

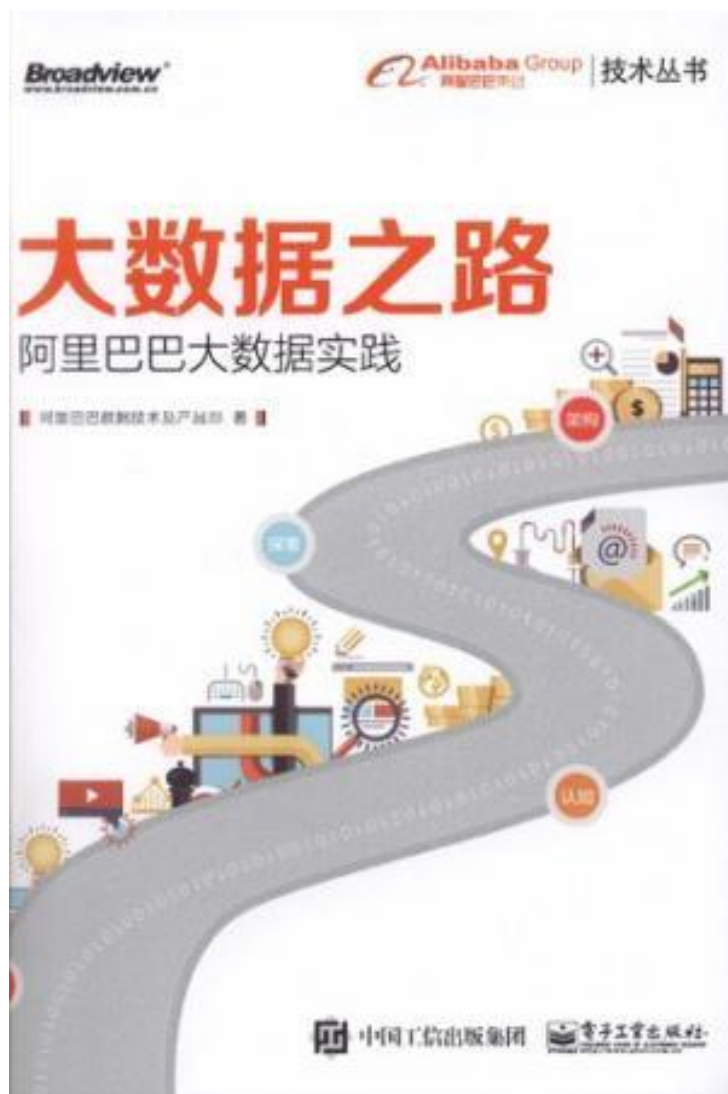


A 1st introduction to **Large Scale Computing** Concepts & Techniques

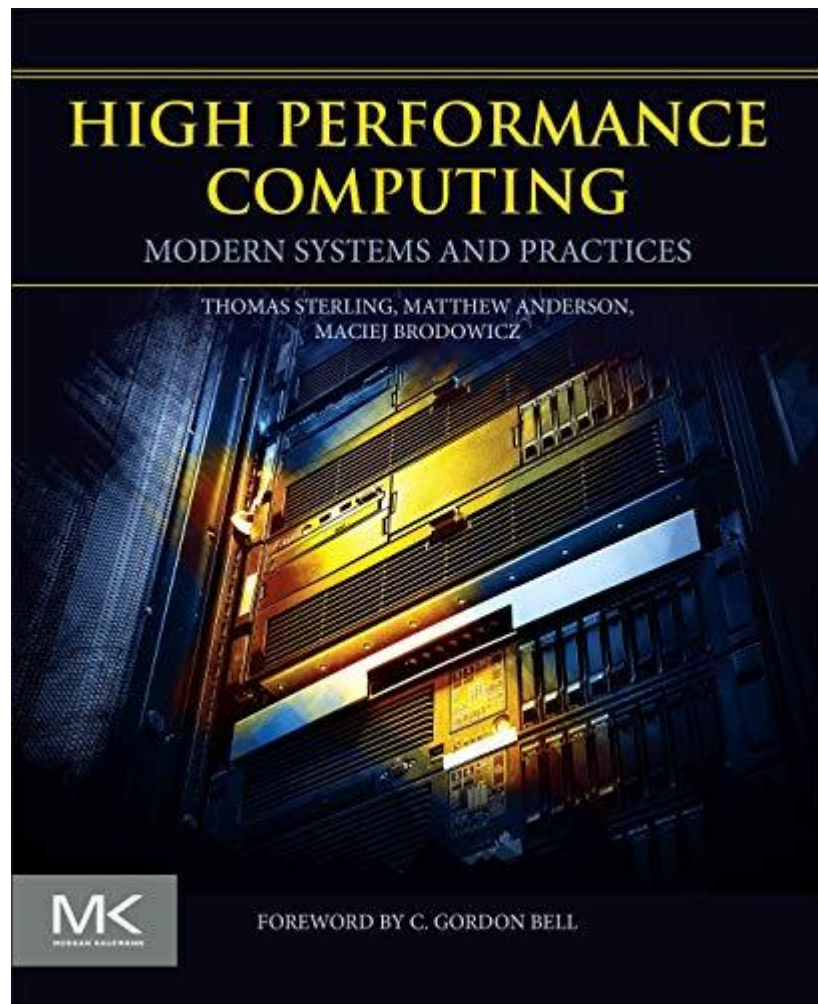
Chapter 1: Introduction



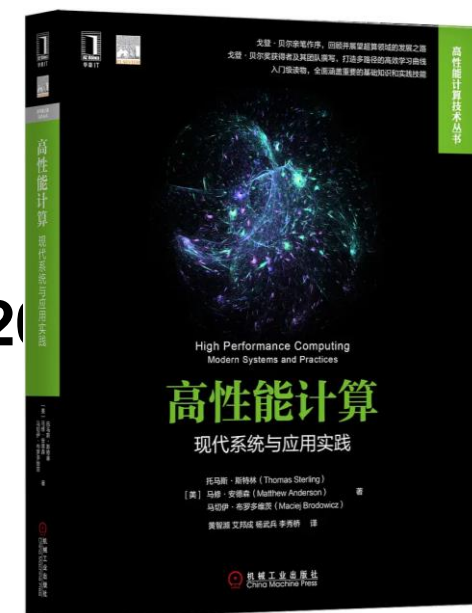
mlinking@126.com
+86 15010255486

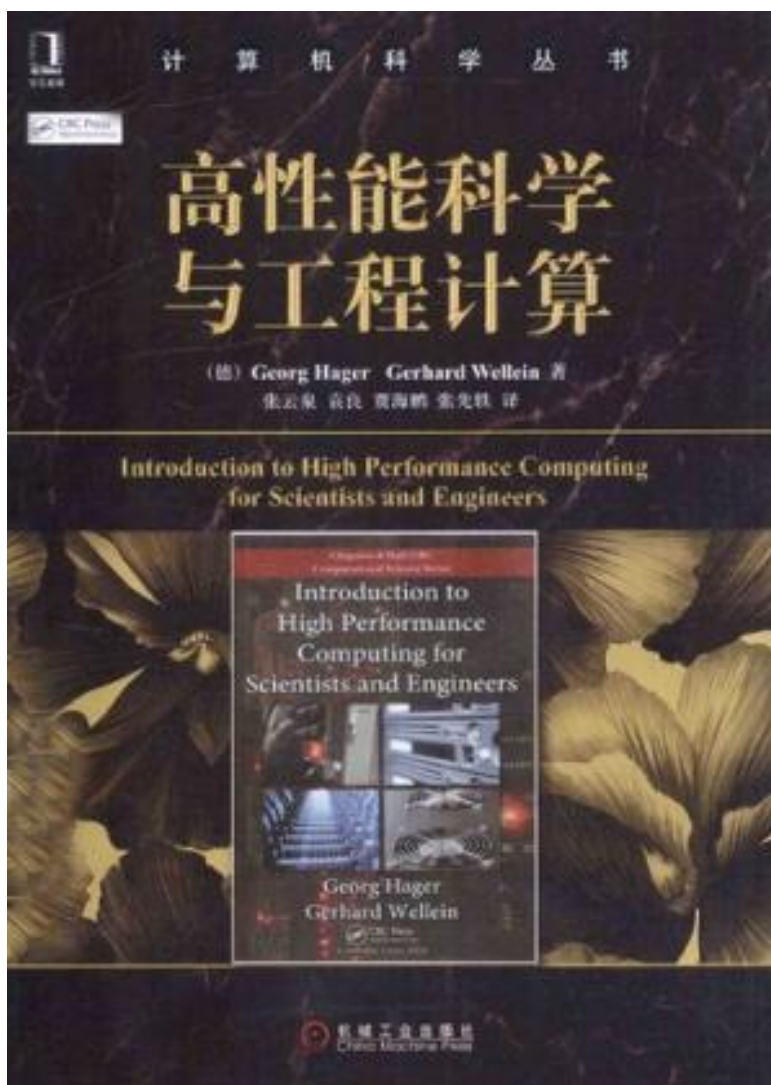


- 大数据之路: 阿里巴巴大数据实践
- 阿里巴巴数据技术及产品部
- 2017
- 电子工业出版社



- ❑ High Performance Computing: Modern Systems and Practices
- ❑ By 作者: Thomas Sterling – Matthew Anderson – Maciej Brodowicz
- ❑ ISBN-10 书号: 012420158X
- ❑ ISBN-13 书号: 9780124201583
- ❑ Edition 版本: 1
- ❑ Release Finelybook 出版日期: 2015
- ❑ pages 页数: (718)

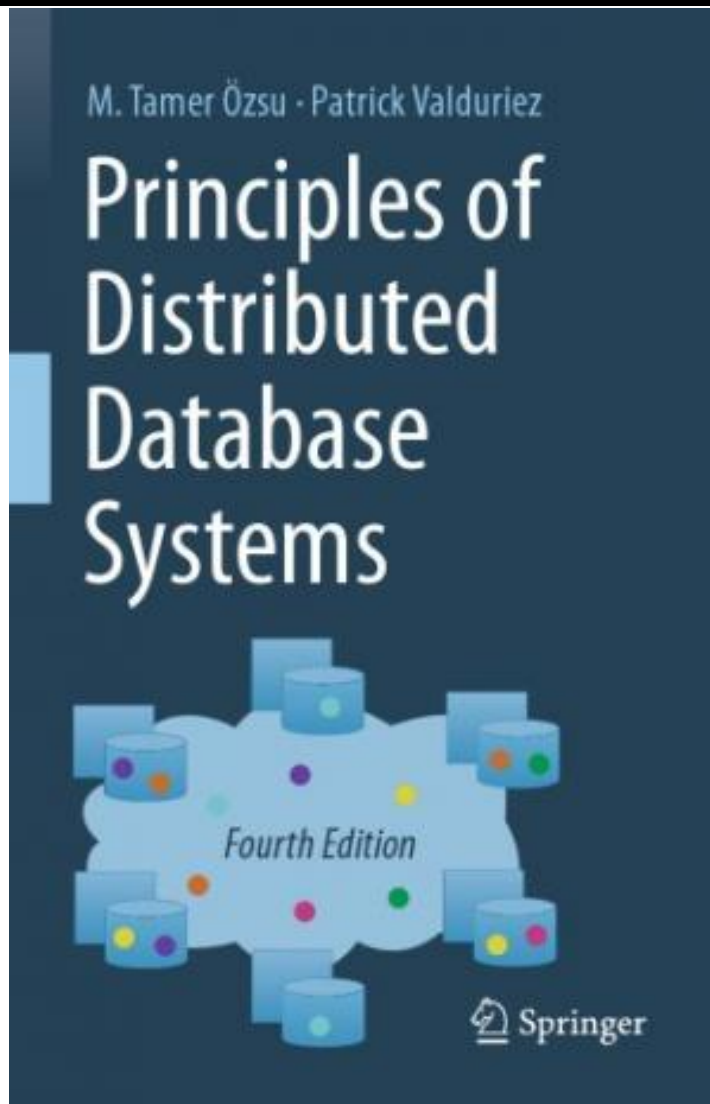




□ 高性能科学与工程计算

□ [德] Georg Hager, [德] Gerhard Wellein

- 《计算机科学丛书：高性能科学与工程计算》从工程实践的角度介绍了高性能计算的相关知识。主要内容包括现代处理器的体系结构、为读者理解当前体系结构和代码中的性能潜力和局限提供了坚实的理论基础。
- 接下来讨论了高性能计算中的关键问题，包括串行优化、并行、OpenMP、MPI、混合程序设计技术。
- 作者根据自身的研究也提出了一些前沿问题的解决方案，如编写有效的C++代码、GPU编程等。



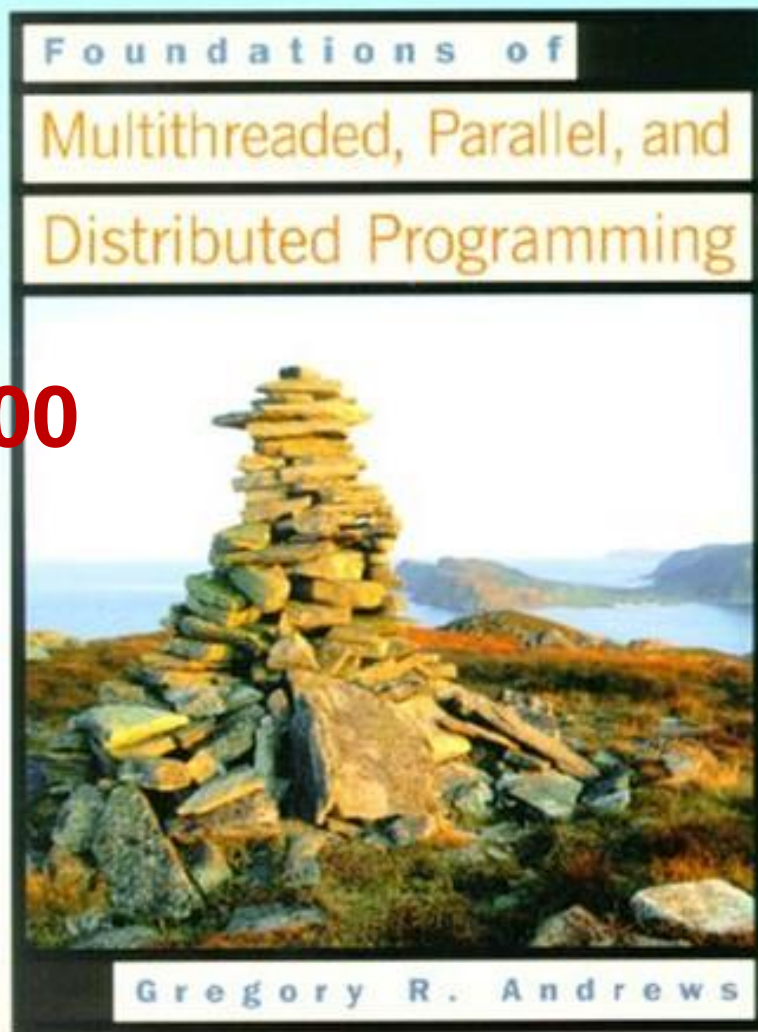
- ❑ Principles Of Distributed Database Systems
- ❑ M. Tamer Özsu, Patrick Valduriez
- ❑ Publisher: Springer
- ❑ Year: 2020

