



“十一五”国家重点图书

数字时代图书馆学情报学研究论丛

# 信息检索新论

An Updated Discussion on Information Retrieval

焦玉英 温有奎 陆伟等 编著



WUHAN UNIVERSITY PRESS

武汉大学出版社

P 数字时代图书馆学情报学研究论丛

“十一五”国家重点图书

# 信息检索新论

An Updated Discussion on Information Retrieval

焦玉英 温有奎 陆伟等 编著



WUHAN UNIVERSITY PRESS  
武汉大学出版社

## 图书在版编目(CIP)数据

信息检索新论/焦玉英,温有奎,陆伟等编著.一武汉:武汉大学出版社,2008.8

“十一五”国家重点图书

数字时代图书馆学情报学研究论丛

ISBN 978-7-307-06390-7

I. 信… II. ①焦… ②温… ③陆…[等] III. 情报检索 IV.  
G252.7

中国版本图书馆 CIP 数据核字(2008)第 093153 号

责任编辑:郭 静

责任校对:程小宜

版式设计:詹锦玲

---

出版发行:武汉大学出版社 (430072 武昌 珞珈山)

(电子邮件:wdp4@whu.edu.cn 网址:www.wdp.com.cn)

印刷:武汉中远印务有限公司

开本:720×980 1/16 印张:30.25 字数:432 千字 插页:2

版次:2008 年 8 月第 1 版 2008 年 8 月第 1 次印刷

ISBN 978-7-307-06390-7/G · 1205 定价:48.00 元

---

版权所有,不得翻印;凡购买我社的图书,如有缺页、倒页、脱页等质量问题,请与当地图书销售部门联系调换。

## 总序

“图书馆学情报学”是我国的习惯用法，是涵盖图书馆学、情报学、档案学、出版发行学等学科的名称。在我国台湾被称为“图书馆与资讯科学”，英文为 Library and Information Science。美国也用 Library and Information Studies 来称谓这一学科。

1807 年，德国学者马丁·施莱廷格 (Martin Schrettinger, 1772 ~ 1851) 首次使用了“图书馆学”这一概念，1808 年他又在《试用图书馆学教科书大全》中建立了以图书馆整理为核心的学科体系，标志着图书馆学学科正式诞生。

自 1887 年美国学者杜威 (Melvil Dewey, 1851 ~ 1931) 在哥伦比亚大学创办世界第一所图书馆学校，1930 年在卡内基基金的资助下芝加哥大学设立第一所图书馆学博士班课程以来，图书馆学开始走进大学殿堂，成为高等教育中的一个专业。

图书馆学教育在美国的兴起带动了全球图书馆教育的发展。1919 年英国在伦敦大学建立了图书馆学院。目前，美国有 56 所美国图书馆学会 (ALA) 认可的图书馆学院，每年招收图书馆与情报学学生 26 000 人左右。

在施莱廷格后的两个世纪，图书馆学科不断变化。特别是在 20 世纪 50 年代以来的冷战期间，美苏军备竞赛，两大阵营形成。苏联卫星上天，美国实施阿波罗计划，科技文献激增。科学家对文献信息的获取变得困难。一门新型学科——情报学应运而生。1963 年美国文献工作学会正式更名为美国情报学会 (ASIS)。大量增设图书馆学与情报学硕士点、博士点。图书馆学课程表中也增加了大

量的情报学课程。

20世纪70年代，计算机技术在图书馆与信息工作中广泛应用，自动化、地区性图书馆网络形成，机读目录广泛应用，国际图联将世界书目控制列为核心计划。图书馆学（Library Science）发展为“图书馆与情报学”（Library and Information Science），后来又进一步演变为“图书馆与情报研究”（Library and Information Studies）。

20世纪80年代高新技术迅速发展，信息时代到来。美国里根政府实施星球大战计划，欧洲实施尤里卡计划等。联机图书馆系统广泛建立，并扩展至世界主要发达国家。商业性联机数据库如ORBIT，DIALOG发展迅速，图书馆与情报职业面临挑战。为适应信息时代要求，国际上图书馆学情报学专业开始调整。国际上有较多大学将图书馆学院易名为图书馆与情报学院或信息研究学院，图书馆学、情报学在硕士、博士层次合二为一。

