UNIVERSITY PRESS**.  
Original English language title from Proprietor's edition of the Work.

Original English language title: **Practical Distributed Processing** by Phillip J. Brooke, Richard F. Paige,  
Copyright © 2009 All Rights Reserved.

This edition has been authorized by Springer-Verlag(Berlin/Heidelberg/New York) for sale in the People's  
Republic of China only and not for export therefrom.

本书中文翻译由 Springer-Verlag 授权给清华大学出版社出版发行。  
北京市版权局著作权合同登记号 图字 01-2009-3651 号

本书封面贴有清华大学出版社激光防伪标签, 无标签者不得销售。  
版权所有, 侵权必究。侵权举报电话: 010-62782989 13701121933

### 图书在版编目(CIP)数据

分布式处理实践 / (美)布鲁克(Brooke, P. J.), (美)佩奇(Paige, R. F.)著; 孙燕等译.  
--北京: 清华大学出版社, 2010.3

书名原文: Practical Distributed Processing  
ISBN 978-7-302-21781-7

I. ①分… II. ①布… ②佩… ③孙… III. ①分布式处理系统 IV. ①TP338

中国版本图书馆 CIP 数据核字(2010)第 001734 号

责任编辑: 龙啟铭

责任校对: 徐俊伟

责任印制: 王秀菊

出版发行: 清华大学出版社

<http://www.tup.com.cn>

社 总 机: 010-62770175

投稿与读者服务: 010-62776969, [c-service@tup.tsinghua.edu.cn](mailto:c-service@tup.tsinghua.edu.cn)

质 量 反 馈: 010-62772015, [zhiliang@tup.tsinghua.edu.cn](mailto:zhiliang@tup.tsinghua.edu.cn)

地 址: 北京清华大学学研大厦 A 座

邮 编: 100084

邮 购: 010-62786544

印 刷 者: 北京季蜂印刷有限公司

装 订 者: 三河市新茂装订有限公司

经 销: 全国新华书店

开 本: 185×230

印 张: 14

字 数: 239 千字

版 次: 2010 年 3 月第 1 版

印 次: 2010 年 3 月第 1 次印刷

印 数: 1~3000

定 价: 29.00 元

---

本书如存在文字不清、漏印、缺页、倒页、脱页等印装质量问题, 请与清华大学出版社出版部联系调  
换。联系电话: 010-62770177 转 3103 产品编号: 029948-01

在生活中,分布式系统应用已经越来越流行和普遍。分布式系统具有类似的模式:使用并发实现资源的共享和高性能;涉及同步,并协调活动时间;需要管理和减缓故障。在实际的分布式处理系统的构建中,不仅要涉及并发的常见内容(如竞争条件、信号量和互斥),还要涉及网络系统的实际实现,并考虑底层的操作系统。

本书从实用的角度讲解分布式处理的整个构建过程,以及在此过程中应用的工具、技术和原理。

根据所讲述的内容,本书大致分为如下几部分:

- 第1部分涉及第1章~第6章,讲解有关分布式处理的基本概念、并发的概念以及模型、操作系统的并发处理、进程通信、协议的使用。
- 第2部分的内容包括第7章~第9章,讲解工程问题,包括安全性和分布式处理的语言,并提供构建分布式处理的示例和案例研究。
- 第3部分包括第10章~第11章,第10章介绍一个游戏的分布式处理部分,从而将前面各章中所学的知识,特别是那些实用的内容,比如BSD套接字、TCP/IP和UDP,与复杂的分布式系统构建过程关联起来。第11章对全书内容加以概括总结。

本书的每一章都有相应的练习题,并在书后提供了参考答案。读者可以通过练习,巩固所学知识。

本书主要由孙燕翻译,参加翻译的还有韦笑、王雷、李志云、李晓春、陈安华、孙宏、赵成璧、侯佳宜、许伟、戴文雅、于樊鹏、刘朋、王嘉佳、李腾、邓卫、邓凡平、陈磊、李建锋、刘延军、魏宇等人。最后,祝广大读者从本书中挖掘出更多的宝藏。

