

# 文献信息数据库 建库技术

戴维民 葛 敏 韩建新  
王兰诚 王绍平 刘永丹 著

北京图书馆出版社

# 文献信息数据库建库技术

戴维民 葛 敏 韩建新 著  
王兰诚 王绍平 刘永丹

北京图书馆出版社

图书在版编目(CIP)数据

文献信息数据库建库技术/戴维民等著. —北京:北京图书馆出版社, 2001. 6

ISBN 7-5013-1777-1

I. 文… II. 戴… III. 文献数据库 - 基本知识 IV. G356.1

中国版本图书馆 CIP 数据核字(2001)第 027180 号

---

书名 文献信息数据库建库技术

著者 戴维民等著

---

出版 北京图书馆出版社 (100034 北京西城区文津街 7 号)

发行 (010)66126153 传真 (010)66174391

E-mail Btsfxb@publicf.nlc.gov.cn

经销 新华书店

印刷 北京双桥咸宁侯印刷厂

---

开本 850×1168 毫米 1/32

印张 11.3125

字数 280(千字)

版次 2001 年 6 月第 1 版 2001 年 6 月第 1 次印刷

印数 1—3000

---

书号 ISBN 7-5013-1777-1/G·486

定价 22.00 元

# 目 录

<b>第1章 绪论</b> .....	(1)
<b>1.1 数据库的概念</b> .....	(1)
<b>1.1.1 数据库的定义</b> .....	(1)
<b>1.1.2 数据库的特点</b> .....	(1)
<b>1.2 数据库的类型</b> .....	(2)
<b>1.2.1 按数据库性质划分</b> .....	(2)
<b>1.2.2 按数据库的结构模型划分</b> .....	(8)
<b>1.2.3 数据库的其它类型</b> .....	(9)
<b>1.3 数据库建设的意义</b> .....	(9)
<b>1.3.1 数据库建设是信息产业的核心</b> .....	(9)
<b>1.3.2 数据库建设是图书馆自动化管理的标志</b> .....	(10)
<b>1.3.3 数据库建设是文献信息网络化的基础</b> .....	(10)
<b>1.3.4 数据库建设是信息资源共享的前提</b> .....	(10)
<b>1.3.5 数据库建设是军队军事训练信息网建设的主体工程</b> .....	(11)
<b>1.4 数据库建库模式</b> .....	(11)
<b>1.4.1 集中模式</b> .....	(11)
<b>1.4.2 分散模式</b> .....	(12)
<b>1.4.3 集中和分散并行模式</b> .....	(13)
<b>1.5 数据库技术和数据库产业的发展</b> .....	(14)
<b>1.5.1 国外数据库技术和数据库产业的发展</b> .....	(14)
<b>1.5.2 我国数据库建设的发展</b> .....	(17)
<b>1.5.3 数据库技术和数据库产业的发展趋势</b> .....	(20)
<b>第2章 文献数据库建库的一般技术</b> .....	(24)

<b>2.1 数据库的设计与系统分析</b>	(24)
2.1.1 数据库的选题论证	(24)
2.1.2 数据库的命名	(24)
2.1.3 数据库系统分析与设计	(25)
2.1.4 数据库建库条件与可行性分析	(28)
<b>2.2 数据准备</b>	(29)
2.2.1 入库文献的选择与确定	(29)
2.2.2 数据加工	(31)
<b>2.3 数据录入</b>	(34)
2.3.1 数据录入方式	(34)
2.3.2 输入数据的具体规定	(34)
2.3.3 输入数据质量要求	(35)
2.3.4 微数据库建设及套录	(35)
<b>2.4 数据库的应用软件开发</b>	(36)
2.4.1 数据库应用软件的开发平台	(36)
2.4.2 数据库应用软件的开发	(37)
2.4.3 数据库应用软件的推广	(37)
2.4.4 数据安全	(38)
<b>2.5 数据库技术规范的编写</b>	(38)
2.5.1 数据库技术规范的基本要求	(38)
2.5.2 数据库技术规范的技术内容	(39)
2.5.3 数据库技术规范的编排格式	(40)
<b>第3章 书目数据库建库技术</b>	(41)
<b>3.1 书目数据库概述</b>	(41)
3.1.1 什么是书目数据库	(41)
3.1.2 书目数据库的发展	(42)
3.1.3 书目数据库的结构	(45)
<b>3.2 机读目录</b>	(46)

