

清华大学出版社

分布式数据库技术

Technology of Distributed Database

第二版

黄海 王志英 周伟红 李斌 编

清华大学出版社

00105737

并行与分布计算技术丛书

TP311.13

109



分布式数据库技术

Technology of Distributed Database

贾 焰 王志英 著
韩伟红 李 霖

国防工业出版社



C0508615

图书在版编目(CIP)数据

分布式数据库技术/贾焰等著. —北京:国防工业出版社,2000.7
(并行与分布计算技术丛书)
ISBN 7-118-02205-5

I . 分… II . 贾… III . 分布式数据库 - 数据库系统
IV . TP311.13

中国版本图书馆 CIP 数据核字(1999)第 52611 号

J6476/34

国防工业出版社出版发行
(北京市海淀区紫竹院南路 23 号)
(邮政编码 100044)
北京怀柔新华印刷厂印刷
新华书店经售

*
开本 787×1092 1/16 印张 16 1/4 350 千字
2000 年 7 月第 1 版 2000 年 7 月北京第 1 次印刷
印数:1—2000 册 定价:30.00 元

(本书如有印装错误,我社负责调换)

致 读 者

本书由国防科技图书出版基金资助出版。

国防科技图书出版工作是国防科技事业的一个重要方面。优秀的国防科技图书既是国防科技成果的一部分,又是国防科技水平的重要标志。为了促进国防科技事业的发展,加强社会主义物质文明和精神文明建设,培养优秀科技人才,确保国防科技优秀图书的出版,国防科工委于1988年初决定每年拨出专款,设立国防科技图书出版基金,成立评审委员会,扶持、审定出版国防科技优秀图书。

国防科技图书出版基金资助的对象是:

1. 学术水平高,内容有创见,在学科上居领先地位的基础科学理论图书;在工程技术理论方面有突破的应用科学专著。
2. 学术思想新颖,内容具体、实用,对国防科技发展具有较大推动作用的专著;密切结合科技现代化和国防现代化需要的高新技术内容的专著。
3. 有重要发展前景和有重大开拓使用价值,密切结合科技现代化和国防现代化需要的新工艺、新材料内容的科技图书。
4. 填补目前我国科技领域空白的薄弱学科和边缘学科的科技图书。
5. 特别有价值的科技论文集、译著等。

国防科技图书出版基金评审委员会在国防科工委的领导下开展工作,负责掌握出版基金的使用方向,评审受理的图书选题,决定资助的图书选题和资助金额,以及决定中断或取消资助等。经评审给予资助的图书,由国防工业出版社列选出版。

国防科技事业已经取得了举世瞩目的成就。国防科技图书承担着记载和弘扬这些成就,积累和传播科技知识的使命。在改革开放的新形势下,国防科工委率先设立出版基金,扶持出版科技图书,这是一项具有深远意义的创举。此举势必促使国防科技图书的出版随着国防科技事业的发展更加兴旺。

设立出版基金是一件新生事物,是对出版工作的一项改革。因而,评审工作需要不断地摸索、认真地总结和及时地改进,这样,才能使有限的基金发挥出巨大的效能。评审工作更需要国防科技工业战线广大科技工作者、专家、教授,以及社会各界朋友的热情支持。

让我们携起手来,为祖国昌盛、科技腾飞、出版繁荣而共同奋斗!

