

自然语言理解的方法和途径

陈永明

中国科学院心理研究所, 北京

摘 要

本文结合汉语处理的问题,叙述和讨论了国外有关计算机理解自然语言的某些主要的方法和途径。作者认为,从语言理解研究的发展过程来看,计算机理解自然语言的研究,已由最初的主要是属于语言学的问题,逐渐成为主要是认知心理学的问题。因而,只有深入了解人类理解语言的实际过程,才能建造出有效的理解自然语言的计算机系统。

一、前 言

理解和使用自然语言,是人类的一个主要的智慧特征。除非让计算机也具有这种能力,否则它就不能称之为智能计算机。因此,自然语言理解的问题,是当代人工智能研究的核心课题。一方面,它试图参照人类处理语言的特点,建立语言理解的计算机模型,给计算机装上语言理解系统,使它能在一定范围内通过自然语言与人类交往;另一方面,更深远的目标,是要搞清人类理解语言的机制,探索人类思维的一般规律,使人自身成为一个更好的语言使用者,并建立更为有效的机器理解语言的模型。

“理解”是指把一种表征转换成另一种表征。这个概念没有多少绝对的意义。把一个句子从英语译成中文,这当然是理解;正确地派定一个句子的结构,这也是理解。因此,理解是有不同层次的。

理解自然语言,这是一种很困难的作业,因为语言是极为灵活、复杂的。语言材料几乎是无限的,总是在发展和扩大;语言的形式是多样的;同一意义可以由不同的形式来表达(同义歧形),同一形式也可表达多种意思(同形歧义);语言使用中,经常出现各种省略,只有依靠知识和推理,才能理解这种省略;而且,即使是句法,也总是有例外的。因此,下面谈到的一些语言理解技术和系统,仅在有限范围内才是有效的。

自然语言有“口头语言”和“书面语言”之分,因此,语言理解问题涉及到两个方面。第一,利用语言的句法、语义以及有关世界的知识,去理解书面语言;第二,利用上述所需的全部信息,再加上语音学的知识,去理解有声语言。通常所说的自然语言理解,则是指理解书面语言而言的。

二、句法分析

理解一个句子,是把词组合起来形成表达一个句子意义的结构。为了做到这一点,必须利用各种信息,包括利用句法信息来对句子进行分析。

对句子进行句法分析,就是把构成句子的一串词转换成一种表示这些词怎么联系起来的结构。如果某些词的序列违反了句法规则,那么,它就会被认为是不能接受的。假如

说有这样一个句子:

Boy the go the to store.

一个英语的句法分析程序,将发现这是错误的。

为了分析一个句子,就需要有一种文法,用来描述特定语言的结构。有了这样的文法,分析程序就可以给每一个符合文法的句子指派一种结构,同时拒绝那些不合文法的句子。

(一) 转换生成文法

这是乔姆斯基(Chomsky, N.)提出的用来分析和说明英语句子的文法。^{[1][2]}它包括短语结构文法(Phrase structure grammar)和转换文法(Transformation grammar)两部分。前者表示句子各成分在层次上的关系,说明句子是由一系列重写规则生成的。图1是短语结构文法的片段。图2是用该文法片段分析一个句子所产生的结构,称之为“分析树”。

S → NP + VP	Art → the/a
NP → Art + NP ₁	N → boy/girl/.....
NP → NP ₁	ADJ → little/big/.....
NP ₁ → ADJs + N	V → hit/tran/.....
ADJs → ADJ + ADJs/E	
VP → V + NP	
VP → V	

