

自动文摘系统EAAS

李小滨 徐越

(中科院软件所)

EAAS: AN AUTOMATIC ABSTRACT SYSTEM

Li Xiaobin and Xu Yue

(The Institute of Software, Academia Sinica)

ABSTRACT

This paper introduces an English automatic abstract system EAAS, describes its architecture, internal representations, algorithms and background knowledge.

摘 要

本文介绍了一个英文自动文摘系统EAAS, 详细描述了其总体结构、各环节的内部表示和算法, 以及背景知识的组织和表示。

§1. 引 言

自动文摘的广泛应用前景是不言而喻的。但自从十年前美国Yale大学首次开展对自动文摘的研究以来, 越来越多的学者从事这方面的研究却至今未见十分成功的报道。因此, 要探讨这样一个具有相当难度的问题, 必须对其性质、背景以及困难进行客观的分析。

自动文摘是一类特殊的自然语言理解问题。语言的层面模型的观点指出, 语言具有三个主要层面: 结构层面、意义层面和功能层面。由于对语言各层面的研究至今尚很不充分, 自动文摘就难免面临诸方面难以逾越的障碍。首先在意义层面上, 由于语言可以有许多比喻性用法, 对其意义进行了不同的引申, 语句里各词的词义不是几个范畴能包括的, 故准确地把握语言的意义十分困难; 其次在功能层面上, 由于语言的功能过于广泛致使歧义问题十分突出。因此, 基于目前的语言研究水平, 只有采取一些避开这些困难的有效对策才能使当前对自动文摘的研究不至于重蹈旧辙。

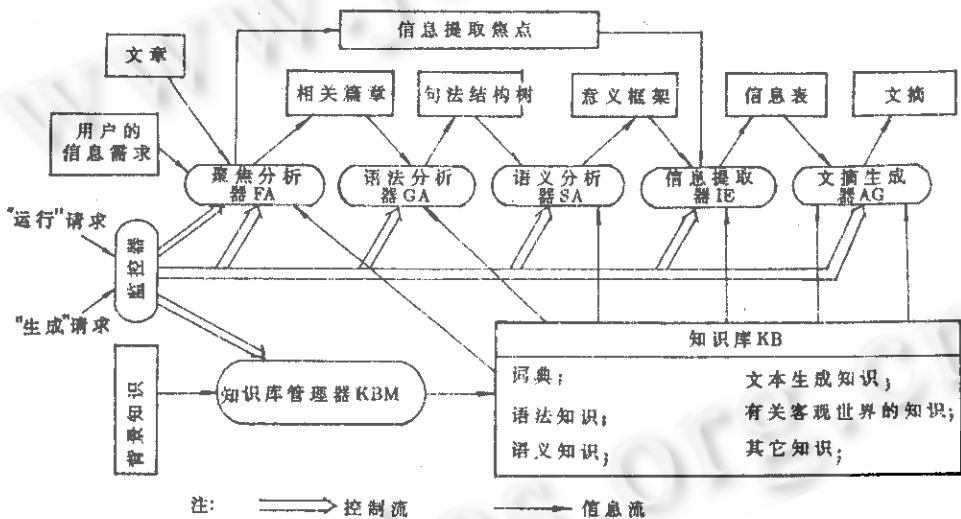
1989年8月25日收到, 1990年3月26日定稿。

对策之一是将自动文摘局限在意义较单一的范围内进行,以避开意义层面上的困难。因特定领域里语言一般比喻、引申的用法较少,大多采用了准确的表达,易于被计算机系统所理解。对策之二是从合适的角度固定语言使用的客观环境之后才进行自动文摘,以避开功能层面上的困难。显然将自动文摘系统固定在文章的读者(即系统用户)的角度上,按用户的需要从文章中摘取信息而不考虑文章中的其它内容,可以避开确定文章主题及处理文章中出现的新概念等目前难以圆满解决的问题,获得更大的成功。

基于上述切实可行的对策,两年来我们在研究自动文摘方面取得了初步的成果:利用基于知识的方法设计并实现了一个实验性的英文自动文摘系统EAAS。下文将介绍EAAS系统总体结构、各主要环节的内部表示和算法思想,以及背景知识的组织和表示。

§2. EAAS系统的总体结构及工作流程

EAAS系统由监控器、聚焦分析器FA、语法分析器GA、语义分析器SA、信息提取器IE、文摘生成器AG、知识库管理器KBM及知识库KB组成(见EAAS系统总体结构图)。