20世纪90年代，全球进入后信息时代——数字时代到来。克林顿政府开始实施国家信息基础设施计划（NII）、全球信息基础设施计划（GII）。新一代互联网投入使用。欧美初步建成信息社会，全球进入无缝信息环境。世贸组织建立和一揽子贸易协定生效，使全球经济一体化并逐步进入知识经济时代。各国继续加强图书馆学、情报学学科调整。图书馆学、情报学学科内容向情报科学汇集。

进入21世纪以来，国际上信息管理学科变化很快。自雪城（SYRACUSE）大学将学院更名为信息研究学院（The School of Information Studies）后，在美国立即出现了iSchool的浪潮。伊利诺依斯大学、华盛顿大学、密歇根大学、匹兹堡大学、加州大学伯克利分校、北卡罗来纳大学等知名大学的图书馆与情报学院宣称自己为iSchool。这些iSchools通过宪章组成I-Schools联盟（ISG）。目前共有20所美国的大学加入联盟（联盟宪章不允许超过25个）。iSchool强调信息、技术与人的关系（relationship between information, technology and people）。iSchool的标准包括：必须有杰出的研究和杰出的博士教育；必须能在科学、企业、教育与文化进

步过程中提供任何形式的信息所需的专业技术；必须能提供信息技术及其应用、信息使用与用户方面的专门知识。2004~2006年的联盟领导委员会协调人是雪城大学信息研究学院的 Raymond von Dran 院长，2006~2007 年将由匹兹堡大学信息学院院长 Ron Larsen 担任。联盟成员的标准主要强调研究即实质性承担研究活动（三年中每年研究支出达到 100 万美元），同时，致力于培养未来的研究者（通常通过研究型的博士点），引领推动信息职业领域。

国际上图书馆与情报学科的发展表现出明显的特征：研究范围由传统的图书馆领域扩大到信息领域（information field），研究视野由实体的图书情报机构扩大到虚拟空间，研究对象由图书文献转向了信息内容。一系列相关学科如图书馆学、情报学、档案学、出版科学、信息管理与系统乃至数字商务汇集于信息科学（Information Sciences）下，从而使图书馆学情报学研究发生了根本的变化。

武汉大学图书馆学科起源于 1920 年美国学者韦棣华女士创办的武昌文华大学图书科，档案专业起源于 1940 年的文华图书馆学专科学校的档案管理科。1978 年武汉大学创办科技情报学专业，后改为情报学专业。1983 年创办图书发行学专业，2002 年创办电子商务专业。1984 年经教育部批准建立武汉大学图书情报学院。2001 年更名为信息管理学院。图书馆学和情报学两个二级学科被国务院学位委员会批准为国家重点学科。“图书馆、情报与档案管理”被国务院学位委员会批准为一级学科博士学位授权点。教育部批准“武汉大学信息资源研究中心”为国家人文社会科学重点研究基地。信息产业部批准成立“国家信息资源管理（武汉）研究基地”。新闻出版总署批准建立“新闻出版总署武汉大学高级出版人才培养基地”。“网络信息资源开发与数字图书馆建设”被国家计委、教育部等批准为“十五”211 重点学科建设项目。建立一级学科博士后流动站。武汉大学信息资源研究创新基地被列为国家“985 二期工程”建设项目。一批院内校级重点研究基地如武汉大学四库学研究所、武汉大学中国科技评价中心、武汉大学政府信息

研究中心、武汉大学数字图书馆研究所、武汉大学出版发行学研究所、武汉大学图书馆学情报学国际合作研究中心也在科研和人才培养中发挥着重要平台作用。

强调一级学科内学科群建设和学科协调发展是武汉大学图书馆与情报学科建设的基本目标。以图书馆学、情报学两个国家重点学科为龙头促进图书馆学、情报学、档案学、信息资源管理、出版发行学等学科的协调发展。

我们深刻认识到信息资源与自然资源、人力资源共同构成支撑现代经济社会发展的资源体系。信息资源是知识经济时代重要的国家战略资源，是实现经济和社会全面、可持续发展的基础条件。对信息资源的拥有、开发和利用水平，是衡量一个国家综合国力和国际竞争力的重要标志之一。消弭信息鸿沟、实现信息公平，是消除贫困、促进经济发展、构建和谐社会的重要条件之一。

