

语言与语言学系列

# 语言的认知研究和计算分析

袁毓林 著

北京大学出版社  
北京

**图书在版编目(CIP)数据**

语言的认知研究和计算分析/袁毓林著. —北京:北京大学出版社,1998.10

ISBN 7-301-03843-7

I. 语… II. 袁… III. 语言-认知-文集 IV. HO

中国版本图书馆 CIP 数据核字(98)第 22129

**书 名:语言的认知研究和计算分析**

著作责任者:袁毓林著

责任编辑:胡双宝

标准书号:ISBN 7-301-03843-7/H·413

出版者:北京大学出版社

地址:北京中关村北京大学校内 100871

网址:<http://cbs.pku.edu.cn/cbs.htm>

电话:出版部 62752015 发行部 62754140 编辑部 62752032

电子信箱:[zpup@pup.pku.edu.cn](mailto:zpup@pup.pku.edu.cn)

排印者:北京大学印刷厂

发行者:北京大学出版社

经销者:新华书店

850毫米×1168毫米 32开本 10.875印张 275千字

1998年10月第一版 1998年10月第一次印刷

定 价:16.00元

## 内 容 提 要

本书是关于语言的认知研究和计算分析的一组论文的结果，中心内容是：尝试从认知科学的角度对语言的结构方式和语义理解的心理机制进行研究并加以计算分析，以探索语言研究怎样为计算机理解自然语言提供恰当的方法和合适的规则。透过本书的各篇文章，读者可以大概地了解到：流行于哲学、心理学、语言学、计算机科学和神经科学中的认知主义思潮，认知语言学的基本观念及其在句法、语义上的各种研究路线；语言知识的形式表示和计算分析的有关方法和步骤，计算机理解自然语言所必需的结构形式线索和所调用的规则与策略；怎样用动词和名词的配价、谓词隐含等概念来刻画语句中不同词项之间的句法、语义连结方式，并用扩散性激活的语义记忆模型和缺省推理这种非单调逻辑来建立一种面向计算的语义理解的微观机制；怎样用心理学上的原型理论来分析汉语的词类问题，怎样用家族相似性的观念来证明根据词在分布上的相似性来聚类的可能性；语言学范畴具有心理现实性，语言信息和生物信息在编码时都遵循双重分节的原理。

本书可供对现代语言学、认知心理学和人工智能(特别是机器学习、自然语言的计算机处理)感兴趣的读者阅读和参考。

## 目 录

序.....	陆俭明(1)
1. 认知科学背景上的语言研究 .....	(1)
2. 语言的认知研究和计算分析 .....	(24)
3. 关于认知语言学的理论思考 .....	(49)
4. 现代汉语二价名词研究 .....	(74)
5. 一价名词的认知研究 .....	(103)
6. 简论计算机理解自然语言 .....	(126)
7. 句法空位和成分提取 .....	(151)
8. 谓词隐含及其句法后果 .....	(166)
9. “者”的语法功能及其历史演变 .....	(194)
10. 词类范畴的家族相似性 .....	(221)
11. 基于原型的汉语词类分析 .....	(247)
12. 语言学范畴的心理现实性 .....	(270)
13. 话题化及相关的语法过程 .....	(279)
14. 语言信息的编码和生物信息的编码之比较 .....	(308)
后 记.....	332



理工科的学生开设语言分析导论、汉语语法分析课,这又促使他学习了计算机科学方面的一些知识,并阅读了一定的有关计算语言学方面的论著。收入本论文集的文章就是他这几年来有关汉语的认知研究的初步成果。

认知语言学的理论背景主要是认知科学,而认知科学本身应该说还是一门新兴的学科。目前,无论在心理学界、哲学界、神经生理学界、计算机科学界,还是语言学界,对于人脑的认知活动的说明都众说纷纭。在这种情况下,语言的认知研究,具体说,应以什么样的认知理论为指导理论,应确立什么样的研究导向,在研究中应以什么为突破口,研究途径是什么,应采用什么样的研究策略,到底怎样界定认知语言学,等等,也都众说纷纭。而袁毓林同志又有他自己的看法。他在本书中(19—20页)认为:

我们倡导语言的认知研究必须以基于信息加工观点的认知心理学为背景,以人工智能的语言信息处理为导向,重视对语言理解的认知加工过程和语义信息处理的微观机制的研究,希望研究结果能在计算机上模拟实现和验证,从而使这种研究途径带有一定的技术色彩。在研究策略上,我们倾向于采用通过一些假设的概念来建立认知模型的方法。具体地说,通过对语言的结构和意义进行考察和分析,形成一些基本的理论概念(如名词配价、谓词隐含、缺省推理等),并在认知科学的指导下把相关的概念组成一个整合的模型,以便:

(1)组织观察并使观察具有意义;

(2)把所观察到的所有要素联系起来,以便作出更为宏观和抽象的假设;

(3)指导以后的观察并预测未来观察到的事件;

(4)通过明确的概念和直观的模式来跟其他研究者交流和对话。

他收入本书的有关从认知角度具体讨论汉语中的某些语法现象的论文就是按他对认知语言研究的认识和思路来写的。对于袁毓林同志关于语言的认知研究的认识和思路,我不想在这里加以评论,大家在阅读本书后可以进行讨论和评判。不过我在这里要指出的

## 序

《语言的认知研究和计算分析》是袁毓林同志的一本论文集，所收的论文都是九十年代写的，从书名可以知道，论文所谈的内容，都是关于语言的认知研究和语言的计算分析的。

我们知道，从人类文明的历史看，人类已经历了两个大的时代——农业时代和工业时代，而现在又进入了一个新的时代——信息时代。这个时代将以进一步开发计算机以便进一步减轻人的体力劳动并逐步减轻人的脑力劳动为主要特点。九十年代开始，更明确提出了建造“信息高速公路”、研制智能计算机的任务。因此，现在比以往任何时候都更迫切需要探索和揭示人脑处理、储存信息的机制和人的语言机制，以及它们之间的相互关系；同时要求我们还需从认知的角度去描写说明语言规则，并加以形式化，使之具有可计算性，从而为计算机处理、理解自然语言提供切实有用的语言研究成果。正是在这种背景下，语言的认知研究和语言的计算分析逐渐成为语言研究中的一个热点。

在中国，语言的认知研究可以说还刚刚起步，从事这方面研究的人还很少。袁毓林同志这几年来一直比较关注认知语言学，在这方面进行了一定的研究。他也有这方面的条件。他在杭州大学攻读硕士学位时师从王惟贤教授，王惟贤教授很强调从逻辑的角度来研究语言，所以袁毓林同志在逻辑方面受到比较好的训练。他在北京大学攻读博士学位时师从朱德熙等教授，这又使他受到了现代语言学，特别是严格的汉语语法研究的训练。他博士研究生毕业、获得博士学位后，应聘到清华大学中文系任教，给计算机系等

是,他关于谓词隐含、缺省推理和扩散性激活的语义记忆模型等理论对解释为什么“小王的帽子”有时可以用“小王的”来指代,而“小王的爸爸”就不能用“小王的”来指代;对解释为什么在汉人的心目中“这房子很大”是说房子的面积大,而“这箱子很大”是说箱子的体积大;这应该说比汉语语法学界先前已有的解释要更深刻些,也更能让人理解。而他根据认知科学里的“典型范畴”(prototypical category)理论所提出的用“词类范畴家族相似性”理论概念来思考、解决汉语词类问题的观点,我觉得这比用“柔性观念”来思考、解决汉语词类问题的看法似乎也更深刻,更具说服力一些。当然在这里我也不能不指出,从本书所收的文章看,袁毓林同志对于汉语的认知研究还是很初步的;本书还是介绍多,理论思考多,而怎么用认知语言学的理论去分析、解释汉语中的一些具体的语言现象还是太少。

语言的计算分析现在越来越受到人们的注意,但运用什么样的计算理论,具体使用什么样的计算方法,如何进行运算,大家的看法也不一致。关于这一方面我完全是门外汉。我只是觉得袁毓林同志所谈的语言的计算分析跟我们国内其他一些人所谈的语言的计算分析不一样。孰是孰非?是不是计算的方面或角度不一样?我不知道。