EAAS系统总体结构图

EAAS系统有三种工作状态: 监控态、生成态、运行态。在监控态下,系统用户(包括知识工程师和文章的读者)可以提出“生成”请求或“运行”请求,在监控器的调度下进入系统的生成态或运行态。在生成态下,知识工程师可通过KBM以人机交互的方式生成KB(将背景知识装入KB)或维护KB(对KB中的知识进行增、删、改等)。在运行态下,文章的读者可以将一些文章和自己的信息需求输入系统,在KB中各类知识的支持下,FA在这些信息需求的驱动下确定文章所属领域和信息提取焦点并筛选出与焦点相关的篇章和段落,GA和SA接着进行语法和语义分析,产生出这些篇章段落的一种内部表示,IE则在此基础上根据焦点进行信息搜索、推理及归纳,得出一个信息表,最后AG将此信息表转换为一篇文摘。

下面将着重介绍运行态下系统各环节的工作。

§3. 聚焦分析

聚焦分析是为了了解用户对信息的需求, 确定输入文章的应用领域和信息提取焦点集 F , 并根据 F 对文章进行筛选, 从而避免处理与 F 无关的内容, 达到提高系统效率的目的。

设用户所需信息项为 $inf_1, inf_2, \dots, inf_n$, 则 $F = \{f_1, f_2, \dots, f_n\}$, $inf_i \in f_i$, f_i 实际上是包含 inf_i 的最小信息范畴(这里称为焦点), $i = 1, 2, \dots, n$ 。

聚焦过程涉及 KB 中以下几类知识:

1. 信息范畴体系 ic

$$ic ::= \{sic_1, sic_2, \dots, sic_m\} \quad sic_k ::= ic \mid \langle \text{终结信息范畴} \rangle \quad k = 1, 2, \dots, m$$

2. 聚焦规则 FR

$$FR ::= P \rightarrow C, cf \quad P ::= ((ca : se) | f) OP \{((ca : se) | f)\}^* | f$$

$$ca ::= \langle \text{逻辑格} \rangle \quad se ::= \langle \text{语义参数} \rangle$$

$$OP ::= \langle \text{操作符} \rangle \quad C ::= f | (ca : se)$$

$$cf ::= \langle \text{可信度} \rangle$$

$$f ::= \langle \text{焦点名} \rangle$$

3. 细节约束框架 CFR

$$CFR ::= \langle \text{焦点名} \rangle \{ (SN \ VA) \}^*$$

$$VA ::= \{ se_1, se_2, \dots, se_m \} | n \{ OP \ se \}^*$$

$$se_i ::= \langle \text{语义参数} \rangle \quad n ::= \langle \text{整数} \rangle$$

$$OP ::= \langle \text{操作符} \rangle$$

$$SN ::= \langle \text{约束语义参数集} \rangle \mid \langle \text{关键词语义参数集} \rangle \mid$$

$$\langle \text{约束条件个数} \rangle \mid \langle \text{约束条件} \rangle$$

聚焦分析器 FA 的工作分为两个阶段:

1. 确定信息焦点: FA 通过交互获得一些待处理的文章和用户用自然语言提出的信息需求 re , 对 re 进行语法语义分析之后得到其意义框架 $remf$ (关于语法分析和语义分析过程见 §4 和 §5), 然后在 $remf$ 之上运用聚焦规则 fr 和信息范畴知识推出焦点集 $F = \{f_1, f_2, \dots, f_n\}$, 最后利用细节约束框架 cfr 产生 f_i 的相应焦点框架 f_i , $i = 1, 2, \dots, n$ 。

2. 筛选输入文章: FA 对输入文章各段落的首句进行语法、语义分析得到其意义框架 smf , 再利用 fr 产生其涉及的焦点集 F_1 。如果 $F \cap F_1$ 为空, 则删除该段落; 否则选择 $F \cap F_1$ 中可信度最高的焦点 fg ($1 \leq g \leq n$) 为该段落的焦点, 并采取关键字匹配法检查段落中的每个句子, 看其是否满足 fg 的焦点框架 f_g 中的约束条件, 从而决定该段落的取舍。

§4. 语法分析

语法分析器 GA 以 KB 中的词典及语法规则为背景对文章中的语句进行语法分析, 确定文章中各词的词义、词的形态、词之间句法上的联系及语句的各语法成分, 并以一棵层次结构的语法树描述之。

