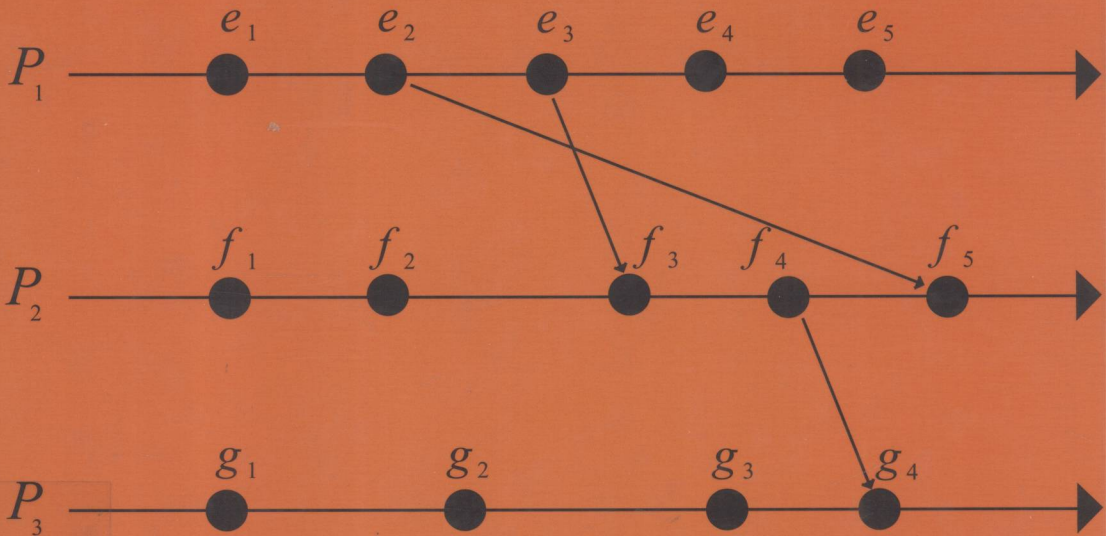# ELEMENTS OF DISTRIBUTED COMPUTING
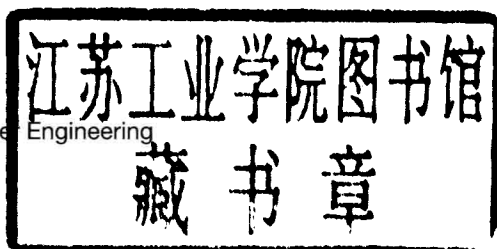


VIJAY K. GARG

# ELEMENTS OF DISTRIBUTED COMPUTING

**VIJAY K. GARG**

Department of Electrical and Computer Engineering
University of Texas at Austin

**WILEY-INTERSCIENCE**

# ELEMENTS OF
# DISTRIBUTED COMPUTING

*To*
*my parents (Shri Saran Lal and Shrimati Laxmi Devi),*
*my wife (Meenakshi),*
*and my children (Sonika and Rohan)*

# Foreword

Distributed computing arises as soon as one has to solve a problem in terms of processes that individually have only a partial knowledge of the several parameters associated with the problem. Thus, distributed computing appears everywhere (even in everyday life!) and is consequently at the heart of lots of applications. The difficulty in designing a distributed algorithm comes from the fact that each of the processes cooperating in the achievement of the common goal cannot have instantaneous knowledge of the current state of the other processes; it can only know their past states. Moreover, because of the possibility of process or communication failures, this knowledge can become very difficult or even impossible to obtain.

Whereas parallel computing is mainly concerned with efficiency, distributed computing addresses uncertainty created by the multiplicity of control flows, the absence of shared memory and global time, and the occurrence of failures. Although it is true that distributed algorithms are often made up of only a few lines, their behaviors can be difficult to understand and their properties hard to state and prove. Hence, distributed computing is not only a fundamental topic but also a challenging topic where simplicity, elegance, and beauty are first-class criteria.