信息资源管理人才培养是学院的基本任务。学院每年为国家培养本科生 260 名，硕士研究生 150 名，博士研究生 55 名左右。学院有一支知识结构和年龄结构合理的优秀学术队伍。这支队伍中有武汉大学人文社会科学资深教授 1 人，博士研究生导师 26 人，国务院政府特殊津贴专家 6 人，教育部新世纪优秀人才支持计划 3 人，武汉大学珞珈特聘教授 2 人。作为实现研究型学院建设目标的一部分，在教学的同时，广大教师承担了大量的科学研究任务。为了推动本学科领域的前进，分享他们的见解，在武汉大学出版社的大力支持下，并报有关部门批准，我们拟出版《数字时代图书馆学情报学研究论丛》（简称《论丛》）。

为了编辑这套丛书，武汉大学邀请了国内外知名学者担任《论丛》的学术顾问，组建了主要由信息管理学院的博士研究生导师担任委员的编辑委员会。

《论丛》拟用 4 年时间出版著作共 20 卷。20 卷著作将分为三个系列：(1) 学科年度进展。主要约请信息管理学院图书馆学系、档案与电子政务学系、信息管理科学系、现代出版系、信息系统与

电子商务系的有关教师和校外专家共同编写本学科的年度研究进展，主要有《图书馆学研究进展》、《情报学研究进展》、《档案学研究进展》、《出版学研究进展》、《信息资源管理学研究进展》；（2）个人学术专著。涉及图书馆、情报与档案管理基本理论研究、信息组织与检索、信息资源管理、信息资源建设与信息服务、文献编纂与出版、数字图书馆与信息系统工程等研究方向；（3）研究报告系列。我院研究人员共承担教育部哲学社会科学研究重大攻关项目、国家社会科学基金重点项目、教育部人文社会科学研究基地重大招标项目、国家自然科学基金项目、国家社会科学基金项目多项。特别是211项目和985项目，围绕数字信息资源开发与管理、数字信息资源服务与保障、信息资源公共获取与知识产权协调管理、数字图书馆关键技术与系统、资源与服务整合、信息构建与知识管理等主题正在进行探索。在信息构建的理论与方法、信息系统与资源整合、元数据知识表达、网络计量与参考、信息服务集成机制、信息资源与服务集成技术、媒体及数字出版、数字内容分销、信息资源的长期保存、商务信息流等关键领域力图实现图书馆学科在数字图书馆领域、情报学科在数字资源管理领域、档案学在数字化政务信息管理领域、出版发行学在数字出版与数字化分销、信息系统科学在集成系统以及数字化商务信息流研究方面取得研究成果。本系列将对部分研究结果进行报告。

丛书的出版是学院广大教师和研究人员辛勤探索的结果，在此，谨向严谨治学、辛勤耕耘的各位著作者表示感谢！对武汉大学出版社的支持表示感谢，对责任编辑严红女士在策划编辑过程中付出的艰辛劳动表示感谢。同时，还望广大读者不吝批评指正，共同推动图书馆学、情报学、档案学、出版发行学和信息资源管理学科的进步！

武汉大学信息管理学院院长 陈传夫  
武汉大学信息资源研究中心主任 马费成  
2006年10月8日

# 前　　言

网络与数字化技术使分布在世界各地主机上的信息资源联为一体,构成了跨时空、跨行业、高效、快速的国际化知识共享信息环境。网上聊天、个人意见与观点展示、即时沟通与互动、邮件传递、在线购物等已成为人们学习、工作、科学研究、教学以及日常生活中不可缺少的组成部分。五花八门、多功能的搜索工具及导引方式、在线帮助等使用户足不出户,不必经过专门培训即可从网上获取自己所需要的信息。

从专业的角度看信息检索,以传统文献为基础的“提问—检索”功能正逐渐萎缩,取而代之的是网上“浏览—查询”。虽然通过网络,信息的主动性得到空前提高,但网上“浏览—查询”检索模式运行至今仍有诸多不尽如人意之处。用户采用“点击”浏览远远未达到预期的查询效果。这就涉及一系列检索理论、技术、方法与策略问题。本书旨在综合检索领域的热门、前沿课题,为关注这一领域的教师、学生以及其他与信息工作有关的人们提高获取网络信息,检索专门知识的效率,开展进一步研究提供参考。

