# 分布式操作系统： 概念与实践

# DISTRIBUTED OPERATING SYSTEMS:

## CONCEPTS AND PRACTICE

### Doreen L. Galli

Pearson
Education

英文版

国外著名高等院校信息科学与技术优秀教材

# 分布式操作系统：概念与实践

# (英文版)

## Distributed Operating Systems:

## Concepts and Practice

Doreen L. Galli

人民邮电出版社 PRENTICE HALL Pearson Education 出版集团

## 版 权 声 明

# 内 容 提 要

　　本书从理论和实践两个方面，阐述了分布计算的主要概念、理论和各种成功实例，主要内容包括：内核、进程间通信、存储管理、基于对象的操作系统、分布式文件系统、事务管理与协调模型、分布进程管理、分布同步、分布计算中的安全性等。选取的实例包括：Amoeba、Clouds、Chorus、CORBA、DCOM、NFS、LDAP、X.500、NFS、RSA、Kerberos 及 Windows 2000 等。

　　本书适于用作计算机科学与技术系本科高年级及研究生分布计算、分布式操作系统等课程的教材或主要参考书，也适用于在相关领域工作的科技工作者。

# 出版说明

  2001 年，教育部印发了《关于"十五"期间普通高等教育教材建设与改革的意见》。该文件明确指出，"九五"期间原国家教委在"抓好重点教材，全面提高质量"方针指导下，调动了各方面的积极性，产生了一大批具有改革特色的新教材。然而随着科学技术的飞速发展，目前高校教材建设工作仍滞后于教学改革的实践，一些教材内容陈旧，不能满足按新的专业目录修订的教学计划和课程设置的需要。为此该文件明确强调，要加强国外教材的引进工作。当前，引进的重点是信息科学与技术和生物科学与技术两大学科的教材。要根据专业（课程）建设的需要，通过深入调查、专家论证，引进国外优秀教材。要注意引进教材的系统配套，加强对引进教材的宣传，促进引进教材的使用和推广。

  邓小平同志早在 1977 年就明确指出："要引进外国教材，吸收外国教材中有益的东西。"随着我国加入 WTO，信息产业的国际竞争将日趋激烈，我们必须尽快培养出大批具有国际竞争能力的高水平信息技术人才。教材是一个很关键的问题，国外的一些优秀教材不但内容新，而且还提供了很多新的研究方法和思考方式。引进国外原版教材，可以促进我国教学水平的提高，提高学生的英语水平和学习能力，保证我们培养出的学生具有国际水准。

  为了贯彻中央"科教兴国"的方针，配合国内高等教育教材建设的需要，人民邮电出版社约请有关专家反复论证，与国外知名的教材出版公司合作，陆续引进一些信息科学与技术优秀教材。第一批教材针对计算机专业的主干核心课程，是国外著名高等院校所采用的教材，教材的作者都是在相关领域享有盛名的专家教授。这些教材内容新，反映了计算机科学技术的最新发展，对全面提高我国信息科学与技术的教学水平必将起到巨大的推动作用。

  出版国外著名高等院校信息科学与技术优秀教材的工作将是一个长期的、坚持不懈的过程，我社网站（www.ptpress.com.cn）上介绍了我们陆续推出的图书的详细情况，敬请关注。希望广大教师和学生将使用中的意见和建议及时反馈给我们，我们将根据您的反馈不断改进我们的工作，推出更多更好的引进版信息科学与技术教材。
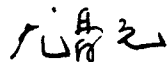
<div align="right">

人民邮电出版社

2001 年 12 月

</div>

# 序 言

分布式系统是基于分布计算技术和网络技术、在网络基础设施之上建立起来的。随着网络技术的出现和发展，从 20 世纪 70 年代中期开始，计算机科学技术界对分布计算基础理论的研究给予了很多关注，在分布进程管理、分布进程通信、数据一致性模型、分布同步、事务处理与协调模型等方面的研究工作都取得了显著成果。20 世纪 80 年代中期，特别是 90 年代初期，网络基础设施在全球范围内的各种企事业单位中纷纷建成，它们的通信带宽和质量得到不断提高，同时，信息获取、处理和存储技术的快速发展，对分布式系统的需求广泛而迫切。在这种背景下，对分布计算技术和系统的研究开发工作的重点也从分布式操作系统转移到了更宽广的范围，出现了 CORBA、DCOM、.net、EJB/J2EE 应用服务器等典型的分布计算标准/规范和平台，近年来对 Web Service 和网格计算（Grid Computing）等方面的研发工作更是方兴未艾。与这样的技术和应用发展相适应，高等院校的计算机及信息处理专业纷纷开设分布操作系统、分布式系统或分布计算技术方面的课程，一些国内外知名的出版社，包括 Prentice Hall 则邀请相关专家编著和出版了若干此方面的教材和参考书，本书就是其中影响比较大的一本。