Chapter 1: Introduction

- ❑ About me
- ❑ Evaluation in this course
- ❑ Why do we have this course, and How do I organize?
 - Modern Society needs more computation power
 - Organize the topics with an example
- ❑ Resources related to this course



2000



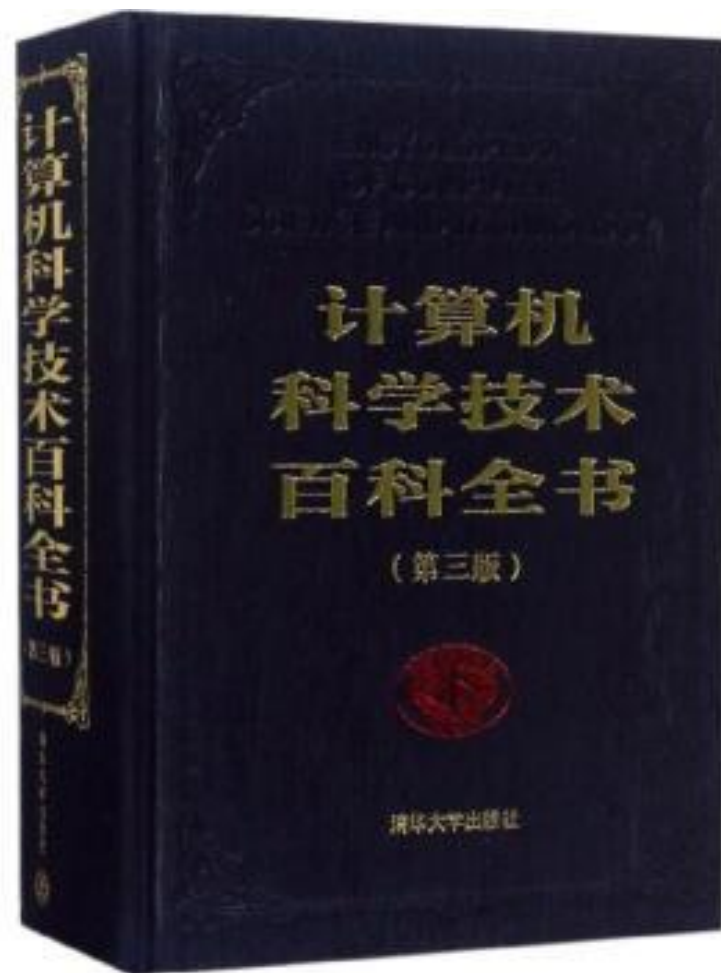
GREGORY R. ANDREWS

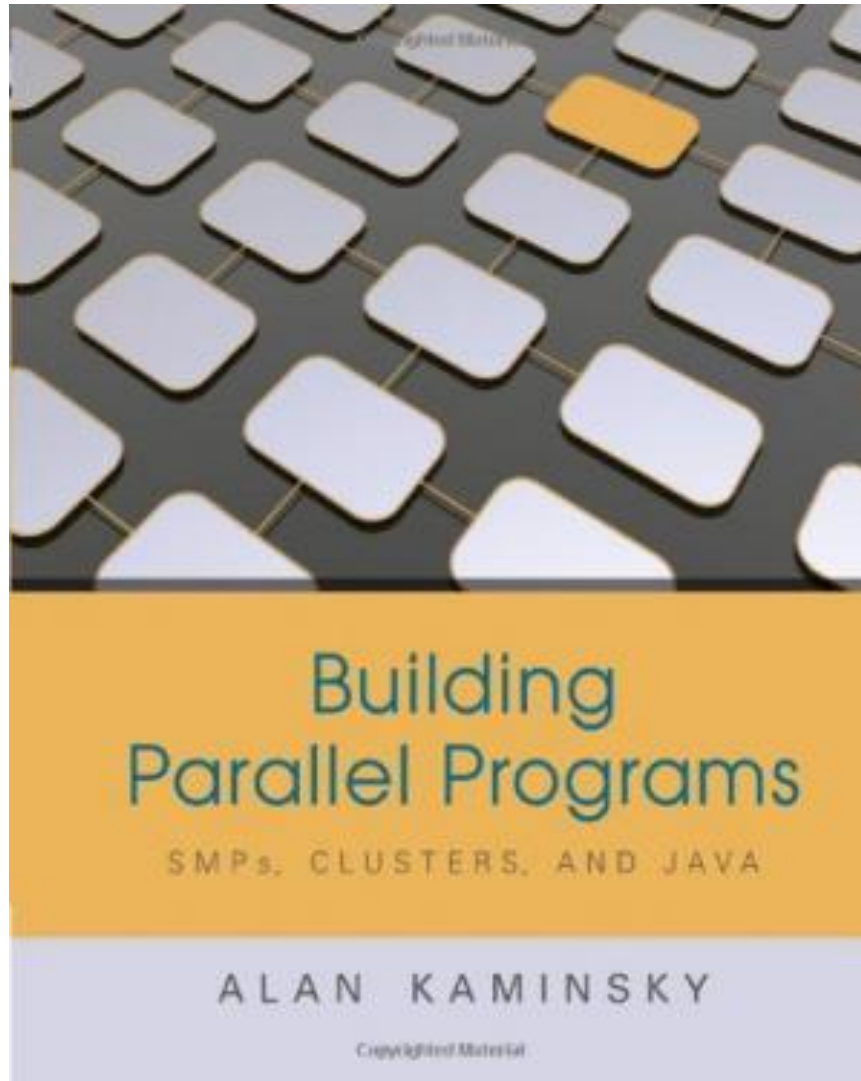
CONCURRENT PROGRAMMING

PRINCIPLES AND PRACTICE

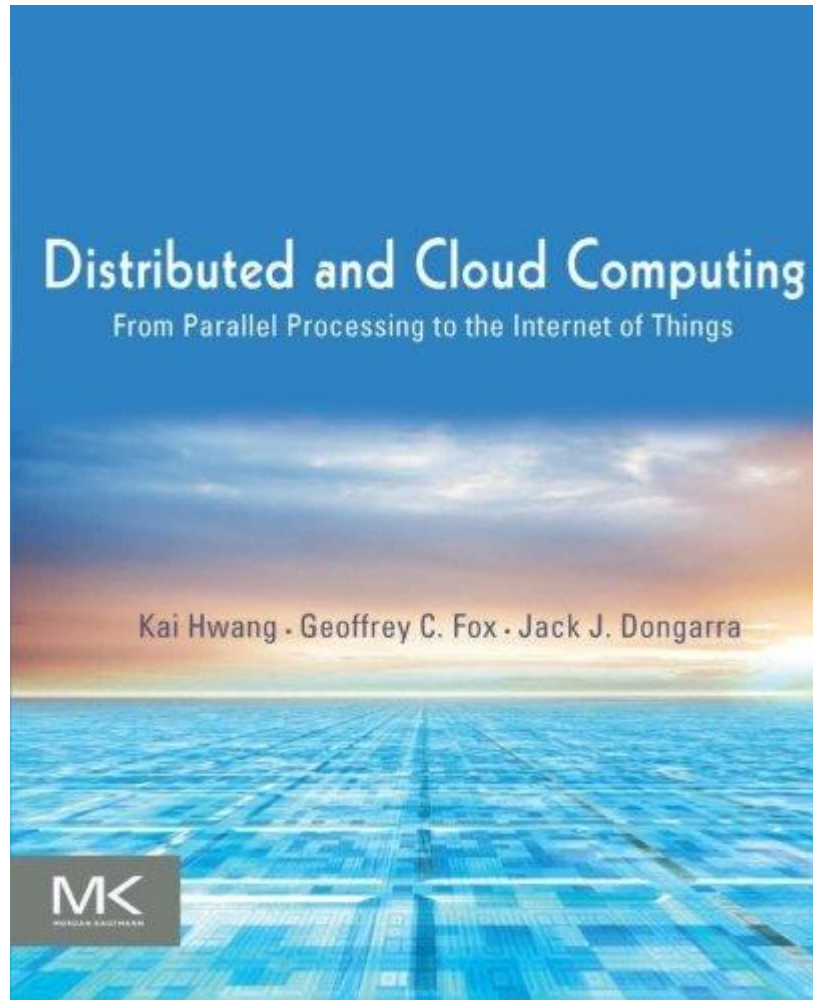


<https://www2.cs.arizona.edu/~greg/>

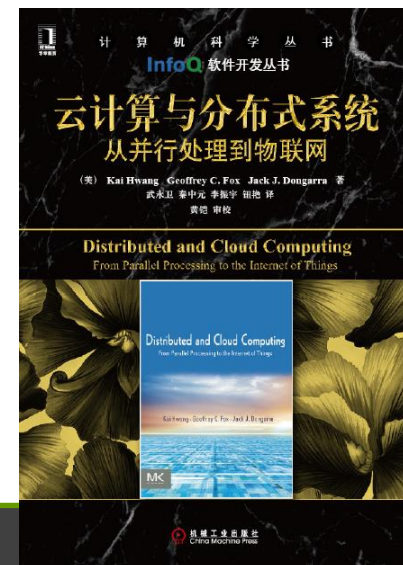


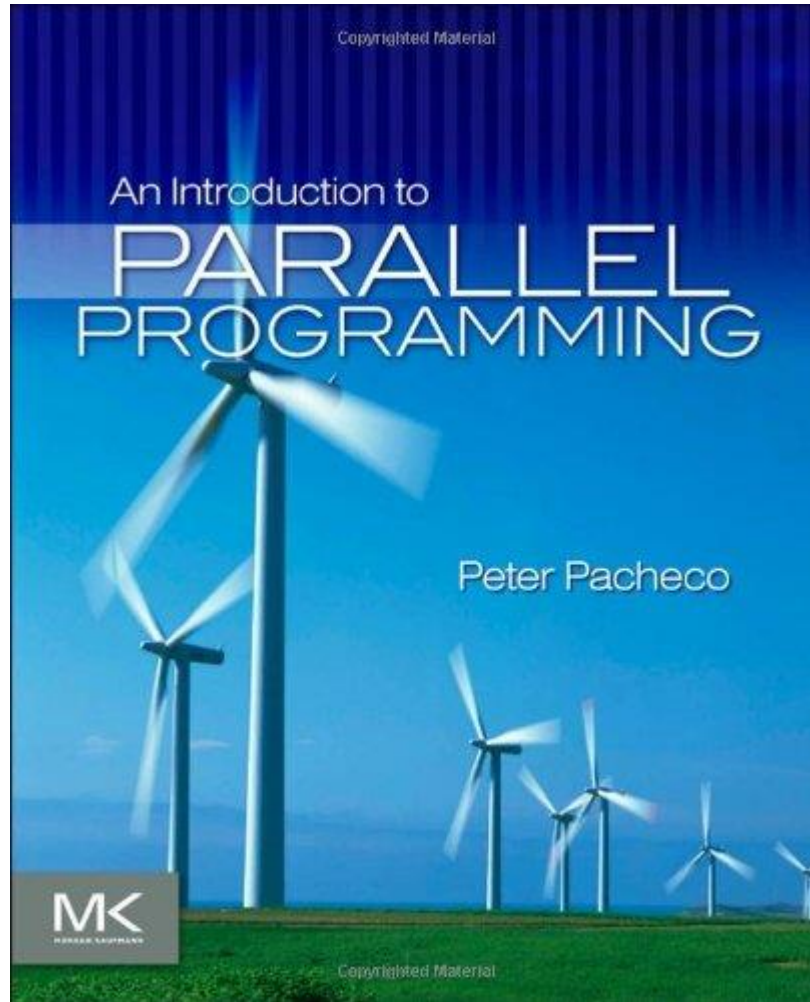


- ❑ Building parallel programs: SMPs, clusters, and Java
- ❑ Alan Kaminsky
- ❑ **2009**

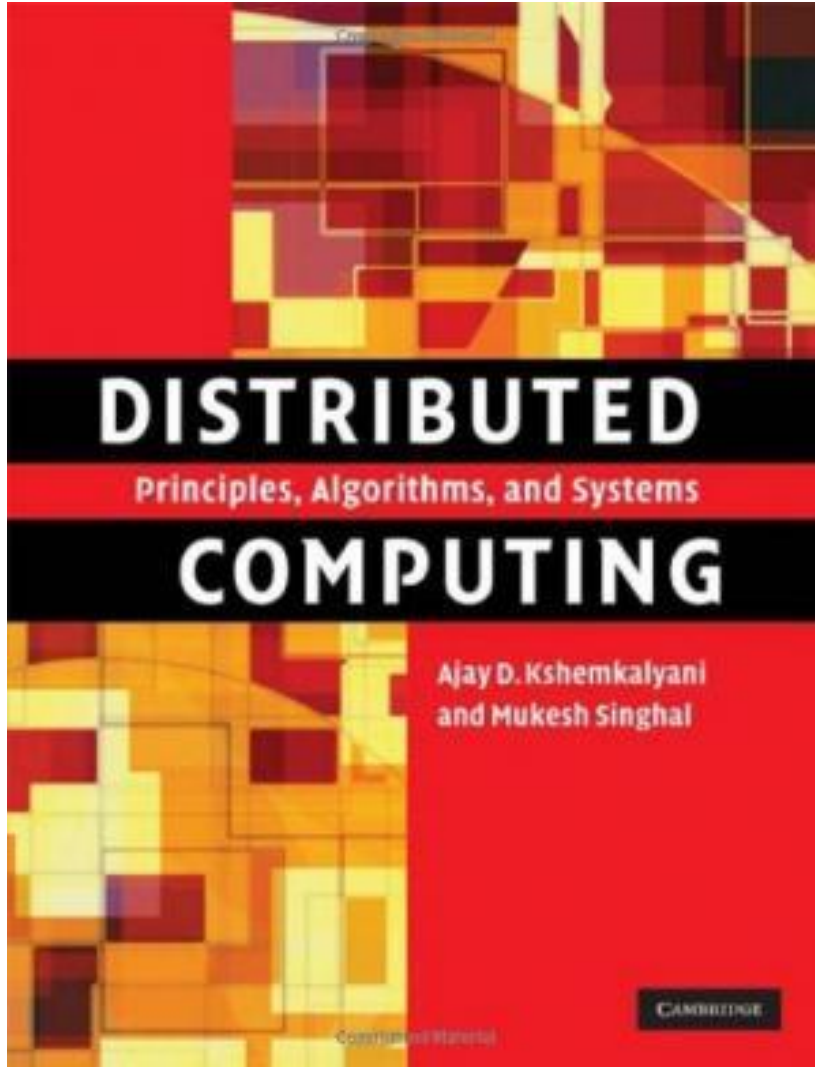


- ❑ **Distributed and Cloud Computing: From Parallel Processing to the Internet of Things**
- ❑ **Kai Hwang, Jack Dongarra, Geoffrey C. Fox**