图1 短语结构文法片段。其中, S = 句子, NP = 名词短语, E = 可缺省, N = 名词, ADJ = 形容词, Art = 冠词。

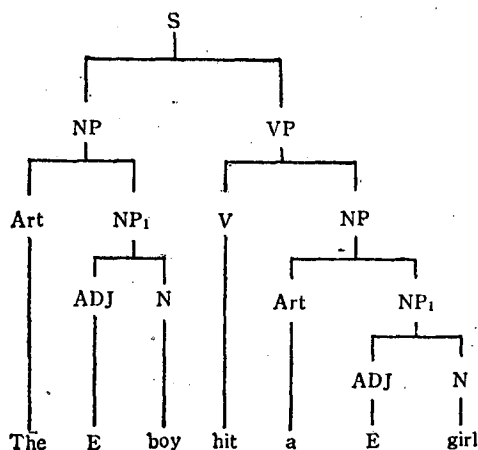


图2 句子“The boy hit a girl”的分析树。

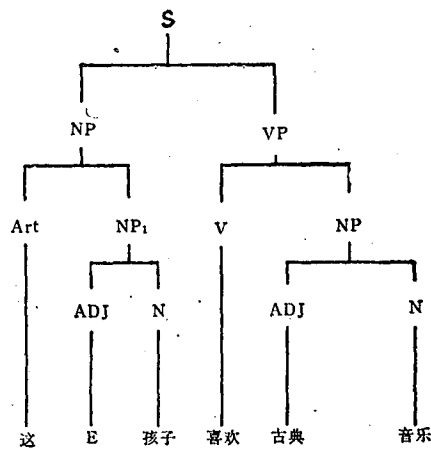


图3 该文法片段用于分析句子“这孩子喜欢古典音乐”。

汉语句子也是由各种短语构成的,它的基本句式的结构与英语是相似的。因此,该文法片段也可用于分析类似的汉语句子,如图3所示。

短语结构文法并没有考虑主动语态和被动语态的问题;也没有提供一种手段,去识别那些具有不同的表层结构但其意义是类似的句子。这表明需要在更高的层次上去分析。对此,乔姆斯基提出了一集新的规则即转换文法,用以解释这种差别和类似性。

由短语结构文法生成的是一种简单的、主动的、陈述的肯定句,称为核心句。转换规则

可对其中某些成份进行重排或替代,增加或删除某些元素。因而,通过转换文法的使用,可产生被动、疑问、否定、乃至复杂句。图4用被动语态的转换作为例子,说明它是怎么工作

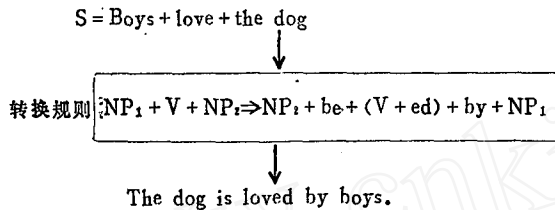
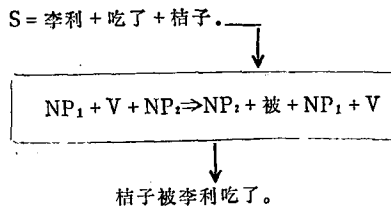


图4 被动语态的转换

的。句子 S 输入给转换规则,经过转换后,生成新的词串,再让助动词与新的主语相一致,从而就生成了句子 S 的被动语态。乔姆斯基的转换文法,对处理汉语来说,在一定程度上也是可用的,只是规则的形式需改变。例如:



乔姆斯基的理论是一个语言学家的途径,而不是心理学的途径。但心理学家对该理论是否有心理学基础的问题,却有兴趣。例如,有人采用“联合对偶学习”的方法作过探讨,结果表明,人对句子的知觉和记忆,受句子短语结构的影响。也就是说,短语结构文法,反映了人们如何组织句子以便于知觉和记忆的特点。另一些实验表明,人对句子的加工,也受“转换次数”的影响。虽然不能肯定地说,人们在实际加工句子时,是象转换文法那样工作的,但是,转换的次数相应于某些需要时间的心理操作,这一点是有人证实了的。

(二) 扩充转移网络

扩充转移网络(Augmented transition networks,简称ATN)是由伍兹(Woods, W. A.)⁽³⁾提出来的。它由两个部分即状态和弧所组成。弧上附有各种标记,指明分析过程的工作性质。图5是用来分析英语的ATN网络片段。如果用它来分析句子“The heavy rain has stoped”,那么,其整个工作过程包括许多步骤。首先,从S开始,因为弧上标明为NP,所以下推到NP子网络,接着,作一个范畴检验,看看第一个词the是否为冠词。检验结果为正,把the放入寄存器,并走向状态6。由状态6伸出的弧指明,要检验下一个词是否是形容词。该检验也为正,并把heavy存放到形容词寄存器。这时,过程仍处在状态6,因此,再作检验,看下一个词是否仍为