本书从理论和实践两个方面阐述了分布计算的主要概念、理论和各种典型实例。在理论方面包括的主要内容是：内核、进程间通信、分布存储管理、基于对象的操作系统、分布式文件系统、分布事务管理和协调模型、分布进程管理、分布同步、分布计算中的安全性等。它们涉及到了分布式系统设计、实现时应当考虑的各主要方面。选取的主要实例是：Amoeba、Clouds、Chorus、Windows 2000、CORBA、DCOM、NFS、LDAP、X.500、RSA、Kerberos 等，它们是当前分布式操作系统、分布式系统的重要典型，涵盖了分布计算主要技术的应用。

本书适于用作计算机科学与技术系及信息类有关系科，本科高年级及研究生分布计算、分布式系统或分布式操作系统课程的教材或主要参考书，在相关领域中工作的科技工作者也会在阅读本书时得到收益。

上海交通大学计算机系

2002 年 5 月

*To my loving husband, Marc and my little ones, Marc Jr. and Steven*

# Acknowledgments

$A$s Sir Isaac Newton once said, "If I have seen further, it is by standing on the shoulders of Giants." To that I would like to add, "and I have always tried to find the tallest giants." To those giants I say thank you. A book such of this could not exist without the help of others, and I would like to thank all of those who have assisted me, including those giants whose contributions have advanced and continue to advance distributed computing as well as the computer science field as a whole. Their contributions are far reaching and have affected this book even though I may have never met or worked with them.

I would like to thank those giants who let me stand upon their shoulders possibly never knowing how much they had given me or that it could lead to this project. I would like to thank Mark Fishman, Ed Gallizzi, and the rest of the faculty at Eckerd College, who first sparked my interest in and love for computers. I would like to thank all of the wonderful faculty and students at the University of Waterloo, including my graduate supervisor, Charlie Colbourn. You taught me more than I ever realized from consistent formatting to advanced research techniques. You were the best supervisor one could have ever hoped to have worked with in graduate school.

I would like to thank all of those involved in the IBM CORDS (Consortium of Research in Distributed Systems) project, one of my greatest source of tall shoulders, including Gopi Attaluri, Michael Bauer, Dexter Bradshaw, Neil Coburn, Mariano Consens, Pat Finnigan, Masum Hasan, James Hong, Paul Larson, Kelly Lyons, Pat Martin, Gerry Neufeld, Jon

1

dents who helped me class test the material. In particular, I would like to extend a special thank you to my students who won the Submit the Most Suggestions - Find the Most Typos contest, John Fritz (the winner with an exceptional talent for proofreading), Moji Mahmoudi, and Haoyuan Chen. The quality of the material is greatly increased due to their due diligence, and I am confident that all readers will benefit from their unique insight as someone learning from this book. I sincerely welcome suggestions and comments from all readers as it is only with readers' insight that continual improvement for future revisions is possible. Readers can contact me through my Prentice Hall Web site, www.prenhall.com/galli.

I would like to thank my family for supporting me during this quest. Thank you Jr. and Steven, who gave up "Mommy time" and cheered me on to the finish line. Most of all, I would like to thank my husband and best friend, Marc. Thank you for your inspiration and the many hours of your invaluable assistance, including the times you let me read the chapters out loud to you in order to identify and correct any awkward wording as well as your assistance in smoothing out the rough spots. You made it fun; and yes, I'm sure it took a lot of patience, too. It is great to have a family that comes first yet appreciates and supports me in endeavors such as this. I'm very fortunate to have been so blessed with each of you. May you always be inspired and have the support necessary to chase and catch all of your dreams as well.

<div align="right">

Doreen L. Galli, Ph.D.
August, 1999

</div>

# Preface

$T$his book examines concepts and practice in distributed computing. It is designed to be useful not only for students but for practitioners and corporate training as well. Over the past decade, computer systems have become increasingly more advanced. Most computers are connected to some type of network on a regular basis. The installation of LANs at smaller businesses is even becoming commonplace. LANs are also being installed in custom homes at an ever-increasing rate. Software technology must keep up and so must our future and current practitioners! At the current pace, it is only a matter time before a working knowledge of distributed systems is mandatory for all computer scientists, because this technology pertains to a majority of all computers and their applications.

## INTENDED AUDIENCE

While the study of standard operating system concepts is extremely important for computer science undergraduates, there is a significant and ever-increasing demand to extend this knowledge in the graduate and fourth-year undergraduate curriculum as well as for the practitioner out in industry. Therefore, there is a great need to study distributed operating

1

systems concepts as well as practical solutions and approaches. This book is intended to meet this need for both students and practitioners.

