

周龙骧 等 著

# 分布式数据库 管理系统 实现技术

数  
据  
库  
从  
书



科学出版社

TP36.3

11.82-2

413393

数据 库 丛 书

分 布 式 数据 库  
管 理 系 统 实 现 技 术

周 龙 骚 等 著



科 学 出 版 社

1998

## 内 容 简 介

本书是在作者十余年来对分布式数据库系统进行研究和开发所取得的成功的实践经验的基础上撰写的。书中系统论述和分析了国内外一些著名的先驱分布式数据库管理系统的工作原理、设计思想、实现方法以及若干典型的概念、方法、算法和技巧。

本书内容共分十一章。主要内容包括：分布式数据库管理系统的体系结构、编程语言、编程语言编译器的设计和实现、编程语言涉网的全局编译和优化、分布事务管理与并发控制机制、分布式数据库管理系统的目录结构及其管理、通讯子系统、恢复子系统、数据执行子系统、分布式数据库系统用户接口的生成和管理等。

本书可作为高等学校计算机软件、计算机通讯、计算机应用等专业的教材和参考书，也可供软件系统、应用系统的设计人员、开发人员和程序员学习和参考。

### 图书在版编目(CIP)数据

分布式数据库管理系统实现技术 / 周龙骥等著。- 北京：  
科学出版社, 1998  
(数据库丛书)  
ISBN 7-03-006533-6

I . 分… II . 周… III . 分布式数据库—数据库管理系统  
IV . TP311.13

中国版本图书馆 CIP 数据核字(98)第 07878 号

科学出版社出版

北京东黄城根北街 16 号

邮政编码：100717

中国科学院印刷厂印刷

新华书店北京发行所发行 各地新华书店经售

\*

1998 年 7 月第一版 开本：787×1092 1/16

1998 年 7 月第一次印刷 印张：12 1/2

印数：1—2 900 字数：276 000

定价：19.00 元

《数据库丛书》是我國数据库專家學者團結協作、合力撰寫的一套系列著作。它比較全面地反映了國際数据库技术的丰富内容與最新发展，和我國数据库科技工作者多年來的主要研究成果，具有較高的理論水平和学术價值。

数据库是計算機科學技術中發展最快的領域之一，也是應用最廣的技术之一。是計算機信息系統与應用系統的構成基礎。相信《数据库丛书》的編輯出版，必將有益於推動我國数据库技术的研究与发展，促進我國数据库技术的普及与提高，加快数据库應用的推廣与深入，為我國社會經濟信息化作出貢獻。

張 故 祥

九九年一月

## 《数据库丛书》编委会

**主 编** 萨师煊

**副主编** 罗晓沛 王 珊

**编 委** 王能斌 施伯乐 郑怀远 童 烨  
唐世渭 周立柱 徐秋元 周龙骧  
徐洁磐 郑振楣 何新贵 马应章  
李建中 张大洋 董继润 瞿兆荣  
张作民 何守才 姚卿达 唐常杰  
冯玉才 尹良瑛 杨冬青 邵佩英  
李昭原 周傲英 于 戈

# 序

数据库是计算机领域发展最快的学科之一,因为它既是一门非常实用的技术,也是一门涉及面广、研究范围宽的学科。因此,它吸引了理论研究、系统研制和应用开发等不同方面众多的学者、专家和技术人才致力于其研究和实践。

数据库系统所管理、存储的数据是各个部门宝贵的信息资源。在信息化时代来临、Internet高速发展的今天,信息资源的经济价值和社会价值越来越明显。建设以数据库为核心的信息系统和应用系统,对于提高企业的效益、改善部门的管理、改进人们的生活均具有实实在在的意义。正因为数据库技术与经济、社会的发展和信息化建设有着密切的关系,这门学科才获得了巨大的源动力和深厚的应用基础。

数据库系统已从第一代网状、层次数据库系统发展到第二代关系数据库系统和第三代以面向对象为主要特征的数据库系统。数据库技术与网络通信技术、面向对象技术、并行计算技术、多媒体技术、人工智能技术等互相渗透,互相结合,成为当前数据库技术发展的主要特征。它使数据库领域中新的技术内容层出不穷,新的学科分支不断涌现,形成了新一代数据库系统的大家族。与传统的数据库相比,当今数据库的整体概念、技术内容、应用领域,甚至某些原理都有了重大的发展和变化。

面对如此丰富的学术内容和技术方法,如此广阔的研究方向和应用领域,从事数据库研究、开发和应用的科技人员,攻读数据库方向的研究生都迫切希望有一套丛书能系统而全面地介绍数据库学科的多个分支和相关领域。