语句的句法结构树 ST 定义如下:

$ST ::= \{ < SU >, < MV >, < IO >, < DO >, < MOD >, < PD > \}$
 $< SU > ::= NP|NC|ST$ $< MV > ::= \{ aux, verb \}$
 $< IO > ::= NP|NC|ST$ $< MOD > ::= PP|ST|adv$
 $< DO > ::= NP|NC|ST$
 $< PD > ::= NP|NC|ST$ $NC ::= \{ NP, conj \}$
 $NP ::= \{ < art >, < HEA >, < DES >, < QUA > \}$
 $< HEA > ::= NP|NC|noun|pron|proper|ppn$
 $< DES > ::= NP|ST|adj|num|noun|inf|ger$
 $< QUA > ::= PP|adv$ $PP ::= \{ Pre, PRO \}$
 $PRO ::= NP|NC$
 $art ::= 'a'|'an'|'the'$ $aux ::= 'do'|'have'|'be'|'can'|'will'|'should'$

注: { } 为集合, < > 为句法成份, 大写字母串表示短语范畴, 小写字母串表示词法范畴, ' ' 为单词。

在语法分析器GA的工作过程中涉及KB中两类知识: 词典(提供了单词的词性、词义及物性、助动性、人称等), ATN文法规则。由于ATN的灵活性、文法描述与修改的方便性使得一些语义知识也能加入其中, 这样, 在利用ATN文法进行文法分析时就增加了语义测试功能, 方便有效地解决了语法分析中的歧义问题。GA采用自顶向下的分析策略和深度优先的搜索方法。由于这种算法已广泛使用, 这里不再赘述。

§5. 语义分析

要有效地从文章里提取出用户所需信息, 文章的内部表示应能反映文章里各语句之间的联系、语句里各成分的语义及它们之间的内在逻辑关系。语义分析器SA的任务就是在KB中有关知识的支持下将语法分析之后产生的一系列语句的结构树转换为文章的一种深层结构——文章的意义框架amf。amf定义如下:

$amf ::= (no_csmf)$ $no_csmf ::= no : csmf|no : csmf; no_csmf$
 $no ::= < 语句序号 >$
 $csmf ::= (rel_ssmf)$ $rel_ssmf ::= rel : ssmf|rel : ssmf; rel_ssmf$
 $ssmf ::= pred pmf$ $rel ::= < 复合句里各简单句的关系 >$
 $pmf ::= (ca_val)$ $pred ::= < 谓词语义参数 >$
 $ca_val ::= ca : val|ca : val; ca_val$
 $ca ::= < 逻辑格关系 > | < 语义范畴 >$
 $val ::= < 语义参数 > | pmf | ssmf$

注: csmf为复合句意义框架; ssmf为简单句意义框架; pmf为短语意义框架。

语义分析以KB中下列知识为背景:

1. 语义参数sep和语义框架sef。在EAS中词义用sep和sef描述之(sef仅用于描述谓词的意义)。

$sef ::= sep : (ca_1 : properties_1; ca_2 : properties_2; \dots; ca_n : properties_n)$
 $ca_i ::= < 逻辑格名 >$ $properties_i ::= < 语法约束 > < 语义约束 >$

2. 语义范畴体系SC

$SC ::= \{ SSC_1, SSC_2, \dots, SSC_m \}$ $SSC_k ::= SC | < 语义参数 >$ $k = 1, 2, \dots, m$

3. 用于处理歧义问题的有关知识

SA在语法分析之后得到的语句结构树ST上工作。对于一ST, SA先由KB查得各词的sep及谓词的sef, 并根据语义范畴体系得到ST中各短语的语义范畴, 然后以sef为模式

建立语句意义框架, 按properties, 递归地自上而下地搜索ST的有关句法成分以确定出语句意义框架csmf中各val, 并确定出csmf的pred(pred=sep), 从而建立了该ST的csmf, 最后在依次建立各ST的csmf的基础上再根据文章中语句的顺序、逻辑关系等建立起文章的意义框架amf。在此过程中的某些二义性问题往往借助意义框架里各成分之间语义制约关系或领域知识来引导解决。

§6. 信息提取

信息提取即按信息焦点集 $F = \{f_1, f_2, \dots, f_n\}$ 从文章的意义框架amf里搜索、推理出有关信息, 并进行归纳, 产生出一个信息表il。il定义如下:

$il ::= (inf_items)$
 $inf_items ::= ic : imf | ic : imf ; inf_items$
 $ic ::= < \text{信息范畴} > \quad imf ::= < \text{信息意义框架} >$

信息提取需借助KB中下列知识:

1. 信息提取特征集ie_p, $ie_p = \{ief_1, ief_2, \dots, ief_n\}$

$ief_i = ic_i : (ca_1 : cond_1; ca_2 : cond_2; \dots; ca_m : cond_m) \quad i = 1, 2, \dots, n$

$ic_i ::= < \text{信息范畴} > \quad ca_j ::= < \text{逻辑格名} > \quad cond_j ::= < \text{语义约束} > \quad j = 1, 2, \dots, m$

2. 信息推理规则ir

$ir ::= p \rightarrow c \quad p ::= ief, sef \rightarrow op$
 $ief ::= < \text{信息提取框架} > \text{ (同上)} \quad sef ::= < \text{语义框架} > \text{ (同 §5)}$
 $op ::= < \text{对意义框架进行修改的某种操作} >$

3. 精炼规则cr

$cr ::= p \rightarrow c \quad p = \{ic_{i1}, ic_{i2}, \dots, ic_{in}\} \quad c = \{ic_{j1}, ic_{j2}, \dots, ic_{jm}\}$

这里 $m < n$ 且 $p \cap c$ 为空, $ic_{i1}, \dots, ic_{in}, ic_{j1}, \dots, ic_{jm}$ 为信息范畴。

信息提取器IE针对每个信息焦点 f_i 依次进行信息搜索和推理。对于一个 f_i , IE先将其相应ie_p中的一系列ief与文章意义框架里的各子框架smf进行匹配(即用ief里各cond检查smf里的相应槽值), 如某smf_i匹配上某ief_k, 则ief_k中的ic_k就和smf_i一道成为信息表il中的一项。由于人类语言的字面意义之外还蕴含着大量信息, 为提取这类信息IE借助ir对某些蕴含着与 f_i 有关信息的smf进行转换, 推导出蕴含信息的显式意义框架smf', 并和ir中ief里的ic一道汇入信息表。在获得与各 f_i 有关的信息后还需利用cr对信息表里的信息进行整理和归并, 最终得到一个较为简洁、与 f_i 信息层次尽可能吻合的信息表il。

§7. 文摘的生成

文摘生成器的功能是由il产生一篇用户易读的英文摘要。鉴于系统的目标是以文摘形式满足用户的信息需求, 我们在文摘生成器AG的设计中避开了自然语言生成中一些复杂的问题(如文摘的风格等), 着重于研究文摘的逻辑性和简明性。

AG在工作过程中用及KB中下列知识:

1. 信息范畴体系ic (见 §3)

2. 信息缺省说明nf, $nf ::= ic : sent$

$ic ::= < \text{信息范畴} > \quad sent ::= < \text{说明语句} >$

3. 文本生成规则gr (类似于一般文法规则, 但其终结符不为单词, 为逻辑格)

文摘生成分为两阶段：段章结构形成阶段和语句生成阶段。段章结构确定了il 里各信息项之间的关系，从而确定了各语句在文摘里的出现顺序，决定了文摘的逻辑性。这里采用段章结构树pst 描述文摘的段章结构。pst 定义如下：

$$\begin{aligned} \text{pst} &::= (\text{nodes}) & \text{nodes} &::= \text{node}|\text{node}, \text{nodes} \\ \text{node} &::= \text{tnode}|\text{fnode} & \text{tnode} &::= \text{ic} : \text{imf}|\text{ic} : \text{pst} \\ \text{fnode} &::= \text{ic} : \text{sent} & \text{imf} &::= \langle \text{信息意义框架} \rangle \\ \text{ic} &::= \langle \text{信息范畴} \rangle & \text{sent} &::= \langle \text{信息缺省说明语句} \rangle \end{aligned}$$

