

术语计算与知识组织研究

宋培彦 ◎ 著



术语计算与知识组织研究

宋培彦 著



科学技术文献出版社
SCIENTIFIC AND TECHNICAL DOCUMENTATION PRESS

· 北京 ·

图书在版编目 (CIP) 数据

术语计算与知识组织研究 / 宋培彦著. —北京：科学技术文献出版社，2018.6
ISBN 978-7-5189-4559-7

I . ①术… II . ①宋… III . ①自然语言处理—研究 IV . ① TP391. 1

中国版本图书馆 CIP 数据核字 (2018) 第 131143 号

术语计算与知识组织研究

策划编辑：周国臻 责任编辑：张 红 责任校对：文 浩 责任出版：张志平

出 版 者 科学技术文献出版社

地 址 北京市复兴路15号 邮编 100038

编 务 部 (010) 58882938, 58882087 (传真)

发 行 部 (010) 58882868, 58882870 (传真)

邮 购 部 (010) 58882873

官 方 网 址 www.stdp.com.cn

发 行 者 科学技术文献出版社发行 全国各地新华书店经销

印 刷 者 北京教图印刷有限公司

版 次 2018 年 6 月第 1 版 2018 年 6 月第 1 次印刷

开 本 710 × 1000 1/16

字 数 172 千

印 张 13

书 号 ISBN 978-7-5189-4559-7

定 价 58.00 元



版权所有 违法必究

购买本社图书，凡字迹不清、缺页、倒页、脱页者，本社发行部负责调换

自序：走向智能化知识服务的“立心之战”

美国对中国芯片技术等领域进行“卡脖子”式的封锁，“缺芯之痛”引起举国关注。然而，看不见、摸不着但却同样重要的知识库，其作用并不逊色于芯片硬件，对智能信息处理和智能化知识服务这两个方面同时具有全局性、根本性和长期性的影响，堪称是不可或缺、不可替代、价值不可估量的“软芯片”。

智能信息处理需要强有力的知识库。这是因为，知识库是计算机智能信息处理的基石，也是未来智能计算机内嵌的核心技术，所以“不可或缺”；知识库本身是一个庞大的知识密集型系统，加上中国语言和文化自身的特殊性，因此很难指望外国人代劳或者从国外购买、坐享其成，而只能主要依靠国人自力更生、自主创新，因此“不可替代”；正是由于其重要性高、技术难度大、研发周期长而市场价值极为显著，因此其价值“不可估量”。离开了知识库，计算机还是擅长数据计算的“超级机器”，但绝不能称为一目千行、勤学善思的“电脑才子”，也与“智能”二字无缘。知识库是衡量“智能”水平高低的标志之一。

智能化知识服务也同样离不开知识库的有力支撑。在大数据环境下，图书情报界正在从传统的图书情报服务转向智能化知识服务，以术语库、叙词表、本体等为代表的组织工具本质上是作为计算机内置使用的知识库，扮演了计算机“大脑”的关键角色。智能检索、关联数据、语义网、知识图谱等智能化知识服务理念，都有赖于知识组织工具（知识库）。研制汉语知识库特别是专业领域知识库，可以为智能推送、精准检索、自动推理等智能化知识服务提供强有力的知识基础。可以说，知

识库建设涵盖了图书情报学（知识组织）、人工智能（自然语言处理）、语言学（术语学）、统计学乃至认知科学等多学科，堪称中国图书情报界开展智能化知识服务的“立心之战”，其科学意义和应用价值不言而喻。

所幸，我们有许多默默无闻、潜心研究的先哲早就深刻意识到了这一点，凭借专业技能和满腔热忱肩负起为电脑“立心”、坚持自主创新的重任，数十年如一日，创立了知网 Hownet、概念层次网络 HNC、《同义词词林》、《汉语主题词表》等许多奠基性、开创性的知识工程。一些富有战略远见的创新型企业，如百度、阿里巴巴、腾讯等 IT 企业及出版业界，也紧紧抓住机遇、投入巨资，纷纷将知识图谱、网络百科、智能检索等基础研究变成落地生根的产品，服务于社会公众并取得良好经济效益。特别是，我国重点研发计划等重大工程着眼长远，将中文信息处理、类脑计算、知识工程等列为国家重大战略项目予以稳定支持，中文知识库这一“软芯片”研究与产业化基本与国际先进水平保持齐头并进，不仅成功避免了被外国技术封锁的被动局面，而且在某些领域已经走进世界舞台中央，从“跟跑”变成“领跑”。本书就是在这样的时代背景下应运而生的，最终希望对我国专业领域智能信息处理和知识服务水平的提升有所裨益。