《数据库丛书》的编写宗旨是把当前数据库学科各个分支的最新学术成果介绍给读者,以促进国内的学术研究;同时,又介绍数据库技术的发展过程,各分支之间的内在联系及在数据库大家族中的位置,以促进数据库和计算机科学的其他领域技术的结合。

本丛书由各分册组成,包括《数据库进展》、《分布式数据库》、《分布式数据库管理系统实现技术》、《并行关系数据库管理系统引论》、《数据仓库技术与联机分析处理》等。本丛书的每一分册涉及数据库学科的一个或几个分支。其中《数据库进展》则与其他分册有所不同,是本丛书的总纲、指南和补充,是给本丛书穿针引线、铺垫基础,从而使丛书成为一个各部分既相互独立又相互联系的整体。

《数据库丛书》是开放的,故丛书的分册将随着数据库学科的发展而不断补充。

本丛书各分册的主编和作者,多是长期从事数据库各分支领域研究工作的专家、学者。他们学术造诣高深,实践经验丰富,书中许多内容是他们长期研究的成果。本丛书不仅反映了国际数据库技术的最新成果和发展方向,也展示了我国数据库工作者的学术成果和研究深度,具有较高的理论水平和学术价值。它的出版是我国数据库学术界的一件大喜事。我向本丛书的所有作者和编委的辛勤工作表示崇高的敬意。

萨师煊

1998年1月

## 前　　言

分布式数据库管理系统(Distributed Data Base Management System, DDBMS)的研究开始于 70 年代中期,当时集中式的关系型数据库系统的研究继层次型和网状型数据库系统之后已趋成熟,相应的产品也已陆续推出。同时,计算机网络包括广域网和局域网亦在逐步推向应用,著名的如美国的 ARPANET,欧洲各国 PTT 的 X. 25 网等。在计算机硬件方面,小型机如 PDP11 及后继的 VAX 机,特别是 80 年代推出的微机由于其相对价廉和量大面广而使得计算机的应用得到了迅速普及。集中式数据库系统的日趋成熟加上计算机网络和小型、微型机两股潮流的汇合,为分布式数据库系统的研究和发展提供了前提条件。

在社会应用领域,对于分布式数据库系统的需求则更是日趋迫切,来势迅猛。诸如银行的通存通兑及划汇,保险业务跨地区的处理,国际洲际民航订票业务的受理,连锁商场的管理,军事上的情报系统、决策指挥系统及地域分布的军事基地的联络和管理,以及城市、企事业、医院、学校、机关的管理等等。众多的应用都牵涉到地理上分布的统一组织的管理,分布式数据处理及其核心分布式数据库管理自然地成为这些应用的技术基础,成为 70 年代末和整个 80 年代计算机科学技术及其应用的主要研究方向之一。

在以上的物质技术发展和应用需求的背景下,各先进国家均不约而同地紧紧抓住分布式数据库系统这个新的发展方向,投下巨资进行研究和开发。典型的如美国国防部委托美国计算机公司(CCA)研制的 SDD-1 系统,美国加利福尼亚大学伯克利分校的分布式 INGRES 系统,美国 IBM 公司的 R\* 系统,法国国家级的涉及一百多个大学、研究所和公司的 SIRIUS 计划,前西德 Stuttgart 大学研制的 POREL 系统和美国 CCA 公司的 ADA-DDM 系统等。这些研制计划规模宏大,耗资惊人,历时以 10 年计。如前西德的 POREL 系统的研制始于 1975 年,至 1987 年结束,其中截至 1983 年已耗资 450 万马克,50 人年,其它研究计划和系统研制的投资也均在此数量级上下。

巨额的投资和历时 10 年的研究,每年数次大型的国际学术会议和各种学术研讨会(workshop)的活跃的研究交流,若干先驱分布式数据库管理系统的研制,使 DDBMS 领域取得了决定性的成就。从理论上,DDBMS 的总结性专著已有数部出版。从实践上,各先驱研究计划均陆续完成,若干原型系统已投入运行或试用。市场上的四大数据库公司 ORACLE,SYBASE,INFORMIX,INGRES 都宣称其产品是分布式 DBMS,它们至少吸收了许多 DDBMS 研究中提出和发展的概念、方法、技术和算法及技巧等。可以说 DDBMS 的研究阶段已经过去,其产品化和商品化的工作正在进行。但是在现有市场 DDBMS 产品中尚未出现真正完全透明的系统,这或许由于在系统开销或复杂性方面还需进一步解决,也或许市场并不愿承受这种开销。另一方面,DDBMS 的研究成果正在各种新的计算机和通讯领域发挥其极为重要的作用,如分布一致性、分布式并发控制等在集体合作工作或写作系统、电视会议系统及在 Internet 中都有着不可或缺的应用。