# 前言

## 概述和目标

本书根据两位作者的讲稿和在分布式处理领域所做的研究而编写。在分布式处理的教学中,我们发现,不仅要涉及并发的常见内容(例如,竞争条件、信号量和互斥),还要涉及网络系统的实际实现,并考虑底层的操作系统。

关注并发原理、操作系统和编程的书有很多,但从实用的角度对分布式处理进行系统介绍的书几乎没有。因此,本书的目的是向读者清楚地展现分布式系统的整个构建过程,并介绍这个过程中应用的工具、技术和原理。该书旨在考虑这一工程过程中的重要问题,包括访问和分析分布式系统各部分时所用的模型及这些模型的实现技术(例如套接字)以及任何实际分布式系统必须考虑的有关协议和安全方面的问题。

## 结构和特点

本书分为3大部分:基础(包括原理和实际实现方面)、工程问题(包括安全性和语言方面)以及示例和案例研究。第3部分是前两部分的综合。

本书第1部分从第1章开始,该章对分布式处理进行了概述。分布式处理和并发的核心概念(像信号量、死锁和命名)是在第2章中讨论的。第3章提供了几种精确的并发模型,帮助读者理解第2章中核心概念背后的基本原理和实施方法;该章还介绍了在构建分布式系统时如何使用这些模型。第4章讨论操作系统以及它们在构建分布式系统中的地位。随后的第5章讲述进程间通信。第6章研究了协议,它们是构建实际分布式系统时要考虑的重要部分;该章还提供了协议示例以及设计协议时会遇到的各种挑战。

本书的第2部分转向工程问题。第7章考虑分布式系统安全性的常规问题,该章从工程的角度来探讨安全性问题。值得注意的是,安全性不仅仅是对通信进行加密,还需要对整个系统进行统筹考虑。第8章介绍语言的选

择,并着重于说明如何应用这些语言来构建分布式系统,以解释前面几章中出现的诸多概念,包括理论和实际两方面。

第3部分着重于示例和案例研究,旨在将前面各章结合起来。第9章通过工程生命周期来选择有关构建分布式系统的案例研究,但没有实现这些系统,目标是演示这些系统的构建过程,注重的是确定需求、甄选协议和进行设计。第10章给出了一个完整的示例:构建一个网络游戏。在该章中,首先确定了技术要求和商业要求,然后给出设计,讨论协议和安全方面,并构思实现,旨在将前面各章的内容有机地结合在一起。最后,第11章总结了本书给出了将来的一些工作、阅读资料和试验。

每一章都含有一个系统的总结,并专门给出了大量的练习题。有些练习题注重的是对当前章中问题的思考,还有一些练习题着重于构建分布式系统的一小部分内容。

本书还包含一个详尽的词汇表。

## 读者对象

本书针对有一定的(顺序)程序设计经验,但没有学过操作系统或并发的相关课程的学生。熟悉C语言编程对阅读和理解本书中的内容是有帮助的,但并不是必须的。那些没有太多编程经验的读者需要再学习有关一些特定编程语言的书籍(请参见参考书目)。

本书中的内容已向大二、大三的本科生传授过。通过扩充一些内容,特别是增加一个大的项目,本书也可以适用于软件工程专业的大四学生。

## 教师注意事项

本书的内容可以用作分布式系统的一学期课程或单元,也可以作为分布式系统要素占重要地位的有关软件工程长期课程的一部分。如果将本书用作一门单独课程,建议大致按照该书的顺序进行学习,当然,也可以在不影响整体方向的前提下,对相关章节的顺序进行调整。老师可以根据学生的背景和兴趣,有针对性地选择第3章中的某些正式模块,但强烈推荐包含状态图,因为它们非常重要,而且在第10章的完整性示例中会用到它们。对于以前已经学过操作系统相关课程的学生,可以跳过第2章中有关信号量的内容,但对于初学者来说,这部分内容是比较有用的。

