

# 高级数据库 原理与技术

● 毛国君 编著

 人民邮电出版社  
POSTS & TELECOM PRESS

# 高级数据库原理与技术

毛国君 编著



人民邮电出版社

## 图书在版编目 (CIP) 数据

高级数据库原理与技术 / 毛国君编著. —北京: 人民邮电出版社, 2004.8

ISBN 7-115-12066-8

I. 高... II. 毛... III. 数据库系统 IV. TP311.13

中国版本图书馆 CIP 数据核字 (2004) 第 085080 号

### 内 容 提 要

随着数据库技术本身的发展和其他新技术的渗透, 当今数据库的整体概念、技术内容、应用领域甚至基本原理都有了重大的发展和变化, 形成了庞大的数据库家族。本书将全面介绍这些新型高级数据库, 包括分布式数据库、并行数据库、Oracle 系统、数据仓库以及面向对象数据库和多媒体数据库等相关技术。

本书共分五篇。第一篇是预备知识, 主要是解决一些读者或学生缺乏必要的分布式系统和数据库基础知识的问题。第二篇全面讲述分布式数据库的原理与技术, 包括分布式数据库的概念、设计、查询优化、并发控制及安全性等。第三篇从理论和应用两个视角, 对数据库中的并行处理技术和 Oracle 数据库管理技术进行深入剖析。第四篇集中阐述数据仓库概念、设计基础、核心技术及它的质量管理等问题。第五篇对其他一些新型数据库技术加以介绍, 包括面向对象数据库、多媒体数据库、工程数据库、科学数据库、模糊数据库、演绎数据库、主动数据库、移动数据库、统计数据库等。

本书可作为计算机专业研究生或高年级本科生教材, 也可以作为从事计算机研究和开发人员的参考资料。同时, 对于高职院校也可以选择部分章节进行讲授。

### 高级数据库原理与技术

- ◆ 编 著 毛国君  
责任编辑 滑 玉
- ◆ 人民邮电出版社出版发行 北京市崇文区夕照寺街 14 号  
邮编 100061 电子函件 315@ptpress.com.cn  
网址 <http://www.ptpress.com.cn>  
读者热线 010-67129259  
北京汉魂图文设计有限公司制作  
北京朝阳展望印刷厂印刷  
新华书店总店北京发行所经销
- ◆ 开本: 787×1092 1/16  
印张: 16.25  
字数: 387 千字 2004 年 8 月第 1 版  
印数: 1-5 000 册 2004 年 8 月北京第 1 次印刷

ISBN 7-115-12066-8/TP · 3827

定价: 22.00 元

本书如有印装质量问题, 请与本社联系 电话: (010) 67129223

# 编者的话

众所周知,数据库技术从20世纪60年代诞生至今,经过30余年时间的发展,已经成为计算机应用的主要技术之一。就数据库本身的应用范围而言,各种规模、各种用途的数据库应用系统已经投入使用,而且新的系统正与日俱增。其应用的广度和深度可能是计算机其他分支所不能比拟的。就数据库技术的研究和发展而言,各种学科与数据库技术的交叉渗透,产生了许多新型数据库技术。20世纪80年代,关系型数据库及其相关的数据模型工具、索引及数据组织技术等日渐成熟,一些关系型数据库管理软件和辅助工具得到不断完善。像Oracle这样的数据库商界的领先公司,已经把一些核心的关系型数据库研究成果和技术成功地应用到产品中。从20世纪80年代中期开始,关系数据库技术和新型技术的结合成为数据库研究和开发的重要标志。例如,扩展关系、面向对象以及对象-关系等新型数据模型被引入到数据库中,而且包括空间、时态、多媒体以及Web等新型的数据成为数据库应用的重要数据源。同时,事务数据库、主动数据库、知识库、办公信息库等技术也得到蓬勃发展。随着网络技术的发展,分布式数据库得到了充分的研究,一些理论和技术问题趋于明朗,成为数据库家族的主流技术之一。近年来,数据仓库作为一种新型的数据存储和处理手段,成为异构数据源集成和管理决策制订的一种有效的技术支撑环境。