本书是作者历时 10 年设计和开发分布式数据库管理系统 C-POREL 的研究和实践

经验的总结,同时也介绍和分析了国内外一些著名的先驱分布式数据库管理系统的工作原理、设计和研制中的创新特色和若干典型的概念、方法、算法和技巧。本书可以使读者一窥 DDBMS 的设计和实现的全貌,学会如何设计和实现一个 DDBMS 系统,并能在各具特色的体系结构、进程结构、算法和技巧中进行选择和取舍、分析和评价。迄今尚未见到国内外有类似本书这样的专著出版。本书可作为大学计算机科学和计算机软件专业的高年级学生和研究生的教材,也可作为研究工作者、系统设计和开发者的参考书。

本书内容分为十一章。第一章是分布式数据库管理系统概论。第二章是 DDBMS 的体系结构,这是 DDBMS 设计的基础,其中包括了与操作系统密切相关的 DDBMS 进程结构的讨论和分析。第三章和第四章介绍 DDBMS 的编程语言及编程语言编译器的设计和实现。第五章讨论编程语言涉网的全局编译和优化。第六章介绍 DDBMS 的事务管理和并发控制机制,它是 DDBMS 的核心成分。第七章描述分布式数据库管理系统的目录结构及其管理,讨论了独具特色的 C-POREL 的分布式目录结构。第八章介绍 DDBMS 的通讯子系统。详述了 DDBMS 的系统通讯的各种概念和算法,特别是对最新的可靠状态控制协议的改进及其实现作了详细的介绍。第九章为恢复子系统,分别叙述了集中式 DBMS 的恢复和分布式 DBMS 的恢复以及其实现技术。第十章为数据执行子系统,它相当于一个集中式数据库管理系统的实现技术,并对在 UNIX 操作系统之上实现 DBMS 的有关算法和技术及若干技巧进行了深入讨论。第十一章是分布式数据库系统用户接口的生成与管理,介绍了 DBMS 用户接口的一般概念及最新发展,对这个国际数据库界一致公认的 DBMS 的一个主要研究方向进行了较全面的介绍和强调。

参与本书撰写的均是分布式数据库管理系统 C-POREL 的设计和研究、开发人员,他们均经历了为时 6 年的设计、研究和研制的全过程,直至系统的集成运行和鉴定,其中部分人员还参加了前期为时 4 年的准备,则历时更以 10 年计,因此本书内容不仅经历过理论上的研究,也经过了系统设计和实现中的推敲和考验,具有较充实的实践基础。本书第一章、第二章、第七章和第十一章由中国科学院数学研究所周龙骥研究员撰写。第三章和第四章由上海大学(原上海科技大学)邵伟民副教授撰写。第五章由华东师范大学顾君忠教授撰写。第六章由中国科学院数学所周为群博士(现任职于北京大规模集成电路测试研究所)撰写。第八章由中国科学院数学所徐建礼博士撰写。第九章由中国科学院数学所彭立军硕士(现任职于 SYBASE 公司)撰写。第十章由中国科学院数学所副研究员柴兴无博士撰写。最后由周龙骥研究员总其成。

本书所涉及的研究工作得到了国家“七五”科技攻关、中国科学院重点项目、国防科工委军事预研项目、电子部电子科学院重点项目、国家科委特别基金和上海市高等教育局科研基金的资助。在 C-POREL 研究、设计和开发中,中国科学院合同局(现应用发展局),中国科学院数学所、中国软件技术公司、上海科技大学(现上海大学)和华东师范大学曾给予了大力支持和帮助,中国科技大学研究生院的罗晓沛教授为本套丛书的出版做了大量组织工作,作者在此对他们表示最诚挚的谢意。