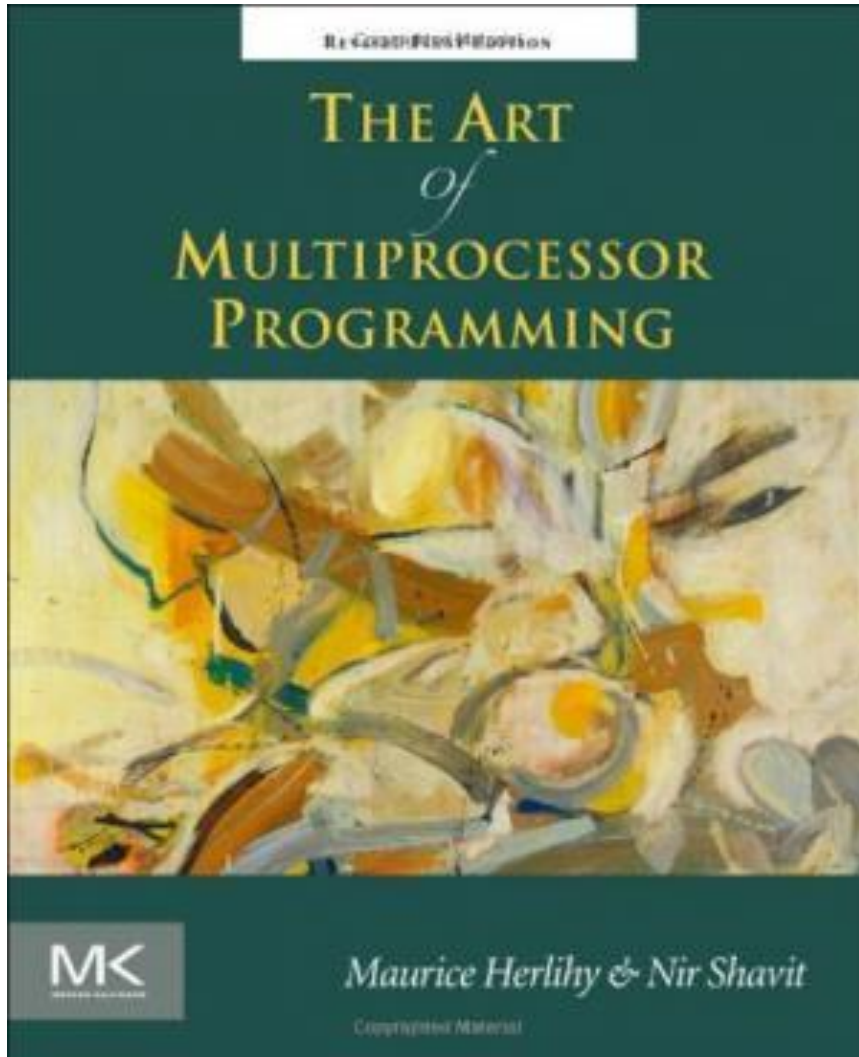




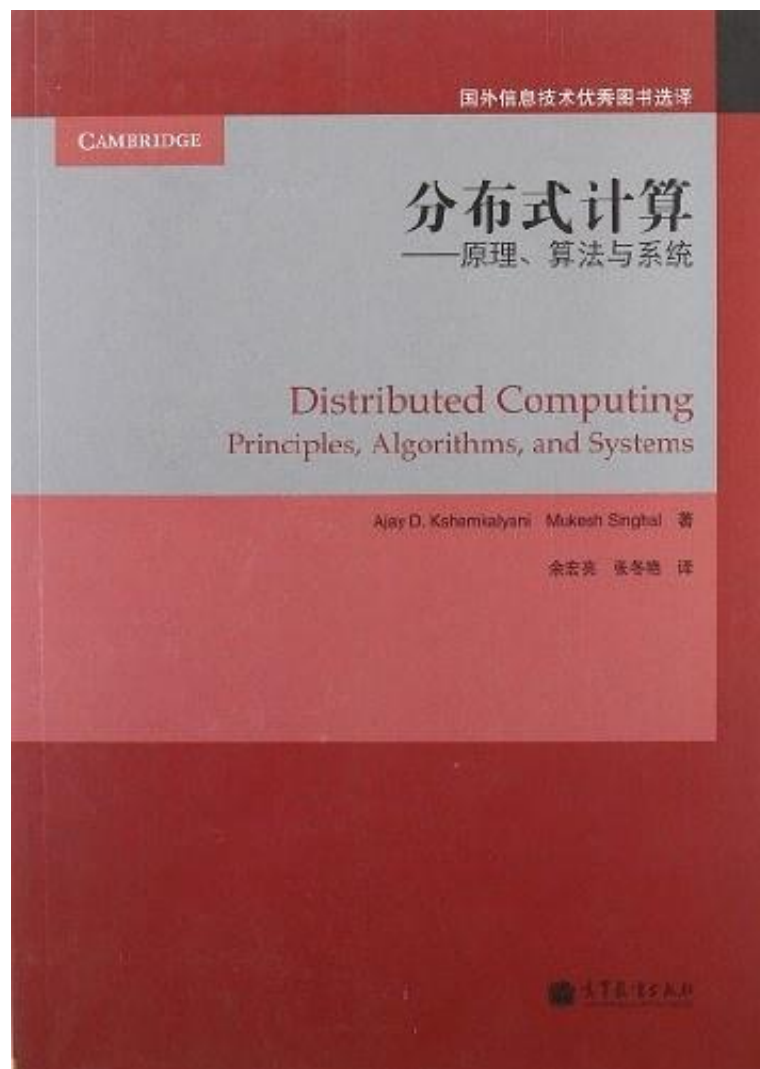
- **An Introduction to Parallel Programming**
- **Peter Pacheco**



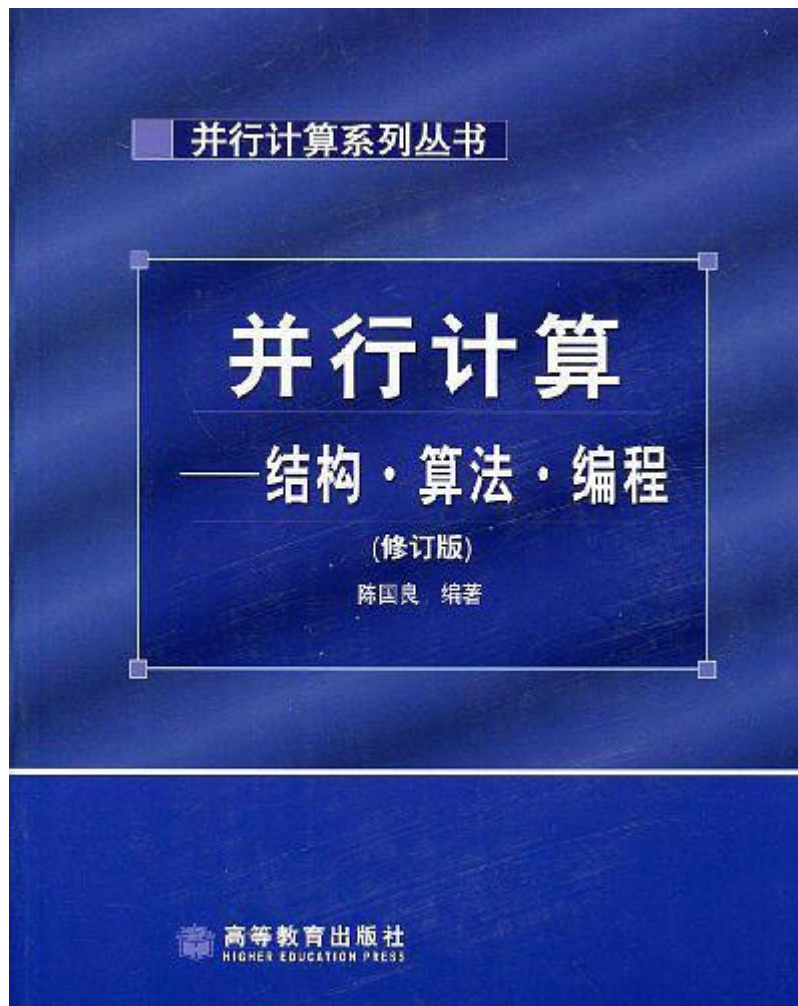
- ❑ Distributed Computing: Principles, Algorithms, and Systems
- ❑ Ajay D. Kshemkalyani, Mukesh Singhal
 - Designing distributed computing systems is a complex process requiring a solid understanding of the design problems and the theoretical and practical aspects of their solutions. This comprehensive 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 issues such as mutual exclusion, deadlock detection, authentication, and failure recovery. Algorithms are carefully selected, lucidly presented, and described



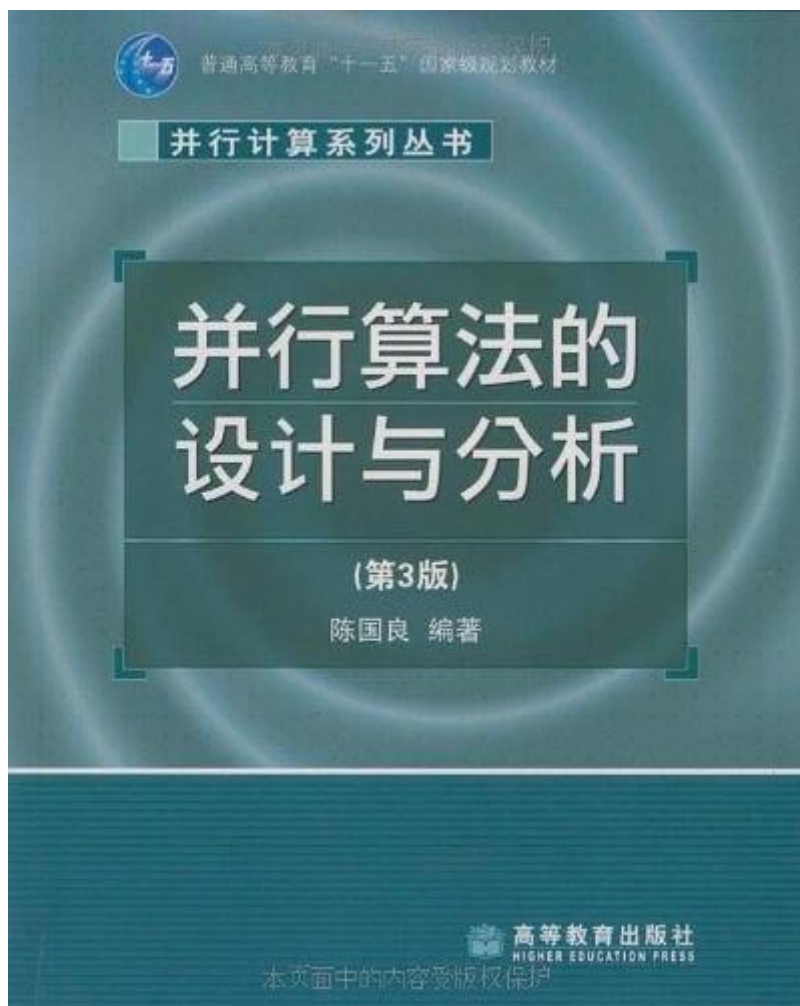
- ❑ The Art of Multiprocessor Programming, Revised Reprint
- ❑ Maurice Herlihy, Nir Shavit
 - This revised edition incorporates much-demanded updates throughout the book, based on feedback and corrections reported from classrooms since 2008
 - Learn the fundamentals of programming multiple threads accessing shared memory
 - Explore mainstream concurrent data structures and the key elements of their design, as well as synchronization techniques from simple locks to transactional memory systems
 - Visit the companion site and download source



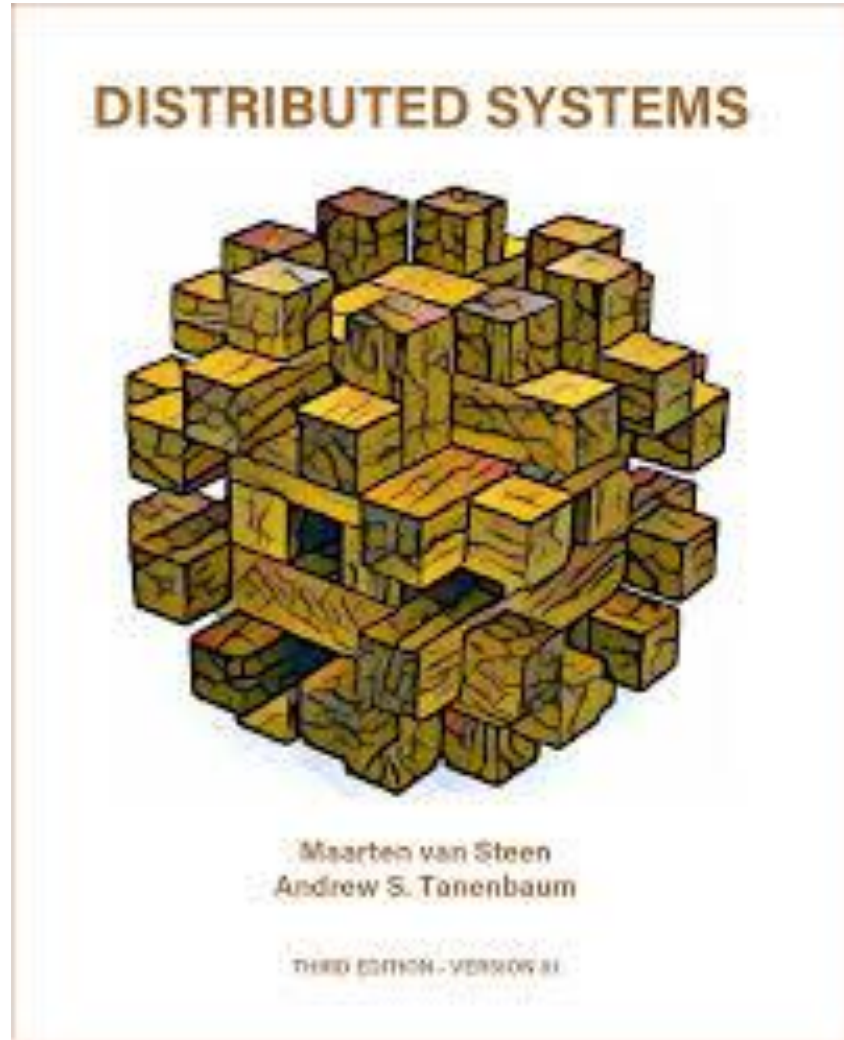
- **分布式计算**
- **作者: Ajay D. Kshemkalyani / Mukesh Singhal**
- **出版社: 高等教育出版社**
 - 副标题: 原理, 算法与系统
 - 原作名: Distributed Computing: Principles, Algorithms, and Systems
- **译者: 余宏亮 / 张冬艳**
- **出版年: 2012-6**
- **定价: 79.00元**
- **ISBN: 9787040324563**



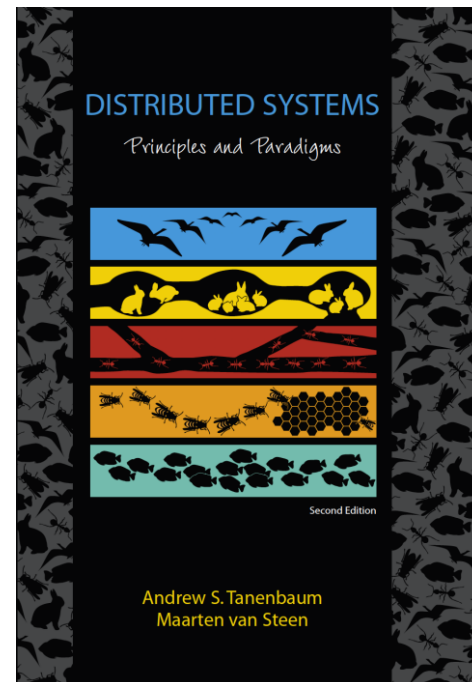
- 并行计算：结构、算法、编程
- 作者：陈国良
- 出版社：高等教育出版社
- 副标题：结构、算法、编程
- 出版年：2003-8-1
- 页数：450
- 定价：36.50
- 装帧：平装
- ISBN：9787040133073

