

PROCEEDINGS OF 1983 INTERNATIONAL CONFERENCE ON CHINESE INFORMATION PROCESSING

中文信息处理国际研讨会 论文集

(2)



中国中文信息研究会和联合国教科文组织联合举办

1983年10月12日至14日于北京

CO-SPONSORED BY CIPSC & UNESCO
BEIJING, PEOPLE'S REPUBLIC OF CHINA
OCTOBER 12-14, 1983



数据加载失败，请稍后重试！

第二集

目 录

汉语语言的理解系统	王锡龙	(1)
结合上下文辅助分词的学习系统	管纪文 谷新英	(11)
汉字字形编码的原理和实践——WBZX汉字编码方案	汉字基本构件	
实用频度表	王永民	(19)
中文计算机系统开发的研究和实践	张寿萱	(32)
微型机汉字输出和汉字报表语言	郑德高	(43)
TRS-80(I)B型文字处理系统	金虎范等	(50)
计算机系统扩充汉字处理功能的通用途径初探	郑秀华	(56)
多用户汉字信息处理的人机对话	朱永平	(63)
高分辨率汉字字形在计算机中的压缩表示及字形点阵的 复原设备	王选 吕之敏 汤玉海 向阳	(71)
ZN9号汉字26键输入码方案	郑易里	(84)
汉语的递归正则模糊语言模型	孙怀民 赵沁平	(94)
关于汉字输入编码评测规则的建议	盛焕烨	(105)
蒙文信息互换标准代码的设计方案	嘎日迪	(108)
海曼公式的一般分析	陈文熙	(118)
指挥自动化与通信拼音字	杜牧平	(125)
汉字结构识别	唐保兴 宋国新	(129)
机器理解汉语的一些探索和设想	李家治 陈永明	(138)
汉语句子的多标记多叉树形图分析法	冯志伟	(144)
“中顿句”和“无顿句”的频度分析	肖申生	(159)
JFY-IV型机器翻译系统设计纲要	刘 健	(164)
计算机辅助科技英语词汇统计结果的初步分析	黄人杰 杨惠中	(174)
汉字信息和信息字典	李金铠	(183)
语言、智能与计算机辅助语言智能学习——兼谈计算 机辅助中文学习系统 SINOCCMS	钱 锋	(194)
廉价的磁笔触式汉字整字输入键盘	龚滨良	(216)

汉语语言的理解系统

吉林大学计算机科学系 王锡龙

摘要：1. 从汉语语句中词的基本要素：词形(X)、词音(Y)、词义(I)和词性(N)出发，给出了理解语言的其他三个信息：词域(U)、前位属性(Q)、和后位属性(H)。从而构成了理解语言的全部信息量(记为W)：

$$W = \{X, Y, N, U, Q, H\}.$$

于是，语句的理解式 L_x 可表示为：

$$L_x = \{V \mid \varphi(W_1, W_2, \dots, W_n)\}.$$

其中，V为语句x中使得条件 $\varphi(W_1, W_2, \dots, W_n)$ 成立的所有词V的集合。

2. 进一步给出了分析语句结构关系(语法和语义)的树型理解图。并给出构造理解图的算法。

3. 最后，给出了限定型语言理解系统的结构框图，并举例说明系统的功能。

由于计算机的广泛应用，迫切需要解决用汉字使用计算机和用计算机处理汉语语言的问题。本着这个目的，我大胆地提出一些看法并给出一些结果，借此与语言学家和计算机科学事业的同行们进行探讨，以便促进我国的语言自动化和计算机的普及与应用。

一、语言的表示和信息量

任何一种语言(自然语言和计算机语言)都是定义在某一个符号集上的无穷多个有限的符号行(串)集。若满足某种语言规范条件的符号行(称为语句)时，则语言就是定义在符号集上的无穷多个语句集合。一般地，可用如下形式表示^[1]：

$$L = \{x \mid \varphi(x)\}.$$

其中，x是代表“使得条件 $\varphi(x)$ 成立的所有语句x的集合”。

事实上，任何有穷多个符号集上的语句x的集合，都是可数无穷多集。语言只是这个集上的子集。而且，并非是可数无穷集^[2]。由于语言的复杂性(语法的和语义的)，我们只能给出某些有穷表示，用来辅助计算机理解(而不是识别)自然语言。