国防科技图书出版基金
评审委员会

国防科技图书出版基金 第三届评审委员会组成人员

名誉主任委员 怀国模

主任委员 黄 宁

副主任委员 殷鹤龄 高景德 陈芳允 曾 锋

秘书 长 崔士义

委 员 于景元 王小谟 尤子平 冯允成

(以姓氏笔划为序) 刘 仁 朱森元 朵英贤 宋家树

杨星豪 吴有生 何庆芝 何国伟

何新贵 张立同 张汝果 张均武

张涵信 陈火旺 范学虹 柯有安

侯正明 莫梧生 崔尔杰

并行与分布计算技术丛书编委会名单

主 编 卢锡城

副 主 编 周兴铭 汪成为 李国杰

主 审 陈火旺

编委委员 施伯乐 康继昌 尤晋元 康立山

沈隆均 李晓梅 朱传奇 王志英

杨学军 杨晓东 李思昆 王怀民

常务秘书 肖 政

总序

并行与分布计算技术是实现高性能计算的重要技术途径。高性能计算机技术是现代科学研究、工程技术开发和大规模数据处理的关键支撑技术。没有高性能计算机，大量复杂问题的计算和事务处理就无法在合理的时间内完成。利用高性能计算机还可以解决一些仅靠理论及实验方法无法解决的问题，分析处理靠传统技术无法应付的海量数据，例如核爆炸模拟、宇宙的形成及演变过程研究、中长期天气预报、石油地质勘探、数字地球和大规模事务处理等。人类对高性能计算能力的需求永无止境。计算速度、存储容量、通信带宽是衡量高性能计算能力的重要技术指标。一些推动人类文明和社会信息化的重大挑战性问题需要百万亿次每秒、千万亿次每秒以上的计算速度，需要万亿字节以上的存储容量和万亿位每秒以上的通信能力。人们已经认识到，高性能计算与网络通信技术是战略技术，是科技创新体系的基础技术，是反映一个国家综合实力的重要标志之一。

由于种种因素的制约及使用计算机方式的改变，在单台计算机系统的性能与功能难以满足应用需求的情况下，由多个计算节点经专门设计的互联网络紧密耦合而构成的大规模并行处理计算机系统，以及由多个自主的高性能计算节点经计算机网络连接组成的分布处理环境，将成为高性能计算机领域的两大重要研究方向。并行计算和分布计算是既有区别又联系密切的两个概念，前者重在发掘计算过程中的并行性，后者则重在有效组织管理各类异构资源，挖掘功能上的并发性。随着基于先进计算机网络的分布并行计算概念的发展，并行计算与分布计算在很多应用领域正面临共同的追求目标和技术挑战，如高效、实用的计算模型、计算方法、并行机制等。

90年代以来，随着高性能计算和网络计算技术的普及，并行与分布计算技术正渗透到现代社会的各个领域，用户通过高性能网络使用各类计算资源提供的服务已成为信息化社会中计算机应用的一种重要形式。事实上，各类基于并行与分布计算的应用系统正在工业、交通、金融、科研、政府、国防等部门支撑着现代社会的高效运行。

20多年来，我国科技人员依靠自己的力量，勇于开拓，奋勇拼搏，研制成功了多种型号序列的高性能计算机系统。国防科学技术大学计算机学院是我国研制高性能计算机系统的重要基地，多种巨型计算机的研制成功及推广应用，打破了国外对我国的技术封锁，缩短了我国同发达国家技术水平上的差距，为推动我国高性能计算机及并行与分布计算技术的发展作出了重要贡献。

为促进并行与分布计算技术领域的研究，国防工业出版社组织国防科技大学计算机学院有关专家、教授撰著了本套丛书。本丛书以并行与分布计算机系统组成为纲，结合作者多年科研与工程实践以及当前研究的热点问题，涵盖了计算机系统结构、计算机网络、系统软件、应用软件、计算方法等多方面内容。本丛书由以下9本专著组成：《并行计算机体系结构技术》论述并行计算机研究和设计的理论及工程问题；《先进计算机网络技术》论

述高性能网络计算的各种关键技术;《并行操作系统原理与技术》论述并行操作系统的机制、基本原理和主要实现技术;《并行编译方法》论述并行编译系统理论、编译器的设计方法;《分布计算——网络化软件新技术》论述分布计算技术的基本概念和关键技术;《分布式数据库技术》论述统一逻辑分布式数据库技术;《数字系统并行 CAD 技术》论述数字系统并行 CAD 的理论和技术;《可扩展并行算法的设计与分析》论述可扩展并行算法的设计与分析的理论和方法;《并行与分布式可视化技术与应用》论述并行与分布式可视化的各种技术及其应用。本丛书既介绍了当前国际上该领域的最新技术发展,又汇集了作者多年的研究成果和工程经验。丛书注重可读性,适合从事该领域工作和学习的科技人员、高等院校高年级学生及研究生作为工作和学习的参考书。

本丛书被列入“九五国家重点图书选题规划”,并获得国防科技出版基金的资助。愿本丛书的出版能为并行与分布计算技术研究园地增添一朵花絮,为以后的研究工作提供有价值的参考。因时间和能力所限,书中不足之处,恳请读者指正。