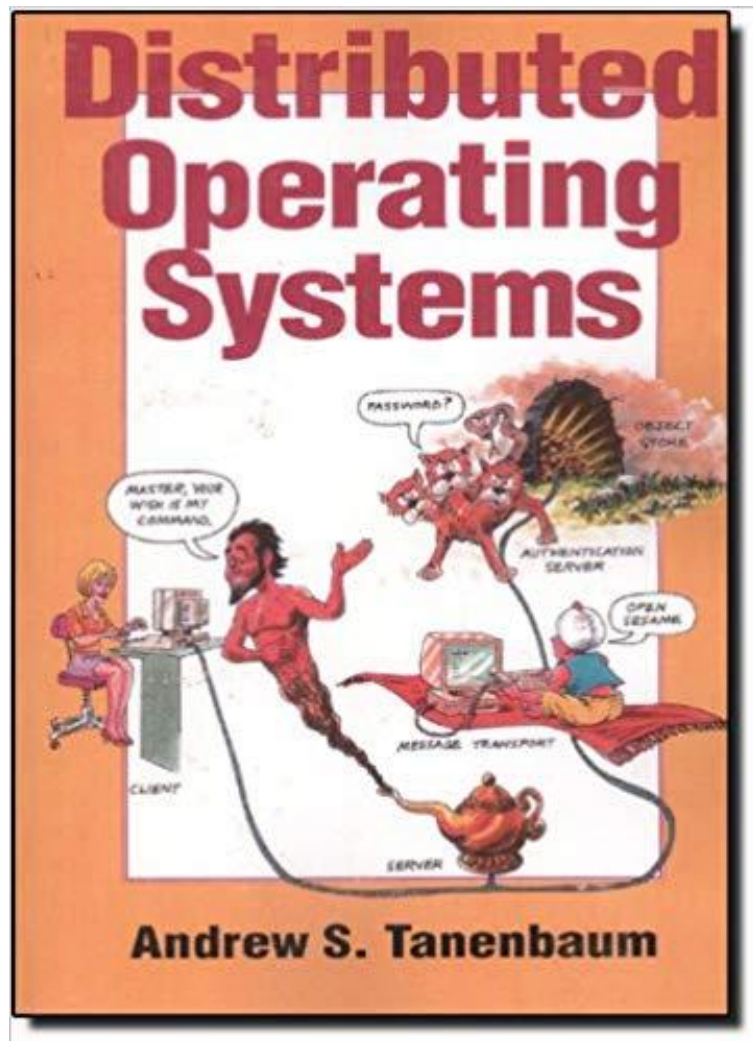


- 并行算法的设计与分析
- 作者: 陈国良
- 出版年: 2009-8
- 页数: 813
- 定价: 66.00元
- 丛书: 并行计算系列丛书
- ISBN: 9787040264364

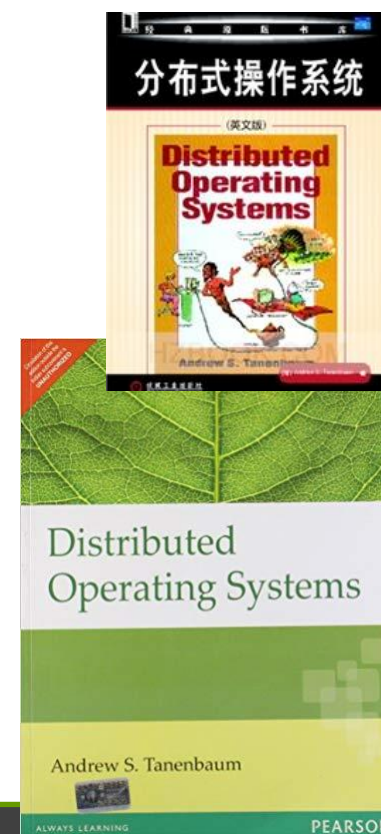
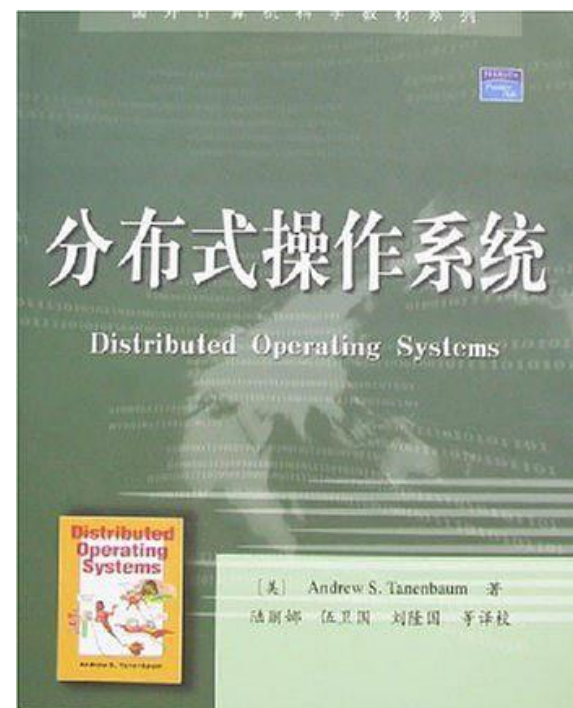


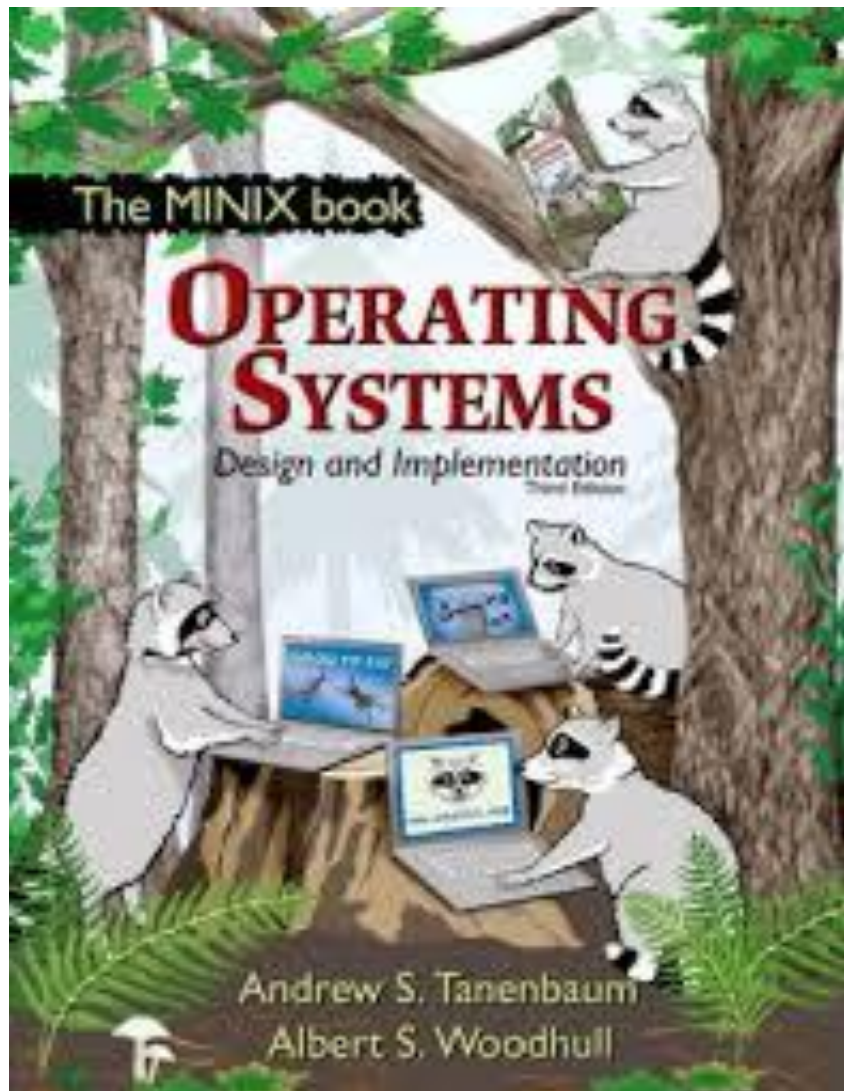
- ❑ **Distributed Systems 3rd edition (2017)**
- ❑ Andrew S Tanenbaum, Maarten Van Steen

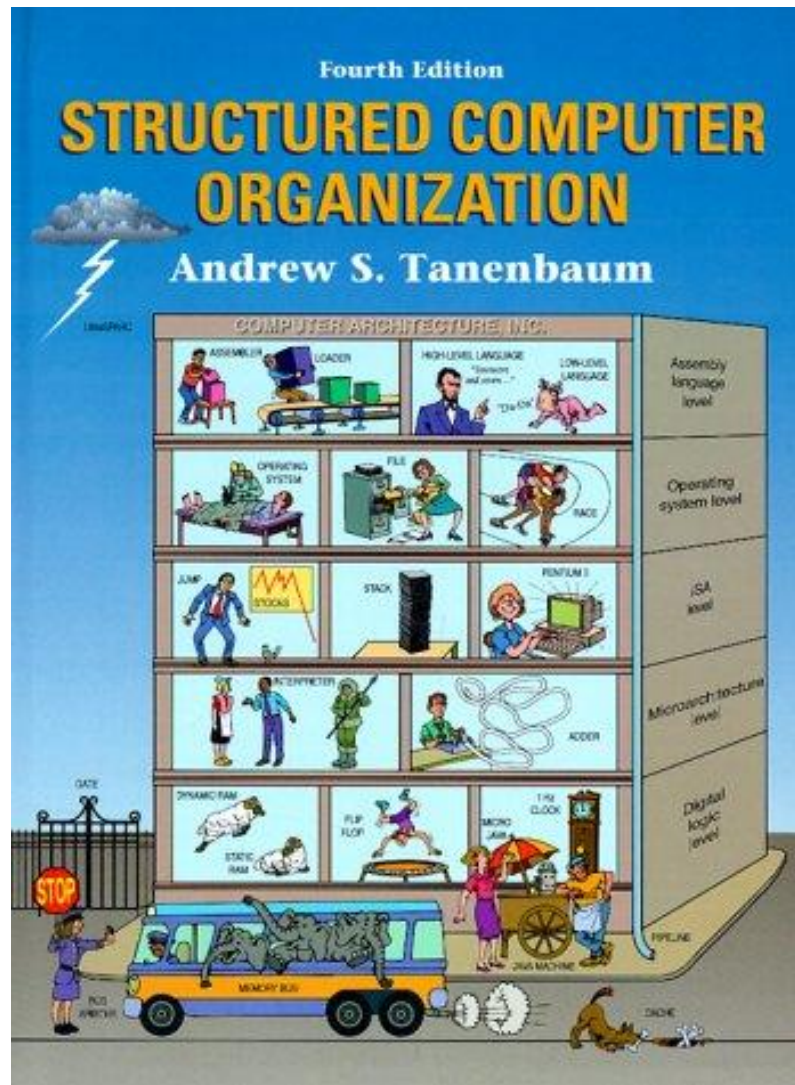




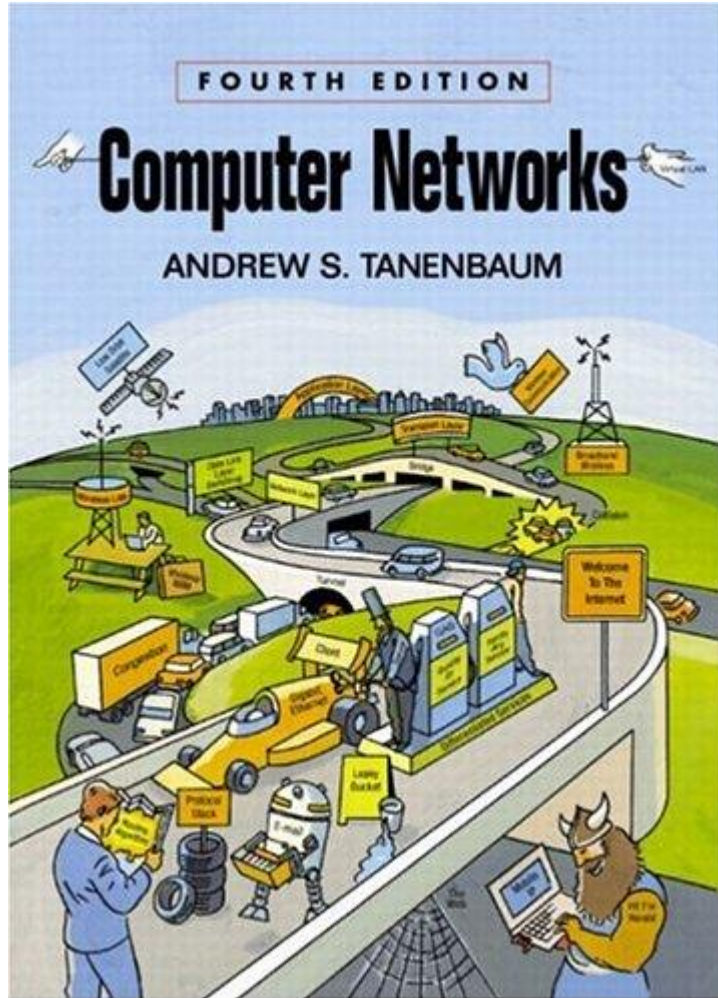
- Distributed Operating Systems
- Andrew S. Tanenbaum



