

王 珊 等 编著

数据仓库 技术与 联机分析处理

数 据 库 从 书

科学出版社

数 据 库 丛 书

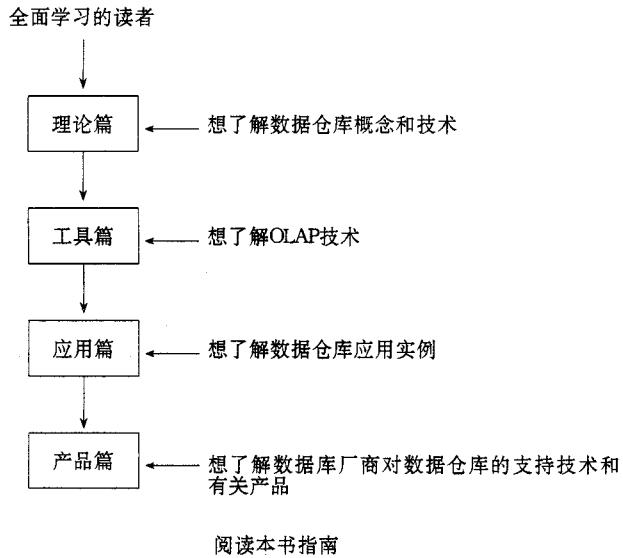
数据仓库技术与联机分析处理

王 珊 等 编著

科 学 出 版 社

1998

中注意做到既使全书各章是一个相互联系的整体，同时又使各篇章能自成一体。读者可以选择其中的某些篇章来阅读。例如，想了解 OLAP 技术的读者可跳过理论篇直接阅读工具篇(如下图所示)。



本书的写作过程也是我组织研究所师生们学习、研讨和实践的过程。两年多来，我们阅读了大量国外的著作、论文，考察和使用了 DBMS 和 OLAP 的有关产品，进行了数十次研讨，写出了数十篇报告和论文，在此基础上进一步充实、整理成为本书。

参加本书研讨和写作的有：王珊、陈红、楼文武、李纪华、周胜、罗立、刘方、张孝、王秋月、李鹏、刘爽、周勇等。特别应指出的是，楼文武在资料整理、内容调整和书稿录入等方面付出了辛勤劳动。

本书收入了中国银行广东省分行建立和应用数据仓库的实例。对国内广大用户来说，他们成功的经验比起国外的应用实例更为亲切，容易借鉴。在此，我向撰写该实例的温巩先生、罗亦文女士、温家明先生等表示衷心感谢。

中国科学院研究生院邵佩英教授审阅了全稿并提出了许多有益的意见，在此向她致以衷心的谢意。

我们在编写本书的过程中，尽可能做到深入浅出，力求概念正确，理论联系实际。由于数据仓库作为一个新的领域，发展非常迅速，加之我们水平有限，故书中一定存在许多不足之处，希望同行和广大读者提出批评和建议。

王 珊

1997 年金秋 于中国人民大学数据与知识工程研究所

内 容 简 介

数据仓库技术(Data Warehousing)和联机分析处理(On-Line Analytical Processing, 简记为 OLAP)是信息领域中近年来迅速兴起的计算机技术。

本书全面而系统地介绍了数据仓库技术和基于数据仓库的 OLAP 应用技术,主要内容包括数据仓库的基本概念、创建技术和方法、数据仓库的体系结构以及投资回报分析(理论篇),数据分析工具、数据分析模型、OLAP 的基本概念、多维数据库、OLAP 的实现技术,以及数据挖掘技术等(工具篇)。本书还在应用篇中给出了数据仓库的若干应用实例,特别是我国自己的应用例子。最后在产品篇中介绍了著名的数据库厂商 Informix, Oracle, Sybase 关于数据仓库的解决方案和相关产品。