1. 符号集

汉语书面语言(社会科学的与自然科学的)材料，一般使用如下符号：

〈字〉：指方块字(目前，计算机采用汉字编码输入)。字是形和音结合的最小单位。计算机只能采用有限多个字。如字典字。所有字典中的字记为Z。

〈字符〉：字母(如拼音字母与英文字母)和有限个符号(如星号“*”，求和号“Σ”等)，记为M。

〈标点〉：语言的停顿符号，即汉语标点符号，记为B。

若汉语语言的符号集用字母“V”表示时，则有：

$$V = \{Z, M, B\}.$$

由于语言的最小单位是词，词才具有形、音、义、性四个基本要素^[3]，才是理解语言的最小单位。若所有的词（系指某词典词）用符号“C”表示，则语言的符号集V*，可表示如下：

$$V^* = \{C, M, B\}.$$

2. 语言的信息量

词构成语句，句构成段（章、节），然后组成文章。由于词是语言的最小单位，句是语言的基本单位，所以，我们还是从词、句的信息讨论入手，从而了解语言的信息量。词的本身要素是词形、词音、词义。词的基本要素是词形、词音、词义和词性。以后分别用X, Y, I, N表示。词性是词在语句中，组词搭配的特性分类，所以，词性对于理解语句尤为重要^[8, 9]。

由于汉语缺少形态变化，在许多语法、词法分析上，都出现了“中间状态”，很难截然分开。这就给计算机的处理带来了很大困难。所以，计算机对汉语语言的处理，至目前为止，只能是限定性的和辅助性的^[4, 5, 6]。然而，几千年来实践表明：汉语语言是最灵活、最丰富、表述最简明的语言之一。比如，同样一篇文章，用汉语的书面语言表达，一般来说，总比用其它语言文字表达，要明了简短。所以，与其说它不发达，不如说它更高级，理解起来更具有智能性。请看下例语句中，标记词的基本要素。

例一

①今天大家都到校园里种花。

zhòng huā

②医生给小朋友种花。

zhòng huā

③这是一株种花。

Zhòng huā

④这种花真好看。

Zhǒng huā

⑤“种花还没嫁人吧？”

Chóng huā

种大嫂关心地问。

Zhòng

⑥种花两家都富起来了。

Chóng huā

(Zhòng)

⑦这钱中花。

Zhōng huā

“种花”为栽培花草之义。这是由于语句环境，特别是“校园里”给定的。

这句的“种花”与①中是同形、同音、异义词。“种花”是取义“种牛痘”方言。

这里与①②中的种花已经是同形、异音、异义、异性词，指留籽儿的花。

此句中，“种”已是量词，与这字结合成指示代词。同时，读音也与①②不同，虽与③同音而义不同。请注意！这里“好”还有hǎo与hào之分。

种字的有两个异体字，即繁体字。它们为异形、异音、同义、同性字。此句都为姓。种花是位女人，因为有“嫁人”一词。当然，如不作标记，种花肯定是词典外词，计算机不会识别出来。种chóng与种zhòng音也不易分清。

因为有后边的“两家”一词得知是姓种的和姓花的两家富起来了。而不是种花草的两家富起来了。若“种花”与“两家”中有“的”字，那就取后一种含义了。

这句中，因为与“钱”字搭配，立即读Zhōng huā，而不会读Zhòng huā。“花”字在此句是动词，指耗费。与①～⑥中的花含义都不同。

通过上例①～⑦可见，词的基本要素都具有多值性。然而，给定一个语句后，一般来说，人是可以很快理解的。汉语中的词具有如下性质：

性质1.有词必有义，有义必可用词来表达^[7]。

性质2.词在语句中的基本要素是客观确定的。

性质1说明了词汇的丰富性、能产性和完备性。性质2说明了语言的可理解性。但是，对于计算机来说，就不那么容易判定和理解了。我们必须限定某些条件，给出计算机理解语言的某种算法，让计算机一步一步的去判定，去理解。通过例一可看出，只从词的基本要素理解语句是不充分的。还需具备一些其它必要信息。下面引进一些新的概念。

词域：词所属的科类别，称为词域，记为U。

词域体现了词在语句中的客观环境。比如：数学，物理，化学，天文，地理，医学，美工，经济，……。词典中，应对词加以适当的域标记，以助理解词义、句义。

词位和属性：设语句x为如下形式：

$$\overbrace{\cdots \cdots V_1 V_2 \cdots \cdots}^{\text{x的前句}} \underbrace{V_i}_{\text{x}} \overbrace{V_{n-1} V_n \cdots \cdots}^{\text{x的后句}}$$

其中， $V \in V^*$ ，i称为词 V_i 的位号，词在语句中的位置称为词位。用i表示。

若 V_i 与 V_j ，当 $i < j$ 时，称 V_i 为 V_j 的前位词符； V_j 为 V_i 的后位词符。每个词都具有一定的性质，称为属性。 V_i 的属性称为 V_j 的前位属性；而 V_j 的属性，称为 V_i 的后位属性^[10]。前位属性和后位属性分别用Q和H表示。位属性提供了参照前后词的意义来判断某一词的条件，表现了词的动态关系。