同时，也应该注意到，我国在知识库方面的研究还存在一些亟待解决的问题。例如，传统中文知识库建设大多依靠专家手工劳动，不同程度地存在着效率偏低、成本高昂、更新迟缓、应用场景不明确等问题；或者过于强调统计模型、深度学习等技术手段，追求“短平快”，造成知识库准确性不足，难以支撑实际应用。因此，本书试图在知识组织的理论框架下，将术语计算技术引入知识库构建中，实现术语计算技术与知识组织理论的结合，探索更加精准化、智能化、开放性的知识库建设与服务新模式。这种结合不仅在理论上是可行的，而且对工程实践和实际应用也是有益的。

知识组织为术语计算提供了较完善的知识表示框架和数据基础。在大数据环境下，以术语为中心，将各类知识组织工具进行融合，形成统一的知识库，进而支撑专业知识的组织、挖掘、管理和服务，不仅可以“复用”叙词表、本体等现有的规范化知识体系，实现知识的继承和发展，而且能够实现对专业知识进行多维度、细粒度的组织、管理和利用，有助于消除散乱的术语造成的“信息孤岛”，并缓解“信息过载”，为用户提供精准的深层次知识服务，提高智能化知识服务的质量和效率。

术语计算为实现知识组织提供了有效技术手段。术语计算（Terminology Computing）是采用自然语言处理技术，对术语所包含的专业知识进行自动挖掘，以较小的颗粒度和可组合性对知识进行表示、推理和计算。通过对术语的挖掘和描述，构建细粒度、动态更新的术语知识库，有助于实现知识的动态聚合，进而以可视化、个性化的方式为用户提供知识服务。本书以自动分类、聚类、词义计算等术语计算技术为依托，形成快速、动态的知识发现与更新方法，从微观层面揭示知识的语义关系，以可视化方式直观展现知识关联，进而为构建词间关系、建立知识关联提供辅助手段，促进大数据时代知识动态更新和精准服务。

实践也充分证明，术语计算与知识组织在智能化知识服务这一根本目标上“殊途同归”，具有很强的可结合性，是构建新型知识库的有效方法。术语计算具有很强的操作性和应用性，在知识组织的框架中才能更好地发挥其应有的作用；现有的知识组织工具需要借助术语计算技术在更大范围内进行有序化组织和关联，才能充分发挥知识组织的价值、满足用户的知识需求。

万丈高楼平地起。以下项目提供的数据资源和应用场景，为本书提供了宝贵的实验基础，并验证了术语计算与知识组织结合的可行性。

(1) 《汉语主题词表（工程技术卷）》编制与科技术语库建设。2010—2014年，笔者作为核心研究人员参加了《汉语主题词表（工程技术卷）》的编制与修订工作。具有完全知识产权的新版《汉语主题词表

(工程技术卷)》核心词(术语)数量达到20余万条;科技术语库收录科技名词术语近400万词,并在语义计算、自动分类、自动标引等方面进行了积极探索,具有较为坚实的数据基础。

(2) 英文超级科技词表与术语库建设。2011—2015年,在国家“十二五”科技支撑计划项目“面向外文文献的知识组织研究与示范应用”支持下,笔者作为术语库负责人,采集了约20万个工程技术领域核心概念、300多万个英文术语,形成了跨语言的知识组织工具,并用于国家科技文献中心NSTL外文文献的标引与检索服务,产生了良好的社会效益。

(3) 术语服务系统建设。2011—2014年,笔者主持了国家社科基金项目“基于知识组织的术语服务研究”,参与研制了术语服务系统、主题自动标引系统等,具备中英文翻译、关系、属性、定义、分类信息等,并实现了对科技术语的语义链接和可视化,在国内一些专业机构开展示范应用。

(4) “国家科技管理信息系统”科技大数据中心建设。2015年以来,围绕“国家科技管理信息系统”建设,笔者主要从事科技大数据中心的设计与建设工作,实现与全国47个省市(含计划单列市、副省级城市、新疆生产建设兵团)的互联互通与共享服务,将国家科研项目、专家、成果等各类科技类数据进行深度组织和管理。术语计算和知识组织对于科技大数据中心建设发挥了重要作用。

(5) 专家推荐系统建设。笔者还主持了2016年度国家社科基金项目“基于知识组织的科研项目评审专家发现研究”,通过知识组织工具的优化实现专家自动推荐,为科研项目评审、决策咨询等提供良好的支持。

