# Distributed Operating Systems

Andrew S. Tanenbaum

# Distributed Operating Systems

## Andrew S. Tanenbaum

*Vrije Universiteit*
*Amsterdam, The Netherlands*

**PEARSON**
Education

# Contents

**10**

**5  DISTRIBUTED FILE SYSTEMS**

# 9  CASE STUDY 3: CHORUS

# 10  CASE STUDY 22: DCE

# Preface

With the publication of *Distributed Operating Systems* I have now completed my trilogy on operating systems. The three volumes of this trilogy are:

- *Operating Systems: Design and Implementation*
- *Distributed Operating Systems*
- *Modern Operating Systems*

The three volumes are not completely independent, however. For schools having a two-course sequence in operating systems (or an undergraduate course plus a graduate course), one possible choice is to use *Operating Systems: Design and Implementation* in the first course and *Distributed Operating Systems* in the second one.

The former book treats the standard principles of single-processor systems, including processes, synchronization, I/O, deadlocks, memory management, file systems, security, and so on. It also illustrates these principles in great detail through the use of MINIX, a UNIX-clone whose source listing is given in an appendix. MINIX is available on diskette from Prentice Hall for the IBM PC (8088 and up), Atari, Amiga, Macintosh, and SPARC processors.

The latter book (this one), covers distributed operating systems in detail, including communication, synchronization, processes, file systems, and memory management, but this time in the context of distributed systems. Four examples of distributed systems are given in great detail: Amoeba, Mach, Chorus, and DCE. Amoeba is available for free to universities for educational use. It runs