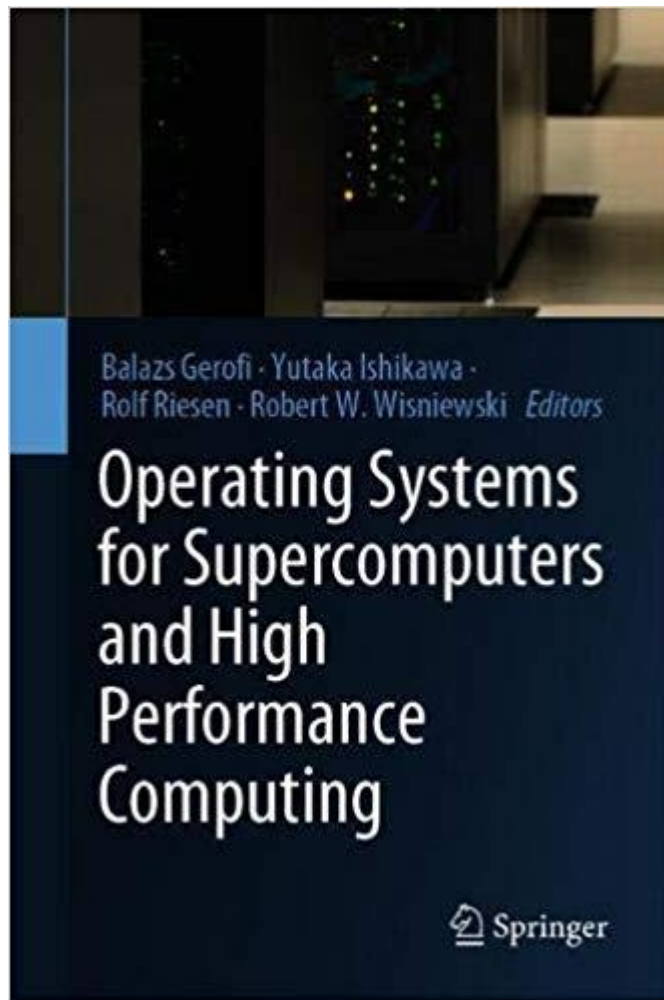




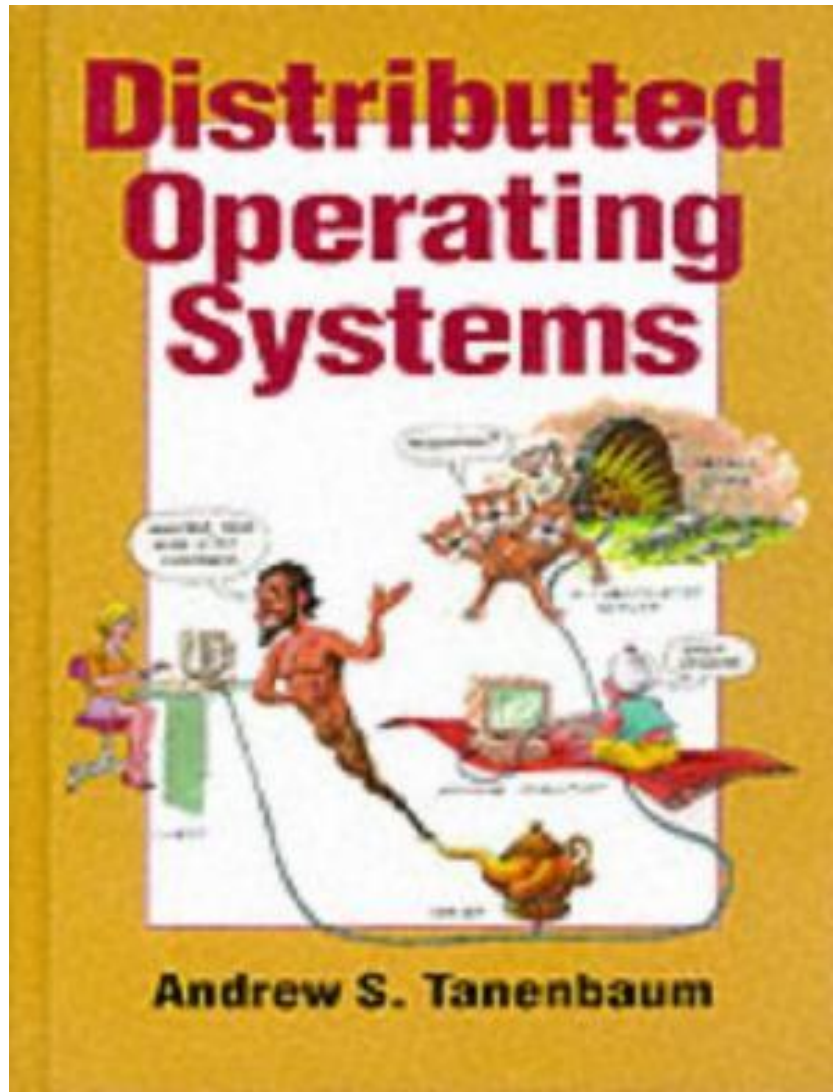
- Structured Computer Organization
- Andrew S. Tanenbaum



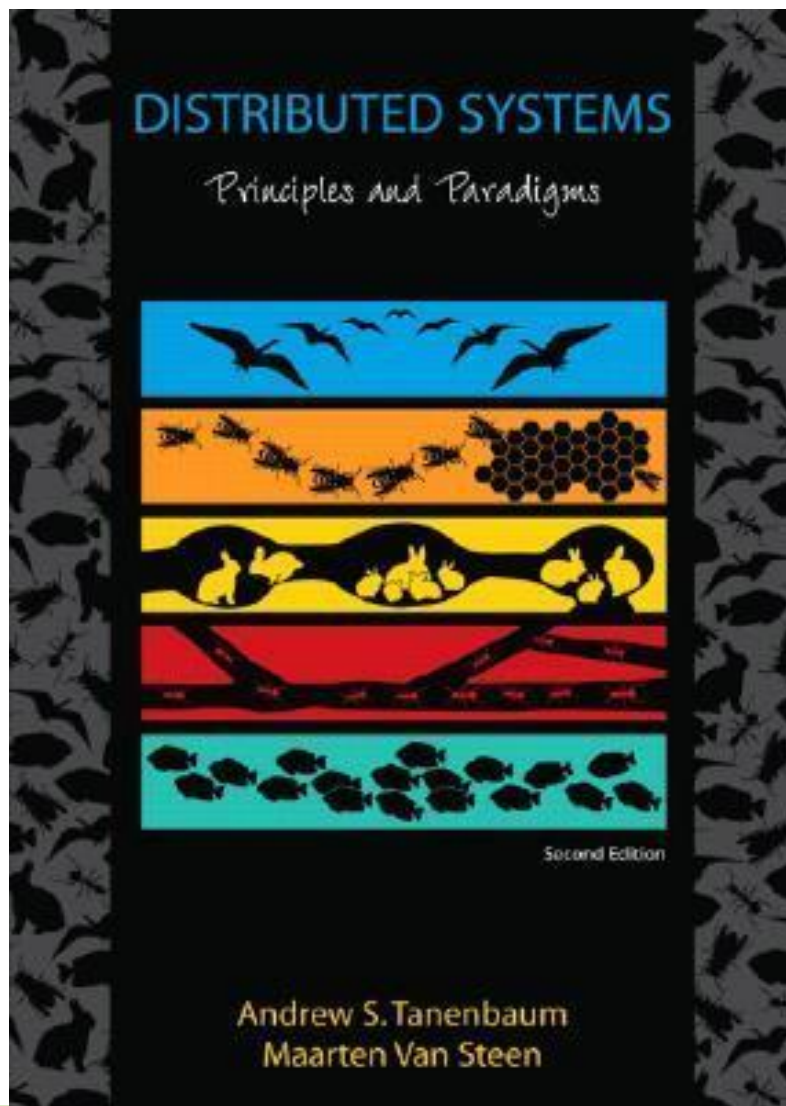
- ❑ **Computer Networks (4th Edition)**
- ❑ **Andrew S. Tanenbaum**



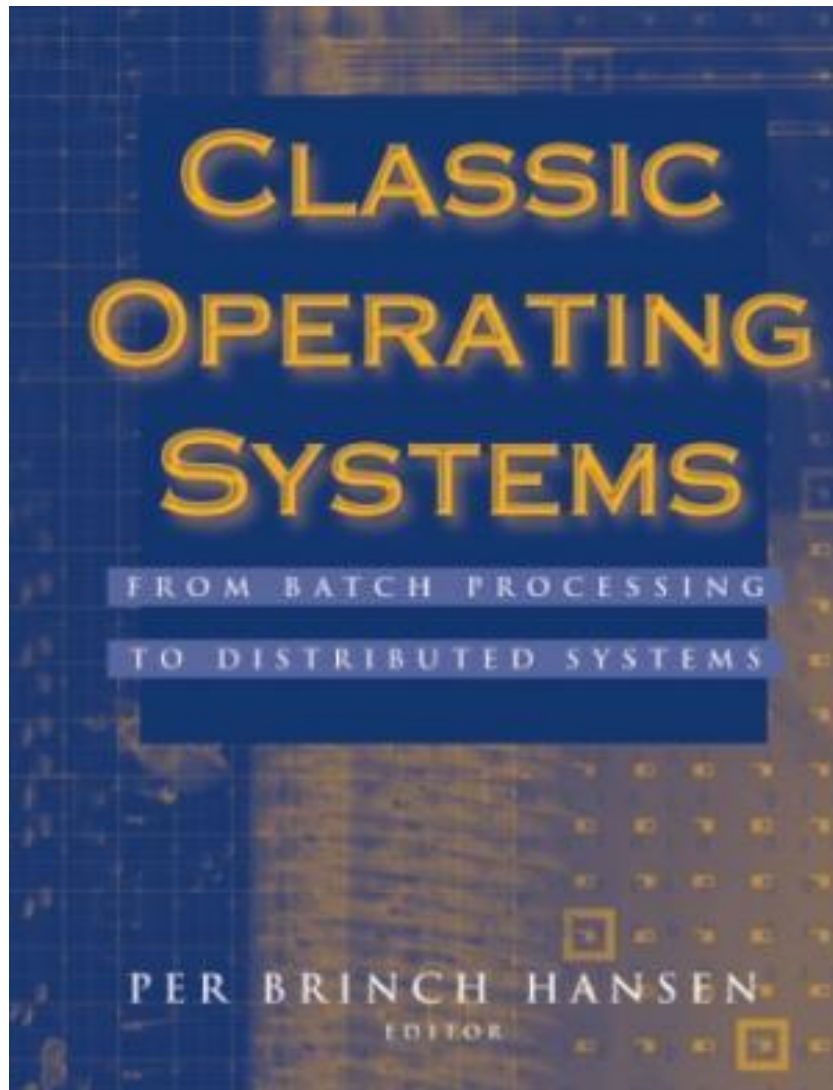
- ❑ Operating Systems for Supercomputers and High Performance Computing (High-Performance Computing Series)
- ❑ By 作者: Balazs Gerofi
- ❑ ISBN-10 书号: 9811366233
- ❑ ISBN-13 书号: 9789811366239
- ❑ Edition 版本: 1st ed. 2019
- ❑ Release Finelybook 出版日期: 2019-10-16
- ❑ pages 页数: (400)



- ❑ Distributed Operating Systems
- ❑ Andrew S. Tanenbaum
- ❑ 1994

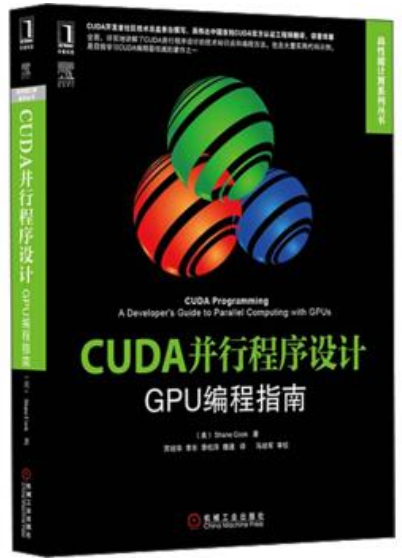
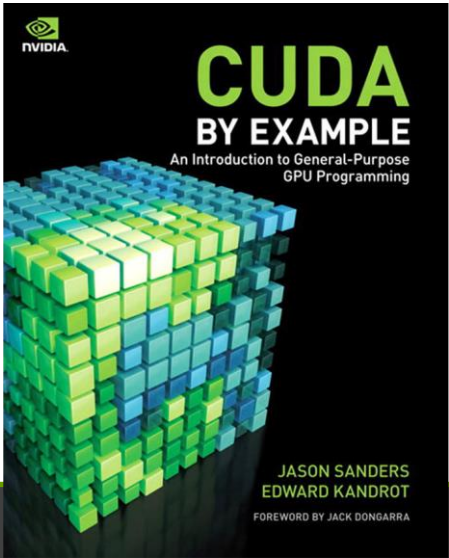
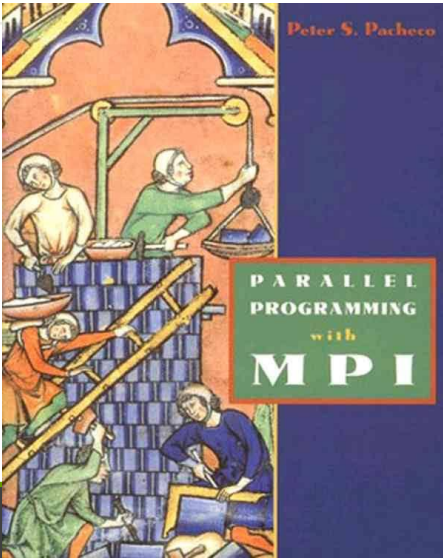
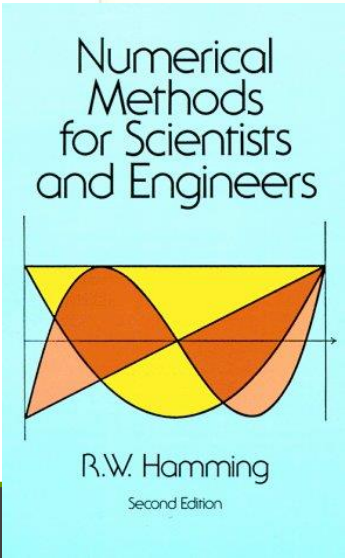
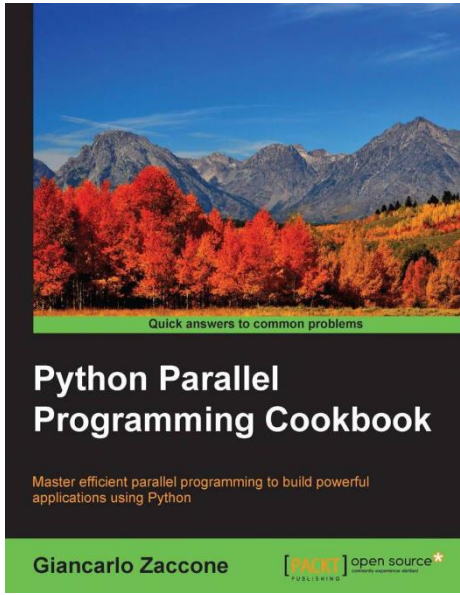
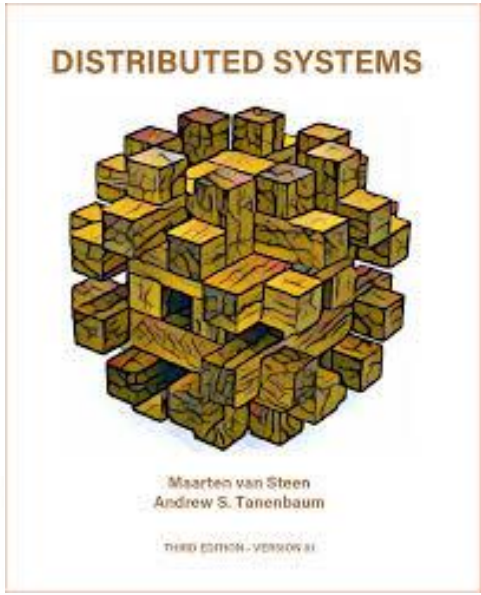
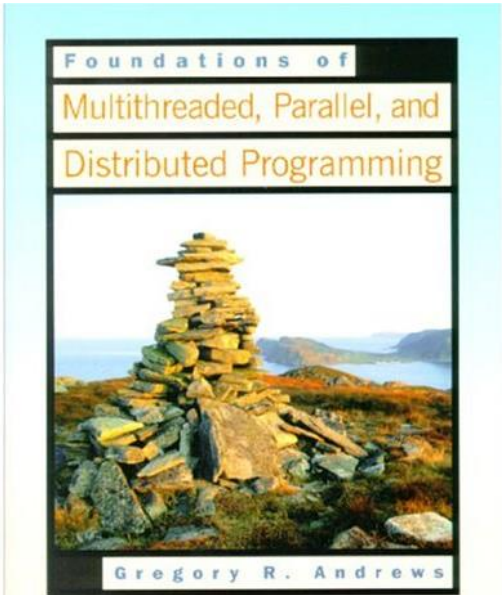