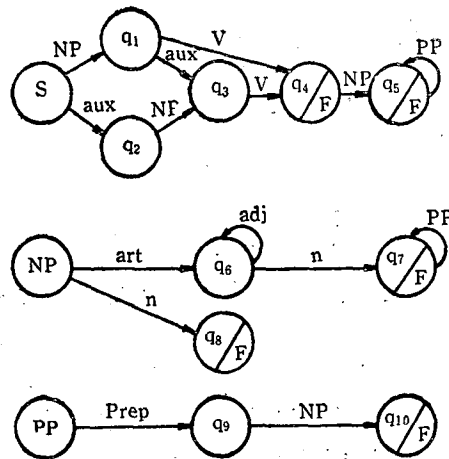


图5 ATN网络片段
q1 = 状态 1, ……; F = 在此状态可结束对句子的分析

形容词。此次范畴检验为负,因此,改向由状态6伸出的另一条弧,并检验下一个词是否名词。名词短语检验完后,根据状态7上弧的指引,分析过程将下推到PP子网络,并检验后面是否有介词。该检验为负。所以,系统把名词短语存入SUBJ寄存器,并返回到主网络的状态1,并进而检验后面是动词还是助动词……。

ATN有其自身的优点。首先,一个子网络尽管在几个地方要使用,但它只需出现一次。例如, NP网络既出现在顶层网络的开始(相当于主语),也出现于末尾(宾语的位置);同时,一个网络可以递归地调用。NP网络可以调用PP网络,而PP网络反过来再可调用NP网络。这对处理有多个介词短语连在一起的句子来说是必要的;其次, ATN使用了任意数量的寄存器,用来存放分析的结果,因此,每分析一步都有相应的记录。而且,这些寄存器的内容,在分析过程中可以改变或彼此交换。比如,若网络用来分析被动语态,当机器发现了被动语态时,那么, SUBJ寄存器的内容可立即转移到OBJ寄存器,而介词后面的宾语可放在SUBJ寄存器。这样,下面两句:

John hit Mary. 和 Mary was hit by John. 将得到同样的解释。

最后,在设计ATN网络时,可以把任何检验(包括语义的)放在弧上,也就是说,可以对句子同时进行句法的和语义的分析。例如,句子“他喝咖啡”, ATN在检验该句动词后面是否是一个名词时,也可检验它是否代表一种“饮料”。因此, ATN为我们提供了把各种知识插入分析系统的方法。事实上, ATN本身并不包含任何一种文法,它只是提供一种机制。人们可以利用这种机制,定义和使用某种文法。这是ATN的一个重要特征。

当然,在使用ATN分析句子时,常需作大量的追踪,因此,运行起来可能消耗较大。尽管如此,它仍是当今很成功的一种分析策略。ATN最早用于LUNAR系统。这是一个关于月球岩石和土壤样本化学分析数据的提取系统^[4]。此后,在许多语言理解系统中都应用了这种分析机制,包括我国范继淹等人研制的关于文学知识的汉语人机对话系统^[6]。

三、语义分析

对一个句子作句法分析,只是理解句子的第一步。有时候,还必须对句子作语义解释。要做到这一点,有两种办法:(1)生成一个完整的句法解释,然后把此结构传送给一个独立的语义解释程序。但这里有一个问题,若不考虑语义信息,通常不可能确定正确的句法解释;(2)把句法的和语义的信息结合起来进行分析。ATN可作为解决这一问题的一种方式。此外,科学家们还提出了各种不同的方法。