作为阅读和描述项目,第9章中的案例研究是非常有价值的,它们是由课堂上的一部分学生提供的。

第10章的内容是分布式系统中较大项目的基础。实际上,我们已经基于此章中的内容来向本科生介绍研究和开发项目以及向硕士生授课。11.2.2小节给出了有关较大型和较长期项目调查研究方面的一些建议。

附录 A 包含许多练习题的概要解答或提示。示例代码,包括第 10 章中讨论的网络游戏,可以从本书的配套网站获得: <http://www.scm.tees.ac.uk/p.j.brooke/dpb/>。将来可能增加对第 10 章中给出的多人游戏的进一步扩充。欢迎大家提出意见和建议,并且可以直接向作者提出。

## 致谢

感谢英国提兹赛德大学(University of Teesside)和英国约克大学(University of York)(以及以前所在的英国普利茅斯大学[University of Plymouth]和加拿大约克大学)的同事们给出的建议、意见和有用的指导。还要感谢本书的审阅者们提出的及时、有用的建议。特别要感谢 Springer 的 Catherine Brett 和 Wayne Wheeler 在编辑和出版过程中所给的帮助以及 Ian Mackie 和 Helen Callaghan 对整个过程的启动。

本书的编写过程中使用了一些软件工具,没有这些工具,编写过程将要艰难得多。在此特别感谢 LATEX 和 Subversion 的开发者和维护者。

Phil 感谢 Christine、Rebecca 和 Alexander,他们在协助编写代码和文字方面付出了大量的精力和时间。

Rich 感谢 Angelika 对自己(和 Phil)编写、重写和测试过程的支持。

Phil 和 Rich 同时感谢审校者对编写错误的指出和改正。当然,对于未指出的错误,责任在我们作者。

Phil Brooke

Rich Paige

# 目 录

<b>第 1 章 什么是分布式处理</b> .....	1
1.1 概述 .....	1
1.2 计算和网络的发展 .....	2
1.3 分布式处理 .....	2
1.4 应用领域 .....	3
1.5 模型 .....	3
1.6 移动代码 .....	4
1.7 分布式系统面临的挑战 .....	4
1.8 本章小结 .....	5
练习题 .....	6
<b>第 2 章 并发的概念</b> .....	7
2.1 概述 .....	7
2.2 并发中的有关结构 .....	8
2.3 命名和寻址 .....	8
2.3.1 名称和地址示例 .....	9
2.3.2 地址映射机制 .....	9
2.4 共享与同步 .....	11
2.4.1 资源分配 .....	11
2.4.2 示例：文件同步 .....	12
2.5 低级同步 .....	12
2.5.1 竞争条件 .....	13
2.5.2 互斥 .....	13
2.5.3 信号量 .....	13
2.5.4 管程 .....	16
2.5.5 会合 .....	18



2.6	定时和实时系统	18
2.7	可靠性	19
2.7.1	故障和失效的类型	19
2.7.2	对故障的响应	20
2.8	服务器类型	20
2.9	簇、负载平衡和网格	21
2.10	本章小结	22
	练习题	22
<b>第3章</b>	<b>并发模型</b>	<b>24</b>
3.1	概述	24
3.2	状态机和自动机	24
3.3	SPIN 和 Promela	25
3.4	进程代数	26
3.4.1	通信顺序进程	27
3.4.2	$\pi$ 演算和灵活性	28
3.5	Linda	30
3.5.1	JavaSpaces	31
3.6	再谈死锁	33
3.7	本章小结	34
	练习题	35
<b>第4章</b>	<b>操作系统中的并发</b>	<b>37</b>
4.1	概述	37
4.2	为什么使用操作系统	37
4.3	进程和线程	38
4.3.1	进程概念	39
4.3.2	CPU 中的用户模式和管理员模式	39
4.3.3	多任务	40
4.3.4	线程和轻量级进程	40
4.4	Linux 中的进程和线程示例	41
4.4.1	Fork	41
4.4.2	Pthreads	43
4.5	Ada 中的任务处理	45