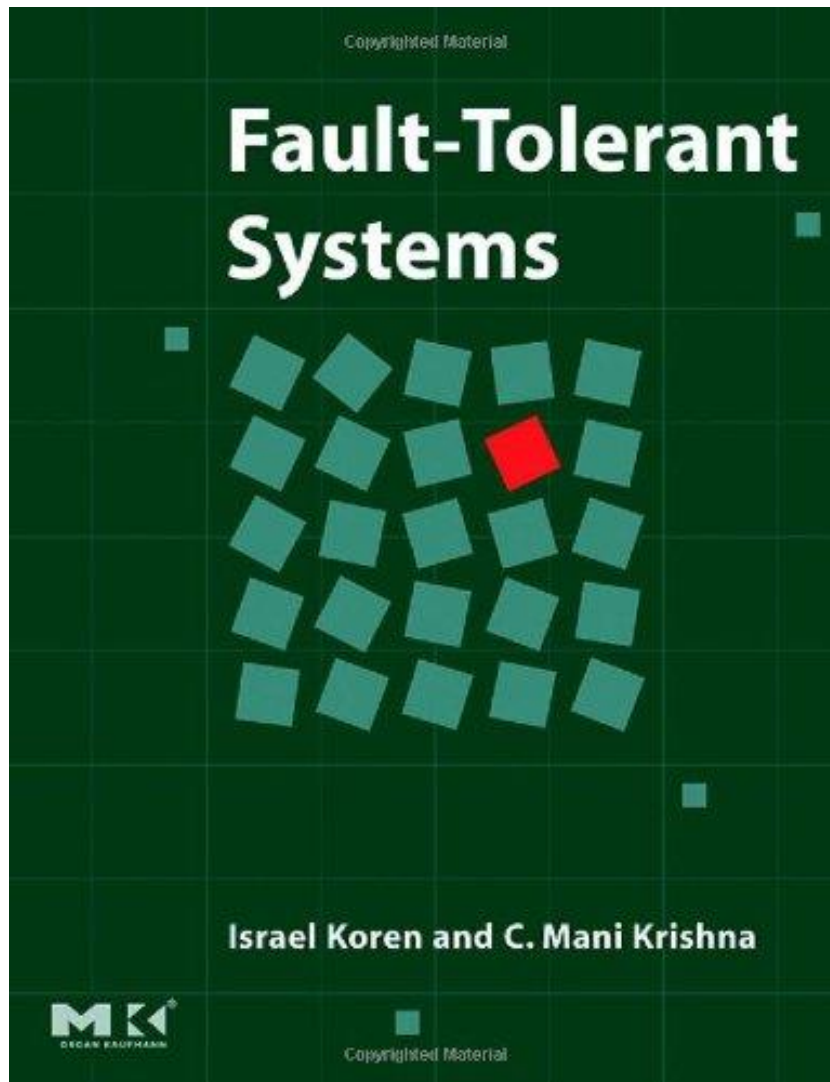
- ❑ Distributed Systems: Principles and Paradigms
- ❑ Andrew S. Tanenbaum, Maarten van Steen
- ❑ 2006



- ❑ Classic Operating Systems: From Batch Processing to Distributed Systems
- ❑ Per Brinch Hansen (auth.), Per Brinch Hansen (eds.)
- ❑ **2001**
- ❑ Springer

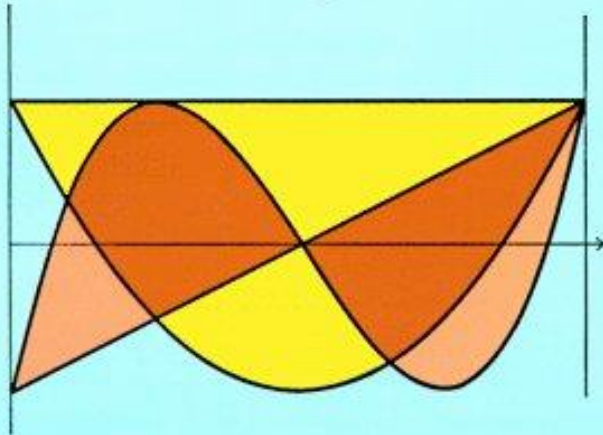
A long list of books, but I prefer





- Fault-Tolerant Systems
- Israel Koren, C. Mani Krishna

Numerical Methods for Scientists and Engineers



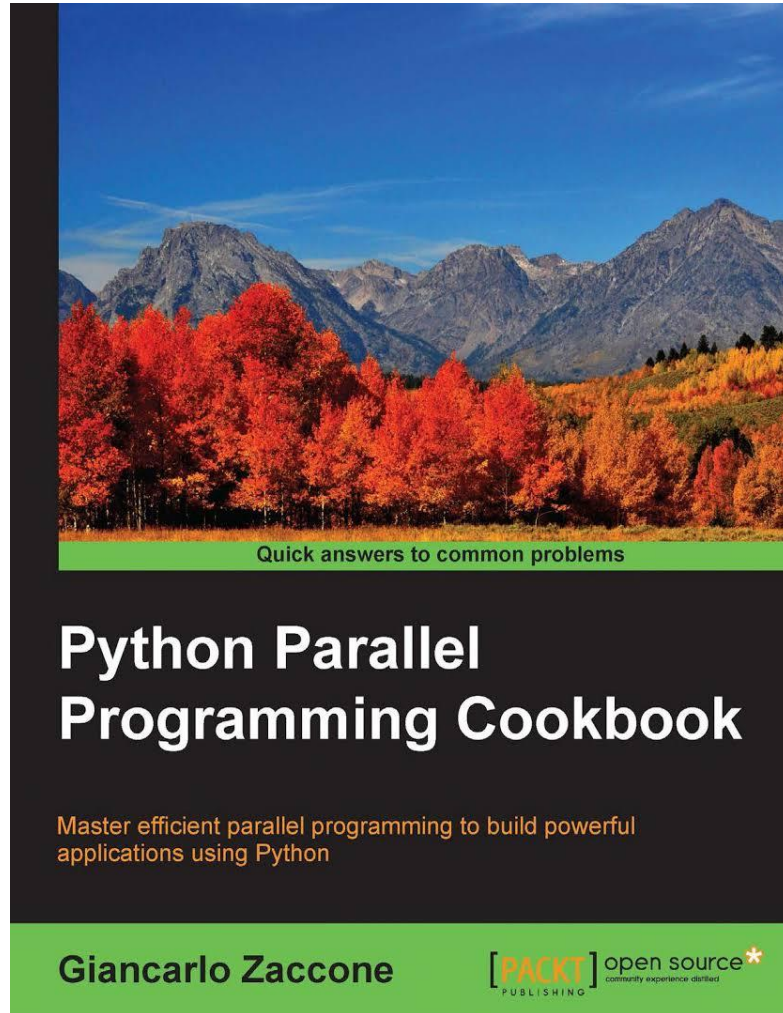
R.W. Hamming

Second Edition

- Numerical methods for scientists and engineers
- Richard Hamming



- 偏微分方程数值解法
- 第二版
- 陆金甫，关治



- ❑ **Python Parallel Programming Cookbook: Master efficient parallel programming to build powerful applications using Python**
- ❑ **Giancarlo Zaccone**

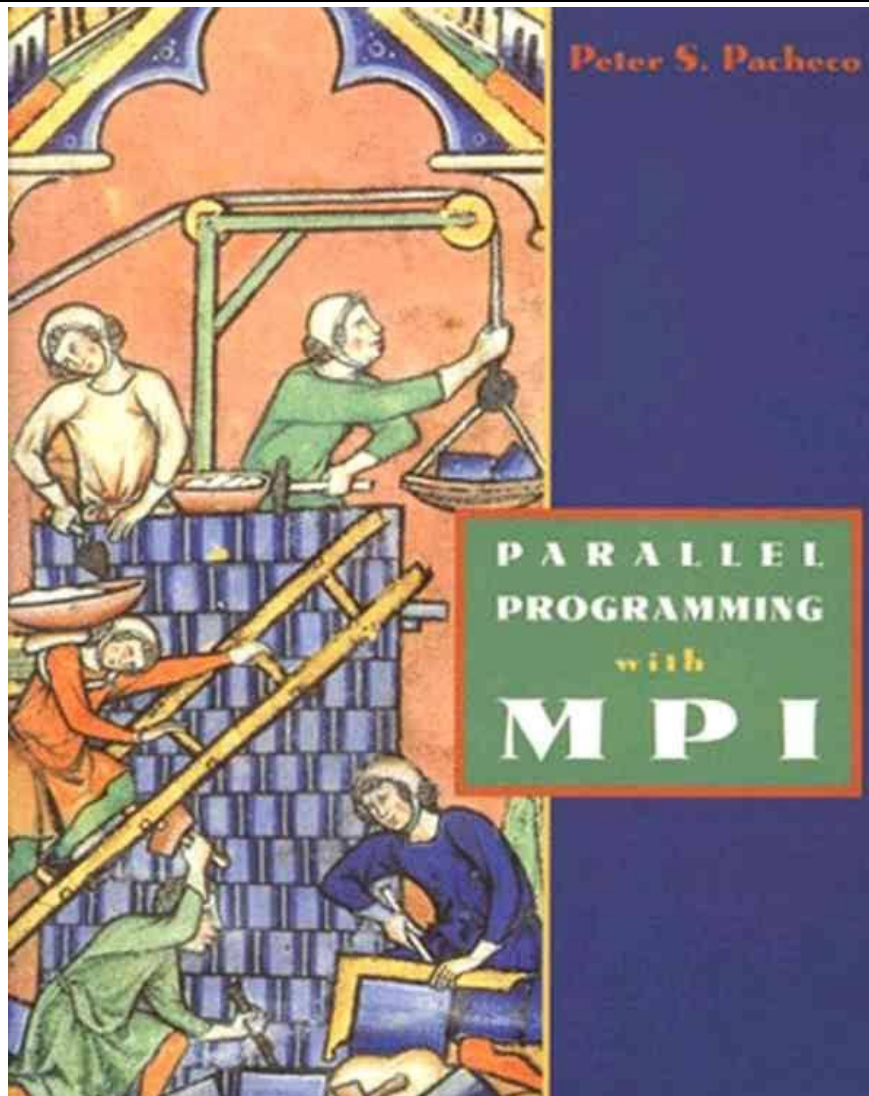
Hands-On GPU Programming with Python and CUDA

Explore high-performance parallel computing with CUDA

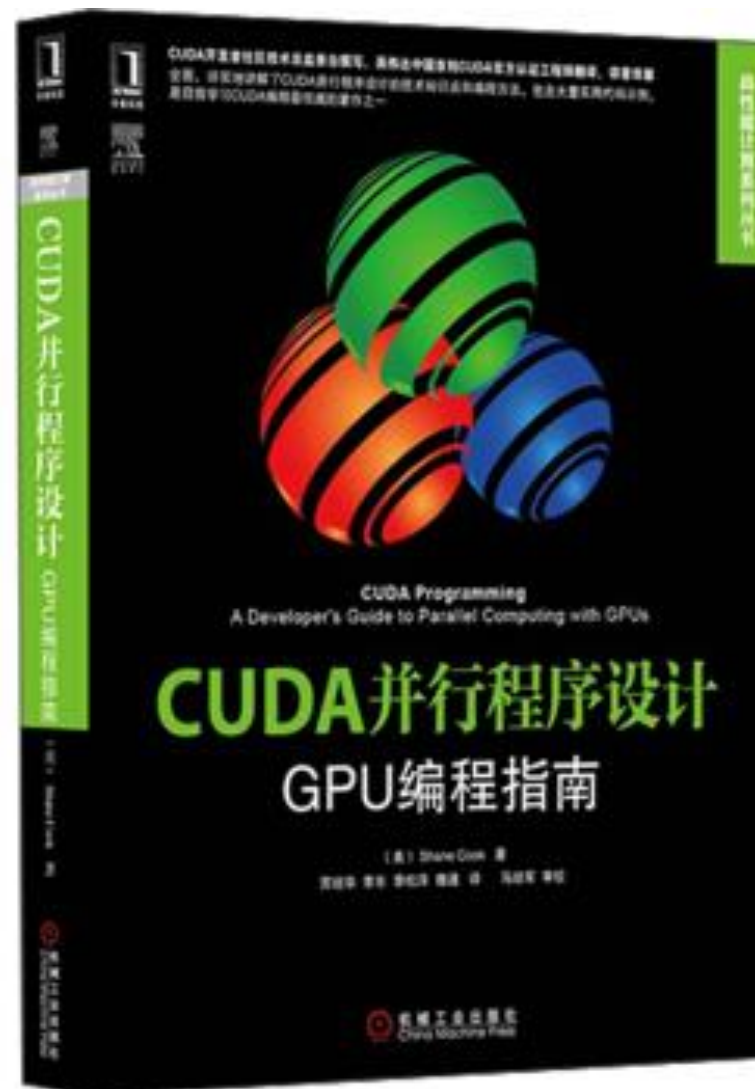
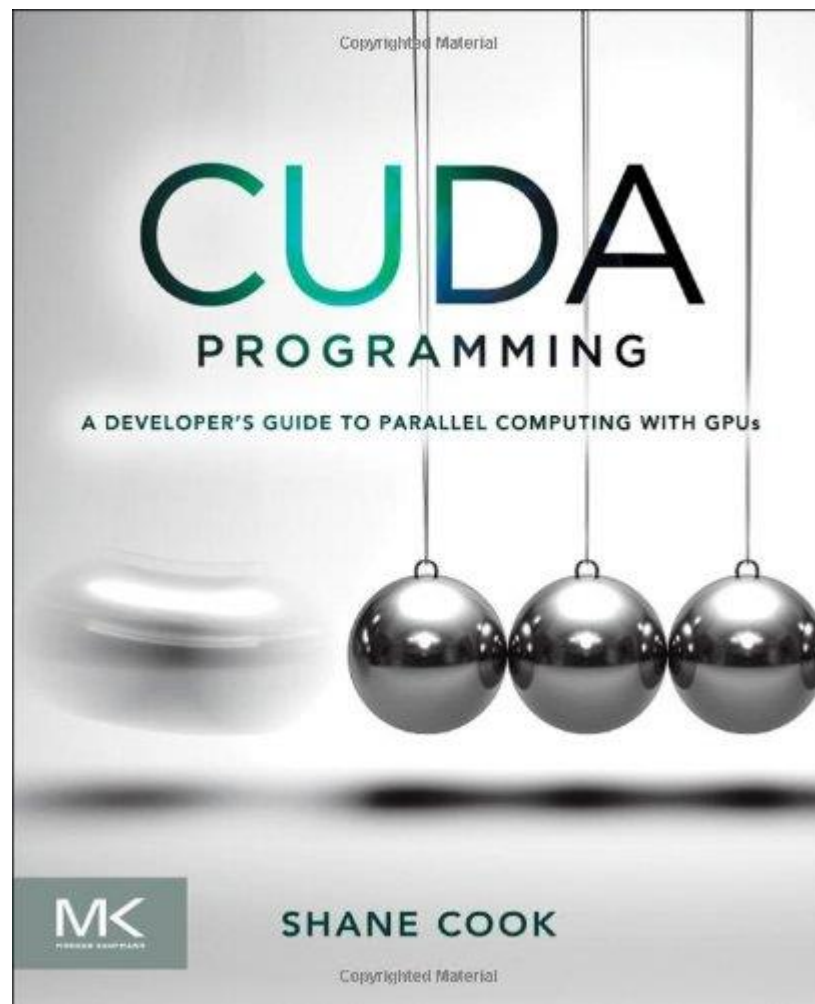


Packt>
www.packtpub.com

Dr. Brian Tuomanen



- Parallel programming with MPI
- Pacheco P. S.