上述研究经历让笔者更加深切地体会到,术语计算与知识组织有天然的紧密联系,具有重要研究价值。本书主要采用实证研究的方法,将术语理论、计算方法、模型与实验等尽可能与叙词表等知识组织工具进

行有效结合和验证，尽量使计算结果具有可重复性，以支撑知识组织的实际应用。同时，以知识组织理论框架为纲，将散乱的术语融合成有机整体，并以智能检索、术语服务、数字出版、科技辅助决策等作为应用场景进行验证，也有助于体现术语计算的效果和价值。理论、方法、技术与应用，在这里相得益彰、并行不悖，构成本书的主线。

本书命名为《术语计算与知识组织研究》，蕴含的目标是希望促进术语计算与知识组织的跨学科结合，寻求智能化知识服务的新突破。笔者主要研究方向是知识组织和自然语言处理，恰好具备图书情报学、计算机和语言学等方面复合型知识结构，加上在知识组织工具建设、术语计算、术语服务等方面具备较好的实际工作经验和科研条件，包括主持国家社科基金项目 2 项、发表学术论文 20 余篇，获得 2 项发明专利、6 项软件著作权等，这些成果均围绕术语计算与知识组织进行，并在一些重要的科技管理信息系统中产生了较好的应用效果，种种机缘巧合，得以管窥其中奥妙。当然，这是一条鲜花与荆棘相伴、机遇与挑战并存、幸福与艰辛同在的羊肠小道，我们不能奢望“毕其功于一役”，试图通过一本书或项目解决所有问题。倘若能够吸引更多的有志者携手前行，开辟更加宽广的科学大道，那将是本书最大的成功。

宋代理学家张载说“为天地立心，为生民立命，为往圣继绝学，为万世开太平”，这是做学问的高层次境界，千百年来激励人心、催人奋进。在大数据和人工智能时代构建高水平知识库，这既是为我国智能信息处理事业“立心”“立命”，有助于提高我国的智能化技术水平和服务能力，也是科技工作者“继绝学”“开太平”的责任担当和价值体现。我们应该与时俱进，在术语计算和知识组织的结合上有更多、更大的作为，推动更强劲的“中国芯”率先巍然屹立在中华大地上。

目 录

第一章 绪论：知识为体，术语为用	1
1 术语：打开知识宝库的一把钥匙	1
1.1 术语与知识组织：如影随形的黄金搭档	1
1.2 大数据时代的知识组织：在继承中创新，在创新中发展	4
2 国内外前沿概览：以他山之石，攻我之美玉	6
2.1 术语计算：知识组织的轻骑兵	6
2.2 叙词表：术语计算的试验田	9
3 本书导读	10
3.1 章节安排	11
3.2 阅读指南	12
第二章 术语知识表示：探究术语背后的奥秘	14
1 术语知识表示模型：让计算机读懂“知识”	14
1.1 术语与知识关系模型：从语言学的视角	14
1.2 术语知识表示框架：让电脑更聪明	17
1.3 小结	22
2 术语知识库构建：将知识变成“记忆”	22
2.1 术语知识库构建模型：“复用”人类知识	23
2.2 术语知识微观结构：知识组织的“蝴蝶效应”	27
2.3 术语知识库：人机两用的知识宝库	31
2.4 小结	33
3 术语 SKOS 形式化描述：助力知识处理驶入快车道	34

3.1 《汉语主题词表》结构框架与术语描述：理性主义的规范控制	34
3.2 术语 SKOS 规范化：与国际接轨的“通行证”	39
第三章 以用户为中心的知识组织：让计算机“善解人意”	43
1 人机知识交互与合作：将“芯”比心，挖掘和预知用户意图	43
1.1 知识组织会话合作原则：人同此心，心同此理	44
1.2 知识组织会话过程：人机对话更智能	46
1.3 小结	50
2 用户自主标注：术语中的“长尾”定律	51
2.1 回归分析算法：用户是如何使用术语的	51
2.2 拟合检验指数：验证术语分布的规律性	52
2.3 术语分布规律分析：长尾也有大用处	55
2.4 小结	59
3 用户交互式术语更新机制：保持知识的“活性”	60
3.1 用户交互机制：用户是知识消费者，也是知识生产者	61
3.2 用户知识交互与协同：让用户更有获得感	62
3.3 小结	65
4 网络百科中的知识组织：汇集众人的智慧	66
4.1 网络百科的知识组织特征：合理规范，顺其自然	66
4.2 网络环境下的术语知识组织层次模型：冲破藩篱，尽情迸发	69
4.3 小结	75
第四章 术语计算技术：从数据中“挖掘”知识	76
1 科技语料库与术语知识抽取：语料为舟，算法为舟	76
1.1 面向知识组织的科技语料库：知其一，也知其二	76
1.2 科技语料库设计与管理：建立“计算”的数据基础	77
1.3 科技语料库应用举例：事实的力量	79

