# Distributed Computing

## Principles, Algorithms, and Systems

Distributed computing deals with all forms of computing, information access, and information exchange across multiple processing platforms connected by computer networks. Design of distributed computing systems is a complex task. It requires a solid understanding of the design issues and an in-depth understanding of the theoretical and practical aspects of their solutions. This comprehensive textbook covers the fundamental principles and models underlying the theory, algorithms, and systems aspects of distributed computing.

Broad and detailed coverage of the theory is balanced with practical systems-related problems such as mutual exclusion, deadlock detection, authentication, and failure recovery. Algorithms are carefully selected, lucidly presented, and described without complex proofs. Simple explanations and illustrations are used to elucidate the algorithms. Emerging topics of significant impact, such as peer-to-peer networks and network security, are also covered.

With state-of-the-art algorithms, numerous illustrations, examples, and homework problems, this textbook is invaluable for advanced undergraduate and graduate students of electrical and computer engineering and computer science. Practitioners in data networking and sensor networks will also find this a valuable resource.

**Ajay D. Kshemkalyani** is a Professor in the Department of Computer Science, at the University of Illinois at Chicago. He was awarded his Ph.D. in Computer and Information Science in 1991 from The Ohio State University. Before moving to academia, he spent several years working on computer networks at IBM Research Triangle Park. In 1999, he received the National Science Foundation's CAREER Award. He is a Senior Member of the IEEE, and his principal areas of research include distributed computing, algorithms, computer networks, and concurrent systems. He currently serves on the editorial board of *Computer Networks*.

**Mukesh Singhal** is Full Professor and Gartner Group Endowed Chair in Network Engineering in the Department of Computer Science at the University of Kentucky. He was awarded his Ph.D. in Computer Science in 1986 from the University of Maryland, College Park. In 2003, he received the IEEE

Technical Achievement Award, and currently serves on the editorial boards for the *IEEE Transactions on Parallel and Distributed Systems* and the *IEEE Transactions on Computers*. He is a Fellow of the IEEE, and his principal areas of research include distributed systems, computer networks, wireless and mobile computing systems, performance evaluation, and computer security.

# Distributed Computing

## Principles, Algorithms, and Systems

Ajay D. Kshemkalyani

University of Illinois at Chicago, Chicago

and

Mukesh Singhal

University of Kentucky, Lexington

CAMBRIDGE
UNIVERSITY PRESS

To my father Shri Digambar and
my mother Shrimati Vimala.

**Ajay D. Kshemkalyani**

To my mother Chandra Prabha Singhal,
my father Brij Mohan Singhal, and my
daughters Meenakshi, Malvika,
and Priyanka.

**Mukesh Singhal**

# Contents

# Preface

## Background

The field of distributed computing covers all aspects of computing and information access across multiple processing elements connected by any form of communication network, whether local or wide-area in the coverage. Since the advent of the Internet in the 1970s, there has been a steady growth of new applications requiring distributed processing. This has been enabled by advances in networking and hardware technology, the falling cost of hardware, and greater end-user awareness. These factors have contributed to making distributed computing a cost-effective, high-performance, and fault-tolerant reality. Around the turn of the millenium, there was an explosive growth in the expansion and efficiency of the Internet, which was matched by increased access to networked resources through the World Wide Web, all across the world. Coupled with an equally dramatic growth in the wireless and mobile networking areas, and the plummeting prices of bandwidth and storage devices, we are witnessing a rapid spurt in distributed applications and an accompanying interest in the field of distributed computing in universities, governments organizations, and private institutions.

Advances in hardware technology have suddenly made sensor networking a reality, and embedded and sensor networks are rapidly becoming an integral part of everyone's life – from the home network with the interconnected gadgets to the automobile communicating by GPS (global positioning system), to the fully networked office with RFID monitoring. In the emerging global village, distributed computing will be the centerpiece of all computing and information access sub-disciplines within computer science. Clearly, this is a very important field. Moreover, this evolving field is characterized by a diverse range of challenges for which the solutions need to have foundations on solid principles.

The field of distributed computing is very important, and there is a huge demand for a good comprehensive book. This book comprehensively covers all important topics in great depth, combining this with a clarity of explanation

and ease of understanding. The book will be particularly valuable to the academic community and the computer industry at large. Writing such a comprehensive book has been a Herculean task and there is a deep sense of satisfaction in knowing that we were able complete it and perform this service to the community.

## Description, approach, and features

The book will focus on the fundamental principles and models underlying all aspects of distributed computing. It will address the principles underlying the theory, algorithms, and systems aspects of distributed computing. The manner of presentation of the algorithms is very clear, explaining the main ideas and the intuition with figures and simple explanations rather than getting entangled in intimidating notations and lengthy and hard-to-follow rigorous proofs of the algorithms. The selection of chapter themes is broad and comprehensive, and the book covers all important topics in depth. The selection of algorithms within each chapter has been done carefully to elucidate new and important techniques of algorithm design. Although the book focuses on foundational aspects and algorithms for distributed computing, it thoroughly addresses all practical systems-like problems (e.g., mutual exclusion, deadlock detection, termination detection, failure recovery, authentication, global state and time, etc.) by presenting the theory behind and algorithms for such problems. The book is written keeping in mind the impact of emerging topics such as *peer-to-peer computing* and *network security* on the foundational aspects of distributed computing.

Each chapter contains figures, examples, exercises, a summary, and references.

## Readership

This book is aimed as a textbook for the following:

- Graduate students and Senior level undergraduate students in computer science and computer engineering.
- Graduate students in electrical engineering and mathematics. As wireless networks, peer-to-peer networks, and mobile computing continue to grow in importance, an increasing number of students from electrical engineering departments will also find this book necessary.
- Practitioners, systems designers/programmers, and consultants in industry and research laboratories will find the book a very useful reference because it contains state-of-the-art algorithms and principles to address various design issues in distributed systems, as well as the latest references.

Hard and soft prerequisites for the use of this book include the following:

- An undergraduate course in algorithms is required.
- Undergraduate courses in operating systems and computer networks would be useful.
- A reasonable familiarity with programming.

We have aimed for a very comprehensive book that will act as a single source for distributed computing models and algorithms. The book has both depth and breadth of coverage of topics, and is characterized by clear and easy explanations. None of the existing textbooks on distributed computing provides all of these features.

## Acknowledgements

This book grew from the notes used in the graduate courses on distributed computing at the Ohio State University, the University of Illinois at Chicago, and at the University of Kentucky. We would like to thank the graduate students at these schools for their contributions to the book in many ways.

The book is based on the published research results of numerous researchers in the field. We have made all efforts to present the material in our own words and have given credit to the original sources of information. We would like to thank all the researchers whose work has been reported in this book. Finally, we would like to thank the staff of Cambridge University Press for providing us with excellent support in the publication of this book.

## Access to resources

The following websites will be maintained for the book. Any errors and comments should be sent to ajayk@cs.uic.edu or singhal@cs.uky.edu. Further information about the book can be obtained from the authors' web pages:

- www.cs.uic.edu/~ajayk/DCS-Book
- www.cs.uky.edu/~singhal/DCS-Book.