本书是学习、掌握和运用数据仓库技术的综合指南,是从事数据库和数据仓库的研究和开发者、设计开发人员,以及需要了解数据仓库实际技术的系统集成人员、系统设计师和有关专业人员的良师益友,也可作为大学高年级学生或研究生相关课程的教材和参考书。

图书在版编目(CIP)数据

数据仓库技术与联机分析处理/王珊 等编著. -北京:
科学出版社,1998.5
(数据库丛书)
ISBN 7-03-006412-7

I. 数… II. 王… III. 数据库系统-信息处理 IV. TP
31 1.13

中国版本图书馆 CIP 数据核字(98)第 02360 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 6 月第一版 开本:787×1092 1/16

1998 年 6 月第一次印刷 印张:16

印数:1—3 100 字数:356 000

定价:30.00 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，对我國数据库科技工作者多年来的主要研究成果，具有較高的理論水平和学术價值。

数据库是计算机科學技术中發展最快的領域之一，也是應用最廣的技术之一。是计算机信息系统与應用系統的構成基础。相信《数据库丛书》的編模出版，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，加快数据库应用的推广与深入，為我國社會經濟信息化作出貢獻。

張致祥

九九年六月

《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 颀
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

近年来,国外信息技术领域中悄然兴起并日益成熟的数据仓库技术引起了国内同仁的广泛重视。1996年7月15日,我们在《计算机世界报》上发表了一组有关数据仓库的文章,引起了学术界、企业界很大反响和浓厚的兴趣。许多同仁来电话、来信,提出了许多问题和要求,希望能结合他们目前的业务,应用数据仓库技术建立决策支持系统,以充分利用现有的数据资源,提取管理决策所需要的信息。在这样的大环境下,本书应运而生。

根据国外数据仓库技术的发展和国内对数据仓库技术的需求,本书全面而系统地阐述了数据仓库的基本概念和方法,介绍了数据仓库作为决策支持系统(DSS)的一种有效而可行的体系化解决方案应包括的三个方面的技术内容:数据仓库技术(Data Warehousing,简记为 DW)、联机分析处理技术(On-Line Analytical Processing,简记为 OLAP)、数据挖掘技术(Data Mining,简记为 DM),以使读者对该领域的理论、技术和应用情况有一个全面的了解。

数据仓库是一种解决问题的方案,而不是可以买到的现成产品。数据仓库以传统的数据库技术作为存储数据和管理资源的基本手段,以统计分析技术作为分析数据和提取信息的有效方法,以人工智能技术作为挖掘知识和发现规律的科学途径。因此,它是诸多学科相互结合、综合应用的技术。本书从数据仓库系统的架构中展示这些技术的作用和相互之间的关系。

本书共十二章,分为理论篇、工具篇、应用篇和产品篇四个部分。

理论篇包括第一至四章,介绍了:数据仓库产生的背景及定义,数据仓库的特征,数据仓库中数据的组织,以及数据库体系化环境;创建数据仓库的方法、模型和步骤;操作数据存储(ODS)的概念、特点,以及 ODS 和数据库、数据仓库的关系,创建 ODS 的两条技术路线等;数据仓库的投资回报定量分析和定性分析,其中列举了 IDC 随机调查的 62 个数据仓库中 3 个有代表性的例子。

工具篇包括第五至七章,论述了数据仓库系统中数据仓库、数据仓库管理系统和数据仓库工具三部分的作用,简要介绍了数据仓库中查询工具、OLAP 分析工具和数据挖掘工具等三类主要的工具(只有通过高效的工具,数据仓库才能真正发挥出数据宝库的作用),介绍了 OLAP 的基本概念、OLAP 产品的 12 条准则和基于多维数据库和关系数据库的两种 OLAP 实现技术,以及数据挖掘的概念、分析方法及其有关技术。