并行与分布计算技术丛书编委会

前　　言

数据库技术从 60 年代中期产生到今天仅仅 30 年,已经经历了第 1 代的网络、层次数据库,第 2 代关系数据库系统。其应用领域也广泛涉及到 CAD/CAM、CIM、CASE、OIS(办公信息系统)、GIS(地理信息系统)、知识库系统和实时系统等。其发展速度之快,使用范围之广是其他技术远远不及的。随着网络技术的飞速发展,以分布式为主要特征的数据库系统正成为该领域的一个新的研究热点和主流方向。

10 年来,作者长期承担大学本科生和研究生的数据库相关课程的授课任务,分析和研究了国外多种的教材和大量的技术资料,为本书的编写奠定了基础。受国家“863”高技术计划、国家“九五”和“十五”预研计划的资助,作者作为课题负责人和主要完成者长期进行数据库系统,特别是分布式数据库系统的研究工作,取得了一系列的科研成果,为本书的撰写积累了素材。传统的分布式数据库专著重点论述的是统一逻辑的分布式数据库技术,随着数据库和网络技术的飞速发展,多数据库、移动数据库和 Web 数据库已成为分布式数据库系统研究的新潮流和新热点。国内已出版了一些与分布式数据库相关的教材和著作,但较全面地反映 90 年代分布式数据库新技术的专著还甚为鲜见。作者编写本书的主要目的是系统分析和总结该领域的最新研究技术和方法,并将它们介绍给读者。

本书共分 10 章。第 1 章概论,简介了分布式数据库系统问题的背景,基本概念和理论,以及发展的历史和方向。第 2 章数据处理——分布与转换,主要论述在分布式环境中数据的划分、分布,以及异构数据库之间的查询转换技术。第 3 章分布式查询的优化,以集中式数据库系统中的查询优化技术为基础,重点论述分布式数据库系统中的查询优化技术,特别是 JOIN 操作的执行方法。第 4 章并发控制,首先介绍了并发控制的基本概念,包括事务和事务冲突,调度表和可串行化等,然后讨论了几种常见并发控制技术。第 5 章恢复,首先介绍了引起系统失败的各种原因,其次讨论了本地恢复协议,最后重点讨论了分布式环境中的恢复协议。第 6 章完整性和安全性,首先介绍集中式系统的完整性和安全技术,并在此基础上重点论述分布式系统中的完整性和安全技术。第 7 章多数据库技术,主要论述多数据库系统的设计原理、体系结构、异构模式消解技术、查询处理技术和事务管理技术。第 8 章多数据库系统实例,主要介绍了 3 个著名的多数据库系统,它们是 UniSQL/M、EDA/SQL 和 Pegasus。第 9 章移动数据库技术,主要论述移动数据库的问题背景,移动相关的技术和移动数据库技术。第 10 章 Web 数据库技术,主要论述 Web 数据库的问题背景,基本的 Web 数据库技术,以及各大数据库公司的 Web 数据库解决方案。

本书的第 1 章由王志英编写,第 2 章至第 6 章由贾焰编写,第 7 章和第 8 章由韩伟红编写,第 9 章和第 10 章由李霖编写,全书由贾焰统一审定。

本书在编著工作中得到了多方面的支持,特别是各级领导和机关的支持,《并行与分

布计算技术》丛书编委会的具体指导和帮助,还有刘艳春、郭歌、夏戈明、杨文波和张文强同学的辛勤工作,在此致以深深的谢意!

由于数据库技术的发展速度快、涉及的知识面广,虽然我们竭尽了全力,但疵漏之处在所难免,欢迎大家批评、指正。

内 容 简 介

本书系统全面地介绍了分布式数据库系统的基本原理和实现技术,充分反映了该领域的最新研究成果。本书的第1章概述了分布式数据库系统问题的背景,基本概念和理论,以及发展的历史和方向;第2~6章论述了统一逻辑分布式数据库技术,主要包括数据分布、查询优化、并发控制、系统恢复和完整性及安全等技术;第7、8章论述了多数据库技术,主要包括多数据库系统的关键技术,以及典型的实用系统;第9章论述移动数据库技术;第10章论述Web数据库技术。

本书可作为相关领域科研工作者的参考书,也可用作计算机和信息技术领域研究生和高年级本科生教材。

This book discusses the basic theory and implementation technology of distributed database, and reflected the newest researches. Chapter 1 summarizes the basic conception, basic theory, developing history and tendency of distributed database. Chapter 2 to Chapter 6 deal with distributed database technology, mainly including data distribution, query optimization, concurrency control, system recovery, integrity and security. Chapter 7 and Chapter 8 deal with technology of multidatabase system, including the key technology of multidatabase, and the typical utility system. Chapter 9 deal with the technology of mobile database. Chapter 10 deal with the technology of Web database.