对大多数读者来说,数据库技术的迅速发展和更新带来的新问题是如何全面而准确地掌握这些新型数据库系统。作者长期从事相关方面的教学工作,也深感在教材选择上的难度。以前研究生的相关课程是“分布式数据库”。一些学校也开设“数据仓库”课程或讲座,对于面向对象数据库等内容基本上涉及很少,或在相关课程提一下。由于没有合适教材可用,不得不通过指定大量参考书或文献来解决。由于目前相关书籍本身就少,而且侧重点不同,内容的完整性和科学性有待商榷。在有限的时间内尽量把数据库的相关知识系统化地介绍给学生,在没有统一教材的前提下给教师带来相当大的负担。同时,对于一些软件工程师或工程硕士、在职硕士进修班等要求提高实践能力的人员来说,也需要在科学的理论(原理)框架下理解和掌握实际系统的应用技术。因此,作者在多年各类教学和软件工程的实践基础上,对积累的素材进行了整理和加工,在统一的理论框架下,系统介绍这些新型数据库的原理和技术。

本书较全面地介绍了数据库家族的经典分支和最新成果。由于本书的内容是按相对独立的篇组织的,教师可以根据情况选用相关篇章来进行教学。对于一般读者,可以根据自己的基础进行选择学习或查阅。本书力求把理论和实践结合起来,不仅选取诸如分布式数据库查询优化、并发控制、数据仓库、多媒体与面向对象数据模型等经典的理论进行了系统化阐述,而且对诸如数据库中的并行处理技术、Oracle数据库管理、数据仓库的数据组织等实用技术进行了全面剖析。在介绍相关实用技术时,力求从理论高度来帮助读者理解相关技术。

本书共分五篇19章来阐述相关内容。第一篇是预备知识,内容包括:分布式系统的概念和原则,数据库基础原理和技术。第二篇全面讲述分布式数据库原理与技术,内容包括:分布式数据库系统概论,分布式数据库设计方法,分布式事务管理与数据库管理系统的概念,

分布式查询处理技术与方法，分布式并发控制的原理与算法，分布式数据库的可靠性和安全性的问题和对策。第三篇讲述数据库并行处理技术与典型数据库管理系统，内容包括：数据库中的并行处理技术，Oracle 系统的主要技术，Oracle 数据库系统的性能优化方法。第四篇讲述数据仓库的概念和技术，内容包括：数据仓库基础概念，数据仓库系统设计的原则和方法，对数据仓库的主要技术，数据仓库的质量管理问题和方法。第五篇是对其他新型数据库系统的介绍，内容包括：面向对象与数据库的结合技术，多媒体数据库的技术特点，对工程数据库和科学数据库等专用数据库系统，并对移动数据库等。

本书可作为计算机专业研究生教材、本科生选修教材，对于高职院校可以选择部分章节进行讲授。本书也可以作为从事计算机研究和开发人员的参考资料。为了保证内容的先进性和深度，对重点内容进行了重点阐述。本书内容相对全面，篇章之间耦合度小。教师可以作为教材来根据学生类型、学时安排等进行内容选取。作为参考书，读者可以根据自己的基础进行选择学习或查阅。

特别感谢我的导师，北京工业大学的刘椿年教授和大连理工大学的杨名生教授，因为本书的许多工作来自于我攻读博士和硕士学位期间的积累。感谢硕士研究生鲁杰和尤春梅，他们从事了部分的文字编排和整理工作。另外要感谢的是北京工业大学参加过我相关课程学习的各类学生（1998~2003年，研究生、硕士班等），他们的许多意见和文字更正，提高了本书的内容编排质量。同时感谢我家人的支持。

作者

2004年7月于北京

# 目 录

## 第一篇 预备知识

<b>第 1 章 分布式系统</b> .....	2
1.1 分布式系统的定义 .....	2
1.2 分布性的刻画 .....	3
1.3 高层操作系统 .....	5
<b>第 2 章 数据库基础知识</b> .....	7
2.1 数据库技术的发展 .....	7
2.2 数据模型 .....	9
2.2.1 概念模型 .....	9
2.2.2 数据模型 .....	10
2.3 数据库系统的基本组成 .....	12
2.3.1 数据库的三级模式设计 .....	12
2.3.2 数据库管理系统 .....	13
2.3.3 数据库系统与计算机应用系统 .....	14
2.4 数据库设计与实例 .....	14
2.5 数据库操作语言 .....	16
2.5.1 数据结构定义功能 .....	17
2.5.2 数据查询功能 .....	18
2.5.3 数据或结构修改功能 .....	19
本篇思考题 .....	20

