

# 基于基频包络的音节聚类研究

刘 涛, 蔡莲红

(清华大学 计算机科学与技术系, 北京 100084)

**摘 要:** 对汉语 TTS 系统的大规模语料库做了基本的韵律参数统计, 分析了音节的韵律特征与其所在的韵律结构位置以及韵律结构边界的关系。进一步, 对有调音节样本集基于基频包络采用 k 中心点算法进行聚类, 通过听辨实验检验了聚类结果, 并分析了音节聚类与其所在韵律结构之间的对应关系。

**关键词:** 文语转换(TTS); 语音合成; 音库裁减; 音节聚类

中图分类号: TP391

文献标识码: A

文章编号: 1000-1220(2004)07-1145-06

## Study of Syllable Clustering Based on Pitch Contour

L U Tao, CA I L ian-hong

(Department of Computer Science and Technology, Tsinghua University, Beijing 100084, China)

**Abstract** In this paper, we analyze statistically the prosodic data in the large corpus for our Chinese TTS system. The relationship between the prosodic features of syllables, the positions and the boundaries of prosodic constituents is discussed. Furthermore, we cluster the syllables according to their pitch contours using the k-medoids algorithm, and evaluate the results through perceptual experiment. In the end, a conclusion of the relationship between the clustered categories and their locations in prosodic structures is presented.

**Key words:** TTS; Speech Synthesis; Corpus Reducing; Syllable Clustering

## 1 引言

近年来汉语语音合成取得了很大进展, 基于大语料库的波形拼接合成系统的可懂度与自然度有了很大提高, 其中语料库的制作越来越受到重视。目前的语料库一般通过设计大量的语料文本, 力求覆盖尽可能多的自然语言现象, 因此语料库的容量日渐庞大。但是无论从应用背景还是语料设计本身考虑, 音库容量都不能无限制增大。另外在某些场合下, 如对于嵌入式系统还需要裁减音库容量。如何有效的控制音库容量, 对音库进行裁减, 是目前需要解决的研究课题。

由于汉语语料库中最基本的单位是有调音节, 所以对语料库的裁减实际上就是对音节样本集的裁减。摩托罗拉中国研究中心的陈芳等人以音节的谱向量和基频作为特征, 采用变长分段量化(variable-length segment quantization)的方法压缩音库<sup>[1]</sup>。中国科学院声学所的孙金城等人以基频作为特征采用 k 均值以及 k 中心点聚类算法裁减音库<sup>[2]</sup>。以上基于基频以及谱参数的聚类可以保证压缩后聚类中的样本在听感上具有一致性, 聚类间的差异则较为明显; 但是还没有研究给出聚类与语境之间的关系分析。目前有很多关于汉语韵律结构的声学特征的研究表明, 音节在韵律边界处会呈现较为一致的声学特性, 如韵律词边界、韵律短语边界与语调短语边界前音节会出现音延现象, 边界后音节则一般出现音高重置现象。而相同声调的音节, 在句中的位置不同, 表现出的声学特

征也不尽相同, 这是因为自然语流中声调要受语调的影响。赵元任先生把汉语声调跟语调之间的关系形象地比喻为“小波浪”与“大波浪”并存叠加的“代数和”关系。如陈述句为下倾语调, 则句首词基频就明显高于句尾词。在语句中, 音节声调的调形基本不变, 但是具体音阶和音域则受语调的调节, 具体来说和音节在语调“大波浪”中的位置以及语句的重音位置密切相关。因此, 研究聚类后音节样本集与语境的关系也十分必要。

本文首先对大容量合成语料库进行了基本的统计分析, 然后以清声母去声音节样本集为实验对象, 以基频包络作为特征, 采用 k 中心点算法进行了聚类实验。实验结果表明: 聚类中的样本不仅在听感上具有一致性, 而且聚类内样本在句中的位置分布以及韵律边界前后的分布也具有一致性。

