

电子计算机档案检索

(上 册)

邱 晓 威 编著

兵器部档案馆编印

一九八五年七月

前 言

档案自动化是一个较新的课题，它涉及电子计算机技术，信息系统、直至档案学理论等许多学科和专门技术。档案自动化的理论还有待于在实践基础上的研究和总结。可以说目前正处在一个实际中迫切需要解决，而理论研究上较为欠缺的时期。本讲义的编写者参与了某些档案自动化实验性系统的设计工作以及参加了与自动化有关的国家标准的制定和起草工作，对档案自动化的理论作了初步探讨，有义务向其他档案工作者提供这方面的经验和知识。

档案自动化知识涉及面广，在本书深入讨论每一种专门技术和理论是很难做到的。为适应档案工作者的实际水平，本讲义有意避开一些较难的数学知识和较专门的计算机技术问题，尽量以概念的说明代替繁琐的论证。对于不承担软件设计的档案工作者来说，讲义中的某些具体算法可以略去不读，但是对于一般的粗略流程图，则建议读通。

本讲义的编写过程中，从王永成、戴璞、孙钢、姜世华等同志处得到很大的帮助和指导，在此一并致谢。

由于时间仓促，编者水平所限，有疏漏不妥之处欢迎指正。

邱 晓 威

一九八四年十月

内 容 提 要

本书是档案电子检索课程的专业教材。书中对档案自动化进行了分析，概括性地讲解和论述了档案自动检索涉及的各类问题。内容共分十章，分别对有关的计算机知识、文献信息的存贮和查找技术、检索系统的设计和与档案工作有关的标准化等问题进行了介绍。考虑到档案工作者急需，本书后面附有一套自动检索算方法图，可直接用于编制编目和检索程序。

本书可作为档案工作者学习档案自动化入门读物，亦可作为档案检索应用软件设计人员和档案标引人员的参考手册。

目 次

上册

第一章	概述	1
§ 1—1	信息与社会的发展	1
一	信息——三大资源之一	1
二	信息技术的进步	5
§ 1—2	新技术与档案	6
一	档案业务的发展对新技术的要求	6
二	新技术对档案的影响	8
§ 1—3	档案自动化现状分析	10
一	档案自动化现状	10
二	现状分析	11
第二章	档案检索使用的电子计算机	13
§ 2—1	电子计算机的发展及应用	13
一	电子计算机发展的几个阶段	13
二	电子计算机的应用	14
三	电子计算机的基本结构	16
四	指令和程序的概念	18
五	电子计算机语言的种类	20

§ 2—2 档案检索系统对计算机的要求.....	2 3
一 计算机运算速度	2 3
二 多道程序和分时系统	2 3
三 多级存贮器系统	2 4
四 多样化的输入输出设备	2 4
五 采用虚拟存贮器	2 4
六 可靠性、可用性和可维修性	2 5
七 灵活的扩充性	2 5
八 字长和指令系统	2 6
§ 2—3 计算机硬件系统的构成	2 6
一 微型计算机系统	2 6
二 集中式多终端系统	2 7
三 计算机网络系统	2 8
§ 2—4 存贮器	3 0
一 内存贮器	3 0
二 外存贮器	3 1
§ 2—5 输入输出设备	3 3
一 输入设备	3 3
二 输出设备	3 4
§ 2—6 计算机的软件系统简介	3 5

一 操作系统	3 5
二 文献信息检索系统软件体系	3 6
第三章 自动化的数学基础	3 9
§ 3—1 数理逻辑简介	3 9
一 命题演算	3 9
二 谓词演算	4 9
§ 3—2 集合论基本知识	6 2
一 概念	6 2
二 集合的表示方法	6 3
三 数理逻辑与集合论中常用符号	6 3
四 集合间的关系	6 4
五 集合的运算	6 6
六 集合运算的若干性质	6 8
七 关系	6 9
八 映射	7 1
§ 3—3 图论基本知识	7 3
§ 3—4 概率论中的一个常用公式	7 6
§ 3—5 离散模糊数学的概念	7 8
一 模糊子集的概念	7 8
二 模糊子集的定义	7 9

三 模糊子集的简单运算	80
第四章 档案检索系统的自动化	83
§ 4—1 手工档案检索系统分析	83
§ 4—2 自动编制专题目录的系统	88
§ 4—3 联机自动检索系统	91
第五章 数据结构与文档、数据库	95
§ 5—1 数据结构简介	95
一 引言	95
二 表格	98
三 链表	99
四 图	101
五 树	103
六 二叉树	105
§ 5—2 文档	110
一 文档概念	110
二 文档的结构	112
§ 5—3 国家标准G B—8206《文献目录 信息交换用磁带格式》介绍	114
一 该标准的用途	114
二 目录信息磁带格式介绍	114