- “高性能计算应用概览”
- 作者: 历军
- 出版社: 清华大学出版社
- 出版年: 2018-6-1
- 定价: 69.8
- 装帧: 平装
- ISBN: 9787302504726

徐明强 编著

Windows HPC Server: Step by Step

微软高性能计算服务器

人民邮电出版社
POSTS & TELECOM PRESS



通俗的语言，精妙的案例，每一个数据分析师都不容错过的最佳读物！

拥抱大数据

新常态下的数据分析典型案例

8个核心的数据分析主题

10个世界级企业的数据分析经验

24个大数据分析算法

横跨15个行业的典型数据分析实例

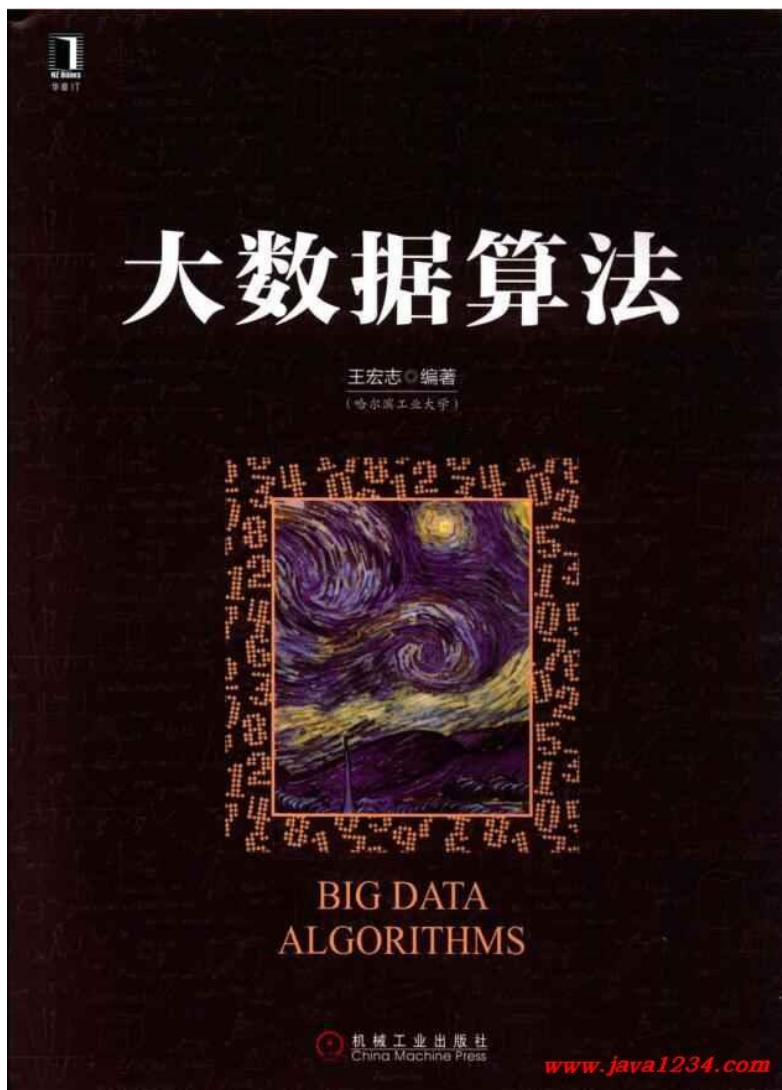
李倩星 王震◎著

//



SPM
商务数据情报
人民邮电出版社





- 大数据算法
- 王宏志
- 哈工大

Project reports @ Autumn 2019

❑ **BD Log analysis**

- Hadoop + Flume + Kafka (Zookeeper) + Spark +

❑ **OpenStack**

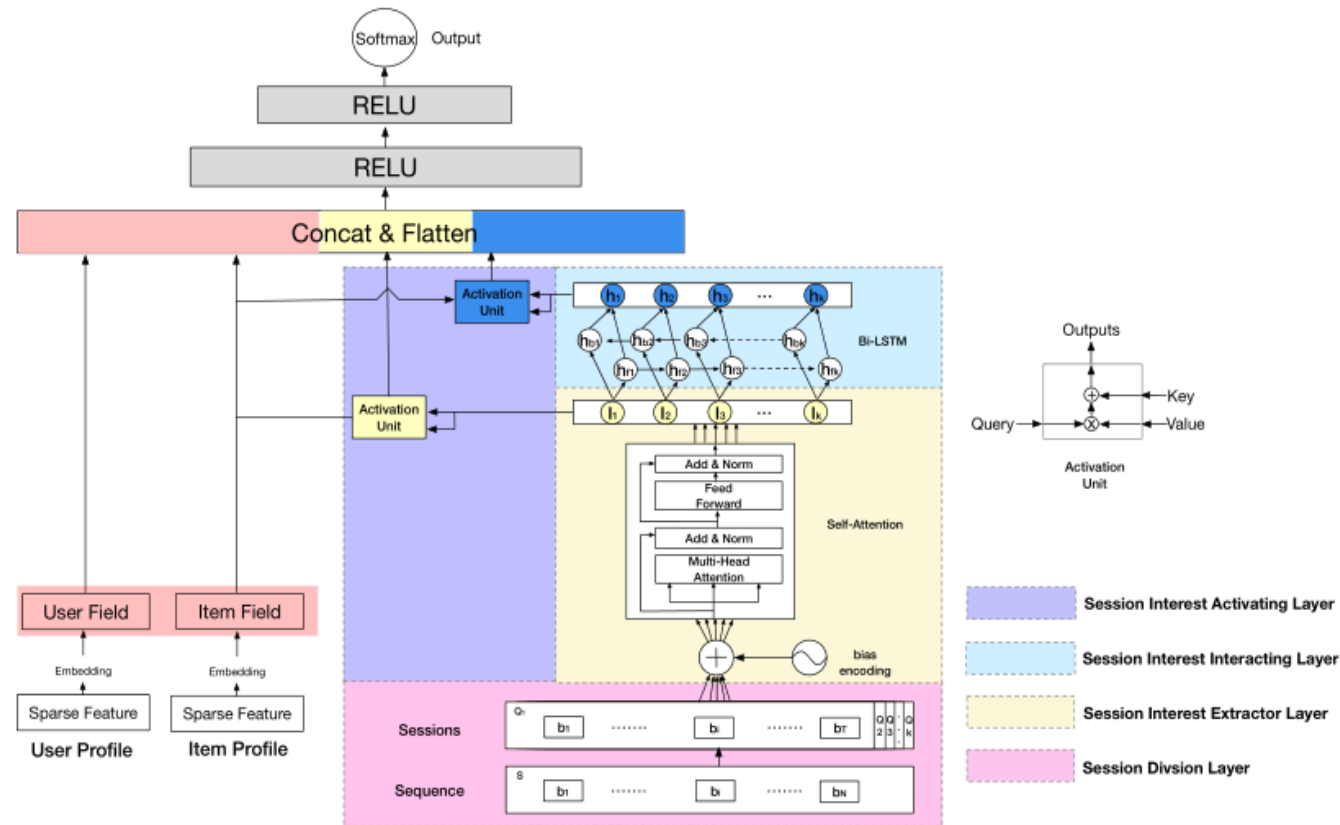
- Manage distributed resources to provide clouding services

❑ **Zookeeper**

- Ensure synchronization for data safety under distributed environment

Paper reports @ Autumn 2019

□ DSIN (Deep Session Interest Network) – LSTM (Long Short-Term Memory)



□ FP-Growth

Frequent Itemset Using FP-Growth (Example)

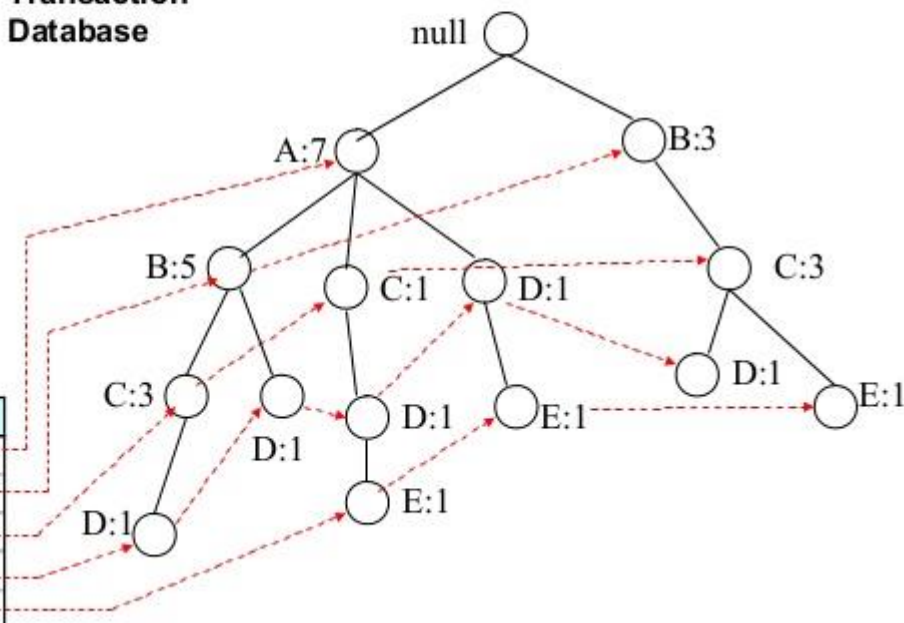
FP Growth Algorithm: FP Tree Mining

TID	Items
1	{A,B}
2	{B,C,D}
3	{A,C,D,E}
4	{A,D,E}
5	{A,B,C}
6	{A,B,C,D}
7	{B,C}
8	{A,B,C}
9	{A,B,D}
10	{B,C,E}

**Transaction
Database**

Header table

Item	Pointer
A	
B	
C	
D	
E	



Suggest for presentation

- ❑ **Always using a simple (vivid) example** to demonstrate the ideas, concepts, algorithms, applications
- ❑ **Distributed and Parallel**
 - 2 keywords you could follow to choose your projects
- ❑ **Try to dive into Design and Implementation**
 - Mechanism (Theory), Data Structure, Algorithms



□ 科学计算与企业级应用的并行优化 (高性能计算技术丛书)

□ 刘文志 著

□ 2015

□ 机械工业出版社

本系列的3本书相互之间有联系，也有其独立性：《**并行算法设计与性能优化**》介绍常见的串行代码优化方法和并行算法的设计；

《**并行编程方法与优化实践**》介绍常见的向量化和并行编程环境及一些实例；《**科学计算与企业级应用的并行优化**》则介绍领域相关的算法与应用的性能优化。



整体而言，本书分为如下几个部分：

·理论基础，本部分主要介绍并行软件和硬件基础，并行算法设计思想以及一些软件优化方法。主要包括第1章、第2章、第3章、第5章。

·代码优化，本部分主要介绍常见的串行代码优化手段（不包括向量化）。主要内容是第4章。

·并行算法设计考量，本部分主要介绍如何设计优良的并行算法并将算法映射到硬件上。主要内容是第6章、第7章、第8章、第9章、第11章和第12章。

·如何将现有的串行代码并行化，主要内容是第10章。

第1章 主要介绍并行化和向量化的相关概念，如并行和向量化的作用、为什么并行化和向量化、并行或向量化面临的现实困难。另外还介绍了一些不写代码也能够利用多核处理器性能的一些方法。

第2章 介绍了现代处理器的特性，如指令级并行、向量化并行、线程级并行、处理器缓存金字塔、虚拟存储器和NUMA（非一致内存访问）。

第3章 介绍了算法性能和程序性能的度量与分析。算法性能分析和度量的主要标准是时间复杂度、空间复杂度和笔者自己提出的实现复杂度。程序性能的度量标准主要有：时间、FLOPS、CPI、指令延迟和吞吐量。用来衡量优化一部分代码

第8章 介绍了并行和析了如何缓解某些缺点的

第9章 介绍了如何化uce、scan和流水线等并读者能够通过模式解决一

第10章 介绍了并行要注意的事项，最后以如

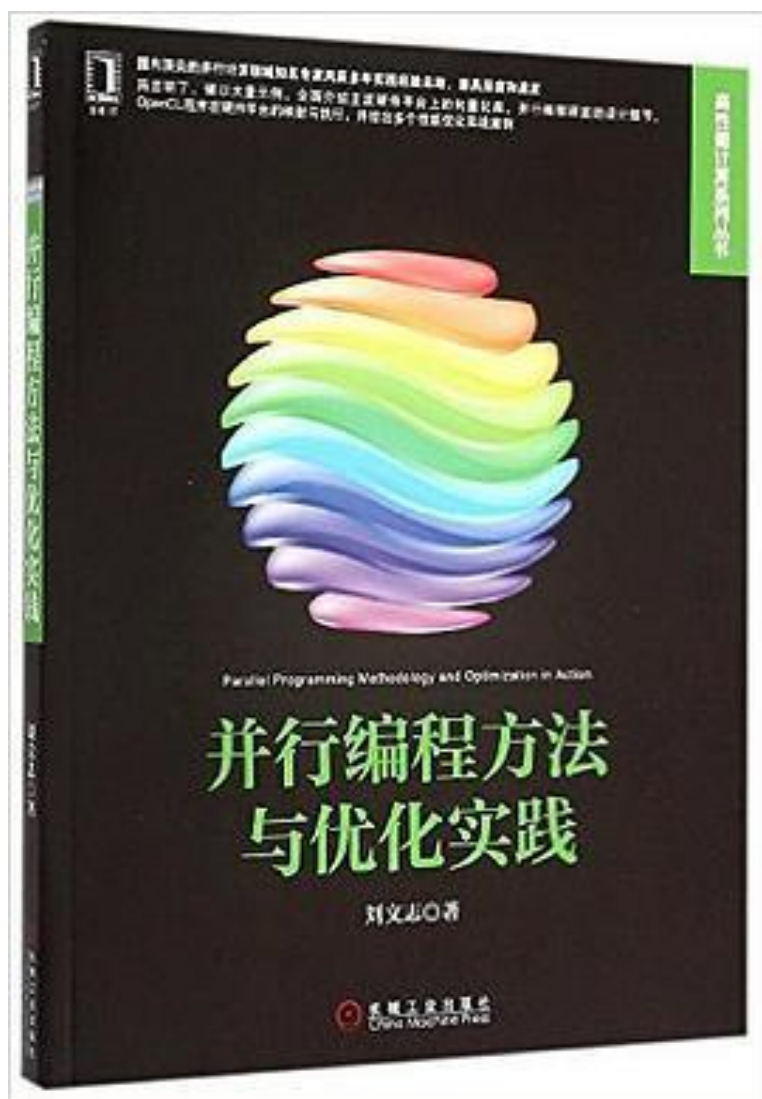
第11章 介绍了常见超级并行模式下如何划分阵乘运算。

第12章 给出了设计

附录A 介绍了整数

第2章 现代处理器特性





- 并行编程方法与优化实践
- 刘文志
- 机械工业出版社
- 2015-06



Introduction

- SIMD architectures can exploit significant data-level parallelism for:
 - matrix-oriented scientific computing
 - media-oriented image and sound processors
- SIMD is more energy efficient than MIMD
 - Only needs to fetch one instruction per data operation
 - Makes SIMD attractive for personal mobile devices
- SIMD allows programmer to continue to think sequentially