应用篇包括第八、九章,概述了数据仓库的应用实例,这些例子是经过精心挑选的,其中有中国银行广东省分行数据仓库的应用和实践。

产品篇包括第十至十二章,介绍了:Informix 公司对数据仓库的支持技术及前端 OLAP 产品 Metacube;Oracle 公司提供的对数据仓库的解决方案及 OLAP 工具 Express 的功能和产品构成;Sybase 公司所支持的数据仓库的多层次体系结构及其 Sybase IQ 数据库服务器,数据仓库设计工具 Warehouse Architect。

为尽可能面向各个层次及不同部门的数据库理论和应用的工作者,我们在编写过程

目 录

理 论 篇

第一章 从数据库到数据仓库	1
1.1 从数据库到数据仓库	1
1.2 什么是数据仓库	4
1.2.1 主题与面向主题	5
1.2.2 数据仓库的其他三个特征	9
1.3 数据仓库中的数据组织.....	10
1.3.1 数据仓库的数据组织结构.....	10
1.3.2 粒度与分割	11
1.3.3 数据仓库的数据组织形式	13
1.3.4 数据仓库的数据追加	14
1.4 数据库体系化环境.....	15
1.4.1 四层体系化环境	15
1.4.2 数据集市	17
1.5 小结.....	18
第二章 数据仓库设计	19
2.1 数据仓库系统设计方法概述.....	19
2.2 数据仓库设计的三级数据模型.....	22
2.2.1 概念模型	22
2.2.2 逻辑模型	22
2.2.3 物理模型	23
2.2.4 高级模型、中级模型和低级模型	23
2.3 提高数据仓库的性能.....	24
2.3.1 粒度划分	25
2.3.2 分割	26
2.3.3 数据仓库物理设计中的其他一些问题	28
2.4 数据仓库中的元数据.....	31
2.5 数据仓库设计步骤.....	32
2.5.1 概念模型设计	33
2.5.2 技术准备工作	36
2.5.3 逻辑模型设计	37
2.5.4 物理模型设计	40
2.5.5 数据仓库的生成	41

2.5.6 数据仓库的使用和维护	42
2.6 小结	44
第三章 操作数据存储(ODS)	46
3.1 什么是 ODS	46
3.1.1 ODS 的定义及特点	46
3.1.2 ODS 的功能和实现机制	47
3.2 DB-ODS-DW 体系结构	51
3.2.1 ODS 与 DW	51
3.2.2 DB-ODS-DW 三层体系结构	52
3.3 创建 ODS	53
3.3.1 ODS 数据模式的形成过程	53
3.3.2 ODS 对数据的控制——获取并传输	55
3.3.3 创建 ODS 的两条技术路线	58
3.4 实例——商场 ODS 系统	60
3.5 小结	62
第四章 数据仓库投资回报分析	63
4.1 概述	63
4.2 数据仓库投资回报的定量分析	64
4.2.1 投资回报的度量标准	64
4.2.2 数据仓库的投资回报率与回报周期	65
4.2.3 数据仓库投资回报分析	66
4.3 数据仓库投资回报的定性分析	68
4.4 数据仓库实现分析	69
4.4.1 建立数据仓库的必要性分析	69
4.4.2 技术选择分析	70
4.4.3 数据仓库实现方法的投资回报分析	70
4.4.4 数据仓库实现目标的投资回报分析	71
4.5 典型企业的投资回报分析	72
4.5.1 美国麻萨诸塞州政府(Commonwealth of Massachusetts)(ROI 44%)	72
4.5.2 荷兰 Interpolis 公司(ROI 568%)	73
4.5.3 美国 Niagara Mohawk 能源公司(ROI 1413%)	74
4.6 小结	76

工 具 篇

第五章 数据仓库工具	77
5.1 数据仓库工具层——数据仓库系统的重要组成部分	77
5.1.1 数据仓库系统的结构	77
5.1.2 数据库系统与数据仓库系统的组成结构的比较	78