This book is for the peoples which work in related fields, and also for graduate students or high grade undergraduate students in computer and information technology fields.

目 录

第1章 概论	1
1.1 问题背景	1
1.1.1 数据分布的需求	1
1.1.2 异构环境中数据集成的需求	2
1.1.3 信息系统集成的需求	3
1.2 数据库技术的回顾	3
1.2.1 数据资源	4
1.2.2 数据库管理系统的体系结构	4
1.2.3 数据库管理系统的组成	5
1.2.4 关系数据库管理系统	5
1.2.5 层次和网络数据库管理系统	9
1.2.6 数据库的设计和规范化	11
1.2.7 查询语言	14
1.3 计算机网络	16
1.3.1 网络的体系结构	16
1.3.2 ISO/OSI 参考标准	17
1.4 分布式数据库系统	18
1.4.1 基于体系结构的分类	18
1.4.2 其他分类	20
第2章 数据处理——分布与转换	22
2.1 数据分布问题	22
2.2 数据分布的例子	23
2.2.1 一个关系的情况	23
2.2.2 多个关系的情况	24
2.3 数据分布问题的语义方法	26
2.4 一种文件分布方法	27
2.5 异构数据库系统的集成	31
2.6 全局数据模式	32
2.7 将网络数据模式转换成关系数据模式	34
2.8 在网络数据库上执行关系查询	37
第3章 分布式查询的优化	39
3.1 查询优化的重要性	40
3.1.1 查询优化的基本方法	40
3.1.2 查询执行的各种途径	41

3.2 等价转换	44
3.3 集中式系统中存取规划的生成和选择.....	45
3.4 联结操作的执行方法.....	48
3.4.1 半联结操作	49
3.4.2 非半联结操作	52
3.5 相关问题的讨论	56
第4章 并发控制	60
4.1 事务	60
4.1.1 基本概念	60
4.1.2 分布式事务	62
4.2 并发事务的冲突	63
4.2.1 丢失更新问题	63
4.2.2 破坏完整性约束问题	64
4.2.3 不一致读问题	65
4.3 调度表与串行性	66
4.3.1 集中式系统的串行性问题	66
4.3.2 分布式系统的可串性问题	68
4.3.3 分布式事务处理	69
4.4 并发控制技术	70
4.4.1 锁方法	70
4.4.2 死锁	72
4.4.3 时戳法	77
4.4.4 基本时戳方法	78
4.4.5 保守时戳方法	80
4.4.6 乐观方法	81
4.5 面向应用的方法和准可串性	84
第5章 恢复	87
5.1 基本概念	87
5.1.1 事务和恢复	87
5.1.2 日志文件	88
5.1.3 检查点	90
5.1.4 数据库的更新问题	91
5.2 引发失败的原因	91
5.2.1 局部事务失败	92
5.2.2 站点失败	92
5.2.3 介质失败	93
5.2.4 网络失败	94
5.3 集中式恢复协议	96
5.3.1 Undo/redo	97
5.3.2 Undo/no-redo	99
5.3.3 No-undo/redo	100

5.3.4 No-undo/no-redo	101
5.4 分布式恢复协议	102
5.4.1 二阶段提交(2PC).....	102
5.4.2 三阶段提交(3PC).....	107
第6章 完整性与安全性	112
6.1 集中式数据库的完整性	112
6.1.1 完整性概念	112
6.1.2 完整性约束	113
6.1.3 关系约束	113
6.1.4 域约束	114
6.1.5 参照完整性约束	115
6.1.6 显式约束	116
6.1.7 静态和动态约束	117
6.2 集中式DBMS的安全	117
6.2.1 数据库的安全问题	117
6.2.2 访问控制策略	119
6.2.3 多级安全	119
6.2.4 SQL中的安全机制	120
6.2.5 统计数据库的安全问题	122
6.3 分布式DBMS的安全	123
6.3.1 认证和授权	123
6.3.2 授权规则的分布控制	123
6.3.3 加密	124
6.3.4 全局视图机制	124
第7章 多数据库系统技术	126
7.1 MDBS的设计原则及其体系结构	126
7.2 异构模式消解	127
7.2.1 异构的LDB模式的例子	128
7.2.2 模式冲突分类	130
7.2.3 冲突消解	131
7.2.4 重命名实体和属性	133
7.2.5 一致化表示	133
7.2.6 属性一致化	135
7.2.7 水平连接	136
7.2.8 垂直连接	139
7.2.9 混合连接及方法冲突	141
7.3 多库系统中的查询处理	142
7.3.1 查询分解	142
7.3.2 查询转换	145
7.3.3 全局查询优化	146
7.3.4 进一步的工作	148
7.4 多库系统中的事务管理	148