《信息检索新论》在保持传统检索的基本理论体系的同时,特别梳理了近年来信息检索领域的最新成果,如信息过滤、网格信息检索、知识元检索、跨语言检索、语义网检索、信息抽取等。许多课题是国家自然科学基金和国家社会科学基金项目的研究成果,有的是国际 TREC 检索会议近年的中心议题。我们希望通过这些问题的研究,向读者揭示 21 世纪信息检索领域的前沿课题和研究方向。

本书由焦玉英组织编著,并撰写前言、第 1 章概论,负责全书的策划、统稿、定稿工作。西安电子科技大学学校级学术带头人温有奎教授专门为本书撰稿,并对全书的内容提出宝贵意见。参加编

写本书的作者，按章节顺序如下：前言与第1章由焦玉英教授负责编写；第2章由刘伟成博士编写；第3章由成全博士编写；第4章由雷雪博士编写；第5章由孙吉红博士编写；第6章由黄传慧博士编写；第7章由余彩霞博士编写；第8章由刘伟成博士编写；第9章由温有奎教授编写；第10章由陆伟博士、副教授编写；第11章由李法运博士、副教授编写；第12章由李进华博士、副教授编写；第13章由曹高辉博士编写；第14章由王娜博士编写；第15章由项英博士编写；第16章由刘颖博士、副研究馆员编写；第17章由武琳博士编写（07JC870005项目部分成果）；第18章由温有奎教授编写。全书约40万字。

博士袁静、硕士龙泉参与了本书的部分资料整理工作。

本书的出版得到武汉大学信息管理学院“985”工程项目——“信息资源管理创新平台”的资助，被评选为“十一五”国家重点图书。感谢信息管理学院院长陈传夫、副院长方卿、系主任唐晓波的大力支持，感谢武汉大学出版社严红、郭静、夏敏玲等同志为本书的出版所付出的辛勤劳动。

本书引用和参考了国内外本领域专家、学者的论著和网站资料，在此一并致谢。水平有限，有不当或疏漏之处，欢迎批评、指正。

焦玉英

2007年12月于武汉大学

# 目 录

前 言 .....	1
1 概论 .....	1
1.1 信息检索的理论研究进展 .....	1
1.1.1 对以传统“提问—检索”模式为核心的文献检索的简单回顾 .....	1
1.1.2 网络环境下信息检索相关理论研究 .....	4
1.2 网络信息检索展望 .....	18
1.2.1 网络信息存储的“动态化” .....	18
1.2.2 人工智能与检索技术的高度融合 .....	18
1.2.3 检索结果相关度的可视化 .....	19
1.2.4 基于内容的检索技术的应用 .....	19
1.2.5 网络信息检索的商业化 .....	19
1.3 语义网检索技术 .....	20
1.3.1 下一代万维网——语义网 .....	20
1.3.2 标记语言提供信息共享的基础 .....	21
1.3.3 RDF 提供处理元数据的基础 .....	22
1.3.4 Ontology 是概念化的规范说明 .....	22
1.3.5 语义查询语言 .....	23
1.3.6 Ontology 在语义信息检索中的使用 .....	24
1.3.7 知识检索 .....	27
1.3.8 语义检索 .....	28

2 信息检索模型理论 .....	31
2.1 国内外有关的研究进展 .....	31
2.2 布尔模型 .....	33
2.2.1 经典布尔模型 .....	33
2.2.2 扩展布尔模型 .....	33
2.3 向量空间模型 .....	36
2.3.1 经典向量空间模型 .....	36
2.3.2 广义向量空间模型 .....	37
2.3.3 潜在语义索引模型 .....	38
2.4 概率模型 .....	40
2.4.1 经典概率模型 .....	40
2.4.2 推理网络模型 .....	41
2.5 逻辑模型 .....	42
2.5.1 信息检索的古典逻辑模型 .....	43
2.5.2 信息检索的非古典模型 .....	43
2.6 统计语言模型 .....	44
2.6.1 N-gram 模型 .....	44
2.6.2 隐马尔可夫模型 .....	45
2.7 结构化文本检索模型 .....	46
2.7.1 基于非重叠链表的模型 .....	47
2.7.2 基于邻接节点的模型 .....	48
2.8 浏览模型 .....	49
2.8.1 平坦浏览模型 .....	49
2.8.2 结构导向浏览模型 .....	49
2.8.3 超文本浏览模型 .....	50
3 信息检索相关性研究 .....	51
3.1 信息检索相关性研究概述 .....	51
3.1.1 相关性的研究与发展 .....	51
3.1.2 相关性的研究学派 .....	55
3.1.3 相关性影响因素分析 .....	59