## 第二篇 分布式数据库原理与技术

<b>第 3 章 分布式数据库系统概论</b> .....	22
3.1 分布式数据库系统的定义 .....	22
3.1.1 分布式数据库系统的发展 .....	22
3.1.2 分布式数据库的定义 .....	23
3.2 分布式数据库管理系统概述 .....	24
3.3 分布式数据库系统的组成 .....	25
3.4 分布式数据库系统的分类 .....	26
3.4.1 紧耦合式 DDBS .....	26
3.4.2 联邦式 DDBS .....	27
3.4.3 组合式 DDBS .....	27

<b>第4章 分布式数据库设计</b> .....	28
4.1 分布式数据库的构成方式 .....	28
4.1.1 单层次分布式数据库 (SL DDB) .....	28
4.1.2 多层次分布式数据库 (ML DDB) .....	28
4.2 分布式数据库的模式结构 .....	28
4.2.1 分布式数据库的模式层次 .....	28
4.2.2 模式间的映射 .....	29
4.2.3 分布式数据库系统参考模型 .....	30
4.3 分布式数据库系统中的透明性 .....	31
4.3.1 分片透明性 .....	31
4.3.2 位置透明性 .....	31
4.3.3 本地透明性 .....	32
4.4 分布式数据库的数据分割方法 .....	33
4.4.1 关系代数介绍 .....	33
4.4.2 数据分割方法 .....	35
4.5 分布式数据库的设计方法 .....	37
4.5.1 分布式数据库设计概述 .....	37
4.5.2 分布式数据库设计的原则 .....	38
4.5.3 分布式数据库的设计方法 .....	39
<b>第5章 分布式事务管理与数据库管理系统</b> .....	42
5.1 分布式事务的定义 .....	42
5.2 事务管理的目标 .....	43
5.3 分布式事务管理的模型 .....	43
5.3.1 主从事务管理模型 .....	43
5.3.2 三角事务管理模型 .....	44
5.3.3 层次事务管理模型 .....	44
5.4 分布式事务的编译与执行 .....	44
5.5 分布式数据库管理系统参考模型 .....	45
<b>第6章 分布式查询处理</b> .....	47
6.1 问题的提出 .....	47
6.2 数据分配与费用 .....	48
6.2.1 数据分配的单位 .....	48
6.2.2 数据分配的费用估计 .....	49
6.3 关系代数的等价变换 .....	51
6.3.1 算符树 .....	51
6.3.2 关系代数的等价变换 .....	52
6.3.3 公共子表达式的问题 .....	53
6.4 把全局查询变换成段查询 .....	53
6.4.1 限定关系的代数学 .....	53

6.4.2	水平分段关系的化简	54
6.4.3	垂直分段的化简	56
6.4.4	分布式分组和聚集函数求值的查询问题	56
6.4.5	关系代数的扩充	57
6.4.6	Group-by 操作的特性	57
6.4.7	参数性查询	58
6.5	基于等价变换的查询优化	59
6.6	基于半连接程序的查询优化	60
6.6.1	半连接程序	60
6.6.2	优化步骤和费用估计	61
<b>第 7 章</b>	<b>分布式并发控制</b>	<b>63</b>
7.1	问题提出与抽象	63
7.1.1	异常情况示例	63
7.1.2	分布式数据库管理系统的抽象	64
7.2	用于并发控制的 DDDBS 抽象结构	65
7.2.1	集中式事务处理模式	65
7.2.2	分布式事务处理模型	66
7.2.3	分布式事务处理模式	67
7.3	分布式并发控制理论	67
7.3.1	无干扰执行与可串行性	67
7.3.2	操作冲突与执行的等价	68
7.3.3	并发控制处理模式	69
7.4	两相封锁并发控制算法	70
7.4.1	基于锁的并发控制基本方法概述	70
7.4.2	两相封锁 (2PL) 算法思想	71
7.4.3	2PL 算法的基本实现方法	72
7.4.4	主副本 2PL 算法	72
7.4.5	表决 2PL 算法	73
7.4.6	集中式 2PL	73
7.5	时间戳并发控制方法	73
7.5.1	时间戳方法的基本实现方法	73
7.5.2	Thomas 写规则	74
7.5.3	多版本 T/O	74
7.5.4	保守的 T/O	75
7.5.5	减少重新启动的启发式方法	76
7.5.6	死锁问题	76
7.6	分布式并发控制算法的性能分析	78
7.6.1	性能评价问题	78
7.6.2	2PL 性能分析	78