由于作者水平有限,书中难免有不妥之处,恳请读者不吝赐教。

作 者

1997 年 12 月

# 目 录

<b>第一章 分布式数据库管理系统概论</b> .....	1
1.1 引论 .....	1
1.2 分布式数据库系统的特征 .....	2
1.3 若干研制计划和原型系统 .....	7
1.4 小结 .....	13
参考文献 .....	14
<b>第二章 分布式数据库管理系统的体系结构</b> .....	16
2.1 DDBMS 体系结构综述 .....	16
2.2 DDBMS 的分层体系结构 .....	17
2.3 DDBMS 的进程结构 .....	19
参考文献 .....	23
<b>第三章 分布式数据库系统的编程语言</b> .....	25
3.1 编程语言的设计要点 .....	25
3.2 RDBL 语言简介 .....	25
3.2.1 数据说明语句 .....	26
3.2.2 数据操作语句 .....	27
3.2.3 数据检查语句 .....	29
3.2.4 分布语句 .....	31
3.2.5 游标语句 .....	31
3.3 R* 对 SQL 语言的扩充 .....	32
参考文献 .....	33
<b>第四章 编程语言编译器的设计和实现</b> .....	34
4.1 编程语言编译器的任务及其体系结构 .....	34
4.1.1 编程语言编译器的任务 .....	34
4.1.2 编程语言编译器的体系结构 .....	34
4.2 预编译法 .....	36
4.3 编程语言的翻译 .....	37
4.3.1 语法分析 .....	37
4.3.2 语义分析 .....	39
4.3.3 局部优化 .....	41
4.3.4 子事务建立 .....	42
参考文献 .....	42
<b>第五章 全局编译的分析和设计</b> .....	43
5.1 全局编译的任务和目标 .....	43
5.1.1 分布透明性 .....	43

5.1.2 数据的全局一致性 .....	44
5.1.3 系统的高性能 .....	45
5.1.4 全局编译的工作流程和结构 .....	45
5.2 编译准备阶段 .....	45
5.3 完整性测试和授权检查 .....	46
5.4 查询优化 .....	46
5.4.1 代数优化 .....	47
5.4.2 分布优化 .....	49
5.4.3 典型的分布优化算法 .....	49
5.5 确定运算执行地点 .....	51
5.6 代码扩充和传递 .....	52
5.7 分布式数据库系统的修改一致性 .....	52
5.8 C-POREL 系统的全局编译 .....	53
参考文献 .....	54
<b>第六章 分布事务管理与并发控制机制 .....</b>	<b>55</b>
6.1 分布事务与分布事务管理 .....	55
6.1.1 分布事务 .....	55
6.1.2 分布事务管理 .....	57
6.2 分布事务处理协议 .....	58
6.2.1 关于分布事务处理协议 .....	58
6.2.2 基本的分布事务处理协议 .....	60
6.2.3 分布事务处理协议的描述工具 .....	62
6.2.4 不阻塞的分布事务处理协议 .....	63
6.3 分布式并发控制的基本方法 .....	66
6.3.1 分布式数据库系统的封锁方法及死锁的预防与检测 .....	67
6.3.2 其它的并发控制方法 .....	68
6.4 分布式数据库管理系统 C-POREL 的事务管理系统 .....	70
6.4.1 TM 的分布事务加工及结束协议 .....	71
6.4.2 TM 的分布式并发控制方法 .....	73
6.4.3 TM 的结构及其实现 .....	74
参考文献 .....	74
<b>第七章 分布式数据库的目录结构及其管理 .....</b>	<b>76</b>
7.1 数据库目录的重要性 .....	76
7.2 若干先驱 DDBMS 的目录体系的回顾与分析 .....	76
7.2.1 SDD-1 .....	76
7.2.2 Distributed INGRES .....	77
7.2.3 POREL .....	77
7.2.4 R <sup>*</sup> .....	77
7.2.5 两类目录体系结构的比较分析 .....	79

7.2.6 SUNDBB .....	80
7.2.7 WDBBS-32 .....	81
7.3 C-POREL 的目录管理策略 .....	82
7.3.1 目录事务 .....	82
7.3.2 C-POREL 的目录结构和目录事务的划分 .....	84
7.3.3 C-POREL 目录的进程结构 .....	85
7.3.4 C-POREL 目录的模块结构和文件组织 .....	85
7.3.5 C-POREL 目录事务的并发控制 .....	87
7.3.6 C-POREL 目录事务的恢复 .....	94
7.3.7 小结 .....	94
参考文献 .....	95
<b>第八章 通讯子系统 .....</b>	<b>96</b>
8.1 引论 .....	96
8.2 若干有代表性的 DDBMS 系统中的通讯子系统 .....	96
8.2.1 SDD-1 .....	97
8.2.2 R* 系统 .....	106
8.2.3 POREL 系统 .....	112
8.2.4 C-POREL 的通讯子系统 CS .....	119
8.2.5 ADA-DDM、分布式 INGRES、SIRIUS-DELTA .....	123
8.3 对通讯子系统的比较和讨论 .....	126
8.4 结束语 .....	128
参考文献 .....	129
<b>第九章 恢复子系统 .....</b>	<b>130</b>
9.1 集中式数据库的恢复 .....	130
9.1.1 事务 .....	130
9.1.2 故障 .....	130
9.1.3 恢复方法 .....	131
9.1.4 恢复方法的选择 .....	132
9.2 分布式数据库的恢复 .....	136
9.2.1 分布式数据库中的故障 .....	136
9.2.2 分布事务的一致性 .....	136
9.2.3 两阶段提交 .....	138
9.2.4 复制技术 .....	141
9.2.5 C-POREL 中恢复的实现 .....	142
参考文献 .....	148
<b>第十章 数据执行子系统 .....</b>	<b>149</b>
10.1 引论 .....	149
10.2 数据执行层的分层体系结构 .....	151
10.3 数据执行层的模块结构 .....	151