5.2 数据分析工具的发展	79
5.2.1 EIS 软件	79
5.2.2 PC 挖掘工具	79
5.2.3 OLAP 服务器	80
5.2.4 面向数据仓库、支持决策应用的数据分析产品	80
5.3 数据分析模型	81
5.3.1 四种分析模型	81
5.3.2 比较	82
5.4 数据仓库工具简介	82
5.4.1 验证型工具	83
5.4.2 发掘型工具	83
第六章 决策支持工具的新进展——联机分析处理(OLAP)	85
6.1 从 OLTP 到 OLAP	85
6.1.1 OLAP 的出现	85
6.1.2 什么是 OLAP	86
6.1.3 OLTP 与 OLAP 的关系及比较	90
6.2 OLAP 的特征及衡量标准	90
6.2.1 Codd 关于 OLAP 产品的十二条评价准则	91
6.2.2 其他厂商对 Codd 的十二条准则的看法	94
6.3 OLAP 实施	94
6.4 基于多维数据库的 OLAP 实现	95
6.4.1 多维数据	95
6.4.2 维的层次关系和类	97
6.4.3 时间序列数据类型	100
6.4.4 多维数据库存储	101
6.4.5 多维数据库存取	102
6.5 基于关系数据库的 OLAP 实现	102
6.6 两种技术间的比较	105
6.6.1 结构	105
6.6.2 数据存储和管理	106
6.6.3 数据存取	107
6.6.4 适应性	107
6.7 OLAP 产品介绍及选择	108
6.7.1 产品介绍	108
6.7.2 产品选择	109
6.8 OLAP 的新发展及在我国的应用展望	111
6.8.1 OLAP 的新发展	111
6.8.2 OLAP 在我国的应用展望	112
第七章 数据挖掘(Data Mining)工具	113

7.1 Data Mining 的技术基础	113
7.1.1 Data Mining 的概念	113
7.1.2 Data Mining 的方法与技术	114
7.1.3 Data Mining 的分析方法	115
7.2 Data Mining 系统的体系结构及运行过程	119
7.2.1 数据挖掘的步骤	119
7.2.2 实例	120
7.3 从技术到实现	121
7.4 Data Mining 与 OLAP 的区别和联系	122
7.5 数据挖掘的应用	124

应 用 篇

第八章 数据仓库应用谈	126
8.1 数据仓库应用概述	126
8.1.1 全局应用	127
8.1.2 复杂分析	127
8.2 数据仓库的应用实例	128
8.2.1 数据仓库解决“蜘蛛网”问题	128
8.2.2 分层决策体系	129
8.2.3 数据抽样分析	132
8.2.4 发挥历史数据的经济效益	133
8.2.5 回扣分析	134
8.3 小结	136
第九章 数据仓库的应用与实践	138
9.1 任务来源	138
9.2 研制过程	138
9.2.1 前期准备工作	138
9.2.2 总体方案的确立	139
9.2.3 数据模型分析与数据库设计	141
9.2.4 应用系统开发	142
9.3 研制成果	143
9.4 作用与效益	144
9.5 结束语	144
附录 中国银行广东省分行 FMIS 系统	144