4.6 本章小结	47
练习题	47
<b>第5章 进程间通信</b>	<b>49</b>
5.1 概述	49
5.2 Linux 中的 Pthreads IPC 示例	50
5.2.1 互斥量和共享内存	50
5.2.2 信号量	52
5.2.3 条件变量	54
5.3 Ada 中的互斥	57
5.4 BSD 套接字	60
5.5 TCP 客户-服务器示例	61
5.5.1 一个简单的 TCP 服务器	61
5.5.2 字符串终止和网络	67
5.5.3 一个简单的 TCP 客户端	68
5.5.4 具有名称查找功能的 TCP 客户端	72
5.6 UDP 客户-服务器示例	72
5.6.1 UDP 服务器	72
5.6.2 UDP 客户端	75
5.7 双向通信	77
5.8 分叉模式的 TCP 服务器	79
5.9 阻塞处理和 select	83
5.9.1 用于双向通信的 select	84
5.9.2 用于多个连接的 select	86
5.10 容错和 IPC 定时处理	87
5.11 本章小结	87
练习题	87
<b>第6章 协议</b>	<b>89</b>
6.1 概述	89
6.2 协议的目的	89
6.3 协议中的有关问题	90
6.3.1 高级和低级协议	90

6.3.2	消息 .....	91
6.3.3	平台依赖 .....	92
6.3.4	容错 .....	93
6.4	定义协议 .....	95
6.4.1	编码 .....	96
6.4.2	表示法 .....	96
6.5	示例: HTTP .....	97
6.6	示例: SMTP .....	98
6.7	示例: 交替位协议 .....	99
6.8	本章小结 .....	101
	练习题 .....	101
<b>第7章</b>	<b>安全性</b> .....	<b>103</b>
7.1	概述 .....	103
7.2	定义、概念和术语 .....	103
7.2.1	风险、威胁和漏洞 .....	103
7.2.2	安全性的目标 .....	104
7.2.3	设计 .....	104
7.3	分布式系统中的安全问题 .....	105
7.4	加密 .....	107
7.4.1	加密示例: 数字签名 .....	108
7.4.2	密钥管理 .....	108
7.4.3	将公开密钥与用户匹配 .....	109
7.5	案例研究: Needham-Schroeder .....	109
7.6	实际问题 .....	110
7.6.1	C 编程 .....	110
7.6.2	Web 应用程序 .....	111
7.6.3	操作系统和网络问题 .....	112
7.6.4	SSL .....	112
7.6.5	使用 SSL .....	113
7.7	本章小结 .....	114
	练习题 .....	115

<b>第 8 章 语言和分布式处理</b> .....	116
8.1 概述 .....	116
8.2 语言的适用性 .....	116
8.3 C 中的分布式处理 .....	117
8.3.1 C 概述 .....	118
8.3.2 调试 C .....	118
8.4 Java 中的分布式处理 .....	119
8.4.1 概述: RMI 模型 .....	119
8.4.2 示例 .....	120
8.4.3 其他方法 .....	123
8.5 Ada 中的分布式处理 .....	123
8.6 Eiffel 和 SCOOP 中的分布式处理 .....	125
8.6.1 SCOOP: Eiffel 的一种并发模型 .....	126
8.6.2 相关工作和原型 .....	128
8.7 语言的比较 .....	129
8.7.1 语言模式 .....	130
8.7.2 类型规则 .....	130
8.7.3 网络支持 .....	131
8.7.4 并发支持 .....	132
8.7.5 进程间通信支持 .....	132
8.8 本章小结 .....	133
练习题 .....	133
<b>第 9 章 构建分布式系统</b> .....	134
9.1 概述 .....	134
9.2 方法 .....	135
9.3 案例分析: 电子邮件 .....	135
9.3.1 典型使用和需求 .....	136
9.3.2 平台和语言要求 .....	136
9.3.3 结构 .....	137
9.3.4 协议和形式 .....	137
9.3.5 示例: 使用 PHP 发送电子邮件 .....	139
9.4 案例分析: 安全外壳 .....	142