3.2 信息检索相关性评价 .....	63
3.2.1 系统相关性评价 .....	63
3.2.2 用户相关性评价 .....	66
3.3 用户相关反馈研究概述 .....	67
3.3.1 用户相关反馈的界定与发展历程 .....	67
3.3.2 用户相关反馈机制的基本原理 .....	69
3.3.3 用户相关反馈的实现类型与主要特征 .....	70
3.4 用户相关反馈技术研究 .....	72
3.4.1 基于向量空间模型的相关反馈技术 .....	72
3.4.2 基于经典概率模型的相关反馈技术 .....	73
3.4.3 基于布尔检索模型的相关反馈技术 .....	75
 4 并行与分布式检索进展 .....	79
4.1 并行与分布式检索的必要性 .....	79
4.2 并行信息检索的基本问题 .....	80
4.2.1 并行计算 .....	80
4.2.2 MIMD 结构的并行检索 .....	81
4.2.3 SIMD 结构的并行检索 .....	85
4.3 分布式信息检索的相关问题 .....	86
4.3.1 分布式检索的体系结构 .....	86
4.3.2 分布式检索的过程 .....	87
4.3.3 分布式检索涉及的标准化问题 .....	101
4.4 并行与分布式检索的发展趋势 .....	102
 5 多媒体检索进展 .....	104
5.1 基于内容的多媒体信息检索 .....	104
5.1.1 多媒体信息检索的关键技术研究 .....	104
5.1.2 基于内容的多媒体信息检索方法 .....	107
5.2 国内外多媒体信息检索最新研究 .....	116
5.2.1 新特征和相似度量 .....	116
5.2.2 新媒体 .....	118

5.2.3 浏览和摘要 .....	118
5.2.4 高性能索引 .....	119
5.2.5 以人为本 .....	120
5.2.6 语义和反馈 .....	121
5.2.7 评价 .....	122
5.3 多媒体信息检索系统及其发展趋势 .....	123
5.3.1 基于内容的多媒体信息检索系统开发概况 .....	123
5.3.2 基于内容的多媒体信息检索技术的发展趋势 .....	126
5.4 多媒体信息检索技术研究面临的挑战 .....	127
 6 智能信息检索进展 .....	129
6.1 智能代理技术的发展 .....	129
6.1.1 人工智能与信息检索 .....	129
6.1.2 智能代理的概念、发展 .....	136
6.1.3 智能代理技术的国内外应用研究 .....	137
6.1.4 智能代理技术在信息检索与服务系统中的应用 .....	138
6.2 智能信息检索与服务系统所涉及的主要理论问题研究 .....	141
6.2.1 基于自然语言处理方面的研究 .....	141
6.2.2 基于智能代理的搜索引擎技术研究 .....	144
6.2.3 智能信息检索服务系统及其运作 .....	147
6.3 智能信息检索与服务系统研究总体特征和趋势 .....	150
6.3.1 我国智能信息检索与服务系统研究的总体特征 .....	150
6.3.2 智能信息检索与服务系统研究趋势 .....	151
 7 光盘与联机检索回顾 .....	153
7.1 光盘检索研究的发展 .....	153
7.1.1 光盘检索技术发展研究 .....	153