10.4	与 TM 的接口 .....	152
10.5	关系代数表达式的执行与优化 .....	153
10.5.1	分布查询的两步优化 .....	154
10.5.2	非代数优化的实现 .....	154
10.5.3	关系运算算法库的实现 .....	157
10.6	单元组接口 .....	157
10.6.1	关系代数运算的完备性 .....	158
10.6.2	将关系代数转化为对元组的操作 .....	159
10.6.3	单元组接口的实现 .....	160
10.7	存取路径管理 .....	161
10.7.1	主键的存取路径 .....	161
10.7.2	查找数据记录集合的存取路径结构 .....	163
10.7.3	通用存取路径 .....	164
10.7.4	存取路径的并发控制 .....	165
10.8	缓冲区管理 .....	166
10.8.1	数据库的外存管理 .....	166
10.8.2	DBS 的系统缓冲区 .....	167
10.8.3	SB 内页查找算法 .....	167
10.8.4	调页算法 .....	167
10.8.5	缓冲区管理子系统的实现 .....	168
10.9	记录管理 .....	169
10.9.1	记录的存储结构 .....	169
10.9.2	记录编址 .....	170
	参考文献 .....	172
<b>第十一章</b>	<b>分布式数据库系统用户接口的生成和管理 .....</b>	<b>173</b>
11.1	引论 .....	173
11.2	智能化输入接口 .....	176
11.3	图形用户接口 (GUI) .....	176
11.4	INFORMIX GUI 工具的分类 .....	177
11.5	INGRES 的数据库产品 .....	179
11.6	Sybase 的多媒体应用开发系统 Gain Momentum .....	180
11.7	集成化数据库前端开发环境 PowerBuilder .....	181
11.8	分布式多媒体数据库管理系统 CDB/M 的用户接口 .....	182
11.9	用户接口软件的一些基本功能 .....	185
11.10	用户接口的评价 .....	185
	参考文献 .....	186

# 第一章 分布式数据库管理系统概论

## 1.1 引 论

分布式数据库系统的研究始于 70 年代中期,迄今已近 20 年了,其基本问题的提出和研究以及国际上具代表性的先驱研究计划的实施和相应原型系统的研制则主要集中在前 10 年。至 80 年代中后期,DDBMS 领域的工作已取得了决定性进展:许多基本问题被提出来并已获解决;提出了一系列新的概念、新的方法和新的技术;一批原型系统已经研制成功并获得了相当的经验;一些产品正在试制或已经推出。总之一句话,此时 DDBMS 的技术已经基本成熟,其产品化的时代已经到来<sup>[1]</sup>。但是原型系统的研制成功距离真正实用的产品系统的推出还有一段相当长的路要走。一般来说,研究工作要比实际系统的应用超前大约 10 年。例如早在 70 年代末 80 年代初就已提出并解决的分布式数据库系统中的两阶段提交协议即 2PC 协议,它在市场上占第一位的数据库管理系统 ORACLE 中直至 1992 年的第六版才予以实现。对于产品系统来说,首要要考虑的是系统的可靠性、安全性、效率以及复杂多变的市场因素。