认知语言学注重分析、描述和解释人类语言构造的认知功能基础,其基本主张是:自然语言是概念化的现实的符号表达,句法结构在相当程度上不是任意的、自主的,而是有自然的动因(motivation),也就是说,句法结构的外在形式常常是由认知上的因素所促动的,与人的身体经验、认知策略乃至文化规约等密切相关。毫无疑问,语言的认知研究为我们研究语言现象提供了一个新的观察角度,启发我们形成新的研究思路。这种研究往往能发前人之所未发,它不仅能帮助我们更好地解释汉语的语法规律,而且也能帮助我们更深入地发掘汉语的语法规律。所以认知语言学作为一

种新的理论,很值得我们重视,要让它跟汉语研究密切结合,进一步发展。在上面我们已经指出,在我国语言的认知研究可以说还刚刚起步,汉语的计算分析也还比较薄弱。这就需要大家,特别是从事汉语研究的学者、从事计算机科学研究的学者以及从事自然语言理解和处理的学者进一步通力合作,群策群力,为尽早使中文信息处理达到预想的实用的目标,适应研究智能计算机的需要,而不断探索,不断努力。袁毓林同志的这本论文集对汉语的认知研究和汉语的计算分析无疑会起到积极的作用。是为序。

陆俭明

1998年5月8日

于北大中关村寓所

# 1. 认知科学背景上的语言研究

本文介绍和评论广泛流行于哲学、心理学、语言学、计算机科学和神经生理学中的认知主义思潮,分析在此基础上形成的认知科学的学科背景和学术源流。其中涉及:(1)哲学上的心脑二元论和物理主义思潮,神经生理学对大脑中神经元之间的通讯原理的研究;(2)智能起源上的天赋论和建构论,以及各自对语言能力的解释力量;(3)基于信息加工观点的认知心理学,以及跟皮亚杰发生心理学的比较;(4)作为人工智能研究和认知心理学的理论基础的物理符号系统假设,以及与之抗衡的神经网络思想;(5)哲学家对认知主义的批判,以及我们的反批判;(6)语言研究上的两种认知主义的操作路线,以及另外一种可能的研究途径。

## § 0 引言:认知主义和结构主义

自乔姆斯基革命以来,当代语言学呈现出范式纷纭、学派林立的景象。如果不计较技术和细节上的巨大差异,勉强概括各种研究路子背后的共同特征;那么,当代一些有影响的语言学流派大都可以归在认知主义这个名目之下。而前此的强调按照语言的内部结构去理解和研究语言的思潮便是结构主义。从学术背景上看,结构主义语言学的哲学基础是广泛流行于当时知识界的逻辑实证主义,心理学基础是摒弃心灵的行为主义心理学。因此,结构主义语言学只注意对现实的话语素材进行切分和描写,拒绝讨论难于摆脱心灵主义纠缠的意义问题。认知主义语言学是对结构主义语言

学及其学术背景的反动,其哲学基础是当代的分析哲学(包括语言哲学),心理学基础是注重内省的各种心灵主义心理学。事实上,认知主义语言学的不同流派在研究取向上是十分不同的。比如,乔姆斯基的生成语法是一种注重形式的认知主义,它深受数学和逻辑原理的形式化研究的影响;R. Langacker 的认知语法是一种注重概念的认知主义,它强调作为语言核心的意义是一种心理现象(因此必须按心理现象来描写);戴浩一的认知功能语法是一种注重语言的信息传递功能的认知主义,它强调语法结构来自对现实的象征(所以语言学必须研究语法结构在观念方面的理据);我们则倡导一种注重可计算性的认知主义,摸索一种基于认知并面向计算的语言研究的途径(toward a cognition-based & computation-oriented approach to linguistic study)。

认知主义是当代广泛流行于哲学、心理学、语言学、神经科学和计算机科学等研究领域的一种思潮,要想对这样一种源远流长又旁涉广泛的学术潮流作出全面而恰当的介绍是不可能的。我们打算紧扣认知、智能、心智、知识、计算和表达等认知主义的核心概念,来讨论它们的学科背景和学术源流,希望能以点带面地展示认知主义的研究纲领和具体成果。