综上所述，词在语句中的信息（记为W）可表示为：

$$W = \{X, Y, I, N, V, Q, H\}$$

共为七个信息量。

3. 语句理解式

现在，我们可以给出语句x的一般理解表达式为：

$$L_x = \{V \mid \Psi(W_1, W_2, \dots, W_n)\}.$$

其中，V是使得条件 $\Psi(W_1, W_2, \dots, W_n)$ 成立的语句x中的全部词符的集合。

同时，我们对一般语言表达式： $L = \{x \mid \Psi(x)\}$ 的理解也就加深了一步。

自然语言是复杂的，不能简单地靠数学运算求解。我们需要提供给计算机处理语言的一个算法，使得计算机按此算法，对任意给定的一个语句x，找出基本符合语句x中全部词符V的信息量W的统一值。这样，就说计算机理解了该语句。并且，可以回答某些需求结果。下面就来讨论这些问题。

二、理解图和算法

语句结构关系的分析，语言学家已经给出了多种分析办法^[11, 12]。我们将根据计算机理解语言的需要和可能，选择和确定一种比较合适的方法。

1. 结构关系的表示

下面，我们通过例句，列出几种分析方法，以便比较，决定取舍。

例句：我们厂最近试制出新产品。

方法一：（我们）厂〔最近〕试制〔出〕新产品。

符号说明：||：主语部分与谓语部分的分界符号；
 |：谓语部分与宾语部分的分界符号；
 =：主语；—：谓语；~w~：宾语；
 ()：定语；〔〕：状语；〈〉：补语^[11]。

符号标记简单，但结构关系不清晰。

方法二（见图1）：

树型结构，层次清楚。但是，没有结构关系。方法一与方法二都是语句的语法成分分析，属层次结构型。

方法三（见图2）：

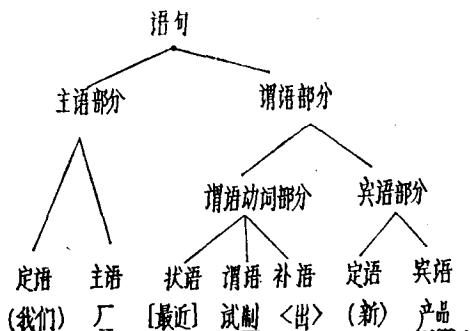


图 1

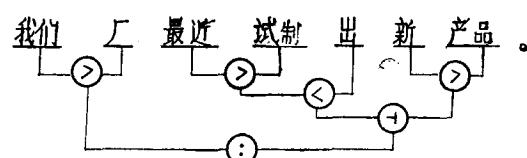


图 2

符号说明：>表示修饰成分在前；<表示修饰成分在后；—符号表示补足成分在后；+符号表示补足成分在前；：是表述关系（即主谓关系）。另外，还有：+表示并列关系等^[12]。

这种方法层次清楚，强调结构关系，即语义关系。

方法四：

<语句> ::= <主语部分> <谓语部分>
 <主语部分> ::= <定语> <主语>
 <谓语部分> ::= <谓语动词部分> <宾语部分>
 <谓语动词部分> ::= <状语> <谓语> <补语>
 <宾语部分> ::= <定语> <宾语>
 <主语> ::= 厂
 <谓语> ::= 试制
 <宾语> ::= 产品
 <定语> ::= 我们
 <定语> ::= 新
 <状语> ::= 最近
 <补语> ::= 出

这种表示为 BNF 表示法（即 Backus-Normal Form 或 Backus-Naur Form 的缩写符）。::= 表示“定义为”。这是计算机语言中常用的一种递归定义。显然，对自然语言是不合适的。比如，按上述定义就会有“新厂最近试制出我们产品”的二义性语句。可以说“我

们的产品”，“我们厂”，而不能说“我们产品”。

本系统将基本上采用方法三。重点放在结构关系（即语义理解）上，兼顾语法成份。对于该例句将产生如下的树型结构形式，称为语句k的理解图。

语句 k： 我们 1 厂 2 最近 3 试制 4 出 5 新 6 产品 7 。 8

理解图见图3。

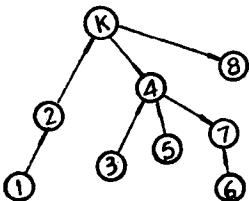


图3 语句k的理解图

其中，k为语句序号。(k)为根节点。根节点至多有一条接收线。体现了语句若有主语，则只有一个，否则为无主句。根节点的发射线至少有一个。根节点的直接子女中，主语节点只有一个向根节点的发射线。可以有多个接收线，做为他的附加成分。谓语节点发射线中，若为一个，则是一个宾语节点；若为两个，则为双宾语结构；它可以有多个接收线，做为附加成分节点。当任一子女节点中，接收线节点词的词序号，大于该节点的词序号时，如⑥→④，则为补足成分，其余皆为修饰、限定成份。其余根节点的直接子女，有且只有一个接收线。宾语节点类似谓语节点，附加成分节点的发射线和接收线各为一条。