7.6.3	T/O 性能分析技术	79
7.6.4	并发控制方法的选择	80
<b>第 8 章</b>	<b>分布式数据库的可靠性和安全性</b>	<b>81</b>
8.1	分布式数据库的可靠性及其含义	81
8.2	分布式数据库系统的故障分析和对策	82
8.2.1	硬件故障及其容错技术	82
8.2.2	软件故障及其容错技术	83
8.2.3	数据的可靠性及其容错技术	84
8.3	分布式可靠性协议	85
8.3.1	可靠性提交协议	85
8.3.2	可靠性终结协议	86
8.3.3	可靠性恢复协议	87
8.4	三阶段提交协议	88
8.5	分布式数据库的安全性及其含义	89
8.6	数据库管理系统的安全级别介绍	90
8.7	分布式数据库的安全机制	92
	本篇思考题	93

### 第三篇 数据库并行处理技术与典型数据库管理系统

<b>第 9 章</b>	<b>数据库中的并行处理技术</b>	<b>96</b>
9.1	数据库系统的应用模式	96
9.2	数据库中并行处理相关问题	97
9.3	多线程并行技术	98
9.4	数据库应用接口	99
9.4.1	数据库连接标准	99
9.4.2	多级分布式 Web 计算模型	100
9.4.3	中间件技术	101
9.5	并行数据库系统的相关技术	102
<b>第 10 章</b>	<b>Oracle 系统</b>	<b>104</b>
10.1	Oracle 数据库系统的基本知识	104
10.1.1	实例与进程概念	104
10.1.2	单进程实例和多进程实例	105
10.1.3	Oracle 后台进程	105
10.1.4	Oracle 内存结构	110
10.1.5	Oracle 的配置问题	116
10.2	Oracle 数据库结构和空间管理	118
10.2.1	Oracle 数据库物理结构及其文件类型	118
10.2.2	Oracle 数据库的逻辑结构	123
10.2.3	数据字典	129

10.2.4	Oracle 模式结构 .....	130
10.3	Oracle 的事务管理 .....	140
10.3.1	事务提交 .....	140
10.3.2	事务回滚 .....	141
10.4	Oracle 的分布处理 .....	141
10.4.1	Oracle 的 C/S 结构与自治性 .....	141
10.4.2	Oracle 的全局数据库名与远程查询 .....	142
10.4.3	Oracle 的透明性 .....	144
10.4.4	Oracle 高级复制技术 .....	145
<b>第 11 章</b>	<b>Oracle 数据库系统的性能优化</b> .....	<b>147</b>
11.1	Oracle 数据库优化问题 .....	147
11.1.1	数据库的系统化优化问题 .....	147
11.1.2	数据库的优化目标与基本过程 .....	148
11.2	Oracle 数据库的逻辑结构设计优化 .....	149
11.3	数据库操作的执行优化 .....	151
11.3.1	SQL 语句的执行计划问题 .....	152
11.3.2	基于规则的优化方法 .....	153
11.3.3	基于代价的优化方法 .....	153
11.3.4	SQL 语句的预处理问题 .....	154
11.3.5	SQL 性能优化的典型方法介绍 .....	155
11.4	Oracle 数据库性能优化和参数调整 .....	164
11.4.1	调整数据库服务器的内存使用性能 .....	165
11.4.2	调整磁盘 I/O .....	166
11.4.3	调整数据库服务器的回滚段 .....	167
11.4.4	调整网络传输与 I/O 代价 .....	168
11.4.5	应用程序的调整 .....	169
11.5	Oracle 系统的初始化参数调整 .....	170
本篇思考题	.....	171

## 第四篇 数据仓库

<b>第 12 章</b>	<b>数据仓库基础</b> .....	<b>174</b>
12.1	数据仓库的概念 .....	174
12.1.1	正确理解数据仓库技术 .....	174
12.1.2	数据仓库的主要特征 .....	176
12.1.3	数据仓库的应用 .....	179
12.2	数据仓库中的数据组织 .....	180
12.2.1	数据组织的层次结构 .....	180
12.2.2	数据分割 .....	181
12.2.3	元数据 .....	182