## § 1 心脑二元论和物理主义思潮

认知科学(cognitive science)是一门以智能系统(包括自然的和人工的)的工作原理为研究对象的新兴学科,它是从哲学、心理学、语言学、计算机科学和神经生理学的交叉领域中发展起来的前沿性学科<sup>①</sup>。在认知科学中,智能(intelligence)是一个核心概念。但是,哲学家对于智能的看法相当分歧。比如,笛卡尔(René Descartes, 1596—1650, 法国哲学家、生理学家、解析几何的创始人)提出了著名的二元论,认为心智(mind)跟身体是两种完全不

同并可以分离的实体；心智的本质是思想，能够脱离身体而存在<sup>②</sup>。与此相对的是物理主义思潮<sup>③</sup>，物理主义(physicalism)强调一切心智事件(精神现象和智力活动)最终都可以还原成大脑的生理过程和物理过程。这种观点把智能看作是人脑这种特殊物质的一种功能，相信心理过程以脑过程为基础并可以还原为脑过程。随之而来的问题是：如何用脑过程这种低层次的现象去说明心理过程这种高层次的现象？

为了回答这个问题，神经生理学从物理、化学角度分析大脑神经元(neuron, 即神经细胞)传递信息的电学和化学过程，证明神经细胞的电学性质与其化学性质是互相联系、互相影响的。神经生物学家设想精神活动可能跟脑子里的神经脉冲(impulse)的图型相关。为了充分评价这一假设的含义，他们研究了神经元是怎样工作的、它们相互之间是怎样通讯的、它们是怎样组成局部的或分布的网络的、它们之间的连接是怎样随经验而改变的。

当代神经科学证明，大脑的结构是非常复杂的。表现为：(1)大脑中至少有  $10^{10}$  个神经细胞，大体相当于银河系中星球的数量级；(2)神经细胞的多样性，有的细胞具有能跟相邻细胞通讯的短轴突、有的细胞具有伸向其他区域的长轴突。当然，神经细胞也有其简单性的一面。功能相似的神经元都集合成柱状或板状而穿过皮层，并且所有的神经元都是以非常相似的方式传递信息的——信息沿着轴突(axon)以动作电位(action potential)的短脉冲的形式传播。具体地说，一个受刺激的神经元产生一种叫动作电位的脉冲而把信号输送给另一个神经元。这些信号像波浪一样沿着神经细胞的单个轴突(即神经纤维)传播，并在突触(dendrite, 即神经细胞之间的接触点)处转化成化学信号。神经元静止时，它的外膜保持大约-70毫伏的电位差(内表面相对于外表面是负的)。静止时，钾离子(K)比钠离子(Na)易于透过膜。细胞受刺激后，膜对钠的渗透性加强，从而导致正电荷的流入而起动一个脉冲，于是膜电位瞬

即逆转。脉冲起始于细胞体和轴突的接合部并引离细胞体。脉冲到达突触前神经元的轴突末端后就诱导释放递质分子。递质通过狭窄的突触间隙而和突触前膜上的受体相结合。这样的结合作用会打开离子通道,而其本身又往往会导致在突触后神经元里产生动作电位<sup>④</sup>。这就是现代神经生理学对神经元之间的通讯原理解释。

总之,神经科学家相信:大脑结构的复杂性足以解释人的记忆、情绪和想像力等心智奥秘,精神活动可以还原为大脑神经网络上的电信号<sup>⑤</sup>。

## § 2 人类智能的天赋论和建构论

在智能的来源问题上,当代心智哲学和心理学中有天赋论(innatism)和建构论(constructivism)两大理论派别。天赋论可以乔姆斯基为代表,他深受笛卡尔哲学的影响,信奉源于柏拉图(Plato)的天赋论。乔姆斯基反对经验论把人类的心智看作是一块“白板”、一切知识都是后天经验习得的想法,他反对行为主义心理学“刺激—反应”的模式。他认为人的智能结构和认知能力是人类这个物种的大脑生物学结构所固有的,这种潜在的结构和能力一旦受到外部诱因的驱动就能被激活,产生出观念和知识。乔姆斯基强调:(1)人具有生成和理解无限多新句子的语言能力,它是人类心智的组成部分,它是人的一种固有的机制,正是这种机制把人类的经验映射为语法;(2)语言中存在成体系的普遍现象(universals),但缺少证据说明它们是经过学习或经验而习得的,可能的解释是作为生物本能的普遍语法(universal grammar)使然。语言作为生物本能具有部分确定的结构(表现为普遍语法),正如生理器官的普遍性对人类来说是确定的一样。乔姆斯基据此证明:人类的所有知识都可以从这个物种天赋的心智特征中推导出来<sup>⑥</sup>。