3.2.1	机读目录格式	.....	(46)
3.2.2	USMARC 格式	.....	(48)
3.2.3	CNMARC 格式	.....	(53)
3.3	<b>计算机编目系统</b>	.....	(57)
3.3.1	计算机编目系统的类型与结构	.....	(57)
3.3.2	计算机编目系统的建立	.....	(59)
3.4	<b>书目数据库的建库过程</b>	.....	(62)
3.4.1	标准源数据的套录(downloading)	.....	(62)
3.4.2	原始编目	.....	(72)
3.4.3	回溯转换	.....	(74)
3.4.4	机读规范档及其他辅助文档的建设	.....	(79)
3.5	<b>书目数据库的检索</b>	.....	(84)
3.5.1	OPAC(联机公共目录系统)	.....	(84)
3.5.2	WebPAC	.....	(87)
3.5.3	Z39.50 检索协议	.....	(88)
3.6	<b>网络信息资源的编目</b>	.....	(90)
3.6.1	从馆藏资源到网络信息资源	.....	(90)
3.6.2	容纳网络信息资源的 MARC 格式	.....	(92)
3.6.3	元数据与 Dublin Core	.....	(93)
3.7	<b>编目专家系统</b>	.....	(94)
3.7.1	OCLC 书名页自动编目计划	.....	(95)
3.7.2	选取检索点的“有向树形图”概念模型	.....	(98)
第 4 章	<b>文摘数据库建库技术</b>	.....	(101)
4.1	<b>文摘、文摘工作与文摘数据库</b>	.....	(101)
4.1.1	文摘概念	.....	(101)
4.1.2	文摘工作	.....	(102)
4.1.3	文摘数据库	.....	(104)
4.2	<b>文摘类型</b>	.....	(105)

4.2.1	报道性文摘(informative abstracts) .....	(105)
4.2.2	指示性文摘(indicative abstracts) .....	(107)
4.2.3	报道/指示性文摘(informative-indicative abstracts) .....	(108)
4.2.4	作者文摘和文摘员文摘 .....	(109)
4.2.5	文章型文摘、电报型文摘和逻辑文摘 .....	(110)
4.2.6	同址文摘和异址文摘 .....	(111)
4.2.7	首次文摘和二次文摘 .....	(111)
4.2.8	其他文摘类型 .....	(112)
4.3	文摘编写技术 .....	(112)
4.3.1	文摘编写的一般原则 .....	(112)
4.3.2	文摘编写的方法步骤 .....	(114)
4.3.3	文摘要素分析 .....	(117)
4.3.4	文摘语言控制 .....	(119)
4.4	文摘刊物的编制 .....	(123)
4.4.1	文摘刊物及其结构 .....	(123)
4.4.2	文摘刊物的编制步骤 .....	(125)
4.4.3	文摘刊物的技术编辑 .....	(127)
4.5	检索期刊条目著录 .....	(131)
4.5.1	检索期刊文献条目类型 .....	(131)
4.5.2	文摘款目的著录项目 .....	(133)
4.5.3	著录格式 .....	(133)
4.6	文摘数据库的结构与功能 .....	(135)
4.6.1	文摘数据库的结构设计 .....	(135)
4.6.2	文摘数据库功能设计 .....	(138)
4.7	检索刊物数据库化及刊库一体化技术 .....	(141)
4.7.1	检索刊物数据库化 .....	(141)
4.7.2	刊库一体化技术 .....	(143)

<b>第 5 章</b>	<b>全文数据库建库技术</b>	(145)
<b>5.1</b>	<b>全文数据库概述</b>	(145)
5.1.1	什么是全文数据库	(145)
5.1.2	全文数据库的发展	(146)
5.1.3	全文数据库的结构	(148)
5.1.4	全文检索系统概述	(149)
<b>5.2</b>	<b>全文数据库的开发步骤与技术</b>	(153)
5.2.1	数据准备	(153)
5.2.2	文本预处理	(157)
5.2.3	数据加载	(158)
5.2.4	数据检索	(159)
5.2.5	数据维护	(159)
<b>5.3</b>	<b>全文数据库的数据结构</b>	(160)
5.3.1	全文数据库的组织	(160)
5.3.2	全文数据库结构的定义	(161)
5.3.3	几种倒排文档结构	(162)
5.3.4	非结构化文档数据库建库	(165)
<b>5.4</b>	<b>全文数据库的标引索引技术</b>	(174)
5.4.1	全文文本的标引处理	(174)
5.4.2	词典法标引	(175)
5.4.3	单汉字标引	(176)
5.4.4	全文数据库的特殊标引	(178)
5.4.5	全文数据库的索引	(180)
<b>5.5</b>	<b>全文数据库的检索技术</b>	(182)
5.5.1	全文检索系统的功能和实现	(182)
5.5.2	位置逻辑检索	(186)
5.5.3	后控词表检索	(187)
5.5.4	期望值与加权检索	(189)