(一)语义文法

语义文法(Semantic grammar)最早是由伯顿(Burton, R. R.)^[6]使用的,他为智能计算机辅助教学系统SOPHIE配制了一个自然语言接口。接着,亨德里克斯(Hendrix, G. G.)^[7]写了更大的语义文法,为美国海军大型数据库管理系统LADDER提供一个自然语言接口LIFER。在语义文法中,句法知识和语义知识结合成一集文法形式的产生式规则。这些规则的选择,既受句法的支配,也受语义功能的支配。图6是LIFER使用的一段简化了的语义文法规则。

可以看到,在语义文法中,没有使用句法范畴(如NP或VP),而是使用了SHIP和COUNTRYS等语义范畴。语义文法的规则,实际上是一种可以用来与输入句进行匹配

S→What is (the) SHIP-PROP of SHIP?
 SHIP-PROP→speed | length | draft | beam | type
 SHIP→SHIP-NAME | the fastest SHIP 2 | SHIP 2 |
 SHIP-NAME→Kennedy | Fox |
 SHIP 2→COUNTRYS SHIP 3 | SHIP 3
 SHIP 3→SHIPTYPE LOC | SHIPTYPE
 SHIPTYPE→carrier | submarine |
 COUNTRYS→American | Russian |
 LOC→ in the Pacific |

图6 语义文法片断

的句式。例如,What is the beam of Fox? 和What is type of the fastest American carrier in the Pacific? 等句子作为输入时,与上述那段语法均可获得匹配。这种语法的主要优点是:(1)当分析完成,表明输入句与某一句式匹配成功时,其结果就可以直接使用,不需要额外的加工;(2)如果某些解释在语义上是没有意义的,那么,语义文法不会把它生成出来。因此,句子的一些歧义可以避免;(3)一些句法上有问题,但不影响语义的句子,语义文法也能生成正确的解释。当然,语义文法没有利用句法的概括作用,因而有时需要较大数量的规则。结果是,机器对句子的分析过程可能是昂贵的。

语义文法一般用于某一特定任务范围的自然语言接口。但是,只要对语义范畴及其定义作适当的改变,就可用于不同的任务范围;例如,把上述语法片段中的语义范畴“船”改成“飞机”,把“飞机的特性”定义为“速度”、“客容量”等,通过这些相应的修改,就可使它适用于民航管理系统的自然语言接口。由于语义文法强调的是句子各成分的语义功能,而不是句法功能,因此,它适用于研制汉语的自然语言理解系统。事实上,国内已有一些自然语言接口,是用语义文法来写的。^[8,9]

(二) 格文法

费尔莫(Fillmore, C.)提出的格文法(Case grammar),^[10,11]是另一种把句法和语义结合起来的分析方法。格文法规则所产生的句子的结构,反映了句子各成分之间的语义关系。假如有这样两个句子:

(1)李利打了王英。(2)王英被李利打了。在这两句中,李利和王英的句法角色是不同的。它们分别在一个句子中是主语,而在另一个句子中却是宾语。但是,在这两个句子中,它们的语义角色是不变的。李利是动作的执行者,而王英则是动作的承受者。格文法对这两个句子的表达是相同的:

(打 (施事 李利)
(与格 王英))

从格文法的角度来看,一个句子是以动词为中心的,其它的词都与动词处于某种关系中,这种关系称之为“格”关系。它表示其它词与动词之间的语义关系。上述两句中,李利是“打”这个动作的执行者,称为“施事格”;王英是该动作的受者,称为“与格”。在某些语言(如俄语和德语)中,也有“格”这个语法概念,但这只是指句子表层的关系。格文法与这种表层格的概念是不同的,它反映了句子深层的语义格关系。

深层的语义格究竟有多少,并无一致的看法。一般而言,有以下十种:

1. 施事格(Agent)——动作的执行者(生物)。
2. 工具格(Instrument)——引起某事件时用到的东西。

3. 与格(Dative)——受动作影响的实体(生物)。
4. 使成格(Factitive)——事件的结果或产物。
5. 处所格(Locative)——事件发生的地点。
6. 源点(Source)——某种东西开始移动的地方。
7. 终点(Goal)——某种东西移向的地方。
8. 受益格(Beneficiary)——事件的受益者(生物)。
9. 时间(Time)——事件发生的时间。
10. 客体格(Object)——被作用或改变的对象(非生物)。