语句k的理解图按着自顶向下，自左至右的顺序读。从主语节点词开始，读取语句的主要成分。然后，按照语句的词序号读取附加成分。

理解图的每个节点词，由计算机按着需求标记信息量的值。所以，它是一个“果树”型。

2. 基本规则

理解图是计算机遵循某些基本规则，执行理解算法产生的。下面给出一部分规则，以便理解算法。

规则1 语句中，几个相邻词，若具有同节点性质时，则称为同节点词。同节点词视为一个词结构。

例如，以下相邻词具有同节点词性质。

- (1) x < 方位词 >，其中，x为一词符，方位词结构。
- (2) x < 助词 >，其中助词中还包括助动词和语气助词。
- (3) x < 量词 >，量词结构。
- (4) xx ，同形重叠词结构，
- (5) < 介词 > x，介词短语。
- (6) x (词缀)，() 表示可为词头，词尾，词嵌。
- (7) < 否定词 > < 动 >。
- (8) < 肯定词 > < 动 >。

规则2 语句中，凡起到修饰、限定、补足成分的词，称为附加成分。附加词，按照附加结构式标记。

例如，以下结构式为附加结构式：

< 形 > → < 名 >，< 副 > → < 动 >，< 形₁ > → < 形₂ >，< 量₁ > → < 量₂ > 等。

规则3 凡对理解图根节点，具有直接子女条件的符号，称为语句的直接子女。

例如：< 叹词 >，< 连词 >，< 标点 > 等。< 名词 >，< 动 >，< 代词 > 具有直接子女性

词。由于代词的指代作用，一般要特殊处理。对人称代词，一般先按直接子女词对待，行再修剪。

规则4 对树型理解图，实施剪枝、嫁接处理的规则，称为剪嫁规则。剪嫁规则分为语法剪嫁，语义剪嫁和综合剪嫁。

(1) “语法剪嫁”例

格式： $\{\text{名}_1\}_{\text{代}_1} \leftarrow \text{动} \rightarrow \{\text{名}_2\}_{\text{代}_2}$

例二 今天 学校 开 大 会。
1 2 3 4 5。 (参看图4)

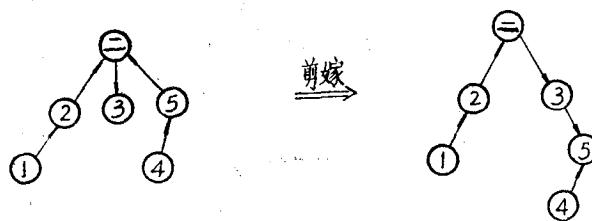


图 4

格式： $\{\text{名}_1\}_{\text{代}_1} \leftarrow \text{动} \rightarrow \{\text{名}_2\}_{\text{代}_2} \{\text{名}_3\}_{\text{代}_3}$

例三 我 问 你 一 个 问 题。
1 2 3 4 5。 (参看图5)

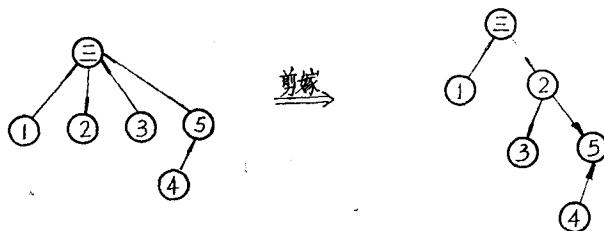


图 5

格式： $\{\text{名}_1\}_{\text{代}_1} \leftarrow \text{动}_1 \rightarrow \{\text{名}_2\}_{\text{代}_2} \leftarrow \text{动}_2$

例四 孩子们 听 了 故 事 很 高 兴。
1 2 3 4 5。 (参看图6)

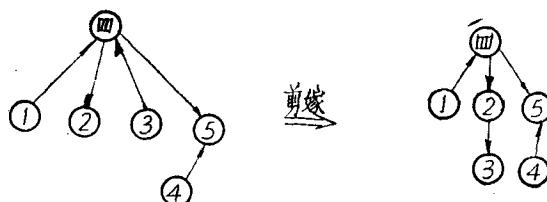


图 6

(2) “语义剪嫁”例

格式: $\{\text{名}_1\} \{\text{名}_2\} < \text{动} >$
 $\{\text{代}_1\} \{\text{代}_2\}$

例五 他 一口 水 都 不喝。 (参看图7)



图 7

格式: $\{\text{名}_1\} < \text{动}_1 > \{\text{名}_2\} < \text{动}_2 > \{\text{名}_3\}$
 $\{\text{代}_1\} < \text{动}_1 > \{\text{代}_2\} < \text{动}_2 > \{\text{代}_3\}$

例六 红红 下 河 洗 衣裳。 (参看图8)



图 8

例七 大家 选 他 当 代表。 (参看图9)