12.2.4	数据装载与追加 .....	182
12.2.5	数据仓库的文件组织形式 .....	183
12.2.6	多维数据模型及其实现 .....	183
12.3	数据仓库系统的体系结构 .....	184
12.3.1	多层的数据仓库环境 .....	184
12.3.2	数据仓库系统的应用体系 .....	184
12.3.3	数据仓库系统的关键部件 .....	186
12.3.4	数据集市 .....	187
<b>第 13 章</b>	<b>数据仓库系统设计</b> .....	<b>189</b>
13.1	数据仓库系统与传统数据库系统设计方法的比较 .....	189
13.2	数据仓库的数据模型 .....	189
13.2.1	星型模式 .....	190
13.2.2	数据仓库的三级数据模型 .....	190
13.3	数据仓库系统的设计和开发 .....	191
13.3.1	数据仓库系统的实现策略 .....	192
13.3.2	数据仓库系统的开发过程 .....	192
13.4	数据仓库解决方案及工具介绍 .....	196
<b>第 14 章</b>	<b>数据仓库的主要技术</b> .....	<b>200</b>
14.1	数据管理技术 .....	200
14.2	数据仓库与 OLAP 技术 .....	202
14.3	数据仓库与 Web 技术 .....	203
14.4	数据仓库与数据挖掘 .....	203
14.4.1	数据仓库和数据挖掘的关系 .....	204
14.4.2	数据挖掘的技术介绍 .....	204
<b>第 15 章</b>	<b>数据仓库的质量管理</b> .....	<b>211</b>
15.1	数据仓库与质量管理 .....	211
15.2	数据仓库系统的层次模式和质量管理 .....	212
15.3	数据仓库系统的组成要素和质量管埋 .....	214
	本篇思考题 .....	217

## 第五篇 其他数据库系统

<b>第 16 章</b>	<b>面向对象与数据库的结合技术</b> .....	<b>220</b>
16.1	面向对象数据库系统的特点 .....	220
16.2	面向对象与数据库技术的结合方法 .....	222
16.2.1	对象-关系数据库 .....	222
16.2.2	面向对象数据库 .....	223
16.2.3	演绎面向对象数据库 .....	224
16.2.4	多种技术相互渗透 .....	225
16.3	面向对象的数据库应用开发工具的发展 .....	226

---

16.4	面向对象与数据库技术结合的产品实例·····	226
<b>第 17 章</b>	<b>多媒体数据库技术</b> ·····	<b>229</b>
17.1	多媒体数据库技术的产生和发展·····	229
17.2	多媒体数据库系统的硬件环境·····	230
17.3	多媒体数据模型·····	230
17.4	多媒体数据库管理系统·····	231
17.5	多媒体数据库的用户接口·····	233
<b>第 18 章</b>	<b>专用数据库系统</b> ·····	<b>234</b>
18.1	工程数据库·····	234
18.2	科学数据库·····	235
18.2.1	科学数据特点和科学数据库的类型·····	235
18.2.2	建立和使用科学数据库·····	236
18.2.3	数字图书馆技术·····	237
<b>第 19 章</b>	<b>其他数据库技术介绍</b> ·····	<b>240</b>
19.1	知识库·····	240
19.2	模糊数据库与演绎数据库·····	240
19.3	主动数据库·····	241
19.4	移动数据库·····	242
19.5	统计数据库·····	242
	本篇思考题·····	243
	<b>主要参考文献</b> ·····	<b>245</b>