In this book, Vijay K. Garg provides a comprehensive and timely treatment of all these challenges. In a clear and consistent way, he addresses the fundamental problems, issues, and solutions that define the core of distributed computing. The book is both versatile and complete. While each chapter (designed to correspond to one lecture) addresses a basic issue of distributed computing, the examination of chapters addressing several facets of the same topic provides the reader with a deep insight into the subject and a good understanding of the subtleties one has to cope with when solving distributed computing problems. Moreover, the exercises accompanying each chapter can help

the reader test her acquaintance with the subject concerned.

The book is addressed to all those who are concerned with algorithms, distributed computing, and distributed systems. It will provide software engineers with provably correct distributed protocols, researchers with a large collection of up-to-date results, teachers with a solid foundation for a course on distributed computing, and students with an introduction to distributed algorithms for distributed computing.

Professor Michel Raynal
IRISA,
Université de Rennes 1
Campus de Beaulieu, Rennes
France

18 February 2002

# Preface

The existing books on distributed computing have mainly been developed as reference books on distributed algorithms. They are more useful to researchers than to students. The present book was designed keeping students in mind. Simplicity has been chosen at the cost of completeness whenever that choice was required.

The topics have been selected very carefully. No book, short of a thick encyclopedia, can cover all topics in distributed systems. I have chosen topics that have practical usefulness as well as aesthetic beauty. Because usefulness and beauty are only in the eyes of the beholder, the selection was necessarily subjective. In my defense, I will only give a brief rationale behind the choice of topics. My rule for selecting topics was simple: I have chosen algorithms over arguments for lower bounds and general mechanisms over special case optimized algorithms—the former choice for usefulness and the latter for beauty. For the sake of usefulness or beauty, I violated my general rule whenever necessary; for example, the result on impossibility of consensus in presence of failures (a negative result) and the algorithm for termination detection (a special case of stable property detection) have been included.

Most proofs in the book are rigorous but not completely formal. I have avoided proofs by contradiction. To quote Royden from his book on Real Analysis: "All students are enjoined in the strongest possible terms to eschew proof by contradiction! There are two reasons for this prohibition: First, such proofs are very often fallacious, the contradiction on the final page arising from an erroneous deduction on an earlier page,... Second, even when correct, such a proof gives little insight into connection between A and B ( antecedent and consequent)."

The book is organized as a large number of small chapters rather than a small number of large chapters. Besides the obvious advantage of providing more occasions for satisfaction gained on finishing a chapter, this organization provides more flexibility to the instructor. Assuming

that there is little dependency across chapters, as is the case for the present book, the instructor can substitute his or her own handouts and papers for some chapters.

The book is designed for a graduate level course on distributed systems. Several examples and exercise problems are included in each chapter to facilitate classroom teaching. My treatment is based on the graduate courses given at The University of Texas at Austin.

The list of known errors and the supplementary material for the book will be maintained on my home page:

<div align="center">

`http://www.ece.utexas.edu/~garg`

</div>

I would like to thank the following people for working with me on various projects discussed in this book. These people include Craig Chase (weak predicates), Om Damani (message logging), Eddy Fromentin (predicate detection), Joydeep Ghosh (global computation), Richard Kilgore (channel predicates), Roger Mitchell (channel predicates), Neeraj Mittal (predicate detection and control, slicing, self-stabilization, distributed shared memory), Venkat Murty (synchronous ordering), Michel Raynal (control flow properties, distributed shared memory), Alper Sen (slicing), Chakarat Skawratonand (vector clocks), Ashis Tarafdar (message logging, predicate control) Alexander Tomlinson (global time, mutual exclusion, relational predicates, control flow properties) and Brian Waldecker (weak and strong predicates). I owe special thanks to Alper Sen for helping me with figures and his meticulous review of an early draft of this book. Ajay Kshemkalyani and Michel Raynal provided valuable feedback on the contents of the book.

I thank the Department of Electrical and Computer Engineering at The University of Texas at Austin, where I was given the opportunity to develop and teach courses on distributed systems. Students in my course gave me very useful feedback. I was supported in part by many grants from the National Science Foundation over the last twelve years. Many of the results reported in this book would not have been discovered by me and my research group without that support. I also thank John Wiley & Sons, Inc. for supporting the project.

Finally, I thank my wife Meenakshi, whose love, support, and understanding kept me inspired in accomplishing this goal.

# Contents