1.4 小结	82
2 跨语言术语自动分类：物以类聚，触类旁通	83
2.1 术语分类推导 - 归并模型：由已知推断未知	84
2.2 术语分类实验：看看结果如何	89
2.3 小结	91
3 术语自动聚类：给知识拍个“快照”	91
3.1 聚类计算模型：分步推进，逐步逼近	92
3.2 两步聚类算法：可计算的语义距离	95
3.3 数据实验：肿瘤领域术语聚类效果	97
3.4 小结	102
4 英文术语同义关系计算：模糊中寻求精确	102
4.1 同义关系：从模糊中划分边界	103
4.2 模糊归并模型：超越 0 和 1	105
4.3 实例分析：将同义转化为量化权重	107
4.4 小结	110
5 中文术语同义关系计算：给术语照照镜子	110
5.1 镜像翻译技术：看看镜子里的中文术语	111
5.2 同义关系推导过程：符号变换的“魔术”	112
5.3 实验结果：简单的方法，可用的结果	114
5.4 小结	116
6 术语知识单元抽取方法：找到知识的“基因”	117
6.1 术语知识单元：以有限应对无限	117
6.2 知识单元抽取算法：在句子的大视野里找知识单元	117
6.3 知识单元抽取实验：构建更细粒度的“知识基因库”	119
6.4 知识多维度聚合：知识的关联不是偶然的	124
6.5 小结	126
第五章 术语示范服务与应用：欲穷千里目，更上一层楼	127
1 术语服务：让知识唾手可得	127

1.1 术语服务基本架构：九尺之台，起于垒土	127
1.2 术语服务模块设计与实现：咫尺屏幕，大千世界	131
1.3 小结	133
2 基于词汇链的辅助标引：编织美丽的语义“项链”	134
2.1 语义网络与词义计算方法：语义地图上的“舞蹈”	135
2.2 词汇链构造算法：寻找文本中的蛛丝马迹	136
2.3 实验分析：找出最耀眼的那颗星	140
2.4 小结	142
3 术语词典辅助出版：为出版插上腾飞的翅膀	142
3.1 术语词典知识组织结构：纲举目张，各入其位	144
3.2 系统总体设计：无纸化的数字出版	147
3.3 小结	149
4 专家推荐：用数据说话，用数据决策	150
4.1 自动推荐技术：发现偶然中的必然	150
4.2 专家推荐应用：找到更权威的“小同行”专家	151
第六章 总结与展望：跨界结合，顺势而为	155
附录 A 国外知识组织协会调研	157
附录 B 汉语主题词表研究热点与发展路径	164
附录 C 国内外典型术语服务系统	178
参考文献	183
后记	189

图表目录

图 1.1 大数据、知识与术语的关系	6
图 2.1 术语与知识关系模型	15
图 2.2 术语知识表示框架	18
图 2.3 术语知识形式化描述举例	20
图 2.4 术语知识库构建模型	23
图 2.5 语义网环境下术语与 KOS 一体化模型	25
图 2.6 术语知识微观结构	27
图 2.7 术语网络化应用：以“纳米材料”为例	33
图 2.8 《汉语主题词表》来源信息结构	35
图 2.9 《汉语主题词表》款目树形结构	36
图 2.10 《汉语主题词表》款目词	36
图 2.11 《汉语主题词表》词属性	37
图 2.12 《汉语主题词表》同义关系	37
图 2.13 《汉语主题词表》上下位及相关关系	38
图 3.1 自然语言向叙词表受控语言的映射模型	47
图 3.2 术语分布计算流程	54
图 3.3 术语词频分布	55
图 3.4 自然标准术语词频幂函数曲线拟合	56
图 3.5 2000—2017 年肿瘤术语幂指数变化趋势	59
图 3.6 用户标签与术语知识库交互模型	62
图 3.7 交互式叙词表更新功能模块	65
图 3.8 主题标引与聚类检索示意界面	65
图 3.9 基于网络百科的知识组织模型	71
图 3.10 基于映射的网络百科微观知识组织	72