产 品 篇

第十章 INFORMIX公司的数据仓库解决方案及其OLAP产品MetaCube技术分

析.....	166
10.1 INFORMIX 数据仓库解决方案	166
10.2 联机事务处理(OLTP)、数据仓库与联机分析处理(OLAP)	168
10.3 Informix 公司数据仓库的数据分析模型——多维模型	170
10.3.1 什么是多维模型	170
10.3.2 多维模型的实现关键——计算中间表的设计	172
10.4 Informix OLAP 产品 MetaCube 介绍及技术分析	177
10.4.1 MetaCube 的技术特色	177
10.4.2 MetaCube Explorer	181
10.4.3 MetaCube Warehouse Manager	182
10.5 MetaCube 使用实例	183
10.5.1 DSS 系统 MetaCube DEMO 的多维模型	183
10.5.2 MetaCube DEMO 的逻辑模型实现	184
10.5.3 通过 MetaCube Explorer 访问 MetaCube DEMO 中的数据	186
10.5.4 利用 MetaCube Optimizer 优化数据仓库	191
10.6 结束语.....	192
第十一章 Oracle 数据仓库解决方案及 OLAP 产品技术分析	193
11.1 ORACLE 数据仓库解决方案	193
11.1.1 数据仓库建模和设计	193
11.1.2 数据抽取	194
11.1.3 数据仓库管理	195
11.1.4 数据分析	196
11.2 Oracle OLAP 产品介绍	197
11.2.1 OLAP 背景	197
11.2.2 OLAP 的两类用户	197
11.2.3 Oracle OLAP 产品系列	198
11.3 Oracle Express Server 技术特色	199
11.3.1 Express Server 结构	199
11.3.2 Express 数据模型	200
11.3.3 Oracle Express 的存储结构	201
11.3.4 Express 多维数据模型的优点	202
11.3.5 Express Server 数据的提取及其与关系数据库的集成	205
11.3.6 SQL 与多维查询的实例	206
11.4 实例.....	207
11.4.1 数据模型定义	208
11.4.2 数据抽取	209
11.4.3 通过 EXPRESS OBJECT 分析	210
11.4.4 总结	213
第十二章 Sybase 的交互式数据仓库解决方案及其特色产品 Sybase IQ	214

12.1 Sybase 的数据仓库三层体系结构	214
12.1.1 多层体系结构的概念与划分	214
12.1.2 三层客户/服务器结构适应数据仓库应用的需要	214
12.2 Sybase 的 QuickStart DataMart 捆绑计划	215
12.2.1 Sybase 数据仓库体系环境	215
12.2.2 数据仓库和数据集市(Data Mart)	217
12.2.3 Sybase 的“WarehouseNOW”策略：Quick Start DataMart	218
12.3 Sybase 特色产品 Sybase IQ 的技术简介	219
12.3.1 Sybase IQ 产品定位	219
12.3.2 Sybase IQ 服务器技术特色	220
12.3.3 Bit-Wise 索引的建立	227
12.4 数据仓库设计工具——PowerDesigner Warehouse Architect 6.0	228
12.4.1 维建模与相关概念	228
12.4.2 Warehouse Architect 功能简介	230
附录 1 设置 Sybase IQ 的基本步骤	231
附录 2 测试比较	231
附录 3 Sybase IQ 的典型用户——美国 MCI 公司的 SOLD 数据仓库	237
参考文献	239

理论篇

第一章 从数据库到数据仓库

随着计算机技术的飞速发展和企业界不断提出新的需求，数据仓库技术应运而生。传统的数据库技术是以单一的数据资源，即数据库为中心，进行从事务处理、批处理到决策分析等各种类型的数据处理工作。然而，不同类型的数据处理有着其不同的处理特点，以单一的数据组织方式进行组织的数据库并不能反映这种差异，满足不了数据处理多样化的`要求。近年来，随着计算机应用，特别是数据库应用的广泛普及，人们对数据处理的这种多层次特点有了更清晰的认识。总结起来，当前的数据处理可以大致地划分为两大类：操作型处理和分析型处理（或信息型处理）。操作型处理也叫事务处理，是指对数据库联机的日常操作，通常是对一个或一组记录的查询和修改，主要是为企业的特定应用服务的，人们关心的是响应时间，数据的安全性和完整性。分析型处理则用于管理人员的决策分析。例如，DSS，EIS 和多维分析等，经常要访问大量的历史数据。两者之间的巨大差异使得操作型处理和分析型处理的分离成为必然。这种分离，划清了数据处理的分析型环境与操作型环境之间的界限，从而由原来的以单一数据库为中心的数据环境发展为一种新环境：体系化环境。