分布式数据库管理系统兴起于 70 年代中期并不是偶然的。一切计算机科学和技术的发展其推动力来自两个方面:一方面是应用需求的刺激,另一方面是硬件环境的发展。在应用方面如全球性的民航订票系统、铁路订票系统、水陆空联运系统、洲际银行的存取和汇兑业务系统、连锁店的管理系统、全国性甚至全球性的保险公司业务系统、上百个军事基地的统一管理系统、跨国公司的管理系统、旅游订票和调度系统等等。这些应用都涉及地理上分布的公司、组织和团体的业务的管理、集中式的数据库系统已无法提供合适的支持。在硬件环境上,计算机及通讯网络则更是突飞猛进地发展。功能强大的计算机和 16 位、32 位的微型机和工作站以及日益广泛装备的公用数据网和局域网,为 DDBMS 的研制提供了一个成熟的实用的环境。在这两股潮流的强力推动下,在集中式 DBMS 成熟技术的基础上,开始了 DDBMS 的迅猛发展。从 70 年代中期起,各先进国家都争先恐后地提出各自的分布式数据库系统研究计划,投资均达百万美元以上,如前西德的 POREL 系统,5 年间斥资 450 万马克。法国的 SIRIUS 计划则是一项国家级规模的项目,涉及全法的许多大学、研究所(如 INRIA)和工业部门。在 DDBMS 领域的几个影响最大的先驱系统中,最重要的首推 CCA 公司(Computer Corporation of America)为美国国防部研制的 SDD-1 系统(System for Distributed Databases)。SDD-1 研制中提出的一些思想、概念和方法如半连接(Semi-Join)、时戳(time stamp)、分布式目录结构、可靠通讯等对 DDBMS 领域的研究和原型系统的开发产生了非常深远的影响。从 70 年代中期到 80 年代中期,DDBMS 的研究如火如荼,国际会议每年都举行多次,发表的文章更是连篇累牍。人们说 80 年代是分布式数据处理及其核心分布式数据库管理的时代不是没有道理的。

## 1.2 分布式数据库系统的特征

分布式数据库管理系统(Distributed Data Base Management System,DDBMS)现在已是随处可见的很普通的术语了,但是它的确切定义到底是什么?它与分散式系统(Decentralized System)有什么区别?市场上见到的声称自己是 DDBMS 的众多产品是真正意义上的分布式数据库管理系统吗?DDBMS 的独有的特征到底是什么?<sup>[2]</sup>在这一小节我们将简要地介绍一下,以利于本书以后各章的展开。

通俗地说,分布式数据库系统是地理上(或物理上)分散而逻辑上集中的数据库系统。管理这样的数据库系统的软件称为分布式数据库管理系统。分布式系统通常由计算机网络(局域网 LAN 或广域网 WAN)联结起来,被联结的逻辑单位(包括硬件如计算机、外部设备等和软件如操作系统 OS、数据管理系统等)称为结点或站点(site)。所谓地理上分散意为各个结点分布在不同的地方,如北京、上海、广州(用广域网联结)或同一大院的不同建筑物,同一建筑物的不同房间(用局域网联结)。所谓逻辑上统一意为网络联结的各结点共同组成单一的数据库。以下是 DDBMS 的一些固有特征。

### (1) 结点透明性(site transparency)

不同结点上的各个用户(全局用户)所面对的是逻辑上统一的同一个分布式数据库。数据分布和事务(transaction)的分布式加工对全局用户透明,亦即每个全局用户都感觉整个数据库就处在他们所在的结点上,好像是一个集中式数据库一样。

### (2) 两种体系结构

存在着两大类 DDBMS 体系结构:同质的(homogeneous)和异质的(heterogeneous)。(由于在我国数据库界流行的翻译方法将其译为同构的和异构的,以下我们有时也沿用这一译法)。所谓同构的意指各结点系统的数据模型(层次型、网状型、关系型、函数型、面向对象型等)是相同的<sup>[3]</sup>,否则称为异构的。分布式数据库(Distributed Data Base, DDB)的体系结构如图 1.1 所示。

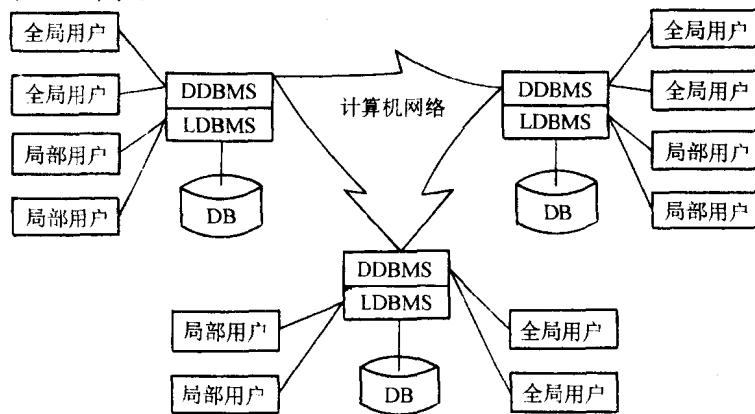


图 1.1 非集中控制的 DDB

注意图 1.1 中各个结点上的 **DDBMS** 和 **LDBMS** 在有些系统中是合二为一的,即合并为 **DDBMS**。此时局部用户(即只存取本地结点上数据的用户)也作为全局用户来看待。