5.5.5 在非结构化数据库中进行全文检索 .....	(190)
5.5.6 中文全文检索软件的选择 .....	(195)
5.6 全文数据库信息处理系统及技术研究 .....	(198)
5.6.1 信息检索系统结构 .....	(198)
5.6.2 TRS全文信息检索系统 .....	(201)
5.6.3 主要研究的课题 .....	(205)
第 6 章 数值数据库和事实数据库建库技术 .....	(209)
6.1 数值数据库和事实数据库的类型 .....	(209)
6.1.1 数值数据库的类型 .....	(209)
6.1.2 事实数据库的类型 .....	(209)
6.2 数值数据库建库技术 .....	(211)
6.2.1 数值数据库的数据与结构 .....	(211)
6.2.2 数值数据库的生成与更新 .....	(218)
6.3 事实数据库建库技术 .....	(219)
6.3.1 事实数据库的结构 .....	(219)
6.3.2 事实数据库的功能 .....	(221)
第 7 章 多媒体数据库建库技术 .....	(223)
7.1 多媒体数据库概述 .....	(223)
7.1.1 多媒体数据的特点 .....	(223)
7.1.2 多媒体数据库研究动态 .....	(225)
7.1.3 多媒体数据库的关键技术 .....	(229)
7.2 多媒体数据库建库过程 .....	(231)
7.2.1 建立数据模型 .....	(232)
7.2.2 多媒体数据库管理系统(MMDBMS)的选定 .....	(233)
7.2.3 多媒体数据的采集与制作 .....	(233)
7.2.4 多媒体数据的录入 .....	(233)
7.3 多媒体数据模型与多媒体数据库管理系统 .....	(234)
7.3.1 建立数据模型 .....	(234)

7.3.2	多媒体数据库管理系统(MMDBMS)的选定	.....	(241)
7.4	多媒体数据的制作与处理	.....	(244)
7.4.1	多媒体数据的分类	.....	(244)
7.4.2	文本数据	.....	(245)
7.4.3	图像数据	.....	(246)
7.4.4	声音数据	.....	(252)
7.4.5	视频数据	.....	(261)
7.4.6	动画	.....	(267)
7.5	多媒体数据的录入与查询	.....	(272)
7.5.1	多媒体数据库与 OLE 技术	.....	(272)
7.5.2	多媒体数据库与 MCI 接口	.....	(276)
7.5.3	前端开发工具实现多媒体数据查询	.....	(279)
7.5.4	基于内容检索技术	.....	(282)
第 8 章	数据库质量控制和评价	.....	(285)
8.1	数据库质量的要素	.....	(285)
8.1.1	数据质量	.....	(285)
8.1.2	系统质量	.....	(287)
8.2	影响数据库质量的因素	.....	(288)
8.2.1	人员素质	.....	(288)
8.2.2	数据源选取	.....	(289)
8.2.3	软件系统	.....	(289)
8.2.4	工作流程	.....	(289)
8.2.5	组织管理	.....	(290)
8.2.6	工作环境	.....	(290)
8.3	数据库建设的标准化控制	.....	(290)
8.3.1	组织工作的标准化	.....	(292)
8.3.2	数据的标准化	.....	(292)
8.3.3	数据库环境的标准化	.....	(297)

8.4	<b>数据库建设的质量控制</b>	(298)
8.4.1	数据库的全面质量控制	(298)
8.4.2	数据库质量控制技术方法	(300)
8.4.3	数据库质量控制的组织保障	(306)
8.5	<b>数据库的评价</b>	(309)
8.5.1	评价指标	(309)
8.5.2	评价方法	(314)
<b>第9章 数据库服务方式和检索方法</b>		(316)
9.1	<b>数据库服务方式</b>	(316)
9.1.1	定题服务	(316)
9.1.2	为特定用户制作微数据库	(317)
9.1.3	联机检索	(318)
9.1.4	网络检索	(319)
9.2	<b>检索方法</b>	(325)
9.2.1	词检索	(325)
9.2.2	布尔逻辑检索	(326)
9.2.3	截词检索	(329)
9.2.4	位置逻辑检索	(331)
9.2.5	加权检索	(335)
9.2.6	字段与范围限定检索	(337)
9.2.7	辅助检索功能	(339)
<b>参考文献</b>		(349)
<b>后记</b>		(352)

# 第1章 絮 论

## 1.1 数据库的概念

### 1.1.1 数据库的定义

现实世界充满着信息。所谓信息泛指反映特定事物的消息、情报或知识，这种消息、情报和知识以可被感受的声音、文字、图像、符号等为表征，并可通过各种方式传播。信息是对客观事物的反映。