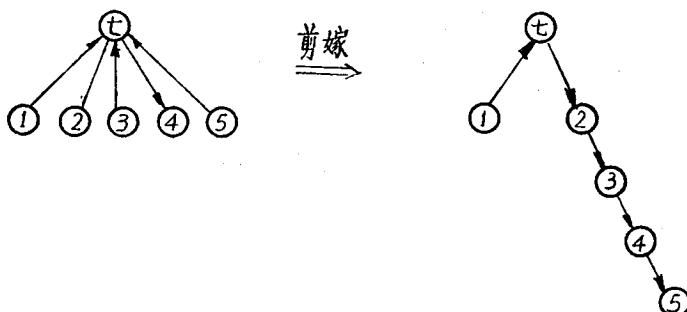


图 9

(3) “综合剪嫁”例

格式: $\{\text{名}_1\} (\text{把字式}) \{\text{名}_2\} < \text{动} >$
 $\{\text{代}_1\}$

例八 你 把房间 收拾 干净。 (参看图10)



图 10

格式: $\{ \text{名}_1 \} \underset{\text{代}_1}{\text{ }} (\text{被动式}) \{ \text{名}_2 \} \underset{\text{代}_2}{\text{ }} < \text{动} >$

例九 敌人 被我们 消灭了。 (参看图11)



图 11

3. 理解算法

算法:

第一步: 取一语句 x_k (k 为语句号)。若语句全部处理完, 则转出口;

第二步: 借助词典和知识库进行分词、定词性、编序号 $[^{13}, ^{14}]$;

第三步: 建立根节点 k 。依次为名词、代词建立接收指示字; 为动词和其它直接子女词符建立发射指示字。

第四步: 归并同节点词和处理代词;

第五步: 按照附加结构式, 建立下一级子女节点;

第六步: 根据剪嫁规则, 实施剪嫁调整;

第七步: 依据需求进行赋值。若取不到值, 则转出等待;

第八步: 问输出结果吗? 是, 则输出; 否则转第一步。

三、系统框图及应用

1. 系统结构

见语言理解的系统结构框图。

共分四部分:

①系统管理部分

控制语言理解系统的管理程序和维护程序, 它在操作系统 和 数据库管理系统 支持下
[15, 16, 17]。

②预处理部分

- 语言文本输入程序；
- 分词处理程序；
- 查阅词典和知识库程序；
- 建立静态库程序等。

所谓静态库，是根据需求，将文本的有关信息量，由词典和知识库中取出，并以某种结构形式存入数据库中。以防止多次扫描词典和知识库，用以提高处理速度。

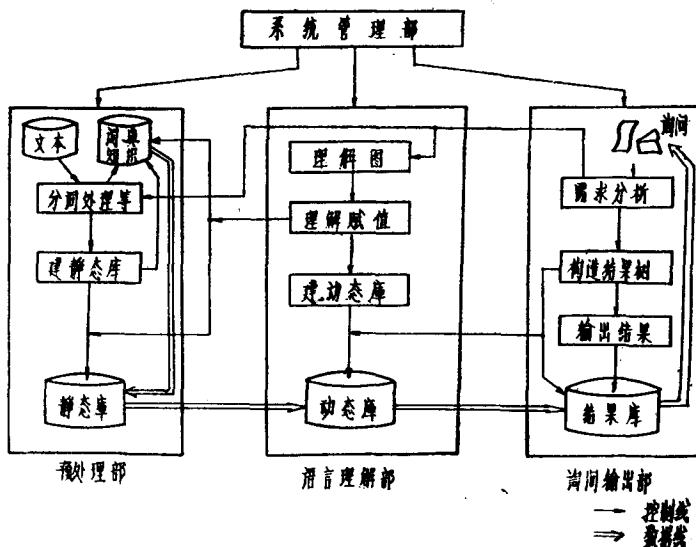


图 12 语言理解的系统结构框图

③语言理解部分

- 语言理解分析程序；
- 建立理解图；
- 赋值程序等。

使赋值后的理解图存于动态库中，以便提供给输出。

④询问输出部分

- 需求分析控制程序；
- 按照需求，从动态库中将结果以输出格式存于结果库中，提供给输出。
- 输出管理程序等。

2. 应用举例

例十 求语句“小毛板着脸不说话”语句中的词音 L(Y) 和第二个词的词性 L(N₂)。计算机对该语句实施处理后，输出格式如下：

L(Y): xiǎo máo bǎn · zhe liǎn bù shuō huà

L(N₂): 板着——动词。

本系统的实现，将对计量语言学、机器翻译、人机对话和计算机辅助教学等有着一定意义。由于篇幅所限，其它应用举例和提问语言省略。

本文曾得到管纪文、孙怀民、张普和谷新英同志的鼓励与指导，在此表示感谢。

参 考 文 献

- [1] John E. Hopcroft Jeffrey D. Ullman, FORMAL LANGUAGES AND THEIR RELATION TO AUTOMATA, Addison-Wesley Publishing Company 1969.
- [2] David Gries, COMPILER CONSTRUCTION FOR DIGITAL COMPUTERS, John Wiley & Sons, Inc.
- [3] “实用词汇新知识”，谭全编著，商务印书馆，1978年3月初版。
- [4] Schank R. C. & Colby K. M. ed, Computer Models of Thought and Language, W. H. Freeman Co. 1973.
- [5] Arens, Y., Using Language and Context in the Analysis of Text, IJCAI 7(1981), P. 25.
- [6] Wilensky, R., A Knowledge-based Approach to Language Processing, A Progress Report, IJCAI 7(1981), p. 52.
- [7] “汉语词义简析”，朱星，湖北人民出版社。
- [8] “拼词法和构形法”，张寿康，湖北人民出版社。
- [9] “词汇知识”，高文达、王立廷编写，山东人民出版社。
- [10] “用属性文法进行编译”，P. Rechenbrg, 浙江大学周炳生译自“Prinzipien des Übersetzerbaus”。
- [11] 高等学校文科教材“现代汉语”，黄伯荣、廖序东主编，甘肃人民出版社。
- [12] “汉语语法分析问题”，吕叔湘著，商务印书馆，1979年6月第一版。
- [13] “结合上下文进行辅助分词的学习系统”，吉林大学管纪文、谷新英，会议论文集。
- [14] “分词词典与知识库”，吉林大学，王锡龙，中国中文信息研究会第二次年会论文。
- [15] ARTIFICIAL INTELLIGENCE, P. H. Winston 1977.
- [16] An Introduction to Database Systems, C. J. DATE.
- [17] Principles of Database Systems, D. Ullman,

结合上下文辅助分词的学习系统

吉林大学 管纪文 谷新英

摘要：随着自然语言自动化处理中有关工作的进展，解决书面汉语分词自动化的问题日趋迫切。

本文提出一个计算机辅助分词系统，它采用按表核对和参照上下文推导相结合的处理方式，可以在一定程度上独立地解决分词问题。但是，为了保证处理结果的正确性，适当的人工干预是必须的。该系统具有学习功能，故其性能可以在使用过程中不断得到改善。

文章阐述了系统的设计思想和工作原理，并给出了“分词处理程序”的算法。

一、引言

本系统的设计立足于现实性、可靠性、通用性、智能性和自完善性。

图1是系统结构概图。

系统的核芯程序是“分词处理程序”，它以词表和知识库为工具，对语言文本实施分词处理，并在适当的情况下接受人工干预。

下面，我们以分词处理过程为主线展开讨论和叙述，逐步地涉及系统结构的各个部分。最后，我们给出“分词处理程序”的具体算法框图。

二、文本的预处理

内容完整、结构正确的一部书、一篇文章、一组语句或者一个词组都可以是系统的处理对象，这里简称为文本。

在把文本送入计算机之前，要人工地进行若干必要的预处理。预处理包括以下内容：

1) 文本线性化

文本应以线性字符串的形式输入计算机，因而对其中的少数例外情形（如底线）要加以特殊处理。

2) 特殊词预标记

特殊词的预标记与系统使用动态词表有关。词表是词的集合，但它还包含一些非词汉字（如：葡、们）和一些非汉字字符（如英文字母、标点符号）。语言的词汇是一个无限集合，任何词表都不可能包容全部词汇。为了使系统具有较好的通用性，同时又不致造成过高的存储和检索代价，我们把词表设计成动态的：一个相对稳定的部分叫作基本词表，保存着一些常见词；另一个比较灵活的部分叫作临时词表，保存着一些用途偏窄的特殊词。

临时词表内容的增生主要依赖于特殊词预标记工作。所谓特殊词预标记，就是由人工事先将文本中那些可能不在词表的词的第一次出现给以特殊标记。系统在处理时若发现果真如此，则把该词加入临时词表。当临时词表增大到一定的限度或对新的处理对象已不适用时，可将其全部或部分地删除，或者转储备用。

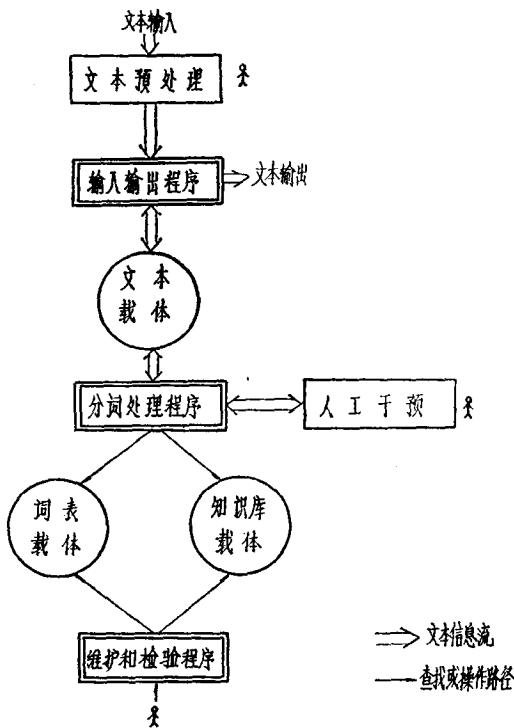


图 1 系统结构概图

特殊词的人工预标记可能有遗漏，这可在随后的人工干预中得到弥补。

3) 含有非汉字字符的词的预标记

后面要谈到，分词处理的第一步是文本切割。由于所有非汉字字符都被列为切割标志，故若不预先标记含有这种字符的词（其例如：一二·九、阿Q、维生素PP），则将造成错误的切割。