三 该标准的特点 119

§ 5—4 数据库 119

一 数据库系统的构成 119

二 数据库的类型 121

三 数据库管理系统 122

下册

第六章 档案检索系统程序设计 129

§ 6—1 程序设计基本知识 129

一 程序和程序设计 129

二 程序设计基本知识 129

三 程序设计语言 132

四 程序(算法框图)的验证 134

五 算法的优化 135

§ 6—2 常用数据处理技术介绍 135

一 堆栈和排队 136

二 紧排顺序存贮结构 143

三 链式存贮结构 149

四 合并与排序 154

五 表查找 157

六 直接存取与 Hashing 函数 161

§ 6—3	自动检索系统的算法	162
一	检索方法和策略	163
二	档案机读目录主文档的建立	166
三	提问式的校验	172
四	源提问加工成目标提问的方法	174
五	顺序文档检索系统	183
六	提问式的波兰倒记法	184
七	倒排文档检索系统	186
第七章	自动检索系统设计	188
§ 7—1	系统工程简介	188
一	系统工程的形成	188
二	系统工程的概况	189
§ 7—2	档案检索系统设计	191
一	档案检索系统的整体性	192
二	自动检索系统研制方式	193
三	档案自动检索系统的可行性调研	194
四	构造系统数学模型	199
五	档案自动化检索系统的线性规划	200
六	方案的决策	205
七	系统实施	205

八 系统评价	2 1 0
第八章 档案自动检索系统实例介绍	2 1 2
§ 8—1 美国档案检索自动化系统实例	2 1 2
§ 8—2 中央档案馆 Z D B J 档案自动编目 检索系统简介	2 1 6
第九章 档案检索系统近期改进与未来发展	2 2 1
§ 9—1 档案自动检索系统的改进	2 2 1
一 系统改进的必要性	2 2 1
二 系统改进的可能性	2 2 1
三 档案自动检索系统的改进原则和方法	2 2 2
§ 9—2 档案自动检索系统的未来发展	2 2 3
一 类自然检索简介	2 2 3
二 智能检索简介	2 2 4
第十章 档案工作与标准化	2 2 6
§ 10—1 标准化及其产生和发展	2 2 6
一 标准的概念	2 2 7
二 标准化	2 3 0
§ 10—2 标准化的原理	2 3 2
一 简化原理	2 3 2
二 统一原理	2 3 4

三 最优化原理	2 3 8
§ 10-3 标准的制订和修订	2 3 9
一 制订、修订标准的原则	2 3 9
二 制订标准的一般程序	2 4 3
三 标准的修订	2 4 5
§ 10-4 标准的贯彻执行	2 4 6
一 贯彻标准的重要性	2 4 6
二 贯彻标准的一般程序	2 4 7
§ 10-5 文献工作标准化	2 4 8
一 文献工作标准化的特点	2 4 8
二 文献工作标准化的组织	2 4 9
§ 10-6 档案标准化	2 4 9
一 档案标准化现状	2 4 9
二 档案检索语言的标准化	2 5 0
附录 1 文献检索自动化常用词汇注释	2 6 0
附录 2 参考文献	2 6 3
附图 1 提问逻辑式语法检查流程图	2 6 4
附图 2 逻辑式的“逆波兰表示”格式检查流程图	2 6 5
附图 3 把提问逻辑式展开为表的流程图	2 6 6

附图 4	主文档记录与提问比较粗框图	2 6 7
附图 5	提问检索词与文档中具有相同标识号 的被检索项的比较流程图	2 6 8
附图 6	提问式转换成波兰式倒记流程图	2 6 9
附图 7	提问的波兰倒记式变成一系列广义 查找指令的流程图例子	2 7 0
附图 8	执行广义查找指令实现倒排文档查找 流程图	2 7 1
附图 9	由顺序文档生成倒排文档流程图	2 7 4
附图 10	分割式排序算法流程图	2 7 5

第一章 概述

§ 1—1 信息与社会的发展

一、信息——三大资源之一

社会的经济、政治、科学、文化等活动与信息有着密切的关系。信息与能源、材料被人们统称为三大资源。

1、信息资源的性质

信息是指对接受者来说预先不知道的报道。信息又是所有容易获得或不易获得的事实和思想的总合，並可在某个时刻供人参考。

知识是人们对于事实和思想的提炼和总结，即在实践活动中获得的对客观事物的认识。知识通常根据实践的不同门类划分成许多学科。

智慧是对知识的综合，是由某些学科的知识升华而成但又超越学科界的理论。可以说信息编织了知识，而知识又升华为智慧。

从使用的角度来分，信息又可以分为二类，一类是直接信息，一类是隐含信息。直接信息是指接收者可以直接使用和吸收的信息，例如各种资料，史料文献各种数据和报导等。隐含信息则是指包含于某一事物中的信息，需要经过分析和考证才可揭示出来。例如一种新式武器，一件发生的故事等等。

现代社会正处于“信息化”的过程之中，在较先进的国家里，从事信息工作的人已达到百分之五十。而且一些经济学家指出，生产率

的提高大部分是靠新知识或新信息。由于信息与物质生产已不可分割，所以信息与物质的关系引起了人们很大的重视。並着重研究了信息这一特殊资源的性质。

能源和材料資源是有形的，而信息是无形的資源。有时信息的载体是有形的，但並不能说信息本身有形。

信息的性质可归纳如下：

(1) 信息是活的。第一，信息之所以可以发挥作用，是由于人脑对它进行观察、记忆、思考、归纳。所以信息是人脑的输入物和输出物。第二，如同生物有生有死一样，信息也有“生”和“死”。“生”是指信息的产生和其价值的存在；“死”是指其价值的消失或信息的失效。