以特定的符号表达的信息称为数据。数据是信息的一种量化表示，数据反映信息，而信息依靠数据来表达。表达信息的符号可以是数字、文字或图形。计算机所能接受、处理和存放的信息须是数字化的信息，因此，必须人为地把信息转换成可以被计算机接受的数据，也就是以二进制形式存储在计算机内并被计算机加工处理的数据。

一定数量的数据存储于特定的存储介质之上，以特定的结构相互联结于一体，形成了数据库。数据库是以一定的组织方式存储在一起，相互关联而独立于应用程序之外，并能为多个用户所共享的数据集合。

### 1.1.2 数据库的特点

#### 1. 有序性

数据库不仅存储数据,而且连同数据之间的逻辑关系一起存储,也就是说数据库中的所有信息不是无序的堆积,而是相互之间存在着一定的逻辑关系,构成一个有序的信息系统。正是这种有序性,使得用户可以方便地查检和利用数据库中的数据。

## 2. 共享性

数据库的数据一经输入,即可多次多重输出,可以为多个用户所共同利用,甚至是同时利用。有效地克服了信息传递的时空障碍。

## 3. 独立性

数据的逻辑组织和数据的物理存储方式与用户的应用程序无关,一旦数据结构改变,与这些数据有关的程序不需要重新编写和调试,可节省大量的开支。这样,也使得同一个数据库可以被不同的应用程序所使用。

## 4. 完整性

数据库所聚集的数据应全面地反映特定主题范围内的相关信息,形成完整的覆盖面,而且通过不间断的数据维护活动,始终保持数据的这种完整和正确。

## 5. 最小冗余性

数据库数据之间多维的逻辑结构使得每条数据可以通过不同的途径得以反映,从而避免了数据的重复存储,最大限度地减少了数据的冗余度,节省空间和数据更新的开支,同时保证了数据的一致性,提高了信息存储利用的效益。

## 1.2 数据库的类型

### 1.2.1 按数据库性质划分

按其提供数据的性质,数据库可以分为文献数据库、数值数据  
2、

库、事实数据库和多媒体数据库。

## 1. 文献数据库

文献数据库是指报道各类文献及文献信息的数据库,包括书目数据库和全文数据库。

### (1) 书目数据库

书目数据库是只存储有关主题领域各类文献资料的书目信息,以二次文献的形式报道文献的数据库,如题录数据库、文摘数据库、引文数据库、期刊目次数据库以及图书馆馆藏目录数据库等。书目数据库以简略的形式向用户报道文献的信息,提供查找、获取文献的线索。这类数据库信息量大,信息密度高,文献报道范围广,数据连续性、累积性强,是用户快速查找文献的有效工具。

书目数据库中的数据来源于期刊论文、会议论文、研究报告、专利文献、学位论文、图书、政府出版物、报纸等各种不同的一次文献,是经过加工、压缩的派生性数据。

书目数据库的数据结构一般比较简单,记录格式也较为固定。因此,其建库速度往往比较快,建设费用相对比较低。

文摘索引数据库和图书馆馆藏目录数据库是最常用的书目数据库。

文摘索引数据库的内容性质与书本式文摘索引相同,主要是简要地通报有关领域特定时期内发表的文章,供人们查阅与检索。它提供确定的文献来源信息,即能准确鉴别相对应的原始文献。但是,它一般不提供原始文献的收藏地点。

图书馆馆藏目录数据库又称为“机读目录”,即 MARC(Machine-Readable Catalogue)。它主要报道特定图书馆实际收藏的各种文献资料的书目信息和存储地址。它既是一般用户查找图书馆资料的工具,又是图书馆业务部门的业务管理工具。它所报道的数据内容详细,除描述文献本身以外,还有许多附加信息,如业务加工信息、管理信息、馆藏信息等,记录格式也比较统一。我国

各类图书馆的馆藏目录数据库都按照《中国机读目录通讯格式》(CNMARC)著录。

## (2) 全文数据库

全文数据库是存储文献全文或其中主要部分,以一次文献的形式直接提供文献的源数据库。通常将经典著作、法律条文及案例记录、重要的科技期刊、新闻报道和百科全书、手册、年鉴,以及图书馆所藏的其他重要文献的全部文字或主要文字转换成计算机可读形式,建成数据库。用户可以从中直接检出所需的原始文献。

与书目数据库相比,全文数据库有许多特点和优点,主要表现在以下方面:

①检索直接。能直接检索出原始文献或解决问题所需的文献资料,不必进行二次检索(即根据检出的书目信息再去查找原文)。

②报道详尽。文献的正文部分或附属部分都可以检索获得和显示,用户可以直接查看到文献正文中的每一段、每一句和每一个词,甚至还可以看到某些边缘性信息。

③传递快速。用户可以通过检索系统迅速地浏览、检索和获得文献原文,不受地理位置的限制。

④标引全面,使用简单。绝大多数全文数据库多利用计算机进行全文自动抽词标引,生成倒排档。能为用户提供标题、著者、关键词、摘要等多重检索途径,使用方便。

⑤用户接口多为菜单驱动型,或采用较简单的检索命令,易学易用。

⑥检索语言多用自然语言,少数用受控语言。检索方法除使用布尔检索以外,位置检索占有相当突出的地位。

而且,随着信息处理技术和储存技术的发展和突破,原先困扰全文数据库发展的存储空间、传递速度、文本转换和信息内容稳定性、可靠性等技术问题已得到基本解决,全文数据库获得了迅速的发展和普及,成为数据库技术和产业的重要发展方向。

## 2. 数值数据库

除了以文献单元为存储和处理单元的文献数据库外,还有一些以数值、事实、术语或图像为存储和处理单元的数据库,即非文献数据库。非文献数据库包括数值数据库、事实数据库、术语数据库和多媒体数据库等。

数值数据库是一种以自然数值形式表示、计算机可读的数据集合,主要记录和提供特定事物的性能、数量特征等信息。其信息报道范围常覆盖某一大类的专业领域。如在军事上,部队的实力、装备、武器的性能等信息;在商业和经济领域中,特定产品的性能特征、价格趋势、国家经济增长率等数值信息;在科技领域,物质的物理化学性质、结构、频谱及实验数据、计算公式等。

数值数据库所记录的数值数据是人们从文献资料中分析提炼出来的,或是从实验、观测和统计工作中直接得到的。数据库生产者将这些数据收集起来,经过核实、检验和加工整理,按一定的方式组织起来,利用计算机进行存储和检索,就成了数值数据库。如果数据库中还含有定义数值或说明这些数据项所必需的文字(文本数据),则又称为文本—数值数据库。与文献数据库相比,数值数据库是人们对信息进行深度加工的产物,它可以直接提供解决问题时所需的数据,是进行各种统计分析、定量研究、管理决策和预测的重要工具。

数值数据库对所收录的数据的可靠性要求比较高,有时还需要列出数据的误差估计、数据来源和实验条件等,以减少用户使用中的误差。数值数据库存储的数据通常成组排列,它们本身并不被检索,但它们与字母或数字形式表达的可检索的关键词或叙词相联结,成为可检索的数据。数值数据库检索的结果是一个或一组特定的数值。

为了满足用户的不同需要,许多数值数据库还常附有特定的检索软件包,以提供数值的运算、统计分析、分类、排序和重组等功能。

能。这就使数值数据库系统除了具备一般检索功能外,还常能提供如下一些功能:

(1)数据运算能力。可以根据现有数据建立函数关系,进行数值的计算,例如可以要求系统以某一年为基准年,重新确定时序的指数;要求系统进行不同计量单位的换算;要求计算时达到某种精确度,等等。

(2)图形处理能力。例如,有的系统能提供分子结构式或晶体结构图、光谱图等,有的可以处理非线性座标等。

(3)报表生成能力。能根据用户的需要自动生成各种报表,或对输出数据进行排序和重新组织,或允许转存数据,建立新的文档。

(4)数据分析能力。例如,为复杂的现象构造仿真模型;或利用各种数值分析技术进行预测、可靠性计算和假设检验;甚至可利用对数据的观察和统计分析产生新的思想等。

(5)试验模拟能力。能够模拟具有数学模式的试验;模拟化学反应过程等。

### 3. 事实数据库

事实数据库是存储有关事物(如人物、机构、事件等)的一般指示性描述的参考数据库。因此又称之为“指示性数据库”或“指南数据库”。其主要用途是供用户查询特定事物的发生时间、地点、过程等简要情况。

事实数据库的类型很多。按信息类型划分,常有收录各种人物传记信息的人物传记数据库,收录各种公司的生产与经营活动信息的公司名录数据库,存储各种基金信息的基金指南数据库,存储各种技术标准或规程信息的技术标准指南数据库,存储各种计算机软件目录信息的软件数据库,存储各种产品和商品信息的产品指南数据库以及面向各行业管理的管理信息数据库,存储各种武器装备的技术性能、制造技术、使用方法的武器装备数据库等。