JEAN BACON & TIM HARRIS

# OPERATING SYSTEMS
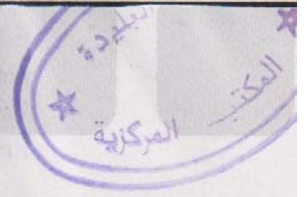
## CONCURRENT AND DISTRIBUTED SOFTWARE DESIGN

ADDISON-WESLEY

# Operating Systems

## Concurrent and Distributed Software Design

# Jean Bacon and Tim Harris

University of Cambridge

▲▲▲ **ADDISON-WESLEY**

*An imprint of* **Pearson Education**

Harlow, England • London • New York • Boston • San Francisco • Toronto • Sydney • Singapore • Hong Kong
Tokyo • Seoul • Taipei • New Delhi • Cape Town • Madrid • Mexico City • Amsterdam • Munich • Paris • Milan

# Brief contents

# Contents

## PART III TRANSACTIONS 485

## 17 Composite operations 487