4) 文体结构的标记（如章、页）是为了跟踪和输出的方便。

三、文本的切割

对于已经存储在计算机上的文本，系统处理的第一步是把它切割成词段（纯粹由词组成的文本片断），然后针对每个词段进行分词。这里需要指出，在算法中，这两个步骤从整体上说是并行实现的。

显然，切割便是粗分。词段切割得越短，对随后的工作越有利。但是，把切割的难度限制在较低而又较为一致的水平上是合适的，复杂的工作应该集中到词段分词中。

基于上述原则，我们选择两类字符作为文本切割标志：一类是全体非汉字字符，一类是一些特殊汉字。其中，后者是具有这种性质的字：它们是某些词的词首或词尾，但不是任何词的其他部分。对于只能作词首的汉字，可从其前面切断文本，故称前位切割标志；对于只能作词尾的词，可从其后面切断文本，故称后位切割标志。但是，由于词汇的可变化性质，认为某些汉字恒具有上述特点是不合适的，于是索性放宽一些，允许这种字有少数的例外情

形。一个有例外汉字切割标志的一个或多个例外词作为一条切割知识保存在知识库中，并由该字的词表条目直接索引。每条知识都是一个逻辑蕴含式而不是等价式。

选择汉字切割标志应综合考虑下列因素：①预选汉字的使用频率；②该字作为前位（或后位）切割标志的例外词个数；③每个例外词的使用频率。一般说来，前者越大，后两者越小的汉字越宜选用。下面开列出在3775个国标一级字中初步选定的574个汉字切割标志。词的确定基本上以〔8〕为依据，选择时参考了王以德等人的工作〔4〕。

汉字切割标志：

(括号中的词是例外词)

1) 前位切割标志 (共327个)

啊 阿 挨 哎 唉 艾 俺 按 航 益 敝 翱 懊 扒 吧 跋(序跋) 罢(也
罢、作罢) 搬 颁 磅 傍 褒 鮑 悲(慈悲) 甭 逆 彼 碧(金碧辉煌) 蔓 陞
卞 扁(侧扁) 憨 濒 摈 丙 秉 玻(毛玻璃) 波 涠 猜 踩 惭(羞惭) 沧
拆 豺 摒 阖 猢 赤(足赤) 抽 踰 稠 雌 捶 簇 崔 挫 傥 殆 逮 怠
耽 拯 捄 佃(租佃) 惕 碉 调 掉 跌 盯 丢 洞 陡 杜 兑 吨(公吨、登
记吨) 跺 讹 扼 遏 琥 诽 吻 丰 烽 舸 辐 弗 辅 釜 阜 讦 肛 搞
各 巍 虱 刷 诡 癸 刻 憨 捍 豁 很 葫 蝴 蒙 焕 幌 恍 蝴 畸 饥(充
饥、朝饥) 激(感激) 讥 汲 伎 既 嘉 佳 艰 兼 缢 僵 骄 傑 狡 揭 皆
睫 梗 纠 圭 拘 犹 沮 撸 摟 扌 咖 凯 慷 磕 棵 颗 挎 夸(浮夸)
匡 岔 僥 垦(坷垃) 嘰 潼 勒(弥勒) 磬 瞭 撂 琳 拎 羚 另 漫 漫
枚 孟 弥 貌 哪 你 悠 柠 拧 凝(冷凝) 扭(歪歪扭扭、别扭) 痒 糯 哟
沤 啪 趴 徘 磐 眇 眇 庞(面庞、脸庞) 搞 抛 呃 烹 瑟 毗 啤 疲 脍
乓 颇 (偏颇、景颇) 菡 沏 琼 涣 痘 壬 妊 仍 冗 熔 揉 蠕 汝 闰 搔 哈 珊
蕃 搀 饮 沁 琼 涣 痘 壬 妊 仍 冗 熔 揉 蠕 汝 闰 搔 哈 珊
煽 赠 稍 苟 呻 甚(过甚、太甚) 孰 曙 甩 摔 拴 谁 朔 撕 伺 忒 艘
虽 裳 他(其他、排他) 它(其它、排它) 她 檀 倘 育 涂 拖 驼 楠 唾 碗
完(玩儿完) 懊 我(自我、忘我、唯我、大我、小我) 牮 羡 姓(复姓、百姓)
酬 旭 绚 衡 阎 佯 耶 椰 噩 颠 矣 屹 臆 肆 窾 瞻 雍 甬 永(隽永)
迂 逾 愉 鸳 曰 酝 哉 再(一再) 乍 辗 崩 肇 遮 这 斡 挣 狰
帧 只(不只) 蜘 秩 皱 拽 卓 琢 苗 啄 啓 兹 滋 仔 棕 摟 最 遵
俎 祚 昨

