

# 连续汉语语音识别中基于归并的音节切分自动机<sup>\*</sup>

张继勇 郑方 杜 术 宋战江 徐明星

(清华大学计算机科学与技术系语音实验室 北京 100084)

E-mail: ZhangJY@sp.cs.tsinghua.edu.cn

**摘要** 文章研究并实现了汉语连续语音中的音节自动切分算法——基于归并的音节切分自动机 (merging-based syllable detection automaton, 简称 MBSDA) 算法。MBSDA 算法利用了包括语音的短时能量、过零率和基音周期在内的多种特征参数, 把特征参数高度相似的相邻帧 (1 帧或若干帧) 的语音信号进行“归并 (merging)”, 形成“归并类似段 (merged similar segment, 简称 MSS)”, 它们被认定属于同一音节的相同状态。这些 MSS 经过一个包含若干状态的“音节切分自动机 (syllable detection automaton, 简称 SDA)”后, 输出音节的切分点。每个确定的切分段中所包含音节个数的范围 (range of syllable number, 简称 RSN) 也由 MBSDA 算法给出。

**关键词** 音节切分, 归并, 音节切分自动机, 韵母特征类段, 音节个数范围。

中图法分类号 TP391

目前, 在非特定人连续语音的识别中, 最广泛使用的声学模型是 HMM 模型及其改进模型, 其中传统的帧同步算法<sup>[1]</sup>或 Viterbi 解码算法<sup>[2]</sup>给出了状态解码序列。借助于一定的词法信息和语言模型, 就可以由声学搜索结果给出最大似然句子输出, 完成句子识别<sup>[3]</sup>。

在这样的方案中, 声学的搜索算法有两个问题不容忽视。(1) 搜索路径的组合爆炸问题; (2) 解码的状态序列错位问题。它们在很大程度上影响了整个识别系统的性能。我们的实验表明, 即使在连续语流中, 在音节切分点已知的情况下, 孤立音节的识别率也能达到理想的水平<sup>[4]</sup>, 因此, 理想的做法是把一段连续语音信号准确地切分至单音节, 然后再加以识别。但是, 由于目前研究条件和水平的局限性, 若不加任何限制就达到这一点是比较困难的。对有些音节, 我们可以给出准确的边界, 而有些音节之间的边界却很难加以区分。本文采取的方案是把完全确定的地方切开, 切不开的地方则给出音节的个数范围 (range of syllable number, 简称 RSN), 从而较好地解决了上述问题。为此, 本文提出了一个基于相似语音帧合并的音节切分自动机 (merging-based syllable detection automaton, 简称 MBSDA) 算法。汉语音节性很强的特点为这种算法提供了很好的理论依据。

## 1 算法基本原理

### 1.1 特征参数的提取

MBSDA 算法使用的主要特征参数是短时帧能量、过零率和基音周期。特征参数的提取和分析以帧为单位。具体的计算请见[5], 在此不再作介绍。

### 1.2 相近语音帧的归并

在同一个音素的发音过程中, 声道会在一定的时间间隔内保持稳定; 而当从一个音素过渡到下一个音素时, 声道会发生变化。因此, 如果连续几帧语音特征没有发生比较大的变化, 我们有理由认为它们是属于同一个音素

\* 作者张继勇, 1977 年生, 硕士, 主要研究领域为语音识别, 信号处理。郑方, 1967 年生, 博士, 副教授, 主要研究领域为语音处理, 语音识别和理解, 信号处理。杜术, 1976 年生, 硕士, 主要研究领域为信号处理, 语音识别。宋战江, 1972 年生, 博士生, 主要研究领域为语音处理, 语音识别和理解, 数字信号处理。徐明星, 1973 年生, 博士生, 主要研究领域为语音识别, 语音信号处理。

本文通讯联系人: 张继勇, 北京 100084, 清华大学计算机科学与技术系语音实验室

本文 1998-10-14 收到原稿, 1998-12-22 收到修改稿

的.基于此,我们提出了对这种特征相近的帧进行“归并(merging)”的概念.

在进行归并之前,我们首先检查语音特征是否发生“转折(transition)”.转折有 I 类转折和 II 类转折两种,分别描述如下.