1.1 从数据库到数据仓库

数据库系统作为数据管理手段，主要用于事务处理。在这些数据库中已经保存了大量的日常业务数据。传统的 DSS 一般是直接建立在这种事务处理环境上的。数据库技术一直力图使自己能胜任从事务处理、批处理到分析处理的各种类型的信息处理任务。尽管数据库在事务处理方面的应用获得了巨大的成功，但它对分析处理的支持一直不能令人满意，尤其是当以业务处理为主的联机事务处理（OLTP）应用与以分析处理为主的 DSS 应用共存于同一个数据库系统中时，这两种类型的处理发生了明显的冲突。人们逐渐认识到，事务处理和分析处理具有极不相同的性质，直接使用事务处理环境来支持 DSS 是行不通的。

具体来说，事务处理环境不适宜 DSS 应用的原因概括起来主要有以下五条：

(1) 事务处理和分析处理的性能特性不同。

在事务处理环境中，用户的行为特点是数据的存取操作频率高而每次操作处理的时间短，因此，系统可以允许多个用户按分时方式使用系统资源，同时保持较短的响应时间，OLTP 是这种环境下的典型应用。

在分析处理环境中，用户的行为模式与此完全不同，某个 DSS 应用程序可能需要连续运行几个小时，从而消耗大量的系统资源。将具有如此不同处理性能的两种应用放在同一个环境中运行显然是不适当的。

(2) 数据集成问题。

DSS 需要集成的数据。全面而正确的数据是有效的分析和决策的首要前提，相关数据收集得越完整，得到的结果就越可靠。因此，DSS 不仅需要整个企业内部各部门的相关数据，还需要企业外部、竞争对手等处的相关数据。

事务处理的目的在于使业务处理自动化，一般只需要与本部门业务有关的当前数据。而对整个企业范围内的集成应用考虑很少。当前绝大部分企业内数据的真正状况是分散而非集成的。造成这种分散的原因有多种，主要有事务处理应用分散、“蜘蛛网”问题、数据不一致问题、外部数据和非结构化数据。

上述问题是事务处理环境所固有的，尽管每个单独的事务处理应用可能是高效的，能产生丰富的细节数据，但这些数据却不能成为一个统一的整体。对于需要集成数据的 DSS 应用来说，必须自己在应用程序中对这些纷杂的数据进行集成。可是，数据集成是一项十分繁杂的工作，都交给应用程序完成会大大增加程序员的负担。并且，每做一次分析，都要进行一次这样的集成，将会导致极低的处理效率。DSS 对数据集成的迫切需要可能是数据仓库技术出现的最重要动因。

① 事务处理应用的分散。

当前企业内部各事务处理应用间实际上几乎都是独立的，之所以出现这种现象有多种原因。有的原因是设计方面的，例如，系统设计人员为了减少系统开发费用和加快开发进度，总是采用简单而“有效”的设计方案，这种“有效”仅指对解决当前面临的问题有效，而不能保证对以后新出现的问题继续有效。有的原因是经济方面的，当经费有限时，企业总是考虑先对关键的业务活动建立应用系统，然后再逐步建立其他业务的信息处理系统。还有的原因是历史、地理方面的，例如，某个大公司由分散在各地的多个子公司组成，企业的兼并，等等。

由于这种事务处理应用分散状况的存在，DSS 应用需要对分散在多个事务处理应用中的相关数据进行集成，以向分析人员提供统一的数据视图。

② “蜘蛛网”问题。

DSS 应用中为了避免与其他用户的冲突和简化用户的数据视图，一种称作“抽取程序”的方法目前被广泛地应用，用户利用抽取程序从文件或数据库中查找有用的数据，然后这些数据被提取出来放入其他文件或数据库中供用户使用。这些经抽取得到的新文件或数据库又被某些用户再进行抽取，这种不加控制的连续抽取最终导致系统内的数据间形成了错综复杂的网状结构，人们形象地称为“蜘蛛网”。企业的规模越大，“蜘蛛网”问题就越严重。