## 2 语音数据库分析

### 2.1 语音数据库介绍

本文的研究对象是汉语 TTS 系统的女声合成音库, 其采样率为 11KHz, 量化精度为 16 位。数据库包含了 3000 多个句子(以陈述句为主), 53909 个音节的音频数据, 覆盖了汉语中的 1283 种有调音节以及多种声学特征的搭配关系。其中阴平、阳平、上声、去声和轻声音节在音库中占的比例为 21%、23%、15%、31% 和 10%。音库中有调音节的样本均来自朗读的自然语句。每个音节带有拼音、语境及基频信息标注。其中

语境信息包括当前音节的位置信息(在句、韵律短语和韵律词中分别所处的位置),前、后音节的拼音、声调信息和音节所在词、所在短语及所在句子的具体信息。基频信息则包含了音节中的最大基频、最小基频、平均基频及波形数据中周期性峰值点的位置标注

在TTS的合成音库中,每个有调音节有若干个样本,组成了它的样本集。其样本的多少和文本语料设计有关,一般来说,汉语中越常用的音节的样本越多。下面是本音库中样本数目最多的5个音节的列表:

音节	声调	样本个数
De	轻声	2343
Shi	去声	953
Yi	去声	575
He	二声	513
Shi	二声	505

从上表可以看出,由于音节在自然语流中的使用频率不同,少部分的音节在音库中占有了大部分的容量。据统计,在目前使用的300M的合成音库中,样本数目在20以上的有调音节有568个,其总样本数为48942个,占据空间259M左右,达音库总容量的86%。这就为裁减音库容量提供了可能

## 2.2 韵律特征的声学参数分析

### 2.2.1 时长统计

经过分析计算,音库中最短的音节为55ms,最长的有668ms,其中时长为100ms到500ms之间的音节占98.68%。所有音节及所有声调的时长分布如表1所示。阴平、阳平、上声、去声音节的时长均值相差不多,去声时长均值最短。轻声音节的时长均值为总均值的67%。

表1 区分声调的时长分布

声调	阳平	阴平	上声	去声	轻声	所有音节	去除轻声
平均时长(ms)	271	273	266	249	170	254	263
方差	63	70	66	65	60	72	67

### 2.2.2 基频统计

我们首先对音库中音节的基频进行了分析计算。按不同声调统计,各音节的基频均值、基频最高点均值(高音点)和基频最低点均值(低音点)呈正态分布。表2给出了统计结果。其中最大基频均值为211Hz,是阳平调;最低基频均值为123Hz,是上声调。总体上看,此数据库的音域较窄,发音较为平稳

表2 声调基频统计分析

声调	阳平	阴平	上声	去声	轻声
基频均值(Hz)	200	162	136	175	142
高音点(Hz)	211	187	156	210	160
低音点(Hz)	186	141	123	139	126

### 2.2.3 韵律参数与语调和节奏的关系分析

音节的时长、基频表现与汉语韵律结构有着密切的联系,一般来说,韵律结构体现了语调、节奏以及重音。语调对音节韵律参数的影响体现在音节所处的位置上,以陈述句为例,句首音节在音高上明显高于句尾音节,在其他韵律结构层次上也表现出这种音高下倾的趋势。汉语的韵律节奏层级比较公认的有韵律词、韵律短语以及语调短语,很多研究指出<sup>[3]</sup>在韵律短语边界前音节会出现较明显的音延现象以及音高重置现象<sup>[4]</sup>。语句的重音则一般会表现为音高的升高以及时长的加长<sup>[5]</sup>。

我们分别对音节在句中所处的不同位置以及韵律结构边界处的韵律参数进行了统计分析,由于缺乏重音标注,在以下的统计中没有考虑重音的影响

#### 2.2.3.1 音节所在韵律结构位置的韵律参数统计

表3 音节所在韵律结构在句中不同位置的韵律参数统计

音节位置	句首词	句中词	句尾词	句首短语	句中短语	句尾短语
时长均值	277	249	246	259	256	242
基频均值	185	167	149	174	166	160
高音点均值	214	190	169	199	190	181
低音点均值	156	145	130	149	144	139