I 类转折.特征发生突然变化.即当前帧的能量(或过零率)大于前一帧能量(或过零率)的  $\alpha$  倍;或当前帧能量(或过零率)的  $\alpha$  倍小于前一帧的能量(或过零率)(这里取  $\alpha=2$ ).

II 类转折.特征发生缓慢变化.即当前帧的前  $T$  帧语音能量(或过零率)的均值与后  $T$  帧语音的能量(或过零率)的均值之间存在着类似于 I 类转折中的变化关系(这里取  $T=3$ ).

如果当前的语音帧发生了上述的 I 类或 II 类转折,则给该帧语音作上“转折标记(transition tag,简称 TT)”.

连续的一帧或几帧没有 TT 标记的语音被归并到同一个归并类似段 MSS(merged similar segment)类段中.类段的一个很重要的性质是,它反映了这段语音中各个音素中最稳定的部分.

### 1.3 音节切分自动机的实现

通过类段并不能直接给出音节的切分边界,我们构造一个音节切分自动机(syllable detection automaton,简称 SDA),由自动机根据其内部状态的转移来确定音节的切分点.SDA 的状态划分为以下几类:静音、噪声、一类声母、二类声母、伪静音、韵母和韵尾.SDA 中的状态分别有如下含义.

静音(SIL)和噪声(NOI).静音是指能量和过零率都是很低的信号,它不包含语音信息.而当环境噪音比较强时,静音转变为噪声.

一类声母(SM1)和二类声母(SM2).对声母特征的研究和统计表明,l,m,n,r 具有类似于韵母的特点,我们将它归入二类声母的状态;剩下的声母的特征与韵母有较明显的区分,我们将它们归入一类声母.

伪静音(Pseudo-SIL).有些音节在发音时,声母和韵母中间有一个能量的“低谷”.比如,音节“kai”,声母“k”和韵母“ai”的能量都比较高,但它们之间的过渡段的能量却很低,类似于静音.如果把它归到静音状态,显然会发生错误.为了对这种情况加以区分,我们在这里引入了伪静音状态.

韵母(YM)和韵尾(YW).韵母的一个很明显的特点是,它具有准周期、较高的能量和适中的过零率.当发音从一个音节转变到下一个音节时,韵母部分能量和过零率都有比较明显的下降,这就是韵尾状态.通过大量的实验,我们得到了各个状态和特征参数之间的大致规律如表 1 所示.自动机各个状态间的转换如图 1 所示.

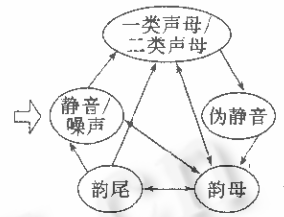


图 1 音节切分自动机的状态转换图

表 1 自动机状态和特征参数的关系

MSS 性质 帧参数	静音	噪声	一类声母	二类声母	伪静音	韵母	韵尾
规整能量	<100	<500	100~1000	100~1000	0~500	>500	100~1000
过零率	<3	10~40	>50	<30	<10	5~80	<30
基音周期	无	无	无	有	无	有	无

### 1.4 SDA 的输出和音节切分

如果我们把静音和噪声都当作一个“伪音节”或“广义音节”来看待,那么根据汉语音节的特点,当表 2 中的任何一个状态转移条件满足时,都表明语音流发生了从一个广义音节到另一个广义音节的转变.这也正是 SDA 需要输出音节切分标记的地方.

表 2 作音节切分标记的条件

状态转移条件	自动机前一状态	自动机的当前状态
条件 1	静音、噪声	声母、韵尾
条件 2	韵母	静音、噪声、声母
条件 3	韵尾	声母、静音、噪声、韵母

### 1.5 切分段所含单音节个数范围估计

切分段音节个数范围(range of syllable number, 简称RSN)的确定有下面几个步骤.