(2) 信息是可扩充的。尽管一部分信息会随着时间流逝而失效，但大部分信息在使用过程中是不断扩充的。这种扩充性体现在现有的信息不可能包含一切事物，另一方面不管怎样广泛收集、深入考察也不可能将反映某一事物的信息完整无缺地得到。

(3) 信息是可浓缩的。信息资源为了便于应用可以被浓缩、综合和归纳。例如许多事物信息反映的规律可总结成一条定理，或者将许多数据信息之间的关系压缩成一个方程式等等。

(4) 信息具有替代能力。信息可以代替资本、工人或其它有形物质。只要考察一下工厂自动化、办公自动化的工作情况就可以认识到

信息的替代能力。

(5) 信息可以以光速传递。电子通讯和通讯技术使得信息按光速进行传递，显然信息的传递速度大大超过了物质的运输速度。

(6) 信息是扩散性的。一般来讲，普通物质的保管和封锁可以做得很严密，而对于信息的保密就困难得多。信息的扩散过程，通常也就是其发挥作用的过程。

(7) 信息是可以分享的。信息一般不能作转手交易，只能作分享交易。其它物质资源在卖给别人、或送给别人后，自己就不再享有它了。而信息传递给别人后，本人则仍然享有它。

关于信息资源的特性，人们会不断有新的发现，并更加广泛而合理地使用它。

2、信息资源的作用

信息资源的影响目前已经极为广泛，并且可以预见对将来的影响更加巨大。有人已经提出了“信息社会”的各种模式。

在信息社会中，图书、情报资料、档案是信息资源的一部分，可以统称为文献信息。图书馆、情报资料馆、档案馆可以看做信息库，或信息源。对这部分信息的重要性，人们是非常清楚的。有人统计，一个人如果上完大学，就意味着生命的四分之一用于与书本打交道。一个科技工作者查阅资料的时间往往占全部工作时间的三分之一。六十年代美国和日本进行的统计，一些化学研究人员业务工作时间分配

如下：查阅文献资料占 50·9%，研究实验占 32·1%，写报告占 9·3%，计划和思考占 7·7%。一个历史学家用于查阅史料的时间比例则更大。这些事实和统计告诉我们文献信息对于社会的发展是举足轻重的。还可以举出许多例子来证明档案资料给政府决策、以及生产、生活带来的效益。档案做为人们社会活动的历史记录，无论现时和将来都是人们必须学习、研究、借鉴的重要信息。

“信息库”，或者“信息源”这样的提法说明了这样一个事实，即历史悠久的档案馆，图书馆和较晚出现的情报资料馆的工作方式和重点，正从文献收集、保管的向心过程转到一种信息发送的离心方式上来。以中央档案馆为例，党中央在起草《关于建国以来若干历史问题的决议》时调查了中央档案三万多卷（件）。近五年来中央档案馆为各方面提供了二十三万卷（件）档案，为 1966 年前提供档案的 1·5 倍。

在人们的社会活动中，由于不能及时获得有关信息，或不注重信息的使用而造成损失的事例是非常多的。例如，50 年代中期，美国曾经耗资 20 万美元动员了大量人力来研究符号逻辑在计算机电路计算方面的应用。事过之后，才得知苏联数学家隆茨 1950 年就做了这方面的工作，他的成果早已赫然在目地存放在图书馆里了。

另一个例子是希腊工程师克利斯多非独自发现的改进粒子加速器的技术。他把结果通知美国加州大学的核实验室后，信件淹没在浩如

烟海的资料之中，杳如黄鹤。同时，许多科学家都大动于研究这项已经被发现的技术。几年后，档案管理人员才在尘封的卷宗中发掘出这封信和这位天才的科学家。

二、信息技术的进步

很早以前人们就以各种方式传递着各种信息。但是限于笨重的载体和落后的传递方式，信息的交流很不方便。随着科学技术的进步，信息技术也发展起来，并且在理论上形成了专门的学科——信息论。

信息论简单地说是利用数学方法，研究信息的计量、传送、变换和存储的一门学科。现代通讯技术、计算技术、自动控制，甚至遗传学都是以信息论作为基础理论。以上这些技术是信息技术的主要支柱。它们使信息的采集、加工、传递进入了自动化阶段。延续数千年在竹、帛、纸上刀刻手写的传统方式已经或正在被自动化处理所取代。而这不过是电子计算机出现后不到四十年发生的变化。信息处理的自动化已发展到了什么规模呢？全面回答这个问题非常困难。下面只以与档案直接有关的文献处理来说明。

进入到 80 年代仅美国、加拿大、欧洲、日本就已有与书目有关的数据库、自然语言数据库 528 个。美国最大的电子计算机化的图书馆网俄亥俄州学院图书馆中心已拥有 3000 台远程终端，一个中心数据库，含有书目记录 700 万条，并以每周 2100 条记录的速度增长。该系统向遍及美国 50 个州的 2200 个成员馆提供服务，