图 4.1 科技语料库在术语发现中的应用流程	79
图 4.2 基于语料库的术语识别	80
图 4.3 通过正则表达式实现术语释义抽取示意	82
图 4.4 跨语言术语自动分类模型	86
图 4.5 术语翻译与分类的对应关系	87
图 4.6 术语按投票机制归并	88
图 4.7 英汉对照法中图分类号归并流程	89
图 4.8 领域知识网络两步聚类法总体流程	93
图 4.9 聚类结果概要	97
图 4.10 关键词对聚类结果的贡献值	98
图 4.11 聚类分析树状图	101
图 4.12 “同形结合”传递原理	104
图 4.13 同义词模糊归并影响因素层次结构示意	108
图 4.14 归并过程实例	109
图 4.15 镜像推导算法流程	113
图 4.16 识别同义词示意	114
图 4.17 术语知识单元抽取流程	119
图 4.18 句法-语义分析树示例	121
图 4.19 知识单元关系存储	124
图 4.20 术语知识单元社会网络分析	125
图 5.1 基于知识组织的术语检索服务体系模型示意	128
图 5.2 术语服务系统模块结构	131
图 5.3 术语服务系统的语义聚合和检索	132
图 5.4 术语服务系统概念关系的可视化	133
图 5.5 词汇链构造流程	137
图 5.6 语义相关系数设置及词汇链可接受程度示意	141
图 5.7 术语词典知识组织结构模型	145
图 5.8 术语词条树形结构示意	147
图 5.9 术语词典辅助编纂系统功能	148
图 5.10 共现知识图谱构建：以“肿瘤学”为例	151

图 5.11 跨层次的知识组织与映射计算框架	152
图 A.1 国际知识组织协会性质与比例	158
图 A.2 国际知识组织协会国别分布	158
图 A.3 时间发展变化状况	159
图 A.4 知识组织相关学科分布	160
图 B.1 论文数量随时间变化趋势	166
图 B.2 期刊分布情况	167
图 B.3 城市发表论文统计	168
图 B.4 城市分布情况	168
图 B.5 研究单位分布	169
图 B.6 各重点单位研究成果时间段分布	170
图 B.7 作者分布的“长尾现象”	171
图 B.8 主题词表领域专家发表论文篇数统计（5 篇以上）	171
图 B.9 单位课题资助情况	172
图 B.10 课题基金资助与论文产出统计	173
图 B.11 “主题词表”“叙词表”共现关系	174
图 B.12 近 10 年来主题词表关键词共现关系网络	175
图 B.13 科研群体合作关系网络	176
表 2.1 《汉语主题词表》与 SKOS 的对应关系举例	41
表 3.1 4 种统计模型的拟合检验和参数评估	57
表 3.2 2000—2017 年肿瘤术语幂指数统计	58
表 3.3 网络百科知识分类比较	67
表 3.4 网络百科主题描述的比较	68
表 4.1 41 458 个术语二级及以上分类效率统计	90
表 4.2 6651 个术语分类准确率	90
表 4.3 g 指数计算过程	94
表 4.4 关键词对聚类结果的重要性	98
表 4.5 关键词聚类归并计算过程	99
表 4.6 同义归并正确率范围分布	109

表 4.7 同义术语推荐实例	115
表 4.8 术语同义关系识别准确率	116
表 4.9 术语释义句实例	120
表 4.10 术语释义句知识单元抽取结果	120
表 4.11 释义知识单元抽取结果准确率	123
表 4.12 术语知识库数据属性	123
表 5.1 词汇链构造结果示例	140
表 A.1 国际知识组织协会成立时间	159

第一章 绪论：知识为体，术语为用

本章导读



本章主要以术语和知识组织的关系为对象，介绍以下内容。

- 术语与知识组织的关系是什么？二者应如何结合？
- 大数据时代，知识组织的现状、趋势和现实需求是什么？为什么术语计算变得更加重要？
- 本书的研究目的和相关背景是什么？如何阅读和使用本书？

1 术语：打开知识宝库的一把钥匙

1.1 术语与知识组织：如影随形的黄金搭档

术语（Terminology）是“在特定专业领域中一般概念的词语指称”^[1]，是科学共同体记录科研成果、进行知识交流的重要交际工具。随着知识的不断积累，术语的数量和所涉及的领域也在逐渐丰富，已经成为人们开展科技交流、学术研究和生产劳动的基础性工具。近年来，随着自然语言处理技术的深入发展和海量文献的迅猛增长，术语学与知识工程、自然语言处理技术等深度结合，人们开始以术语为基础，从工程化的角度进行专业领域知识的表示和挖掘，形成可供计算机利用的专业领域知识库，这已经成为术语研究的重要课题，并引起人工智能领域的关注^[2]。

知识组织（Knowledge Organization）是“以知识为对象的诸如整理、加工、表示、控制等一系列组织化过程及其方法”^[3]，本质上对知识进