# 第一篇 预备知识

---

信息是人类宝贵的财富。随着计算机的发展和应用的深入，人们组织和利用信息的方式和力度得到了根本的改变。毫无疑问，数据库技术已经成为现代社会利用丰富信息的重要手段。

追溯数据库技术的发展历史，不难看出人们在一直追求更高效、更聪明地利用数据和信息的方法和技术。20世纪60年代，为了适应信息的电子化要求，信息技术一直从简单的文件处理系统向有效的数据库系统变革。到了70年代，数据库系统的三个主要模式——层次、网络和关系型数据库的研究和开发取得重要进展。80年代，关系型数据库及其相关的数据模型工具、索引及数据组织技术被广泛采用。80年代中期开始，关系数据库技术和新型技术的结合成为数据库研究和开发的重要标志。从数据模型上看，诸如扩展关系、面向对象、对象-关系以及演绎模型等被应用到数据库系统中。从应用的数据类型上看，包括空间、时态、多媒体以及 Web 等新型数据成为数据库应用的重要数据源。同时，事务数据库、主动数据库、知识库、办公信息库等技术也得到蓬勃发展。从数据的分布角度看，分布式数据库及其透明性、并发控制、并行处理等成为必须面对的课题。进入90年代，分布式数据库理论上趋于成熟，分布式数据库技术得到了广泛应用。数据仓库作为一种新型的数据存储和处理手段，成为异构数据源集成和管理决策制订的一种有效的技术支撑环境。

经过几十年的努力发展，数据库理论和技术得到了充分的发展。特别是近几年数据库和其他技术结合得越来越紧密，形成了一个庞大的数据库家族。要全面地理解和掌握这些数据库技术，的确需要相关的技术准备。而我们很多人由于没有系统地进行计算机相关理论的学习，因此带来一定难度。例如，作者曾多次为研究生开设本课程，发现许多研究生在本科期间并没有系统学习诸如分布式系统、数据库等相关课程。他们有的本科学习的不是计算机专业，有的因为毕业较早也缺乏相应的训练。为此，我们编制本篇，希望通过数据库主要理论和方法的集中介绍，使读者能较好地接受后面的内容。在这里我们将介绍分布式系统和集中式数据库技术的主要知识。读者可以根据情况选用。

# 第1章 分布式系统

分布式系统（DS）或分布式处理系统代表了计算机发展的方向之一，它处理的并行性、节点的自制性、容错可靠性、可扩展性、高可用性以及良好的性能价格比等优势，已经吸引专家学者和商业厂家来为此付出了艰辛的劳动。作为重要的理论铺垫，本节将介绍分布式系统的相关知识。

## 1.1 分布式系统的定义

分布式系统所涉及的技术和概念很多，其实很难用简短的方式给它一个精确的定义；而且随着它本身以及应用它的相关技术的发展，它的含义也在发展中。尽管如此，我们还是挑选了几个比较有代表意义的描述，希望读者能从中了解分布式系统的精髓。

美国电工电子学会下属的计算机学会给出的分布式系统描述为：“包含多个相连的处理资源，这些资源能在系统的控制下，对单一问题进行合作，而且最少依赖集中过程、数据或硬件。”简言之，这样的描述揭示了分布式系统的精髓——硬件的重复性、处理的统一性和控制的分布性。

英国国家科学研究委员会下属的计算机学会给出的分布式系统描述为：“包含多个独立的但又交互作用的计算机，它们可以对公共问题进行合作。这个系统的特点是包含多个控制路径，它们执行一个程序的不同部分而且又相互作用。”这也从另一个侧面揭示分布式系统多机下的自治性、合作性和并行性。

P. H. Enslow 总结了分布式系统的五个基本准则，从中可以对分布式系统有个更透彻的认识。

### 1. 资源的重复性

分布式系统是一个多节点的系统。这些节点可以是计算机或处理器，可以同构或异构。这里的资源从更广泛的观点可以看作是硬件、软件以及数据等。所谓资源的重复性是指分布式系统中硬件、软件以及数据的冗余配置。

### 2. 物理上的分布性

从硬件上看，不同的计算机或处理器以节点形式相对独立地分布。随着概念的发展，其实这里的分布性也可以从不同层次来理解。从软件上看，每个节点都可以有全局相关的系统程序、局部系统程序以及应用程序；从数据上看，它可能是一个数据分布的系统。

### 3. 高层操作系统（或分布式操作系统）

高层操作系统是相对于局部操作系统而言的。它负责对系统的分布性资源进行统一的控制。事实上，在一群独立工作的计算机和分布式系统之间的主要区别就在于系统软件上。高层操作系统使一个简单的硬件堆积转变为一个统一协调的工作系统。

### 4. 系统的透明性

用户可以在任何节点请求系统服务，但是他们不必关心具体提供服务的物理节点。理想的分布式系统应该为最终用户屏蔽具体的系统实现细节，用户像使用单一的集中式系统一样来使用分布式系统。透明性是分布式系统的灵魂，实现不同层次的透明性是分布式系统必须解决的关键问题之一。