在第一阶段中AG 以信息范畴体系的知识为背景，依次考查各焦点和il。对于一焦点 f_i ，首先建立起一个ic 值为 f_i 的结点 node_i 。然后查il，如il 中存在与其相关的信息项 $(\text{ic}_j : \text{imf}_j)$ $j = 1, 2, \dots, m$ ，则据此建立起结点 $\text{node}_1, \text{node}_2, \dots, \text{node}_m$ ，并按 ic_j 之间的隶属关系 rel_1 和平行关系 rel_2 建立这些结点间的联系(如 $\text{ic}_j \text{rel}_1 \text{ic}_k$ 则 node_j 和 node_k 为父子关系；如 $\text{ic}_j \text{rel}_2 \text{ic}_k$ 则 node_j 和 node_k 为兄弟关系)，从而构造起一棵子树 spst ，作为 node_i 的子结点；如果il 中不存在与 f_i 相关的信息项，则在 node_i 中放入该项信息的缺省说明。最后将各焦点对应的子树以兄弟关系连接起来就形成了文摘的pst。第二阶段AG 以深度优先的策略依次输入pst 上各项信息。如该信息结点有信息，就用gr 匹配结点中的imf，按匹配成功的句型依次输出imf 中的各成分；如该信息结点无信息，则输出结点的信息缺省说明。在此阶段AG 选择句型以简单明了为宗旨，在上下句的衔接上还考虑了代词的运用。

§ 8. EAAS 系统的实现及例示

EAAS 系统已在micro-vax 机上用C 语言实现，目前应用于一个就业机会介绍领域。系统可接受杂志Computer 中Career Opportunities 专栏的文章，按用户需求提取信息，以文摘形式输出。

下面给出两个小例子：

输入文章：

Texas A&M University

Applications are invited for faculty positions as Lecturer of Computer Science. Candidates from all areas of computer science will be considered. Experienced teachers who feel that their career would be invigorated by exposure to the computer science research community are invited to apply. Lecturers are not eligible for tenure.

The department of Computer Science has 15 full-time equivalent senior faculty members and additional teaching staff of 8 Lecturers. There is a tradition of quality instruction at the BS, MS and PhD levels. The university has initiated a major commitment to develop excellence in Intelligent Systems, Artificial Intelligence, Software Technology, Simulation and other areas of computer science. This commitment is in support of the rapid build up of microelectronics and computer technology in Texas.

Computer Science at Texas A&M university is located in the College of Engineering. With nearly 10000 students, the College of Engineering is the nation's largest. Computer Science is expanding to become one of the larger and more comprehensive departments in the country. The university has the internal resources to achieve this goal.

Ability in teaching undergraduate is essential. An MS in Computer Science is required with a PhD or comparable experience preferred. Applicants should submit a resume and three references to Bruce H. McCormick, Head, Computer Science Department, Texas A&M University,

College station, Texas 77843.

Texas A&M university is an equal opportunity, affirmative action employer.

例1:

What information do you want to know?

> I would like to know the information about education, speciality which the employee should possess.

> I want to get some information about the employer and pays.

ABSTRACT

Texas A&M University requires an MS in Computer Science with a PhD or comparable experience preferred.

It will consider candidates from all areas of computer science.

Texas A&M university is an equal opportunity and affirmative action employer.

The salary which will be provided by employer does not be mentioned.

例2:

What information do you want to know?

> I want to know the position offered and requirement for the position.

> Please tell me the working environment.

ABSTRACT

Texas A&M University invites applications for faculty positions as Lecturer of Computer Science. Lecturers are not eligible for tenure.

Ability in teaching undergraduate is essential.

The department of Computer Science has 15 full-time equivalent senior faculty members and additional teaching staff of 8 Lecturers. The university has initiated a major commitment to develop excellence in Intelligent Systems, Artificial Intelligence, Software Technology, Simulation and other areas of computer science. There is a tradition of quality instruction at the BS, MS and PhD levels. This commitment supports the rapid build of microelectronics and computer technology in Texas. Computer Science at Texas A&M university is located in the College of Engineering. The College of Engineering is the nation's largest with nearly 10000 students. Computer Science is expanding to become one of the larger and more comprehensive departments in the country.

注: 符号“>”后为用户输入的信息要求

致谢: 在EAAS 系统的设计过程中我们得到了马希文教授的悉心指教, 在此向他表示深深的谢意。

参考文献

- [1] Winograd, T, Language as Cognitive Process, vol, 1.
- [2] De Jong, G. F, Skimming Stories in Read Time: A Experiment in Integrated Understanding, Yale University, Dept. Computer Science, Research Report 158 (May 1979).
- [3] D. Roesner, Schemata for Understanding of Argumentation in Newspaper Texts, Universität Stuttgart.
- [4] J. I. Tait, Generating Summaries Using a Script-based Language Analyser, Acorn Computers, Cambridge.