建构论可以瑞士心理学家和哲学家皮亚杰(Jean Piaget, 1896—1980)为代表。皮亚杰深受从洛克(John Locke)、休谟(David Hume)到斯宾塞(H. Spencer)的古典经验论的影响,在几十年的生物学、哲学、逻辑学和心理学(特别是儿童实验心理学)研究的基础上提出:认知的结构既不是在客体中预先形成的,因为这些客体总是被同化(assimilate)到那些超越于客体之上的逻辑数学框架中去;也不是在必须不断地进行重新组织的主体中预先形成的。认识起因于主客体之间的相互作用,有机体自我调节(self-regulation)的适应机制使客体被同化到主体的图式(scheme)中。图式来源于动作,作用于新客体的活动反复多次并通过一般化而形成图式<sup>⑦</sup>。

皮亚杰认为,跟所有的生物的发展一样,在智力的成长过程中起作用的主要有两条原则:适应和组织。适应(adaptation)指我们的心理越来越有效地对环境的需要作出反应,它有同化和顺应两个双向建构的过程,其中起协调作用的机制是调节(accommodation)。同化指认知主体将客体纳入主体的图式之中,顺应指认知主体调整原有的图式并创立新的图式以适应新的客体。调节则指认知主体如何控制同化和顺应的双重建构以达到平衡,即图式之间的分化和协调问题。组织指我们的心理以越来越整合的方式来结构或组织,其中最简单的水平是图式,它是能作用于物质的某种行动(物质的或心理的)的心理表征(psychological representation)。对于新生儿来说,吸吮、抓握和看都是图式,它们是新生儿得以认识世界的方式——通过作用于世界的动作来认识世界。这种在遗传的反射图式的基础上通过同化和顺应而建立起来的图式叫动作图式,这是在感觉运动的基础上形成的认知的第一图式。动作图式一经形成便会在同化、顺应和调节的不断循环中走向完善。接着,儿童在用某种方式和手段取得某种结果的过程中,逐步建构起有目的控制的图式——目的图式,这便是智能的萌芽。在此基础

上,儿童便会形成一系列能应付新情况的经验图式,并进一步发展成为抽象思维——形式运算图式。这些图式有规律地变得更加整合和协调,最终产生出成人的心理<sup>⑥</sup>。因此,智能的发展就是人的认知结构连续地建构和再建构的过程,是一个从建构初级图式到高级图式的实践的过程。

综上所述,乔姆斯基是预先决定论者(predeterminist),强调作为人类最基本的智能要素的语言能力的认知结构是天赋的,是在人的器官中甚至基因中就早已编制好了的(preprogramming),后天的发育和环境因素只不过是促使这种结构成熟而已。相反,皮亚杰是相互作用论者(interactionist),强调认知结构是后天建构的,智能、知识只能来自后天的活动、操作和实践。但是,乔姆斯基的天赋论在生物学和心理学上都是很难完全确证或否证的。因为人种突变从生物学上看是一种随机的、不好解释的现象。为什么这种随机的突变会给人提供了一种表达力很强的语言能力呢?怎么能把一种合理的语言结构归因于这种随机的突变呢?不过,皮亚杰源于感觉运动的智能建构也说明不了普遍语法不因个体和语种而异这一事实<sup>⑥</sup>。

### § 3 信息加工观点和认知心理学

在智能的研究上,认知心理学采取计算机模拟以及其他实验方法来探索人类心智的一般的工作原理。认知心理学认为人的动机在表现为意识之前要经历一系列的转变、改造和处理,就像通讯过程需要将信号进行编码和译码一样。因此,可以把人看作是一个信息加工系统,人与环境和其他认知主体之间有着信息交流的关系。信息在进入认知主体之后又会经历被转换、简约、合成、存储、重建、再现和使用等加工过程。通过考察信息加工的链条,人们就可以对系统的行为作出预言,并分析这种行为与环境的关系。这