### 5. 协作的自治性

理想的分布式系统应该是每个节点都是一个完整的处理系统，表现出强大的自治性。同时这些节点又是合作的，可以在统一控制下，协作地完成更复杂的工作。

归纳上面的观点，我们可以认为分布式系统是一个多节点的、处理或数据分布的、在统一一下提高综合处理能力的协作体。

## 1.2 分布性的刻画

为了帮助读者更好地理解分布式系统的概念，我们引入三维空间来表征处理系统的分布特性。它们分别是硬件构成、控制方式和数据分布。

### 1. 硬件构成层次

我们根据大多数系统的硬件构成方式，以处理器为核心划分如下层次。

A. 单 CPU：单一控制器、单一运算器、单一内存。

B. 多执行部件：单一控制器、多运算器或存储器。

C. 具有专用功能部件：单一通用控制器、多运算器或存储器、配有通道/输入输出处理机/向量运算部件/辅助数学运算部件等专用功能部件。

D. 多处理机：多个控制器、多运算器或存储器、单一输入输出系统。

E. 多计算机：多台通用计算机。

显然，上面的硬件分布性是递增的。从分布式系统的观点看，随着硬件分布性的增强，系统才更有希望接近分布式系统的目标。上面的 D 和 E 是构成理想分布式系统的硬件环境。

### 2. 控制方式层次

分布式系统是多节点系统，如何控制这些节点协同工作是关键的问题。因此，一个系统的控制方式是刻画分布式系统的重要方面。在多节点的系统中，控制方式是多种多样的，主要的控制方式归纳如下。

A. 单个控制点：物理上的或概念上的一个控制节点。

B. 固定主从关系：有一个节点是主节点，其他的是从节点。当然也可以构成多级主从关系。这种主从关系是预先规定好的，不能修改。

C. 动态主从关系：可以通过程序修改主从关系。

D. 多个控制节点独立工作：例如，多个计算机最多是在 I/O 级别上交换信息。

E. 多个控制点在某个层次上（如任务分割）协同工作。

F. 多个同构控制点完全协同工作。

G. 多个同构或异构控制点完全协同工作。

上面的控制方式中，F 和 G 是分布式系统所追求的。如果放宽标准的话，E 也可以认为是分布式系统的控制形式。

### 3. 数据分布层次

分布式系统的数据是分布的，每个节点都可能存储本地或全局的数据。根据数据的分布性特点，我们选取下面的层次进行分析。

A. 集中式数据库：在文件及目录上只有单一的拷贝，减少数据存储冗余是它追求的目标之一。

B. 文件分布但中央集中式目录：没有本地目录，所有的访问都必须通过这个中央目录来完成。

C. 重复的数据库拷贝：在每个节点都有一份完整的数据拷贝。

D. 主节点存放完整数据，其他节点存放所需的数据或数据分片。

E. 主节点存放数据分布图或目录，其他节点存放所需的数据或数据分片。

F. 所有节点都存放最需要的数据或数据分片，而且任何节点都能形成对其他节点的访问。

当然理想的分布式系统的数据分布是 F。如果放宽标准的话，C、D 和 E 也可以认为是分布式系统的数据组织形式。

为了更清楚地理解分布式系统的特性，我们根据如上描述，以表征分布性的三维图来刻画（如图 1-1）。在图 1-1 中，深颜色的部分是我们所追求的理想分布式系统的分布表征空间。