7.4.1 全局事务管理所面临的问题	149
7.4.2 全局可串行化	151
7.4.3 原子性和持久性	152
第8章 多数据库系统实例	154
8.1 UniSQL/M 系统	154
8.1.1 LDB 模式的例子	155
8.1.2 GDB 实体的定义	156
8.1.3 面向对象技术的应用	159
8.1.4 模式变化	163
8.2 EDA/SQL 系统	169
8.2.1 操作过程概述	170
8.2.2 系统结构	170
8.2.3 客户/服务器操作	172
8.3 Pegasus:一个异构的信息管理系统	174
8.3.1 数据模型及语言	174
8.3.2 系统概述	175
8.3.3 外部数据的引入	177
8.3.4 数据一致化	178
8.3.5 查询处理	179
第9章 移动数据库技术	180
9.1 移动数据库技术概述	180
9.1.1 无线网络技术简介	180
9.1.2 移动计算环境的体系结构	181
9.1.3 典型移动数据库应用	182
9.1.4 移动数据库系统分类	184
9.1.5 移动数据管理与分布数据管理的关系	184
9.2 移动数据库的关键技术	185
9.2.1 复制与缓存	186
9.2.2 数据广播	192
9.2.3 移动查询处理	198
9.2.4 移动事务处理	201
9.2.5 Agent 技术	204
9.2.6 其他技术	208
第10章 Web 数据库技术	210
10.1 基本概念	210
10.1.1 World Wide Web	210
10.1.2 客户/服务器数据库系统	211
10.1.3 Web 数据库网关	213
10.2 Web 数据库网关的实现技术	214
10.2.1 Web 数据库网关分类	214
10.2.2 CGI 执行程序	214
10.2.3 CGI 应用服务器	216

10.2.4 使用服务器 API	217
10.2.5 专有服务器	218
10.2.6 外部查看器	218
10.2.7 浏览器扩展	219
10.3 Web 数据库产品简介	220
10.3.1 Sybase web.sql	220
10.3.2 Sybase jConnect	225
10.3.3 IBM DB2 WWW 连接	230
10.3.4 Oracle WebServer	231
参考文献	232

Contents

Chapter 1 Introduction	1
1.1 Background	1
1.1.1 The Needing of Data Distribute	1
1.1.2 The Needing of Data Integration in Heterogeneous Environment	2
1.1.3 The Needing of Information System Integration	3
1.2 Overview of Database Technology	3
1.2.1 Data Resource	4
1.2.2 DBMS Architecture	4
1.2.3 Components of a DBMS	5
1.2.4 Relation Database Management System	5
1.2.5 Hierarchical and Network DBMSs	9
1.2.6 Database Design and Normalization	11
1.2.7 Query Languages	14
1.3 Computer Networks	16
1.3.1 Network Architecture	16
1.3.2 ISO/OSI Reference Model	17
1.4 Distributed Database System	18
1.4.1 Classification of DDBMSs by Architecture	18
1.4.2 Other Classification	20
Chapter 2 Data Handling —Distribution and Transformation	22
2.1 Data Distribution Problem	22
2.2 Some Examples of Data Distribution	23
2.2.1 Single – Relation Case	23
2.2.2 Multi – Relation Case	24
2.3 A Semantic Approach to Data Distribution Problem	26
2.4 A Distribution Approach to File	27
2.5 The Integration of Heterogeneous Database Systems	31
2.6 The Global Data Schema	32
2.7 Translation From Network Data Schema to Relation Data Schema	34
2.8 How to Execute Relation Query in Network Database	37
Chapter 3 Distributed Query Optimization	39
3.1 The Importance of Query Optimization	40
3.1.1 Basic types of Query Optimization	40
3.1.2 Variations in Ways of Executing Query	41