统计表明,句首词(句中第一个韵律词,下同)以及句首短语(句中第一个韵律短语)在时长上长于句尾词(句中最后一个韵律词,下同)和句尾短语,在基频上高于句尾词和句尾短语,反映了句调对音节韵律参数的影响

#### 2.2.3.2 韵律结构边界前后音节的韵律参数统计

表4 音节处不同韵律结构边界时韵律参数统计

韵律结构	句边界		韵律短语边界		韵律词边界	
	后音节	前音节	后音节	前音节	后音节	前音节
时长均值(ms)	309	245	266	312	251	242
基频均值(Hz)	204	140	177	148	174	160
高音点均值(Hz)	234	159	201	171	198	184
低音点均值(Hz)	171	121	152	130	150	140

表4中韵律短语边界前音节的时长明显长于全体音节的时长均值,体现了韵律结构边界前音节的延长效应,对于句边界及韵律词边界的时长,则稍短于全体时长均值,这是由于受到轻声音节影响的缘故,对于韵律词边界前音节,在去除轻声音节后,时长均值为266ms。研究中还发现,韵律结构边界等级越高,边界前后音节时长差距、基频差越显著

#### 2.2.3.3 去声音节韵律参数与韵律结构边界关系图

为了更直观的观察韵律参数与语调以及节奏的关系,我们以所有去声音节为例,给出了韵律参数以及韵律结构边界的关系图,图中可以看出音高表现为首高尾低的趋势,显示出音库中陈述句呈现的下倾语调,韵律短语以及韵律词也同样呈现下倾趋势。其中句尾、韵律短语和词尾音节的调域下限接近,这与沈炯关于音域的论述一致:“音域下限的延伸反映节奏单元的完整性”,“音域下限呼应关系反映的是节奏单元之

间的层次关系<sup>[4]</sup>。

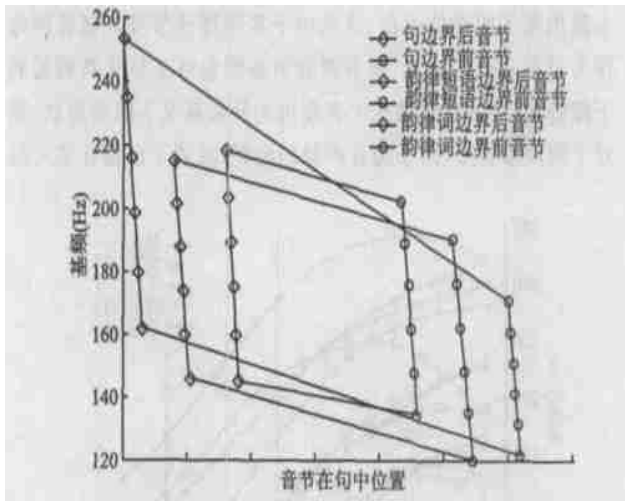


图 1 去声音节韵律参数与韵律结构边界关系图

Fig 1 Relation diagram of prosodic parameters and prosodic boundary of the fourth tone

通过对语音库的基本统计及与语调和节奏的关系分析, 我们可以得出这样的结论, 汉语音节的韵律参数受语调及节奏的影响, 在相似的位置以及韵律结构边界, 会表现出相似的韵律特征。那么当我们以韵律参数作为聚类特征, 得到的聚类结果也应该与音节的位置以及所处的韵律结构边界有密切的关系。

### 3 聚类实验

为了让实验结果更具有统计意义, 用所有去声音节作为实验对象。而去声音节由于声母类型的不同在基频包络上的表现又不尽相同, 因此, 又把去声音节分为清声母、零声母以及浊声母去声音节。