SIMD Variations

- Vector architectures
- SIMD extensions
 - MMX: multimedia extensions (1997)
 - SSE: streaming SIMD extensions
 - AVX: advanced vector extensions
- Graphics Processor Units (GPUs)
 - Considered as SIMD accelerators

Vector Supercomputers

- In 70-80s, Supercomputer \equiv Vector Supercomputer
- Definition of supercomputer



Evolution of Intel Vector Instructions

- **MMX (1996, Pentium)**
 - *CPU-based MPEG decoding*
 - Integers only, 64-bit divided into 2 x 32 to 8 x 8
 - Phased out with SSE4
- **SSE (1999, Pentium III)**
 - *CPU-based 3D graphics*
 - 4-way float operations, single precision
 - 8 new 128 bit Register, 100+ instructions
- **SSE2 (2001, Pentium 4)**
 - *High-performance computing*
 - Adds 2-way float ops, double-precision; same registers as 4-way single-precision
 - Integer SSE instructions make MMX obsolete
- **SSE3 (2004, Pentium 4E Prescott)**
 - *Scientific computing*
 - New 2-way and 4-way vector instructions for complex arithmetic
- **SSSE3 (2006, Core Duo)**
 - Minor advancement over SSE3
- **SSE4 (2007, Core2 Duo Penryn)**
 - *Modern codecs, cryptography*
 - New integer instructions
 - Better support for unaligned data, super shuffle engine



基本概念

- ADDVV.D: add two vectors
- ADDVS.D: add vector to a scalar
- LV/SV: vector load and vector store from address
- Vector code example:

# C code	# Scalar Code	# Vector Code
for (i=0; i<64; i++)	LI R4, 64	LI VLR, 64
C[i] = A[i] + B[i];	loop:	LV V1, R1
	L.D F0, 0(R1)	LV V2, R2
	L.D F2, 0(R2)	ADDV.D V3, V1, V2
	ADD.D F4, F2, F0	SV V3, R3
	S.D F4, 0(R3)	
	DADDIU R1, 8	
	DADDIU R2, 8	
	DADDIU R3, 8	
	DSUBIU R4, 1	
	BNEZ R4, loop	

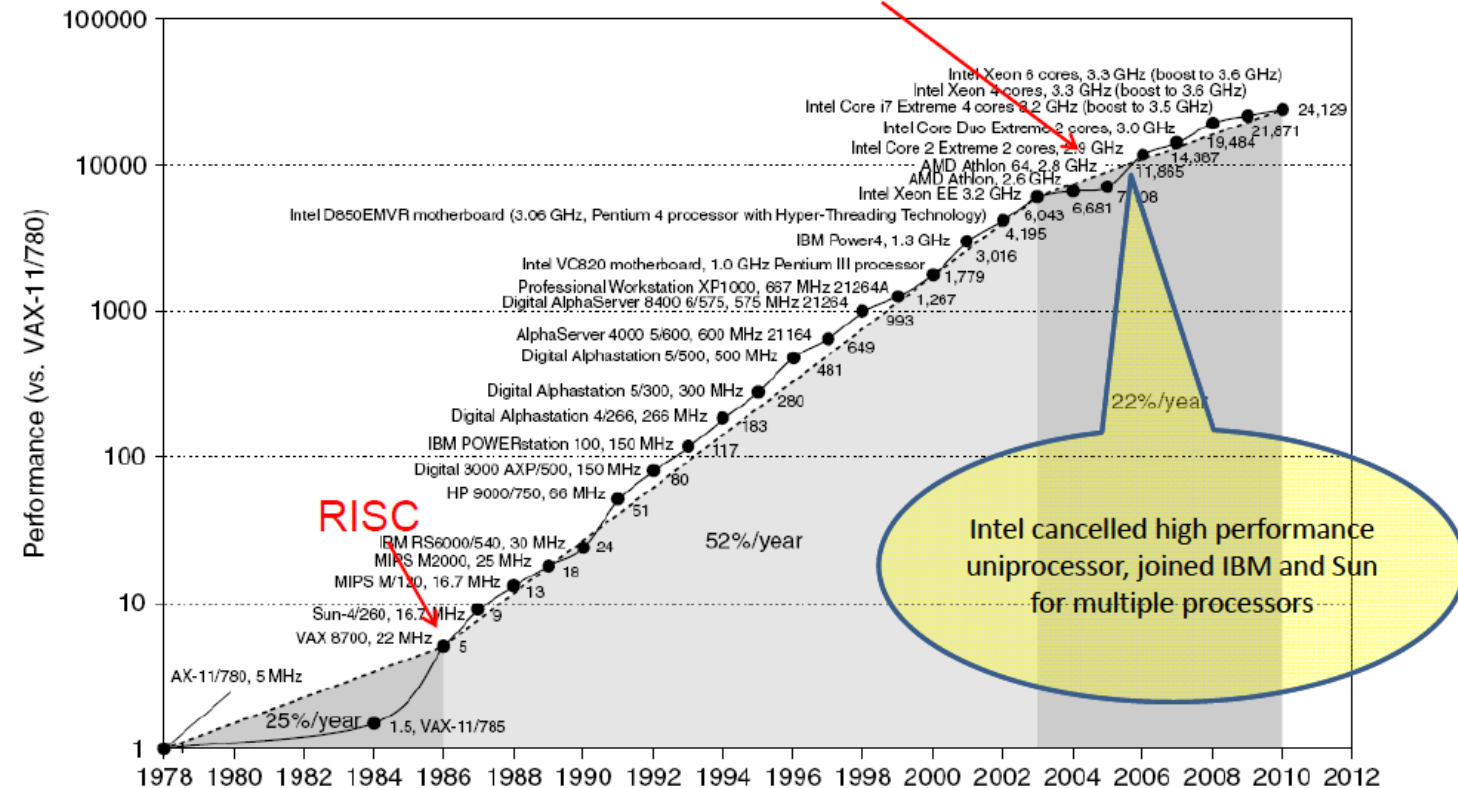
向量处理器系统 (Vector Processor System, VPS), 是采用先行控制和重叠操作技术、运算流水线、交叉访问的并行存
行时还不能充分发挥并行处理潜力。向量运算很适合于流水线计
大程度上克服通常流水线计算机中指令处理量太大、存储访问不
结构的潜力, 显著提高运算速度。

向量运算是一种较简单的并行计算，适用面很广，机器实现迅速发展。TI ASC(1972年)和CDC STAR-100 (1973年) 是世界上第一台巨型机，其中大多数是向量机。中国于1983年研制成功的每秒

<https://baike.baidu.com/item/白>

Single Processor Performance

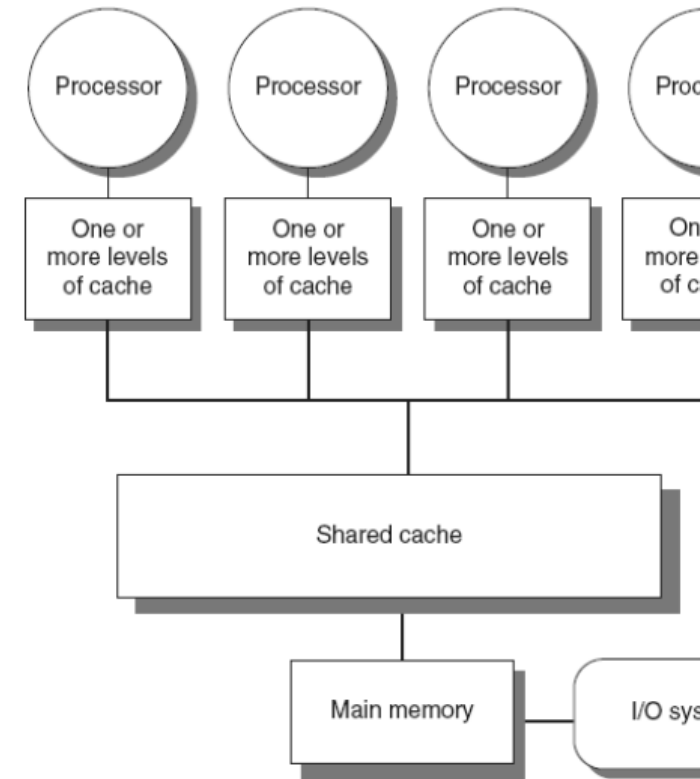
Move to multi-processor



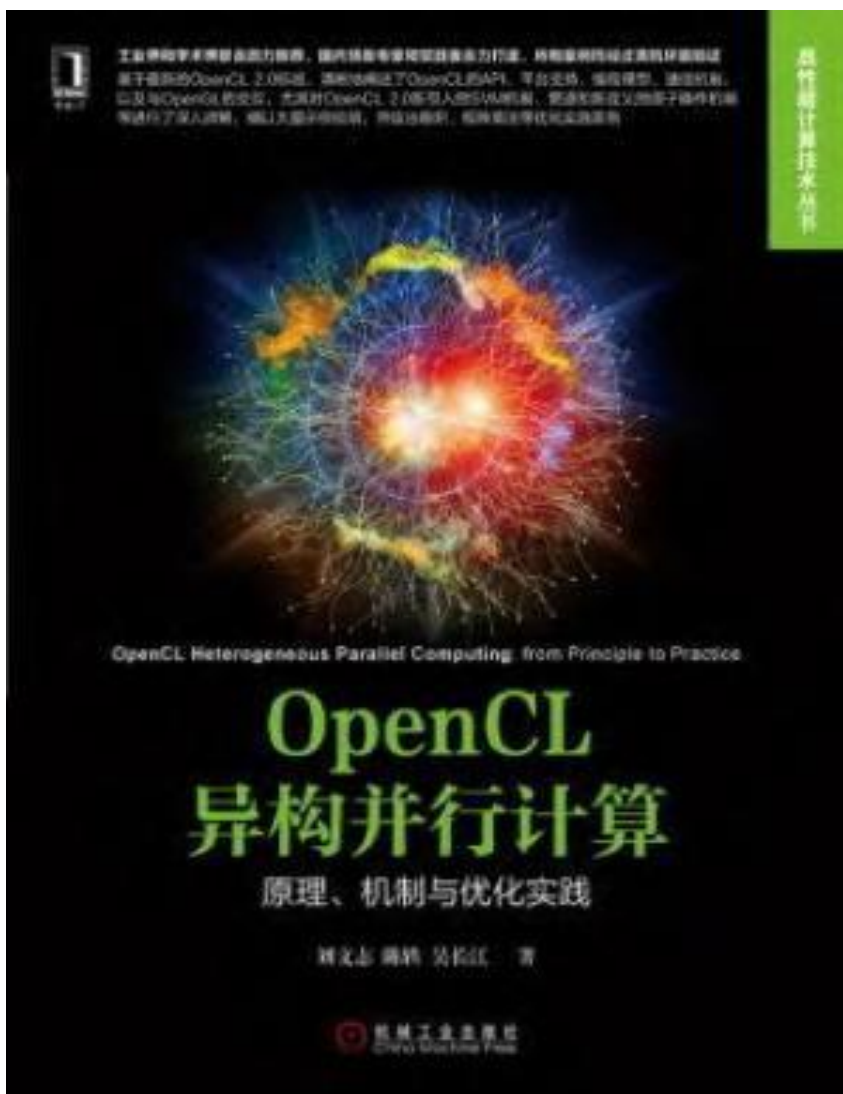
Introduction

Shared-Memory Multiprocessor

- SMP, Symmetric multiprocessing

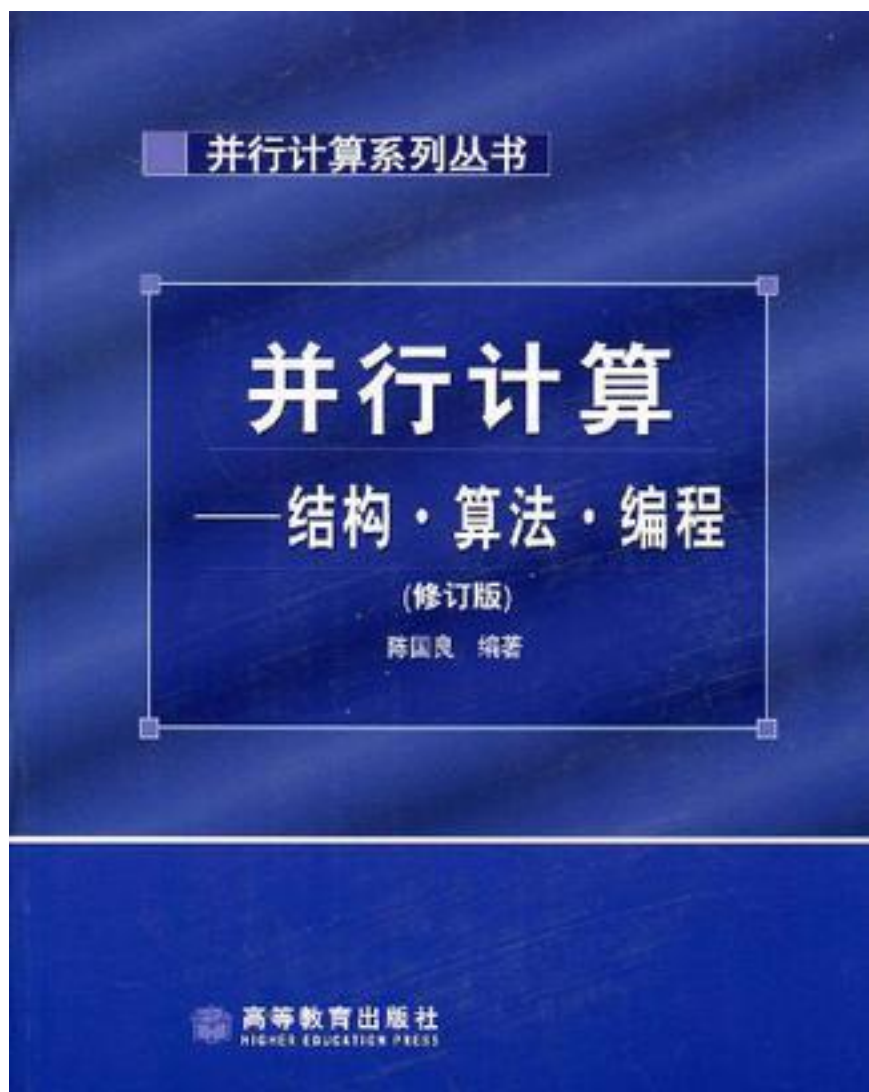


2 Models for



□ OpenCL异构并行计算：原理、机制与优化实践

□ 刘文志 陈轶 吴长江



- 并行计算-结构·算法·编程
- 陈国良编著
- 北京:高等教育出版社
- 2011



普通高等教育“十一五”国家级规划教材

并行计算系列丛书

并行算法的设计与分析

(第3版)

陈国良 编著

高等教育出版社

□ 并行算法的设计与分析

□ 陈国良

□ 2009

□ 高等教育出版社