9.4.1	典型使用和需求	142
9.4.2	平台要求	143
9.4.3	结构	143
9.4.4	协议	144
9.5	案例分析：版本控制和同步	145
9.5.1	典型使用和需求	146
9.5.2	平台要求	146
9.5.3	结构	146
9.5.4	协议	147
9.6	案例分析：Web 应用程序	148
9.7	本章小结	149
	练习题	149
<b>第 10 章</b>	<b>案例分析：一个网络游戏</b>	<b>151</b>
10.1	动机和组织	151
10.2	大概结构和基本需求	152
10.3	分析和设计	152
10.3.1	大纲用例	153
10.3.2	详细设计问题	155
10.3.3	安全性	156
10.4	协议	157
10.4.1	协议消息	157
10.4.2	客户端登录	157
10.4.3	地图服务器启动和关闭	158
10.4.4	UDP 消息	159
10.4.5	协议备注	161
10.4.6	数据视图	161
10.5	实现	162
10.5.1	管理服务器	163
10.5.2	地图服务器	164
10.5.3	玩家客户端	164
10.5.4	运行示例	165
10.6	测试	165
10.7	本章小结	166

练习题	166
<b>第 11 章 结束</b>	<b>168</b>
11.1 小结	168
11.2 建议	169
11.2.1 将来方向	171
11.2.2 有趣的项目	171
<b>附录 A 练习题：提示和注解</b>	<b>173</b>
<b>附录 B 关于示例代码</b>	<b>197</b>
<b>参考文献</b>	<b>198</b>
<b>词汇表</b>	<b>202</b>

## 什么是分布式处理

### 本章主要内容

- 分布式处理概述
- 计算和网络的发展
- 常见应用领域(数据库、文件系统、服务)
- 模型
- 移动代码

### 1.1 概 述

我们一起来思考以下日常事件。假设用户站在自助银行提款机前,准备取一些钱。先是插入银行卡,输入有关信息(例如密码),选择所需的交易(例如,取款并打印凭条),接着选择要取的金额,几秒钟后,用户会收到钱和交易凭条。这是一个简单的日常事务,许多人都依赖这一过程,但现场背后的处理过程是非常复杂的。用户借助通信基础设施、计算机软件和安全机制,及时获得所需响应(以及所取现金)。

上面的银行系统就是一种分布式系统。这种系统包括网络计算机上的多个组件(可以是硬件、软件或两者都是)。这些组件通过消息进行通信,并通过消息协调活动和传递结果,例如用户所取的现金。

当今基于计算机的系统通常都是分布式系统。在日常生活中,有许多分布式系统的实例:

- Internet 就是分布式系统的一种典型实例:它将许多独立的计算机连接起来,并定义了运行在这些计算机上的软件进行通信的标准协议。标准应用程序,例如 Web 浏览器、文件程序、对等文件共享应用程序和计算机游戏,都使用 Internet 进行操作。
- 当我们去一个旅行社预定假日行程,比如坐飞机旅行,旅行社会使用订票系统来查找两座城市之间的可用航班。这种系统称为SABRE[81],它搜索不同航空公

司的分布式数据库,找出不同类型的机票,并为顾客预订机票。整个过程用时很短。

- Search for Extraterrestrial Intelligence(SETI)是一种大型分布式系统,它利用个人计算机(连接到 Internet 并运行 SETI@home 软件)分析来自 Arecibo 射电望远镜的数据,查找某些模式,这些模式可以确定收到的是来自地球以外的射电信号,这些信号说明地球外有可能存在高智能的生命。这种分布式程序特别有意思的地方是它利用闲置的时钟周期,也就是闲置的计算资源。

本书关注的是构建分布式系统过程中会遇到的实际问题,特别是着重介绍构建实用分布式系统时需要的基本理论、原理和技术,本书较后的章节(例如第 9 章)中,通过构建一些有趣的重要系统,举例说明如何将这些理论、原理和技术结合起来。