在一个局部结点上的系统配置如图 1.2 所示。第一级划分是基于各结点系统的数据模型的。如果同构的 DDB 亦即具有相同数据模型的 DDB，其各结点系统所基于的 OS 和 COMPUTER 也相同，则称为完全同构的，否则称为非完全同构的。这是第二级划分。第三级是对于完全同构的 DDB，其各结点系统之间也可能存在差异，例如数据的表示不同，有的用 ASCII 码，有的用 EBCDIC 码。整数、浮点数表示的标度和精度也可能不同。语种亦可能各异，有的为中文，有的为英文。数据的单位和记法也可能不同，电视机有的用厘米而有的用英寸。奶粉在结点 a 用克，而在结点 b 用磅。工资有的只计基本工资和补贴，有的则还要加上奖金、年终分红甚至实物。

对异构的 DDB，除了可能有上述的数据表示、单位和记法的不同之外，还有以下方面的不同：

- 1) 数据结构的不同。例如层次型和网状型会包含存取路径结构如 DBTG 中的系结构 (set)，而关系型则不会包含。
- 2) 完整性断言不同。它通常在数据模型中或用户程序中指明，故各结点间可能差别很大。
- 3) 数据语言不同。关系语言的语句和 DBTG 语言的语句不可能一一对应，例如 DBTG 中沿系结构的查找语句在关系 DB 中不存在对应物。

对于异构 DDB 的处理，一个常用的方法是建立一个公共的全局模式或正则模式，并建立全局模式和各结点局部模式之间的模式转换。这种转换程序的总数为  $2N$  ( $N$  为不同的局部模式的数目)。如果不建立全局模式而实行不同局部模式之间的直接转换，则转换程序的总数将增至  $N(N-1)$ 。

研究和实践表明这样的转换程序是较复杂的，它相当于一个编译程序或一个解释程序。近几年发展起来的联邦数据库和模式合并(schema merging)，其有效的成果有限。只有在定义域(domain)一级其结果较为严格(域的等价、包含、被包含、相交、不相交等)，对模式合并尚需人工交互式操作，开销亦过于庞大。

#### (3) 结点自主性(site autonomy)

在图 1.1 中我们看到每个结点上既有全局用户也有局部用户，前者只与 DDBMS 打交道而后者只与 LDBMS 打交道。前者只知道 DDB，不知有 LDB 的存在，后者则只知道 LDB，不知道有 DDB 的存在。这表明在分布式数据库系统中存在着两级控制：全局控制和局部控制。不同的系统中根据不同的设计原则，这两级控制的程度各不相同。一种极端的情况是不允许结点的自主性，只允许全局事务、DDBMS 和 LDBMS 合并为一个 DDBMS。LDB 被看作是 DDB 的一个逻辑组成部分，即使只存取本地结点数据的局部事务也被看成是一种特殊的全局事务，例如 C-POREL 系统。其它情况则是允许某种程度的结点自主性，结点上允许纯局部用户，由 LDBMS 管理，只存取 LDB。结点还可决定哪些数据专属于本地结点，哪些数据可以提供出来供 DDB 公用，例如 R\* 系统。

#### (4) 目录结构

大多数 DDBMS 都支持全局目录，这是一种面向数据对象的目录结构，它对全 DDB

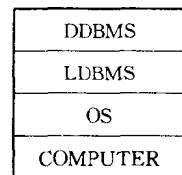


图 1.2 一个局部结点的系统配置

进行全局控制<sup>[4]</sup>。另一种是分布式目录结构,各结点将它愿提供给 DDB 共享的数据通知其它结点,然后其它结点上的各用户将要用的数据记入他自己的目录中,这是一种面向用户的目录结构<sup>[4,5]</sup>。在目录的构造上,多数系统将目录数据和用户关系一律看待,但有的则建立独立的目录系统<sup>[6]</sup>。

#### (5) 数据分片(fragment)

关系可能太大而不够灵活,因此需将关系分割称之为关系分片(relation fragment)。其又分为水平分片和垂直分片。水平分片由限制运算将关系分割为若干个子关系(例如学生关系按系分片)。大多数 DDBMS 的子关系互不交叠,但 VDN 系统允许交叠。垂直分片采用投影运算将关系分割为一组组的属性组,当然应保证无损连接。

#### (6) 副本(replica)

DDB 支持副本的目的有二:一是为了提高效率;二是为了数据安全即保证数据的可用性。一个数据对象在不同结点或同一结点有若干份拷贝,它们对用户是透明的。有了副本对保持数据更新的一致性带来复杂性,许多论文讨论了它们的实现算法。新版的 ORACLE 系统、INGRES 系统和 SYBASE 系统等均装备了专门设施以处理这一至关重要的问题。