## OBJECTIVE

The objective of this book is to describe in detail each major aspect of distributed operating systems from a conceptual and practical viewpoint. Thus, it includes relevant examples of real operating systems to reinforce the concepts and to illustrate the decisions that must be made by distributed system designers. Operating systems such as Amoeba, Clouds and Chorus (the base technology for JavaOS) are utilized as examples throughout the book. In addition, the case study on Windows 2000™ provides an example of a real commercial solution. Technologies such as CORBA, DCOM, NFS, LDAP, X.500, Kerberos, RSA, DES, SSH, and NTP are also included to demonstrate real-life solutions to various aspects of distributed computing. In addition, a simple client/server application is included in the appendix that demonstrate key distributed computing programming concepts such as the use of INET sockets, pthreads, and synchronization via mutex operations.

In summary, this book focuses on the concepts, theory and practice in distributed systems. It is designed to be useful for practitioners, fourth year undergraduate as well as graduate level students and assumes that the reader has taken a basic operating system course. It is hoped that this book will prove to be invaluable not only for those already active in industry who wish to update and enhance one's knowledge base but also for future reference for those who have used it as a course text.

## ORGANIZATION AND PEDAGOGICAL FEATURES

This book is divided into two parts. The first part, Chapter 1-6, presents the base foundation for distributed computing. The second part, Chapter 7-11, expands on these topics and delves more heavily into advanced distributed operating system topics. The pedagogical features included in this book are the following.

1.  Detail Boxes to further enhance understanding. These boxes contain informa tion such as complex algorithms and more in depth examples.

2.  More than 150 figures and tables to help illustrate concepts.

3.  A case study of Windows 2000™ to demonstrate a real life commercial solutions.

4.  Project oriented exercises (those with italicized numbers) to provide "hands on" experience.

5.  Exercises that build upon concepts covered in earlier chapters.

6. Reference pointers to relevant sources including:

    A. *overview* sources for further in-depth study,

    B. *research* papers, and

    C. *'core'* web & *ftp* sites.

7. A simplified distributed application program to demonstrate key distributing programming concepts.

8. Comprehensive glossary of terms (**boldfaced** words appear in the glossary) to provide a centralized location for key definitions.

9. Complete list of acronyms to aid readability and provide a centralized location for easy reference.

10. Chapter summaries.

11. Comprehensive index, primary references in **bold**.

12. Book website located at www.prenhall.com/galli.

## SUGGESTIONS FOR INSTRUCTORS

This book is designed to provide maximum flexibility to instructors and has pedagogical features inherent within the text to allow you to customize the coverage to best meet the needs of your class and your institution's mission statement. In preparing this book, the only assumption made is that a basic introductory to operating systems course has been taken by the reader. Select topics that may be included in an introductory operating system course but are sometimes omitted, covered lightly, often not grasped or may have been forgotten but nonetheless are key to distributed operating systems, are included where appropriate. This material need not be presented in the classroom but is included in the book so that you can be assured that the students have the basis necessary for the more advanced distributed topics. Below are suggestions on how this book may be used for those requiring additional practical emphasis as well as for those desiring additional research emphasis. A graduate course desiring to add both types of emphasis may wish to use suggestions from both categories. Additional information may be available at the author's Prentice Hall website, www.prenhall.com/galli.

## Adding Practical Emphasis

The following are a few suggestions for adding practical emphasis to a course utilizing this text.

1.  Have the students, either individually or as a group complete one or more
    of the 'Project Exercises', those indicated by an italicized exercise number
    at the end of relevant chapters. Additional practical experience may be
    achieved if their design and implementation is orally presented to the class.

2.  Cover all Detail Boxes related to real-life implementations.

3.  Spend class time covering the Windows 2000™ Case study.

4.  Create an individual or group project working with the distributed features
    of Windows 2000™.

5.  Have the students expand or change the Surgical Scheduling Program. This
    may be as simple as changing the type of interprocess communication em-
    ployed or as complex as creating another program utilizing the same dis-
    tributed concepts.

## Adding Research Emphasis

The following are a few suggestions for adding a research emphasis to a course utilizing this
book.

1.  Have the students, either individually or as a group, prepare a paper on a
    topic relevant to distributed operating systems. Reference papers cited at
    the end of each chapter should serve as good starting points. These projects
    may include an oral presentation.

2.  Present lecture material from the relevant RFCs or research papers cited at
    the end of each chapter that are available on the web and include it the list
    of required reading for the students.

3.  Have the students seek the relevant RFCs or research papers cited at the
    end of each chapter that are available on the web and prepare a summary.

4.  Select a subset of the reference papers cited at the end of each chapter and
    create a spiral bound accompaniment to be used in conjunction throughout
    the course with the book. A large number of bookstores at research institu-
    tions have the ability to perform the copyright clearing necessary for this
    purpose.

# Contents