## 1.2 计算和网络的发展

传统计算是单个机器运行由程序员和管理员提供的单机任务。在 20 世纪 70 年代,程序员开始认识到将计算机组成网络带来的价值(以及灵活性)。将计算机组成网络就是将机器通过通信基础设施(例如以太网电缆或无线网络)连接在一起。Tanenbaum 在谈及从集中式的大型机到分布式的而功能依然强大的小型机的计算机发展过程时,将这称为“计算机和通信的融合(merging of computers and communications)”[76]。

分布式系统和软件利用计算机网络:网络提供部分基础结构,供分布式系统在其上执行(在第 4 章中,我们将看到,操作系统提供部分的其他关键基础设施)。此外,计算机网络对用户来说,通常就是一组能够相互通信的计算设备。分布式系统通常提供一致的接口,使得用户可以将一组计算设备看作单个的统一实体。

虽然本书着重于介绍分布式系统,但是,阐明计算机网络的某些方面也是很重要的,特别是网络基础设施是如何用来实现和支持强大可靠的分布式系统的。

## 1.3 分布式处理

分布式系统有 3 项基本特点,这 3 项特点也是本书余下内容的关键章节。

- **并发**: 在实际的分布式系统中,多个线程同时运行。比如,当你的网络浏览器正在从某个网络服务器下载内容时,可能也在确定如何将该内容显示在你的计算机屏幕上。同时,网络浏览器也正在从磁盘加载内容,以将它们通过网络发送给你。并发的复杂性的关键在于处理对资源的共享访问。例如,如果你和你的朋友同时要访问网上的某个文件,就需要有适当的机制来确保解决访问该文件时发生的冲突,譬如,你和你的朋友都同时要编辑文件。关于用来确保共享资源正确一致的



基本原理和实用方法,第 2 章中有详细的讨论。

- **在时间上同步:** 在实际分布式程序中,组件和线程通常必须同步并协调它们的活动。例如,前面简略提到的 SABRE 系统,它允许旅行社为顾客查找和预订机票,显然,该系统必须有时间观念,使得如果伦敦的一家旅行社预定了格林尼治标准时间 09:00 的一张机票,虽然在巴黎该张机票的中部欧洲时间为 10:05,但巴黎的某家旅行社就不能预订这张机票。这个简单的例子说明活动的同步和协调通常依赖于活动发生的时间。当然,时钟不能做到完全精确,就像我们的手表一样,网络上同步的多个时钟也一样存在偏差。在分布式系统中,处理时间问题非常关键,后面的章节对此进行了简要的讨论。
- **失效:** 在依赖于多个独立组件的分布式系统中,不幸的事实是这些组件会出现失效(例如,软件中的错误、硬件故障或外界的错误输入都可以导致出现失效),而且,整个系统对这些组件如何以及何时会失效几乎没有控制。因此,分布式系统设计者的责任是预见组件会出现失效,并纠正整个系统中的潜在失效。在许多领域中,拥有一个能够处理组件失效的分布式系统(也就是拥有一个可靠的分布式系统)非常重要,例如前面讨论的银行提款机系统。有关那些可以用来改善整体可靠性的技术,包括分布式系统的精确模型和冗余,本书会一直贯穿讨论。与失效有着特殊相关性的是安全,这部分内容将在第 7 章中详细讨论。

## 1.4 应用领域

我们在前面提到了几个分布式系统的例子。如今,分布式系统随处可见,而且通常使用在银行业(例如进行交易和清算账目)、电子商务、计算机游戏、嵌入式系统、生物信息(例如生物模型的分析)、移动设备(比如个人数字助理、手机和数码相机)和许多其他领域。

## 1.5 模 型

虽然当今使用的分布式系统各式各样(而且我们已经讨论了一些重要的例子),但我们可以将这些系统归于几类模型。

- **客户-服务器模型**(也称为双层模型)最普遍,通常来说,使用的中央服务器和客户机可以是一个或多个,这些服务器维护共享信息,比如某个数据库,客户机则根据需要访问共享信息。实际上,服务器提供统一的公开界面,通过该界面提供服务,而客户机是根据需要使用这些服务。最为广泛的客户-服务器模型的应用是文件传输协议(File Transfer Protocol,FTP),它用于在两台计算设备之间传输文件。