样,认知心理学为研究人的行为提供了一个全面而一贯的工作模型。由于认知心理学把人类的心智过程看作是一个信息加工(information processing)过程,因而可以利用信息加工的概念和方法来研究人类内在的认知过程。比如,认知心理学运用信息和编码等概念来研究记忆,用存储、检索、控制、缓冲器和网络图等概念对记忆的结构和动态过程进行比喻性的描写,使记忆等一些过去比较含糊的概念转化为可操作的术语,从而为建立形式化的认知模型奠定了基础<sup>⑩</sup>。

认知心理学的理论目标是要解释人在完成认知活动时是怎样进行信息加工的,研究内容包括:(1)人是如何从世界中获取信息的;(2)这些信息是怎样表征的,即以什么样的表达形式存储在人的大脑中;(3)这些信息是怎样转化为知识的;(4)这些知识在大脑中是怎样表征的;(5)这些知识在人解决问题时是怎样被调用的;(6)大脑中的各种信息和知识是怎样指导人的注意和行为的<sup>⑪</sup>。其中,(1)讨论感觉信号的检测和加工、模式识别、知觉广度和选择性注意等问题;(2)一(4)讨论记忆的结构和过程、记忆模型、短时记忆及其容量、长时记忆的组织 and 信息回收、记忆组织及其记忆模型等问题;(5)(6)讨论在概念形成、问题求解、语言理解等高级认知活动中人怎样调用知识和信息加工的策略等问题。可见,认知心理学一方面受惠于计算机科学,另一方面又为计算机模拟人类智能提供了可操作的认知模型。

通过上面的介绍可以看出,基于信息加工观点的认知心理学跟皮亚杰的认知观点在方法论上有着根本性的差别。皮亚杰的理论基本上是一种认知一元论,即假定所有的认知发展都能由单一的、统一的原则集合来加以解释。他的图式发展理论、图式的整合和调节、更加一般的适应和组织原则,形成了一个连贯又统一的理论体系。这种理论坚持建立科学理论的一个崇高的原则:试图用少数一般原理来解释极其多样的智力过程。在理论体系的建构方面

也符合奥卡姆剃刀(Occam's Razor)原则:除非必要,决不在理论体系中增加假设的实在的数目。与此相反,认知的信息加工理论基本上是一种带有实用主义(Pragmatism)色彩的认知多元论。心理被看作是形形色色的个别过程的集合,这些过程不一定遵循相同的规则,也不仅仅受一组发展原理的支配。注意可以离开记忆而单独发展,记忆也并不总是跟知识相联系的。信息加工理论没有皮亚杰理论朴素和高雅,但可以敏感地反映人类思维和智能的多样性,也比较适合于用计算机模拟和在程序上实现。皮亚杰理论试图明确描述认知发展赖以出现的机制,正是在这一点上信息加工理论表现出了明显的弱点:不能找出变化的机制,这在很大程度上是由于它遵循多元化这种方法论所致。必须指出,皮亚杰理论也有其方法论弱点:缺少可靠的实验基础,一批实验设计的合理性值得怀疑<sup>⑨</sup>。有人指出,在皮亚杰的实验中儿童所提的问题是不公平的或是引起误解的;这些问题对于儿童大都是不能理解的,因而引起了荒谬的回答<sup>⑩</sup>。

#### § 4 物理符号系统假设和神经网络思想

为了使计算机更加聪明,而不是只能按照预先为它编制好的程序进行操作,计算机科学家也开始对人类智能进行研究,希望设计出跟人一样能在新情况中作出恰当反应的机器。这种研制具有创新能力的、受程序控制的装置的努力,终于造就了计算机科学中的一门新学科——人工智能(artificial intelligence)。人工智能研究怎样让计算机来完成表现出人类智能的任务,让机器模拟人脑所从事的推理、学习、理解和规划等思维活动,解决需要人类专家才能处理的复杂问题(如:医疗诊断、地质解释、气象预报、运输调度、管理决策和语言理解等)。那么,什么是人类的自然智能呢?人工智能研究者大致同意人类智能至少要包括以下四种能力<sup>⑪</sup>: