

当代语理论丛刊 Contemporary Linguistic Theory Series



Introduction to Computational Linguistics

计算语言学导论

翁富良 王野翊 著

Weng Fu-Liang Wang Ye-Yi

主编 黄德超 编辑 徐

Quet Editors James Huang De-Bao Xu

中国社会科学出版社

当代语言学理论丛书

主编 黄正德 许德宝

计算语言学导论
Introduction to
Computational Linguistics

翁富良 王野翊 著

Weng Fu-Liang Wang Ye-Yi

中国社会科学出版社

(京)新登字 030 号

图书在版编目(CIP)数据

计算语言学导论/翁富良,王野翊著. —北京:中国社会科学出版社,1998.9

(当代语言学理论丛书)

ISBN 7-5004-2080-3

I. 计… II. ①翁…②王… III. 统计语言学 IV. H087

中国版本图书馆CIP 数据核字(98)第 01118 号

0872 / 18

中国社会科学出版社出版发行

(北京鼓楼西大街甲 158 号)

北京奥隆印刷厂印刷 新华书店经销

1998 年 9 月第 1 版 1998 年 9 月第 1 次印刷

开本: 850×1168 毫米 1/32 印张: 7 插页: 2

字数: 178 千字 印数: 1—3000 册

定价: 16.00 元

序　　言

语言学自乔姆斯基以来，对认知科学、心理学、医学、电子计算机以及人工智能等学科都产生了巨大的影响，成为人文科学的带头学科。只要在国外走一走，就会发现几乎所有的大学都设有语言学系或语言学专业。语言学理论不但对语言学系的学生至关重要，而且也是心理系、教育系、社会学系、认知学理论乃至计算机系的学生必修的基础理论课。乔姆斯基的语言学理论为什么对人文科学和社会科学的影响如此之大？他的什么变革使本来默默无闻的语言学（理论）一跃而成为认知科学、心理学、电子计算机以及人工智能等学科的奠基理论？这不是一句话能说清楚的。要回答这个问题，得从现代语言学的立足点说起，系统介绍现代语言学的基本理论和研究方法、研究对象、研究范围以及研究结果等。不说清楚这些问题，现代语言学在人文科学中的带头作用和对社会科学的巨大影响也就无法说清楚。有系统有深度地介绍现代语言学理论，这就是我们这套丛书的编写目的。

要系统介绍现代语言学，各种理论的来龙去脉都得交代清楚，某种理论的发生、发展、不同阶段以及各个流派之间的关系都要说清楚。不能只把一种理论搬来，不管它的过去和与其他理论的联系，那样会让人不知所云。在系统介绍的同时，也要把各种理论的最新研究成果写进去，并评价其优劣不同以及对现代语言学研究的贡献等，做到有深度、有系统，这是我们介绍的第一个原则。介绍的起点一般是以乔姆斯基与哈利的《英语语音系统》

(1968) 为始,介绍的终点就是今天,介绍时以八九十年代发展起来的语言学理论为主,所以这套书叫作《当代语言学理论丛书》。

要介绍现代语言学不容易。台湾、新加坡、香港等地的学者有很好的经验。他们介绍的特点就是把现代语言学理论与汉语的研究结合起来。这样理解起来方便得多,效果也就比较好。单纯介绍,不谈在汉语中的应用,结果理论还是死的东西。我们这套丛书也本着这一点,在选材和编写上都强调在汉语中的应用,尽量用汉语说明。汉语与某种理论不相关的时候,才用其他语言中的例子。这是我们介绍的第二个原则。

我们的第三个原则是以介绍美国语言学理论为主。美国是现代语言学研究的中心,也是生成语言学的发源地。要介绍现代语言学就离不开这个发源地。所以从选材上来讲,我们以美国语言学系研究生(博士和硕士)的必修课为标准,包括语言学史、句法学、音系学、语义学、心理语言学、社会语言学、历史语言学、语言获得理论、计算机语言学与人工智能等。有些新兴学科和边缘学科就放在主要学科中介绍。比如神经语言学归入了心理语言学,音系与句法的交叉研究归入了音系学,语义和句法的交叉研究归入了语义学等。

应该指出,有些学者一直在致力于现代语言学的介绍工作,比如黑龙江大学、上海复旦大学、天津师范大学的学者等。我们希望这套丛书能与他们的研究结合起来,起到使国内外语言学研究接轨的作用。

《当代语言学理论丛书》的编写开始于1993年,由著名句法学家黄正德教授全面负责,许德宝协助作主编工作。编委大都是在美国读语言学博士而且有教授语言学经验的学者,一般是在讲义的基础上增删整理成书。但即使是如此,也都得付出很多的劳动。我们也请了在美国教授多年的语言学家、汉学家和有在国内外介

绍现代语言学经验的学者作为顾问，帮助我们把这一套丛书出好。在此向他们谨致谢意。我们还得感谢中国社会科学出版社对这套丛书的大力支持，特别是责任编辑及其他有关同志的辛苦工作，不然这套丛书也不能和读者见面，在此也一并致以谢意。

《当代语言学理论丛书》编委会
1996年7月于纽约

当代语言学理论丛书
Contemporary Linguistic Theory Series

主 编
Chief Editors

黄正德(James Huang) 许德宝(De Bao Xu)
尔湾加州大学 纽约汉弥尔顿文理学院
University of California, Irvine Hamilton College, New York

编辑委员会
Editorial Board

(按拼音顺序)

包智明(新加坡国立大学)
Zhiming Bao(National University of Singapore)

端木三(密西根大学)
Duanmu San(University of Michigan)

冯建明(新罕布什尔达木斯文理学院)
Jianming Feng(Dartmouth College, New Hampshire)

胡明亮(缅因鲍登文理学院)
Hui Ming Hu(Bowdoin College, Maine)

Mingliang Hu(Bowdoin College,Maine)

蒋严(香港理工大学)

Yan Jiang(Polytechnic of Hong Kong)

靳洪刚(纽约汉弥尔顿文理学院)

Hong Gang Jin(Hamilton College,New York)

李亚飞(威斯康星大学)

Yafei Li(University of Wisconsin,Madison)

陆丙甫(南加州大学)

Bingfu Lu(University of South California)

潘海华(香港中文大学)

Haihua Pan(Chinese University of Hong Kong)

沈榕秋(伯克莱加州大学)

Rongqiu Shen(University of California,Berkeley)

石定栩(香港理工大学)

Dingxu Shi(Polytechnic of Hong Kong)

侍建国(新加坡国立大学)

Jianguo Shi(National University,Singapore)

宋国明(劳伦斯大学)

Kuo-ming Sung(Lawrence University,Wisconsin)

陶红印(新加坡国立大学)

Hongyin Tao(National University of Singapore)

王野翊(卡内基·梅隆大学计算机学院)

Ye-Yi Wang(Language Technology Institute,School of
Computer Science,Carnegie Mellon University)

翁富良(国家语音技术研究室)

Fuliang Weng (Speech Technologies and Research Lab,
California)

吴建慧(伊利诺大学)

Mary Wu(University of Illinois ,Champaign—Urbana)

谢天蔚(戴维斯加州大学)

Tianwei Xie(University of California ,Davis)

徐大明(新加坡南洋理工大学)

Daming Xu(Nanyang Technological University ,Singapore)

许德宝(纽约汉弥尔顿文理学院)

De Bao Xu(Hamilton College ,New York)

张敏(新加坡国立大学)

Min Zhang(National University of Singapore)

顾问编辑委员会 Advisory Editorial Board

(按拼音顺序)

陈渊泉(圣地亚哥加州大学)

Matthew Chen(University of California, San Diego)

戴浩一(俄亥俄州立大学)

James H.-Y. Tai(Ohio State University)

丁邦新(伯克莱加州大学)

Pang-Hsin Ting(University of California, Berkeley)

李艳慧(南加州大学)

Audrey Li(University of South California)

汤廷池(台湾国立清华大学)

Ting-Chi Tang(Taiwan Qinghua University)

王士元(伯克莱加州大学)

Williams S.-Y. Wang(University of California, Berkeley)

薛凤生(俄亥俄州立大学)

Feng-Sheng Hsueh(Ohio State University)

徐烈炯(香港城市理工学院)

Liejiong Xu(City Polytechnic of Hong Kong)

郑锦全(伊利诺大学)

Chin-Chuan Cheng(University of Illinois, Champaign-Urbana)

作者介绍

翁富良 斯坦福国际研究所语音技术和研究实验室研究工程师。1984年毕业于复旦大学计算机科学系。在1984—1989年间,师从吴立德教授进行模式识别和自然语言理解方面的研究。1989年,赴卡内基—梅隆大学机器翻译中心继续自然语言理解的研究。1993年,在新墨西哥州立大学获硕士学位。自1994年起,在斯坦福国际研究所的语音技术和研究实验室从事语言、语音模型研究。先后单独或与同事合作在一些专业杂志和会议上发表论文20余篇,曾获1986年国家教委科技进步一等奖,第三届中国国家自然科学四等奖。

王野翊 1985年于上海交通大学计算机科学与工程系获学士学位。后师从上海交通大学孙永强教授和中科院数学研究所陆汝钤研究员进行自然语言处理的研究,并于1988年获上海交通大学计算机科学与工程系硕士学位。1992年于美国卡内基—梅隆大学(Carnegie Mellon University)获计算语言学硕士学位,现为卡内基—梅隆大学计算机科学学院语言技术研究所(Language Technologies Institute, School of Computer Science)博士候选人。主要研究课题包括统计学机器翻译、语言模型、语言学习、神经网络。

前　　言

语言是反映人的思维的最重要的一面镜子,又是人与人之间交流的最重要的媒介。对语言的研究,是一个经久不衰的古老课题。几千年来,中外学者从语言与思维、语言与现实的关系等不同角度,在词源、注释、分类、语法等各个方面开展了广泛的研究。近一两百年来,西方学者在逻辑学、数学和分析哲学方面的成果,大大推动了语言形式化的研究。而随着计算机科学的发展,建立语言的形式化计算模型成为语言学的重要课题。计算语言学作为以形式化的计算模型来分析、理解和处理语言的科学也就应运而生。而信息革命的展开,更使计算语言学的研究达到了一个空前的程度。

如果说纸的发明对人类文明的继承光大有着巨大作用的话,那么语音和语言技术的发展,将对人类各语种之间的交流,各文化体系间的促进与提高至关重要。语言和语音技术为语言文本和会话的检查、理解、合成、翻译、重组,提供了有效的自动化工具,使得靠人工进行的信息交流和信息处理能够逐步地为具有智能的语言技术所取代。在信息革命使世界日新月异的今天,计算语言学的发展成为我们是否能够跟上世界潮流的一个重要因素。这一点应该激起广大研究人员的高度紧迫感,同时也应该获得科研基金组织和工商界有识之士的重视,对计算语言学的研究给予长期的大力支持。

本书的作者希望此书的出版能够对国内计算语言学的发展起到一定的促进作用。由于篇幅和作者水平的限制,我们不可能面

面俱到地覆盖整个领域。我们在选材时一方面注重本领域的基础性的经典工作,希望读者阅读理解后能够举一反三,用于解决实际问题;另一方面我们侧重于介绍一些当前国际计算语言学界的研究重心,希望有关的研究人员能够站在该领域的前沿。

本书的对象是大学计算机专业、数理统计专业及语言学专业的高年级学生或研究生,与计算语言学有关的科研人员,以及其他有兴趣的读者。由于计算语言学的综合性特点,如有条件,作者建议组织多学科的兴趣小组,相互交流,共同提高。

本书第2章第1节、第3章、第4章、第5章、第6章、第7章由翁富良撰写,第1章、第2章第2、3节、第8章、第9章、第10章由王野翊撰写。

目 录

第一章 计算语言学简介	(1)
第一节 计算语言学是一门边缘科学	(1)
第二节 计算语言学研究的基本问题	(3)
第三节 计算语言学研究的基本方法	(4)
一、理性主义和经验主义：计算语言学研究	
方法的哲学分野	(4)
二、计算语言学研究方法	(5)
第四节 计算语言学的应用	(8)
第二章 预备知识	(10)
第一节 离散数学基础	(10)
一、集合及相关的概念	(10)
二、图及相关的概念	(11)
三、字符串及相关的概念	(12)
四、栈及相关的概念	(15)
五、序及相关的概念	(15)
第二节 概率统计理论基础	(17)
第三节 信息论基础	(26)
第三章 形式语言及自动机	(34)
第一节 形式语言和自动机的直观意义	(34)
第二节 形式语言和自动机的定义	(35)
一、形式语言的定义	(35)

二、自动机的定义	(39)
第四章 语法学理论和表示形式	(46)
第一节 GB 理论	(46)
第二节 词汇功能语法	(49)
第三节 广义词组结构语法	(51)
第四节 树连接语法	(53)
第五节 链语法	(55)
第五章 语言的识别与分析	(57)
第一节 有限状态语法的识别和分析算法	(58)
第二节 上下文无关语法的识别和分析算法	(58)
一、移进—归约法	(58)
二、由底向上的图表法	(64)
三、欧雷算法	(69)
四、GLR 算法	(70)
五、链语法的识别算法	(82)
第三节 其他类型的分析器	(85)
一、基于原则的分析方法	(85)
二、基于归一的分析方法	(87)
第六章 计算语义方面的一些工作	(91)
第一节 语义理论简介	(91)
一、词的指称作为意义	(91)
二、心理图象，大脑图象或思想作为意义	(92)
三、说话者的意图作为意义	(92)
四、过程语义	(93)
五、词汇分解学派	(93)
六、条件真理模型	(94)
七、情景语义学	(94)
八、语义网络	(95)

九、模态逻辑	(95)
第二节 计算语义学的一些代表工作	(96)
一、概念依赖理论	(96)
二、选择限制学说	(99)
三、指代化解	(101)
四、计算语义学的一些其他方面	(103)
第七章 容错分析	(106)
第一节 基于关键词或中心词的方法	(106)
第二节 省略不识词的方法	(107)
第三节 元规则方法	(108)
第四节 同化法	(109)
第八章 概率语法	(116)
第一节 Ngram	(116)
一、减值法 (Discounting)	(118)
二、删除插值法 (Deleted Interpolation)	(121)
三、基于词分类的 Ngram	(122)
第二节 隐马尔柯夫模型	(122)
一、马尔柯夫模型	(122)
二、隐马尔柯夫模型	(124)
三、向前算法	(126)
四、韦特比算法	(130)
五、向前向后算法	(132)
第三节 概率上下文无关文法	(136)
一、向内算法	(137)
二、韦特比算法	(139)
三、向内向外算法	(140)
第九章 语言学习	(145)
第一节 词分类	(145)

第二节 词法学习	(148)
一、语法框架	(148)
二、词汇选择 (Lexical Selection)	(152)
第三节 语法学学习	(155)
一、有限状态自动机的机器学习	(155)
二、语法推导的理论问题	(159)
三、贝叶斯推理在语法推导中的应用	(161)
第十章 当前计算语言学的研究	(166)
第一节 统计学机器翻译	(166)
一、IBM 统计学机器翻译	(167)
二、参数训练	(168)
三、源语言搜索	(169)
第二节 词类标识 (Part – of – Speech Tagging)	(170)
一、隐马尔柯夫模型词类标识	(171)
二、基于规则的词类标识	(172)
第三节 歧义化解 (Disambiguation)	(174)
一、基于结构的语法歧义化解	(175)
二、统计学语法歧义化解	(177)
三、词汇歧义化解 (Lexical Disambiguation)	(180)
附录 A 汉英术语对照	(185)
附录 B 有关计算语言学的重要期刊和会议	(194)
附录 C 参考文献	(196)