对于第  $n$  个切分段, 统计出其中所包含的韵母特征类段(vowel feature segment, 简称VFS)的个数, 记为  $C_{VFS(n)}$ . 当前切分段的平均音节长度(average syllable length, 简称ASL)值, 如果不考虑初始的情况, 可按下式来计算(以帧为单位).

$$ASL_n = \sum_{i=1}^{-M(1)} l_n^{(1)}(i) / M, \tag{1}$$

其中  $l_n^{(1)}(-i) = l^{(1)}(n-i)$ , 表示第  $n-i$  个“音节估计个数为1”的切分段长度(以帧为单位), 也即当前位置以前第  $i$  个“音节估计个数为1”的切分段长度. 上式用当前位置以前的  $M$  个“音节估计个数为1”的切分段长度来估计当前段的ASL值. 这里,  $M$  的选择必须合适, 如果太小, 则估计对音节长度的局部变化太敏感, 缺乏抗干扰能力; 如果太大, 音节个数的估计缺少对语速变化的跟随特性. 通常, 我们选  $M=10$ .

在给出RSN之前, 我们先利用ASL值定义一个上限参考值

$$C_{R(n)} = \lfloor L_n / ASL_n \rfloor, \tag{2}$$

其中  $L_n$  表示待估计的切分段长度.

在求得待估计切分段的  $C_{VFS}$  值和  $C_R$  值之后, 则该切分段所包含的音节个数的上限确定为

$$C_{\max(n)} = \begin{cases} C_{R(n)} - 1, & C_{VFS(n)} < C_{R(n)} - 1; \\ C_{R(n)} + 1, & C_{VFS(n)} > C_{R(n)} + 1; \\ C_{VFS(n)}, & \text{其他.} \end{cases} \tag{3}$$

利用本文给出的切分算法, 一般地,  $C_{\max}$  值不超过3, 大部分情况下为1.

## 2 实验及评价

### 2.1 实验数据

在实验中我们使用了两批测试数据. 第一批测试数据采用了863语音数据库. 这批测试数据是在安静的办公室环境下采集的, 基本上没有噪声的干扰. 我们从中抽取了5男5女的语音数据, 每个录音者各取20句. 这样就组成了一个男声和女声各100句的测试集. 为了对比噪声对切分结果的影响, 我们在环境噪声较强的房间里采集了第二批测试数据. 另外, 还故意加入了说话者的“吹气”、“咳嗽”等干扰信号. 这批数据共有100句, 都是男声. 两批数据的采样频率均为16KHz, 量化精度为16bit. 录音者的语速为每分钟150~180个字. 帧长为16ms(256个样本点).

### 2.2 音节切分算法的评价方案

切点正确率  $p_d$  和个数范围估计正确率  $p_n$  采用如下的公式来计算.

$$p_d = \frac{\text{SDA 给出正确切点个数}}{\text{SDA 给出总的切点个数}} \times 100\%, \quad p_n = \frac{\text{RSN 估计正确的切分段数}}{\text{SDA 给出的总切分段数}} \times 100\%. \tag{4}$$

总的切分正确率  $p_c$  由  $p_d$  和  $p_n$  共同决定. 为了反映在切分算法中能切出的音节占总音节的比例, 我们引入了切出率  $p_{out}$  的概念. 分别定义它们如下.

$$p_c = (p_d + p_n) / 2, \quad p_{out} = \frac{\text{SDA 切出的音节个数}}{\text{实际的音节总个数}} \times 100\%. \tag{5}$$

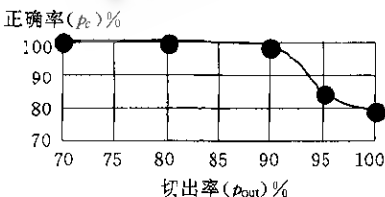


图2 MBSDA切分算法的ROC曲线

### 2.3 实验结果

根据上面的统计公式, 适当调整切分算法的参数, 可以得到切分算法在不同操作点下的ROC(receiver operating characteristics)曲线, 如图2所示, 在该曲线中, MBSDA算法的  $p_c$  是  $p_{out}$  的函数. 我们可以根据不同的切出率需求来调整MBSDA参数, 从而得到ROC曲线上的不同操作点. 但在一般情况下, 一个实际使用的切分算法要求  $p_c$  接近100%.

从ROC曲线上我们可以看到, 如果刻意追求切出率或者正确率都会导致算法整体性能的下降. 在允许发生

一定错误的情况下,当切出率大致为 90% 时,系统整体的性能最理想.表 3 中列出了其中一次实验统计的结果.