## 目 录

7.1.2 光盘检索系统的应用研究 .....	155
7.1.3 光盘检索系统的分析与评价 .....	156
7.1.4 光盘检索服务的应用模式分析 .....	158
7.2 联机检索研究的发展 .....	160
7.2.1 联机检索技术发展研究 .....	160
7.2.2 联机检索系统的分析与评价 .....	163
7.2.3 联机检索系统的应用研究 .....	164
7.2.4 联机检索服务的应用模式分析 .....	165
7.3 联机检索与光盘检索的比较研究 .....	167
7.3.1 传统联机检索的优劣势和发展方向 .....	167
7.3.2 光盘检索的优劣势和发展趋势 .....	169
7.3.3 两者的比较 .....	171
 8 跨语言检索理论与实践 .....	174
8.1 跨语言信息检索研究概况 .....	174
8.1.1 跨语言信息检索研究的发展历程 .....	174
8.1.2 跨语言信息检索研究的主要内容 .....	176
8.2 跨语言信息检索的基本框架 .....	178
8.3 跨语言信息检索的类型和技术 .....	180
8.3.1 基于翻译方法的分类 .....	180
8.3.2 基于翻译工具的分类 .....	183
8.3.3 基于检索媒体的分类 .....	186
8.4 跨语言信息系统评价研究 .....	187
8.4.1 跨语言信息检索评价模型 .....	187
8.4.2 效率评价指标 .....	188
8.4.3 现有测试平台运行状况分析 .....	189
8.4.4 跨语言信息检索测试集 .....	190
8.5 跨语言信息检索研究的主要应用领域 .....	192
8.5.1 在数字图书馆中的应用 .....	192
8.5.2 在科学中的应用 .....	193
8.5.3 在电子商务中的应用 .....	194

8.5.4 在跨文化交流中的应用 .....	194
9 知识元检索 .....	196
9.1 知识检索单元的变化 .....	196
9.1.1 文献单元检索向知识元检索发展 .....	199
9.1.2 语义网成为知识元检索的有效工具 .....	201
9.2. 知识元的表示模式 .....	204
9.2.1 知识元的知识原子 .....	204
9.2.2 知识元的知识因子 .....	205
9.2.3 知识元的知识因子的一元运算 .....	206
9.2.4 知识元的知识项 .....	207
9.2.5 知识元表达式 .....	208
9.3 知识元链接理论 .....	209
9.3.1 文献处理技术的进展 .....	209
9.3.2 知识元结构定义 .....	210
9.3.3 知识元链接框架 .....	211
9.4 知识元的变换性 .....	212
9.4.1 Brookes 方程与理解 .....	212
9.4.2 信息与知识谱的变换性 .....	213
9.4.3 知识谱分析 .....	214
9.5 知识元本体推理模型 .....	216
9.5.1 知识元检索模型 .....	216
9.5.2 本体信息层上知识推理 .....	218
10 信息抽取理论的进展 .....	220
10.1 信息抽取概况 .....	220
10.1.1 信息抽取的定义 .....	220
10.1.2 信息抽取的背景 .....	221
10.1.3 信息抽取相关领域 .....	222
10.2 信息抽取应用现状 .....	223
10.2.1 国外信息抽取应用现状 .....	223

10.2.2 国内信息抽取应用现状	225
10.3 信息抽取技术的分类与比较	226
10.3.1 基于规则的信息抽取	227
10.3.2 基于概率的信息抽取	230
10.3.3 多策略混合方法	231
10.3.4 信息抽取技术的比较	232
10.4 信息抽取技术发展趋势	234
10.4.1 增强系统的适应性	234
10.4.2 以 XML 的格式输出信息	235
10.4.3 面向多种用户需求	235
10.4.4 面向语义的信息抽取	236
10.4.5 与信息检索相结合	237
11 信息过滤与服务研究	238
11.1 信息过滤概述	238
11.1.1 信息过滤的产生和发展	238
11.1.2 信息过滤的类型	239
11.1.3 信息过滤的特点	240
11.1.4 信息过滤的应用领域	241
11.2 信息过滤模型及信息过滤系统	242
11.2.1 信息过滤模型及系统概述	242
11.2.2 基于内容的信息过滤模型及系统	243
11.2.3 协作过滤模型及系统	244
11.2.4 经济过滤模型及系统	246
11.2.5 基于内容和协作过滤相结合的过滤推荐 模型及系统	247
11.3 信息过滤技术	247
11.3.1 信息过滤流程	247
11.3.2 用户模型构建技术	248
11.3.3 信息表征技术	254