#### (7) 优化

由于 DDBMS 的极端复杂性,其系统开销非常可观。为了提供能够被用户接受的效率,优化是很重要的一环。优化分为局部优化和全局优化。局部优化采用的技术就是通常集中式 DBMS 的代数优化和非代数优化。全局优化则涉及通讯费用,包括发送一条消息的费用,消息长度(通讯量)等。对于远程网,通讯费用占主要地位。对于局域网,则通讯、CPU 和 I/O 等开销均应考虑。在设计优化算法的目标函数时可以有两种选择:一种是以总时间作为优化准则,另一种是以响应时间作为优化准则。对于后者可以充分计及在网络各结点上并行执行的收益。

关于优化的时机问题存在着两种不同的策略。一种是静态优化即编译时优化,一种是动态优化即执行时优化。中间生成的关系的尺寸大小是优化时的重要依据,故应研究好的算法以便给出尽可能接近实际的结果。对于静态优化,中间结果的尺寸估计失误将会在后续的估算中积累和传播,从而大大背离优化的初衷。静态优化的优点是简单,系统开销小并可充分发掘在各结点并行执行的好处;缺点是不易有好的方法去估计中间关系的尺寸,因而难于得到最优的执行计划。动态优化的方法可在执行中(一般是当两个关系连接时或考虑更多的关系运算时)估算中间结果的尺寸,以决定送往哪个结点。这样将减少错误估计效果的累积和传播,其缺点是要做更多的估算,系统开销大且不能充分利用并行执行的好处。为了取长补短,有的算法采用混合策略,只有当中间结果的尺寸与静态估算的结果相差太远时才采用动态方法进行估算。

根据所采用的优化方法即可制订出查询计划。查询计划的产生过程存在着不同的方法。一种是集中式的,即查询计划在递交事务的结点上作出。一种是分布式的,由所有参与加工该事务的结点共同决定查询计划。还有一种半集中式方法,其主要决定由递交事务的结点作出,次要的决定可由参与加工的其它结点作出。

关于优化的范围或作用域。通常的优化器只考虑对一个查询语句进行优化。进一步的考虑是可一次对多个查询语句进行优化,一个用户其后继的查询可以利用前面的查询结果。对于多个用户,可考虑他们的平均响应时间,考虑负载对于响应时间的影响等。

表 1.1 中列出一些分布式查询算法。

表 1.1 一些论文及系统中的分布式查询算法

论文及系统	优化方式	目标函数	优化因素	网络拓扑	是否用Semi-Join	分片否	有否用副本	算法复杂性
SDD-1	静态	总时间	消息长度	一般	是	否	否	多项式 (局部最优)
Distributed INGRES	动态	响应时间或总时间	消息长度, CPU	一般或广播	否	是	否	多项式 (局部最优)
R*	静态	总时间	消息个数, 消息长度, CPU,I/O	一般	否	否	否	动态规划
P. M. G. Apers, A. R. Hevner, S. B. Yao	静态	响应时间或总时间	消息个数, 消息长度	一般	是	否	否	多项式 (接近最优)
W. W. Chu, P. Hurley	静态	总时间	消息长度, CPU	一般	否	否	是	指数 (全局最优)
UNITY	静态	响应时间	消息长度, I/O	星形	是	否	否	
POLYPHEME	动态	总时间	消息长度	广播	否	否	否	动态规划
MERMAID	混合	总时间	消息个数, 消息长度	一般	是	是	是	简单且有效

### (8) 数据模型

绝大多数 DDBMS 都是关系型的。关系特别适宜于分布。E. F. Codd 曾谈到由于分布式数据库系统的发展,显示出层次型和网状型的不足<sup>[7]</sup>。道理很简单,关系的结构非常简明清晰,易于作水平和垂直分片和恢复。层次和网状型由于包含了导航式的存取路径结构,如 CODASYL 中的“系”(set),使得相互联系繁复纠缠,难于进行数据分布和相应的运算。

### (9) 并发控制

已经提出了封锁、时戳和乐观方法三种并发控制技术<sup>[8]</sup>。大多数 DDBMS 都采用封锁方法,SDD-1 系统采用时戳方法,C-POREL 采用时戳和封锁混合方法。尚未见到实现乐观方法的系统。

### (10) 死锁

存在死锁预防、死锁避免、死锁检测和恢复以及超时(time-out)等方法。死锁检测可在集中结点进行,也可在一组结点分布式检测。

### (11) 解释和编译

绑定(binding time)的早晚涉及系统的效率。编译时的早绑定可以提高执行效率,特