虽然网上的任意两个节点的数据可能归根结底是从一个原始库中抽取出来的，但其数据没有统一的时间基准，抽取算法各不相同，抽取级别也不相同，并且可能参考不同的外部数据。因而对同一问题的分析，不同节点却会产生不同甚至截然相反的结果。这当然使决策者无从下手。

③ 数据不一致问题。

前述的应用分散和“蜘蛛网”等多个问题，导致了多个应用间的数据不一致。这些数据不一致的形式是多种多样的：

同一字段在不同应用中具有不同的数据类型。例如，字段 Sex 在 A 应用中的值为

“M/F”，在 B 应用中的值为“0/1”，在 C 应用中又为“Male/Female”。

同一字段在不同应用中具有不同的名字。例如，A 应用中的字段 balance 在 B 应用中名称为 bal，在 C 应用中又变成了 currbal。

同名字段，不同含义。例如，字段 weight 在 A 应用中表示人的体重，在 B 应用中表示汽车的重量，等等。

为了将这些不一致的数据集成起来，必须对它们进行转换后才能供分析之用。数据的不一致是多种多样的，对每种情况都必须专门处理，因此，这是一项很繁重的工作。

④外部数据和非结构化数据。

在决策中经常用到外部数据，这部分数据不是由事务处理系统产生的，而是来自于其他外部数据源。例如，权威性刊物发布的统计数据、业界的技术报告、市场比较和分析报告、股票行情等，这些数据通常都是非结构化数据。在事务处理系统中，由于没有对外部数据进行统一管理，用到这些数据的 DSS 应用必须自行集成。

(3) 数据动态集成问题。

由于每次分析都进行数据集成的开销太大，一些应用仅在开始对所需数据进行了集成，以后就一直以这部分集成的数据作为分析的基础，不再与数据源发生联系，我们称这种方式的集成成为静态集成。静态集成的最大缺点在于，如果在数据集成后数据源中数据发生了改变，这些变化将不能反映给决策者，导致决策者使用的是过时的数据。对于决策者来说，虽然并不要求随时准确地探知系统内的任何数据变化，但也不希望他所分析的是几个月以前的情况。因此，集成数据必须以一定的周期（例如 24 小时）进行刷新，我们称其为动态集成。显然，事务处理系统不具备动态集成的能力。

(4) 历史数据问题。

事务处理一般只需要当前数据，在数据库中一般也只存储短期数据，且不同数据的保存期限也不一样，即使有一些历史数据保存下来了，也被束之高阁，未得到充分利用。但对于决策分析而言，历史数据是相当重要的，许多分析方法必须以大量的历史数据为依托。没有对历史数据的详细分析，是难以把握企业的发展趋势的。

通过(2)、(3)、(4)所述可见，DSS 对数据在空间和时间的广度上都有了更高的要求，而事务处理环境难以满足这些要求。

(5) 数据的综合问题。

在事务处理系统中积累了大量的细节数据，一般而言，DSS 并不对这些细节数据进行分析。这主要有两个原因，一是细节数据数量太大，会严重影响分析的效率；二是太多的细节数据不利于分析人员将注意力集中于有用的信息上。因此，在分析前，往往需要对细节数据进行不同程度的综合。而事务处理系统不具备这种综合能力，根据规范化理论，这种综合还往往因为是一种数据冗余而加以限制。

以上这些问题表明，在事务型环境中直接构建分析型应用是一种失败的尝试。数据仓库本质上是对这些问题的回答。但是数据仓库的主要驱动力并不是过去的缺点，而是市场商业经营行为的改变，市场竞争要求捕获和分析事务级的业务数据。建立在事务处理环境上的分析系统无法达到这一要求。要提高分析和决策的效率和有效性，分析型处理及其数据必须与操作型处理及其数据相分离。必须把分析型数据从事务处理环境中提取出来，按照 DSS 处理的需要进行重新组织，建立单独的分析处理环境，数据仓库正