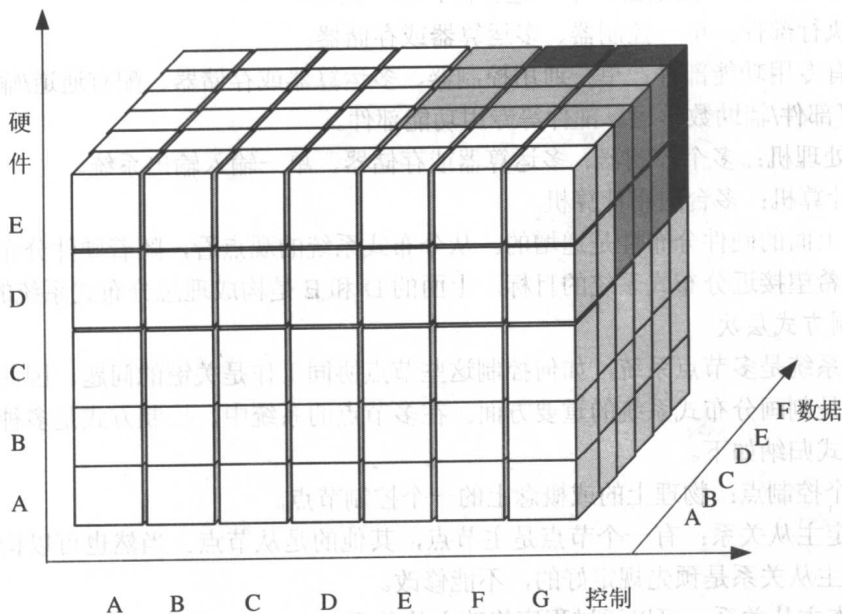


图 1-1 分布特征示意

图 1-1 可以帮助我们直观上区分哪些系统是理想的分布式系统，哪些是准分布式系统，哪些只是使用了部分分布式思想或技术的系统，哪些根本就不是分布式系统。

有这样一些系统，或许是商家在概念上炒作，或许是认识上的偏差或不全面，很容易被联想成分布式系统。

下面我们就几个典型的非分布式系统及其特征加以分析。



### 1. 系统中引入专用处理部件

在系统中引入专用处理部件的系统大多是一台主机和若干用于固定任务的专用部件。像引入通道的计算机、控制通信的前端处理机、向量乘法器、快速傅里叶变换器等属于这种。尽管这些使用很普遍的系统在结构上具有多处理机特性，但是那些专用部件只能完成某些特定的任务，在地位上和主机是不平等的。

### 2. 系统中主从关系明显

在硬件上，系统的各处理部件有不同的地位，有些只是被动地接受分配的任务。在软件控制上，它们执行的不是协作性协议，而是按主从模式工作，严重违反了分布式系统的协作自治性原则。上面介绍的在系统中引入专用处理部件就属于这种情况。另外，具有智能终端的系统也属这类（尽管广告总是把它和分布式处理联系起来）。这样的系统一般由一台处理机和若干终端组成。终端能通过运行处理机上的程序提供智能式输入、共享文件、远程作业提交、与主机通信，但是这类系统的控制是集中式的，终端不能完全根据自己的负载和能力决定任务的实施。

### 3. 简单网络互连结构

以一个具有多个计算机组成的网络互连结构为例，它们有很好的自治性，而且能交换信息。但是，这种协作是有限的。它们无法因为硬件故障来重新分配任务，即使是双机或多机备份系统，离真正的分布式系统仍有距离，因为它们很难同时合作解决一个大问题。

## 1.3 高层操作系统

高层操作系统是分布式系统的一个关键问题。理想的高层操作系统是非层次的，即内部没有任何的主从关系。加之系统的自治性要求，使得分布式系统的控制问题的难度加大。相对于集中式系统，分布式系统必须面对如下问题。

### 1. 不完整系统状态信息

在集中式系统中，总是假设操作系统是在完整而准确的系统状态信息下工作的。但是，在分布式系统中，这种假设是可望而不可及的。这是因为获得这些完整而准确的系统状态信息的代价太大了，根本无法在用户可忍受的范围内收集和整理出这些状态信息并及时得到利用。

### 2. 时间延迟

在集中式系统中，操作系统可以及时请求状态信息而保证被询问的部件在稳定的状态下作出决定。但是，在分布式系统中，由于自治性的局部处理和多任务的交叉作用使发生时间的滞后是必然的。这种时间延迟也会带来系统状态信息的不准确。

### 3. 通信的代价

分布式系统的通信是复杂的，特别是如何使通信的代价降到用户可以忍耐的层次。在单处理机系统中，我们可以使用信号灯、标志、加锁等来解决同步处理问题。但是，在分布式系统中，这些方法会大量消耗时间和降低系统的吞吐能力。解决合理的通信代价问题是分布式操作系统设计的重要目标。

### 4. 负载均衡

在集中式系统中，各部件的任务明确。但是分布式系统是多机或多处理机协同工作的系