格文法认为,任何特定动词所要求的格,形成一个有序集,称为“格框架”。例如,动词“开”至少要求有一个“客体格”(即O格),此外,也可能有“工具格”(即I格)和“施事格”(即A格)。所以,它的格框架是:

动词“开”的格框架: [……O (I) (A)]

例句: 门开了。[O]

李利开了门。[O A]

李利用钥匙开了门。[O I A]

每个动词都有相应的格框架。这样,机器在理解句子时,一旦动词确定下来,就可以根据该动词提供的格框架期望句中将出现什么成分。因此,使用格文法的语言理解系统,是“期望驱动”的系统。格文法反映了人类加工句子的特点。人类在理解和记忆前述的两个句子时,总是力图发现和记住“谁打谁”这一句子的深层意义,而其表层形式往往是被忽略的。人类对句子的加工是“期望驱动”的。因此,格文法可以说是语言理解的心理学模型。

格文法强调的是句子各成分之间的语义格关系,较少受不同种语言特点的制约,因此是可用于处理汉语的。至于格文法如何用于处理汉语的问题,有些学者已在研究。^[12]

(三) 概念依赖

概念依赖(Conceptual dependency)是由山克(Schank, R. C.)首先提出的一种表征自然语言句子意义的理论,简称CD理论。^{[13], [14]}他认为,人的记忆中存储着各种概念内容。概念内容有其完善的结构,是由概念及其相互之间的依赖关系构成的。人理解语言的过程,就是把语句编码成某种相应的概念内容。因此,要使计算机模拟语言理解过程,就必须研究概念内容的性质,以及由语句映射(即转换)为概念内容的规则。

CD理论与通常的语义网络不同。后者只提供了一种表征意义的层次结构,而前者既提供了一种表征意义的结构,也提供了一集专用的原始概念。图7给出了句子的CD表达的一个简单例子。在图中,箭头表示依赖方向;双箭头表示动作者与动作之间的双向关

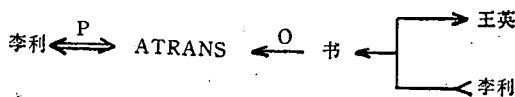


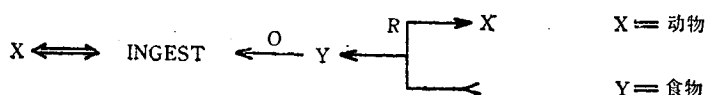
图7 句子“李利给了王英一本书”的CD表达

系;P表示过去时态;O表示动作的对象;R表示接受的关系;ATRANS是CD理论专门使用的一个原始概念,它表示所有权的转移。

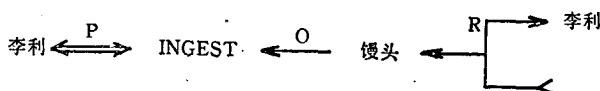
CD理论提出了十一个原始概念。除了上述ATRANS外,还有PTRANS(表示客体

物理位置的转移)、PROPEL(物理力作用于客体,如“推”)、INGEST(动物摄入客体,如“吃”)、EXPEL(排出某种东西,如“叫”)、MTRANS(心理信息的转移,如“告诉”)、MBUILD(从原有信息中构造出新信息,如“作决定”)、以及SPEAK、MOVE、GRASP、ATTEND等。

把一个句子分析成概念依赖的表达,与使用格文法的分析过程相类似,也是由“期望”来驱动的。这些期望是根据句子的主要成分来设置的。但是,在CD理论中,对动词是在较低的层次上作表征的。所以,它有更大的概括程度及更高的预测能力。机器在把一个句子映射为它的CD表达时,分三个步骤。首先,由句法处理程序抽取句子的主要动词和名词,并确定动词的句法范畴。接着,由概念处理程序利用一种“动词—动作词典”来进一步加工。该词典中有各种词条,描述一个动词可能出现在其中的概念结构。例如,动词“吃”的词条内容是:



假如输入的句子是“李利吃了馒头”,由句法处理程序的指引,在词典中找出这一词条,概念处理程序就把句中各个成分插入概念结构的相应的空槽中。这样,就可得到该句的CD表达:



CD理论在较高的语义水平上提供了分析句子的手段。句法只起指引的作用,无须对句子作详细的分析。山克一再强调,他的理论反映了人类理解语言的过程,因此,是语言理解的一个心理学模型。

除了上述三种语义分析方法以外,语义网络也是语言理解中值得采用的一种重要手段。语义网络理论是一种心理学模型。它认为人脑中的知识,就是一些概念及概念间的联系。它们错综复杂地相互交织在一起,形成了网一样的结构。人类理解语言的过程,就是对网络的搜索、比较过程。我国第一个汉语理解系统,就是采用语义网络理论建立的。^[16]

四、理解篇章

前面谈的只涉及理解单个句子的一些方法和途径。为了理解构成某种情节的一组句子,或一段课文,就不单要求把某种结构指派给一个句子,而且还要求发现句子之间的联系。要识别句子之间的种种联系,就需要大量的关于世界的知识。因此,对于理解程序来说,其成功的关键,乃是组织和利用知识的方式。在这一方面,也已提出了一些方法。例如,利用焦点(focus)和分区语义网络(Partitioned semantic net)相结合的方法来理解物理客体事件的情节,^[16]利用目标结构(目标和计划相结合)的方法来理解人及其活动的故事;以及在理解中采用“脚本”(Script)的办法^[17]限于篇幅,下面仅谈一下利用脚本来理解一段故事。

在人的生活经验里,一些事件经常是有序地联在一起发生的。它们之间的联系似乎已成了一种模式。脚本就是对经常在一起出现的一系列事件的固定的记忆结构,它记录了经验的典型实例,如去餐馆吃饭,到剧院看戏。脚本所包含的特定情境下的信息较丰富,所以能进行较大幅度的推理。故事理解程序SAM(Script Applier Mechanism)的功能证明,在自然语言理解中利用脚本这种知识表达形式来进行推理,是方便和有效的。^[18]图8给出了SAM系统使用的一个“餐馆脚本”所包含的信息。

利用脚本达到理解的过程分三个步骤。首先,调用一个概念依赖分析程序,把故事中的句子映射为CD表达;然后,从记忆中选取合适的脚本;最后,利用脚本推出故事中没有说明的事件。假设输入下面一段课文:

“李利走进一个餐馆。他点了一份牛排。他付了钱,并离开了餐馆。”

人物:顾客、服务员……。

道具:餐馆、餐桌、菜单、食物、饭钱……。

事件:

- | | |
|-------------|--------------|
| 1. 顾客走进餐馆; | 2. 顾客在餐桌前坐下; |
| 3. 服务员拿来菜单; | 4. 顾客点食物; |
| 5. 服务员端来食物; | 6. 顾客吃食物; |
| 7. 顾客付饭钱; | 8. 顾客离开餐馆。 |

图8 SAM系统内部的“餐馆脚本”(略有删节)。

然后问SAM系统:李利吃了什么? 故事中并没有说明这一点,但SAM系统能够回答这个问题。它根据脚本中的事件1,认出了“李利走进一个餐馆”这句话,同时,李利被指定为顾客的角色;根据事件4,它认出了“他点了一份牛排”这句话,同时把牛排指定为食物的角色,等等。这样,脚本中的顾客就是李利,而食物就是牛排。所以,SAM就可以把脚本中的事件6具体化为“李利吃了牛排”,并以此作为上述问题的回答。显然,这一回答是合理的。

从这里可以看到,第一,SAM所经历的推理机制与人在思考类似问题时的推理机制是相似的;第二,理解一段有情节的故事,不但需要具有关于语言的知识,而且还需要大量有关世界的知识。