表 3 实验统计结果

测试数据库		结果(%)		切点正确率	RSN 正确率	总正确率	切出率
				( $p_a$ )	( $p_n$ )	( $p_e$ )	( $p_{out}$ )
第 1 批数据	男声			99.2	98.3	98.8	90.3
	女声			99.7	98.5	99.1	92.1
第 2 批数据	男声			98.9	96.2	97.6	89.8
	女声						
平均				99.3	97.7	98.5	90.7

从实验的结果可以看出,总体来讲,切分的效果很不错.在男声和女生的对比中可以发现,女生的切分效果要比男生好,原因主要是由于数据库中的女声发音比男声发音要清晰.此外,通过实验结果还可以看出,噪声对切分的结果有一定的影响,但算法的整体性能没有很大幅度的下降,这反映了本算法对噪声有较好的鲁棒性.

### 3 结 论

无论是从算法的可行性还是算法的复杂度来看,基于语音切分的汉语语音识别都是一种行之有效的方案.如何能保持本算法现有的切分正确率;提高算法的音节切出率,将是我们今后努力的方向.可以看出,在此算法的基础上,如果再加入语音信号的一些其他特征参数,如倒谱参数,将会提高算法的切出率.当切出率达到很高的值时,实际上也就是达到了切分至单音节的目標.此外,由于本切分算法的自动机已经给出了每一帧的状态,因此对该算法加以适当的修改即可实现连续语音中音节的声/韵切分.

#### 参考文献

- 1 Lee C H, Rabiner L R. A frame synchronous network search algorithm for connected word recognition. *IEEE Transactions on ASSP*, 1989, 37(11):1649~1658
- 2 Zheng Fang, Chai Hai-xin, Shi Zhi-jie *et al.* A real-world speech recognition system based on CDCPMs. In: *Proceedings of the International Conference on Computer Processing of Oriental Languages (LCCPOL'97)*. Hong Kong, 1997. 204~207
- 3 郑方,牟晓隆,徐明星等.一个语文学转换文本编辑器的实现.见:王承发,张凯等编.第 5 届全国人机语音通讯学术会议(NCMMSC'98)会议论文集.哈尔滨:哈尔滨工业大学出版社,1998. 280~285  
(Zheng Fang, Mu Xiao-long, Xu Ming-xing *et al.* The implementation of a speech-to-text transition editor. In: Wang Cheng-fa, Zhang Kai *et al.* eds. *Proceedings of the '98 National Conference on Man-Machine Speech Communication*. Harbin: Harbin Institute of Technology Press, 1998. 280~285)
- 4 郑方,吴文虎,方隼棠. CDCPM 及其在语音识别中的应用. *软件学报*, 1996, 7(863 专刊):69~75  
(Zheng Fang, Wu Wen-hu, Fang Di-tang. CDCPM with its application to speech recognition. *Journal of Software*, 1996, 7(special issue of 863):69~75)
- 5 杨行峻,迟惠生. *语音信号数字处理*. 北京:电子工业出版社,1995  
(Yang Xing-jun, Chi Hui-sheng. *Speech Signal Digital Processing*. Beijing: Publishing House of Electronics Industry, 1995)

### Merging-based Syllable Detection Automaton in Continuous Chinese Speech Recognition

ZHANG Ji-yong ZHENG Fang DU Shu SONG Zhan-jiang XU Ming-xing

(Speech Laboratory Department of Computer Science and Technology Tsinghua University Beijing 100084)

**Abstract** In this paper, an automatic syllable detection method namely merging-based syllable detection automaton (MBSDA) is studied and implemented. The MBSDA uses a variety of features including the frame energy, the zero crossing rate and the fundamental frequency to merge similar consecutive frames (one or several frames) into one merged similar segment (MSS). The frames in the same MSS are treated as frames of the same state of a phonetic. These MSSs are passed into a syllable detection automaton (SDA) to give the syllable detection results. In addition, the MBSDA gives the range of syllable number (RNS) of each definite detection segment.

**Key words** Syllable detection, merging, syllable detection automaton, vowel feature segment, range of syllable number.