# Scalable Parallel Computing

## Technology, Architecture, Programming

# 可扩展
# 并行计算

## 技术、结构与编程

（英文版）

## Kai Hwang　Zhiwei Xu 著

机械工业出版社
China Machine Press

**Mc Graw Hill**

**WCB**
*McGraw-Hill*

# 可扩展并行计算

## 技术、结构与编程

（英文版）

# Scalable Parallel Computing

## Technology, Architecture, Programming

Kai Hwang　Zhiwei Xu　著

凡购本书，如有缺页、倒页、脱页，由本社发行部调换

# Trademark Notice

ANSI Standards X3T11, X3T9, and X3H5 are trademarks of the American National Standards Institute.

ALLICACHE and KSR-1 are trademarks of Kendall Square Research.

X Window, Alewife, J-Machine, and *T are trademarks of Massachusetts Institute of Technology.

CM-2, CM-5, C*, *Lisp, CM Fortran, and CMOST are trademarks of Thinking Machines Corporation.

HP/Convex Exemplar is a trademarks of Hewlett-Packard Company and Convex Computer Systems.

Cray 1, Cray 2, Y-MP, C90, CRAFT, T3D, T3E, UNICOS are trademarks of Cray Research, Inc.

Dash and SPLASH are trademarks of Stanford University.

Ethernet is a trademark of Xerox Corporation.

Gigabit Ethernet is trademark of the Gigabit Ethernet Alliance.

SCI, POSIX Threads, and Futurebus+ are trademarks of the Institute of Electrical and Electronic Engineering, Inc.

Illiac IV, Cedar, and Parafrase are trademarks of University of Illinois.

Intel 80x86, i960, i860, SHV, Paragon, Intel Hypercube, Pentium, Pentium Pro, and Pentium-II are trademarks of Intel Corporation.

MIPS R2000, R4000, and R10000 are trademarks of SGI/MIPS Technology, Inc.

Mach/OS, C.mmp, and Cm* are trademarks of Carnegie-Mellon University.

NYU Ultracomputer is a trademark of New York University.

SX/4 is a trademark of NEC Information Systems.

Symmetry and NUMA-Q are trademarks of Sequent Computers.

Gigaplane, Gigaplane-XB, S-Bus, TSO, PSO, NFS, SPARC, Ultra Enterprise, SPARCcluster, SPARCstation, and Solaris MC are trademarks of Sun Microsystems, Inc.

Tera is a trademark of Tera Computer Systems.

Cosmic Cube and Mosaic C are trademarks of California Institute of Technology.

UNIX 4.3 BSD, NOW, GLUnix, IRAM, xFS, and 4.4BSD are trademarks of the University of California at Berkeley.

DEC, DECsafe, VAX, VAXCluster, VMS, Digital Unix, TruCluster, Digital NT Cluster, Alpha 21164, and Alpha AXP are trademarks of Digital Equipment Corporation.

Apollo Domain, HP, Series 9000, PA-RISC, and UX are trademarks of the Hewlett-Packard Company.

IBM PC, AIX, ESSL, 801, RP3, EUI, HACMP, IBM, LoadLeveller, NetView, POWER, PowerPC, RS/6000, S/370, S/390, SP2, SP, and Sysplex are trademarks of International Business Machines Corporation.

Lifekeeper, NCR and Teradata are trademarks of National Cash Register.

M680x0 is trademark of Motorola Corporation.

OSF/1, DCE, and DFS are trademarks of the Open Software Foundation, Inc.

Oracle Parallel Server is a trademark of Oracle Corporation.

Express is a trademark of Parasoft Corporation.

Load Sharing Facility is a trademark of Platform Computing Corporation.

Origin 2000, Power C, POWER CHALLENGEarray, and Cellular IRIX are trademarks of Silicon Graphics Computer Systems.

Sybase is a trademark of Sybase, Inc.

UNIX and SVR5 are trademarks of X/Open Company, Ltd.

Microsoft, Wolfpack, and Windows NT are trademarks of Microsoft Corporation.

OpenMP is a trademark of OpenMP Standards Evaluation Board.

Myrinet is a trademark of Myricom, Inc.

# About the Authors

**Kai Hwang** presently holds a Chair Professor of Computer Engineering at the University of Hong Kong (HKU), while taking a research leave from the University of Southern California (USC). He has been engaged in higher education and computer research for 26 years, after earning the Ph.D. in Electrical Engineering and Computer Science from the University of California at Berkeley. His work on this book started at the USC and was mostly completed at the HKU.

An IEEE Fellow, he has published extensively in the areas of computer architecture, digital arithmetic, parallel processing, and distributed computing. He is the founding Editor-in-Chief of the *Journal of Parallel and Distributed Computing*. He has chaired the international conferences: ICPP86, ARITH-7, IPPS96, ICAPP 96, and HPCA-4 in 1998. He has received several achievement awards for outstanding research and academic contributions to the field of parallel computing.

He has lectured worldwide and performed consulting and advisory work for US National Academy of Sciences, MIT Lincoln Laboratory, IBM Fishkill, TriTech in Singapore, Fujitsu and ETL in Japan, GMD in Germany, CERN School of Computing, and Academia Sinica in China. Presently, he leads a research group at HKU in developing an ATM-based multicomputer cluster for high-performance computing and distributed multimedia, Intranet, and Internet applications.

**Zhiwei Xu** is a Professor and Chief Architect at the National Center for Intelligent Computing Systems (NCIC), Chinese Academy of Sciences, Beijing, China. He is also a Honorary Research Fellow of HKU. He received a Ph.D. in Computer Engineering from the University of Southern California. He participated in the STAP/MPP benchmark projects led by Dr. Hwang at USC and HKU during the past 5 years.

He has taught at the Rutgers University and New York Polytechnic University. He has published in the areas of parallel languages, pipelined vector processing, and benchmark evaluation of massively parallel processors. Presently, he leads a design group at NCIC in building a series of cluster-based superservers in China. His current research interest lies mainly in network-based cluster computing and the software environments for parallel programming.

# Dedication

To my family, for their love and support -- k.h.

To my high-school teacher,
Ms. Zhang Dahua, who taught me the
serenity of science at a tumultuous time -- z.x.

# Foreword

## Michael J. Flynn, Stanford University

Parallel processors are the future of computer design. The replicability of technology has proven an ever-increasing impetus towards parallel processor implementations. The demand for increasing performance in important applications is another persuasive argument. However, building efficient, scalable parallel processors has proven to be a very difficult problem.

For decades, it seemed that large-scale ($n > 100$ processors) parallel processors could be realized efficiently with just a little bit more effort and research. This has not proven so. It has been difficult to find the prerequisite parallelism in programs, and when this parallelism was found, it did not translate into a speedup of program execution proportional to the number of processing elements involved. This has been especially true when $n$ exceeds 100.

One important problem in the realization of efficient parallel processors is the computational model that underlies the program model and language used to represent applications. This model was developed for and efficiently suits the single processor. Because our linguistic notions of sequential specification of actions are familiar, they are readily mapped onto a serial computational and programming model. Attempts to remap these representations onto the parallel model have proved so far quite inefficient.

There are three legs that support the understanding of efficient parallel processing: *computational models, underlying alternatives,* and the *programming paradigms.* This book is uniquely concerned about all three issues. Metrics for performance provide a quantitative basis for understanding basic models of computation.

Professors Hwang and Xu comprehensively analyze hardware models of processors, interconnection networks, and serial wide area networks giving a complete picture of the hardware state of the art. Their coverage extends from the hardware ILP (instructional level parallelism) to forms of parallelism achieved by NOW (networks of workstations). Their systems viewpoint integrates the computational model, the hardware, and the programming model in a review of major parallel processor implementation efforts that are currently underway.

It is in this integration of hardware and software that this book is most valuable. It is only by understanding the mapping from application onto system and into processor configuration that we can expect the consistent progress in building efficient scalable parallel processor systems.

# Preface

## Oh, A Digital World!

This book covers scalable architecture and parallel programming of multiprocessors, multicomputers, and network-based cluster platforms. Digital technology has made the computer industry. Now, digital technology is making another wave of fundamental impact to telecommunications and information industries. Converting everything to digital is the key to future success in a highly automated society.

The crossbreeding of technologies demands a new generation of computers that can adapt to scalable, parallel, and distributed computing. These changes in computer and information technologies has prompted computer professionals to study the material presented in this text. The ultimate goal is to become ready for new challenges in the 21st Century.

## A Glance at the Book

The book consists of 14 chapters presented in four parts. We provide a balanced coverage of four aspects: *principles, technology, architecture,* and *programming*:

- In Part I, three chapters cover scalable computer platforms and models, basics of parallel programming, and parallel performance metrics.

- Part II assesses commodity microprocessors, distributed cache and memory architecture, switched interconnects, Gigabit networks, and communications.

- Part III covers *symmetric multiprocessors* (SMP) and *cache-coherent, nonuniform memory-access* (CC-NUMA) machines, *clusters of workstations* (COW), and *massively parallel processors* (MPP).

- Part IV presents parallel languages, programming models with emphasis on Unix programming environments, message passing, data parallelism, and the use of PVM, MPI, Fortran 90, and HPF on scalable computers.

## A Trilogy on Computer Systems

Over a period of 15 years, Hwang and his associates have produced a trilogy on computer systems, all published by McGraw-Hill Book Company.

- *Computer Architecture and Parallel Processing,*
  by K. Hwang and F. A. Briggs (1983)

- *Advanced Computer Architecture*:
  *Parallelism, Scalability, Programmability,* by K. Hwang (1993)

- *Scalable Parallel Computing*:
  *Technology, Architecture, Programming,* by K. Hwang and Z. Xu (1998)

XVI

More than 90% of the topics treated in this book are based on new technological advances and research development within the past 5 years. This book is newly written, not a revision of Hwang's previous books. Unique features are highlighted below:

## Hot Chips and Interconnects

We assess commodity microprocessors and hot chips for building scalable multiprocessors and multicomputer clusters. Distributed cache/memory and Gigabit networks are studied along with latency hiding mechanisms. In particular, we study multiprocessor buses and crossbar switches, SAN (*System Area Network*), and LAN (*Local Area Network*) such as Gigabit Ethernet, SCI (*Scalable Coherence Interface*), and ATM (*Asynchronous Transfer Mode*) networks.

## Scalable Platforms and Clusters

We focus on scalable architectures, fast messaging mechanisms, latency hiding, distributed shared memory, cache coherence protocols, and memory consistency models. We cover software extensions for higher availability, single system image, failure recovery systems, and job management in clusters of computers.

Case studies include the HP/Convex Exemplar, Cray T3D/T3E, IBM SP2, Digital TruCluster, Microsoft Wolfpack, Sun Ultra Enterprise 10000, SGI Origin 2000, Sequent NUMA-Q, Intel/Sandia ASCI Option Red. We discuss the lessons learned from Stanford Dash, Berkeley NOW, Princeton SHRIMP, and Rice TreadMarks.

## Parallel Software Environments

This book devotes more than half of the material to software tools and parallel programming systems. In the shared-memory approach, we study the ANSI X3H5, Pthreads, SGI Power C; OpenMP, and C// language. We study Solaris MC and LSF (*Load Sharing Facility*) for availability, single system image, and cluster job management.

For parallel programming, we study data-parallel, message-passing, shared-memory, and implicit paradigms. We study MPI (*Message-Passing Interface*), PVM (*Parallel Virtual Machine*), Fortran 90, and HPF (*High Performance Fortran*) for explicit parallelism; and languages and compliers for implicit parallelism.

## Benchmark-based Evaluation

This book benefits from our benchmarking experience with six scalable computers: namely, the SP2 at Maui High Performance Computer Center, the T3D/T3E and Paragon at the San Diego Supercomputer Center, the T3D at Cray Eagan Data Center, the SP2, SGI server, and Pearl cluster at the University of Hong Kong.

Collective MPI communications in various machine platforms are evaluated with firsthand benchmark results. We reveal architectural implications from NASA parallel NAS and USC/HKU STAP benchmark results. These benchmark performance results are evaluated along with scalability analysis over machine sizes and problem sizes.

## Web Resources

In a rapidly changing world, any computer book becomes obsolete in a few years. We have strived to prolong the life cycle of this book by selecting practical topics and discussing fundamental issues that can last over generations of computer systems. Examples and quantitative data are drawn from real designs or benchmark experiments.

We have complied at the end of the book an extensive *Web Resources List*, linking to thousands of home pages of computer companies, research projects, information technology centers, and major application groups across academic, business, and government sectors. An on-line home site of this list is also maintained at HKU. See *Guide for Instructors/Students* to access our home pages.

## Acknowledgments

We thank the professional reviews of our draft manuscript by six leading experts in this field. Their suggestions are very useful in revising the manuscript to its present form and contents. We appreciate Choming Wang for preparing the indices and the help from Dr. Cho-Li Wang in maintaining the book's Web site at HKU.

Intellectual exchanges with Dharma Agrawal, Jean-Loup Baer, Gordon Bell, David Culler, Jack Dongarra, Michael Flynn, Ian Foster, Jeffrey Fox, Mark Franklin, Wolfgang Giloi, Allan Gottlieb, Anoop Gupta, John Hennessy, Ken Kennedy, Duncan Lawrie, Charles Leiserson, Kai Li, Guojie Li, Lionel Ni, David Patterson, Gregory Pfister, John Rice, Sartaj Sahni, Chuck Seitz, Bruce Shriver, H. J. Siegel, Burton Smith, Daniel Tabak, H. C. Torng, and Ben Wah are always inspiring and appreciated.

We appreciate the sponsorship from McGraw-Hill editors, Eric Munson, Lynn Cox, and Betsy Jones. The production work from Richard DeVitto, Francis Owen, and Nina Kreiden and assistance from Polly Leung of HKU are gratefully acknowledged.

During the courses of writing this book, research funding supports from MIT Lincoln Laboratory, Hong Kong Research Grants Council, and the University of Hong Kong are appreciated. In particular, the excellent facilities and environment provided by HKU make the writing of this book a very pleasant undertakin,

## Points of Feedback

For all technical contacts, suggestions, corrections, or exchanges of information, readers and university instructors are advised to contact either author via Email:

        **kaihwang@cs.hku.hk**    **zxu@apple.ncic.ac.cn**

We appreciate your feedback and hope you enjoy reading the book.

<div align="right">

**Kai Hwang**
**Zhiwei Xu**

</div>

November 15, 1997
Hong Kong

# Guide to Instructors/Readers

This book is designed as a standard text for classroom adoption in Computer Science or Computer Engineering curriculum at college/university levels. Suitable courses include: Computer Architecture, Parallel Processing, Distributed Computing, Concurrent Programming, Network-based Computing, Computer Engineering, etc.

**Flowchart for Reading**   The flowchart shows the logical flow of the 14 chapters in this book. The four parts are indicated on the side labels.There are two chapters in theory and modeling (slashed boxes), six unshaded boxes for hardware and architecture chapters, and six shaded boxes for software and programming chapters.



**Course Offerings**

Suggested below are five possible course offerings in adopting this book for use in an one-semester course with 45 hours of lectures:

- **Computer Architecture**: For hardware-oriented students in Electrical and Computer Engineering programs, cover Chapters 1, 3-6, and 8-11.

- **Parallel Programming**: For programming and software-oriented students in Computer Science programs, Chapters 1, 2, 7-10, and 12-14 are suitable.

- **Parallel Processing:** For mixed students from Computer Science/Engineering and Electrical Engineering programs, cover Chapters 1, 2, 4-6, 8-10, and 12.

- **Distributed Computing:** For mixed CS and EE students, Chapters 1, 2, 5-7, 9-10, and 13-14 are suitable for one-semester use.

- **Computer Engineering:** An advanced course in computer technology and software system design. Chapters 1, 2, 4-7, and 8-10 should be covered.

All readers should start with the material in chapter 1. Engineering students may read more of the technology and architecture chapters on the left side and software-oriented students on the right subtree of the flowchart. Logically, the reading of this book should flow from the top to the bottom chapters shown by the arrows.

The six hardware and architecture chapters form the core of the course in *Computer Architecture* or in *Computer Engineering*. The six chapters in software and programming form the core for the *Parallel Programming* course. The *Parallel Processing* course covers all scalable parallel systems with more emphasis on shared-memory multiprocessing. The *Distributed Computing* course emphasizes message-passing systems and network-based cluster computing.

## Instructor's Manual

An *Instructor's Manual* is available to proven instructors without charge. The Manual provides solutions to selected homework problems, viewing-graph masters, some sample tests, and topics suggested for term projects.

Instructors can request the Manual by writing to: Lynn Cox, McGraw-Hill College Division, 55 Francisco Street, Suite 200, San Franscico, CA. 94133-2117, U.S.A. or submit a written request to Fax No. 415-989-7702.

## Web Site Access

As an updated reference, the book is also intended for self-study use by system designers, academic researchers, application programmers, system analysts, resource managers, solution providers, and computer professionals in general. To avoid becoming obsolete, readers are invited to visit our Web site for updated WWW links.

**http://www.cs.hku.hk/~kaihwang/book98.html**

This Web site is updated dynamically. If you want your organization or your projects to be added into the list. Contact Dr. C. L. Wang of the University of Hong Kong by Email: **clwang@cs.hku.hk.**

# Table of Contents