五、结束语

上面叙述了语言理解的一些主要方法和途径。从这里我们可以看出,虽然最初它们被提出是为了处理英语的,但对于汉语理解来说,在一定程度上它们也是可借鉴和适用的。当然,汉语有它自身的特点:词的形态变化不发达,主要依靠词序来表示其结构关系,词序有时又是比较自由的,常见省略与简缩形式等等。因此,汉语的理解更依赖于语义和语用情境;知识和推理在汉语理解中起着更重要的作用。因此,把这些方法和途径用于理解汉语时,必然要结合汉语的特点,参照我国人民的思维习惯。同时也可以看出,语言理解的研究不单是一个语言学的问题,而更主要的是认知心理学应该加以研究的一个问题。语言理解是人类最重要的认知过程。那些著名的语言理解的研究者,常常希望或强调自己的语言理解模型是符合人的语言理解过程的,是一个心理学模型。显然,这是因为人是最巧妙和最有效的语言使用者和理解者。解过程,才能建造出有效的理解自然语言的计算机系统。

参 考 文 献

- [1] Chomsky, N., *Syntactic structure*. Mouton, The Hague, 1957.
- [2] Chomsky, N., *Aspects of the theory of syntax*. Cambridge Mass. MIT Press.
- [3] Woods, W. A., Transition network grammar for natural language analysis. *Communications of the ACM*, 1970, Vol. 13, P. 591—606.
- [4] Woods, W. A., Progress in natural language understanding: An application to lunar geology. in *Proc. AFIPS conference 1973*, 42.
- [5] 范继淹、徐志敏: RJD-80型汉语人机对话系统的语法分析。《中国语文》1982年,第3期。
- [6] Burton, R.R., *Semantic grammar: An engineering technique for constructing natural language understanding system*. Technical Report 3453, Bolt Beranek and Newman, 1976.
- [7] Hendrix, G. G, et al, Developing a natural language interface to complex data. *ACM Transactions on Database System*, Vol. 3, 1978.
- [8] 陈群秀、黄昌宁等: 计算机的汉语通用接口。中文信息处理国际会议论文集, 1987, 第2卷, P. 104—109.
- [9] 吴智君、符启缨: 在数据库上建立实用的自然语言接口, 华中工学院学报, 1986年, 第2期, P.167—172.
- [10] Fillmore, C., The case for case. in *Universals in Linguistic Theory*, 1968, 译文见《语言学译丛》第二辑, 社会科学出版社, 1980年。
- [11] 杨成凯: Fillmore的格文法理论。《国外语言学》, 1986年, 1、2、3期。
- [12] 鲁川、梁镇韩: 计算机对汉语的理解和生成。自然语言理解学会第二次学术会议论文, 1986年。
- [13] Schank, R. C., Identification of conceptualizations underlying natural language. in *Computer Models of Thought and Language*, R. C. Schank and K. M. Collby (Eds.), 1973.
- [14] Schank, R. C., *Conceptual Gnformation Processing*. 1975.
- [15] 李家治、郭荣江、陈永明: 机器理解汉语—实验I。《心理学报》, 1982年, 第1期, P. 29—40.
- [16] Grosz, B. J., The representation and use of focus in a system for understanding dialogs. in *IJCAI 5*, 1977.
- [17] Schank, R. C. and Abelson, R. P., *Scripts, Plans, Goals and Understanding*, 1977.
- [18] Cullingford, R., SAM. in *Inside Computer Understanding*, edited by R. C. Schank and C. K. Riesbeck, 1981.

THE APPROACHES OF NATURAL LANGUAGE UNDERSTANDING

Chen Yong-ming

Institute of psychology, Academic Sinica

Abstract

This paper has described and discussed some main approaches of computer understanding natural language in relation with Chinese language processing. The author has pointed out that: (1) although these approaches are proposed for processing English, they are also suitable for processing Chinese language to certain degree; Of course, Chinese language has its own feature. Understanding Chinese depends more on semantic and pragmatic knowledge (2) Natural language understanding not only is a problem of linguistics but mainly is a problem that cognitive psychology should study. There is a close relation between cognitive psychology and computer understanding natural language. The actual process of understanding natural language by human must be explored in more detail in order to develop more powerful computer system of understanding natural language.