2) 后位切割标志 (共247个)

唉 蔴 隘 坝(坝埽 坝子) 般(般配) 锅 狍 甭 蔽(蔽塞) 痢 敝(敝
屣) 柄(柄子) 箔 舶(舶来品) 脖 泊(泊位) 埠(埠头) 怖 眇 踩(踩水)
舱(舱位) 蹤(蹬蹭) 砧(砧儿) 僮 擎(擎肘) 漱(漱底) 忱 騞 匙(匙子)
畴(畴昔) 踏 础 捶 囗 痒 粹 殆 悚蹈(蹈海、蹈袭) 涤 甸(甸子) 睹 妒
(妒忌) 盾(盾牌) 刨(刨斧石) 饵 伐 阀(阀阅、阀门) 肀 啡 吠 氛(氛
围) 峰 幅 袱 弗 甫 釜 脯 傅 阜 缚 咻 漑 搞 墉(埂子) 菇 刷 豢
涸 很 乎 肅 壶 噤 猥 徊 瘴 簍 凤 卉 海 坡 簍(箕踞) 姬(姬鼠)
绩 辑(辑录) 己(己方、己任) 悸 际(际遇) 件 润 椒(椒盐) 窑(窑肥)

皆 结 诫 烬 茎 睛 炙 咎 疾 灶 倦(倦游) 偷 眇 犁 澜 览 佬 蕃
偏 例 哩 倆 梁 辄 瞑 捺 贲 颉 龄 溜 窟 篓 戮 侷 虑 率 峦 李
仑 纶 淚 寐 眇 们 棂 觅 盔 饭 摰 拧 泞 你(你们) 哟 帕 琶 湛 畔
呸 坪(坪坝) 颇 珀 泣(泣诉) 讫 嵌 呛 桅 撚(撬杠) 沁 擎 権 刃(刃
具) 纩(纫佩) 蓉 蕊 涕 衫 裳 梢(梢头) 媚 蚀 驶(驾驶员) 孰 墅 帅
拴 谁(谁边、谁个) 艘 撒 他(他们、他人、他日、他乡) 它(它们) 她(她们)
蹋 汤 膛 淌 惕(惕厉、惕励) 翁 我(我们) 隙(隙地) 暇 翔(翔实) 眉
墟 讶 呀 堪 殽 舂 漾(漾奶) 矣 亦 谊 谆 邑 吻 泳 佑 孟(孟
兰盆) 眇 域 呼 喻 誉 取(取手) 垣 曰 哉 蚊 燥 盍 绽 彪 蛰 撇(撇
口) 帧 址 帜 炙 岑 州 洲 沽 粥 宙 拽 缀(缀文) 兹 淳 浑 摱
纂

非汉字字符和汉字切割标志的词表条目中有相应的标记，处理程序据此（有时还要涉及一次知识库查找）判断是否应在文本中的一个字符的前面或后面切断文本。

四、词段的分词

词段的分词是整个处理工作的核心。

为叙述方便，下面继续引进几个概念。

词段中的任意一串字符若是词汇中的词（不管它在文本的语义下是不是词），则称之为该词段的形式词。一个形式词若在文本的语义下是词，则称之为实际词。

在下例的词段P中，所有的形式词下都划以底线，共八个，其中三个划以曲底线的是实际词。

P：发展 中 国 家

将一个词段P按词表L分成以形式词为单位，称为按L分解P，其结果叫作L下P的一个形式分解。

显然，在特定的词表下，一个词段可以有零个或多个形式分解。例如对于上例的P，若词表L中包含全部八个形式词，则P有六个形式分解；若L中不包含“发”，则P只有三个形式分解；若“发”和“发展”都不在L中，则P无形式分解。

若一个词段的某个形式分解的所有分解单位都是实际词，则此形式分解叫作正确分解。

在一个词表下，一个词段可能有也可能没有正确分解。若有，则应是唯一的。对于上例P，下列分解是正确分解：

发展 中 国 家

显而易见，对于任意文本，我们总可以通过扩充词表，使对该文本来说，它的任意词段都存在正确分解。对于这样的词表，我们说它对于这个文本是充分的。于是，我们有下述命题。

命题：若对于文本T词表L是充分的，则对于T的任意词段P，至少存在一个形式分解；当P的形式分解只有一个时，该形式分解就是正确分解。

我们对一个词段进行分词，实际上就是要求取它的正确分解。本系统采取穷尽所有形式分解而求正确分解的方法，以保证处理结果的正确性。正确分解的判定首先依据知识库中的