#### 3.1 数据预处理

对于每个音节, 通过标注得到其基频包络后, 进行了 3 点中值平滑处理, 但是音节之间基频序列的长度 (即音节中有基频段波形峰值点的个数) 并不相等, 还需要进行归一化处理。为了避免长度差异太大对聚类结果的影响, 对样本集进行了预分类处理: 首先, 统计得到基频序列长度的均值  $\mu$ , 然后按照基频序列长度的范围  $(0, 0.7\mu]$ ,  $(0.7\mu, 1.2\mu]$ ,  $(1.2\mu, \infty)$  把基频序列分为短、中、长三类。在第一类中, 设基频值采样点为 10 个, 也就是取基频序列中时间平均的 10 个点的基频值作为聚类原始数据; 第二及第三类则分别为 15 和 18 个基频采样点。经过平滑以及重新采样, 最终得到所有音节的等长基频序列作为实验数据, 基频序列样本的相异度则通过其基频序列向量的欧式距离来度量。这里需要说明的是, 基频序列的长度一般与音节时长成正比, 因为韵母在汉语音节时长中占主要部分, 而一般情况下, 只有韵母部分才有基频特性, 所以按基频序列长度分为的 3 类, 也对应了音节时长的短、中、长三类。表 5 给出了清声母去声音节预分类后的音节个数与时长均值统计。

表 5 实验数据统计

Table 5 Statistic of experimental data	
平均基频序列长度 $\mu$ (点数)	24
音节个数	8669
总体时长均值	251
$(0, 0.7\mu]$ 音节个数/时长 (ms)	1926 / 209
$(0.7\mu, 1.2\mu]$ 音节个数/时长 (ms)	3963 / 253
$(1.2\mu, \infty)$ 音节个数/时长 (ms)	2414 / 319

#### 3.2 聚类算法描述

本文采用使用最为广泛的基于划分的方法对音节样本集进行聚类, 其中最为著名也最为常用的划分方法是  $k$ -均值与  $k$ -中心点算法<sup>[6]</sup>。由于  $k$ -均值算法对于孤立点是敏感的, 一个有极大值的样本可能会相当程度扭曲聚类的分布, 所以本文采用了基于代表性样本的  $k$ -中心点算法。 $k$ -中心点算法采用类最靠近中心的样本而非计算出的均值作为新的类中心, 它仍然基于最小化所有样本与中心样本之间的相异度之和的原则来执行。其基本策略是: 首先为每个类随意选择一个代表点即类中心; 剩余的样本根据其与该代表样本的距离分配给最近的一个类; 然后反复的用非代表样本代替代表样本, 以改进聚类的质量。聚类结果的质量用聚类总体误差平方和  $J_e$  来衡量,  $J_e$  定义如下:

$$J_e = \sum_{i=1}^k \sum_{y \in \Gamma_i} |y - m_i|^2$$

其中:  $k$  为类数,  $y$  为第  $i$  个聚类  $\Gamma_i$  中的样本,  $N_i$  是聚类  $\Gamma_i$  中的样本数目,  $m_i$  即为代表 (中心) 样本, 它是聚类  $\Gamma_i$  中使

$\sum_{y \in \Gamma_i} |y - m_i|^2$  最小化的样本。 $k$ -中心点算法具体描述如下:

输入: 结果类的数目  $k$ , 包含  $n$  个对象的数据集合

输出:  $k$  个聚类, 使所有样本与其聚类中心点的相异度总和最小

方法:

- 1) 随机选择  $k$  个样本作为初始的中心点
- 2) repeat
- 3) 指派每个剩余的对象给离它最近的中心点所代表的聚类
- 4) 随机的选择一个非中心点对象  $O_{random}$ .
- 5) 计算用  $O_{random}$  代替  $O_j$  后的  $J_e$ .
- 6) if  $J_e < \min J_e$ , 用  $O_{random}$  代替  $O_j$  形成新的  $k$  个中心点集合.
- 7) until  $J_e$  不再发生变化

它具有如下特点:

1. 能有效的处理中等规模的数据集;
2. 对孤立点不敏感;
3. 算法的执行结果和样本的顺序、初始代表点的选择有关
4. 计算复杂度较高为  $O(k(n-k)^2)$ .

此算法需指定聚类数目, 这里我们取聚类数目为 5, 即基频值最高、次高、中、次低和最低的 5 类。

### 4 实验结果及分析

#### 4.1 去声音节聚类结果 (类中心基频曲线)

图 2、3、4 分别为清声母、零声母和浊声母去声音节的聚

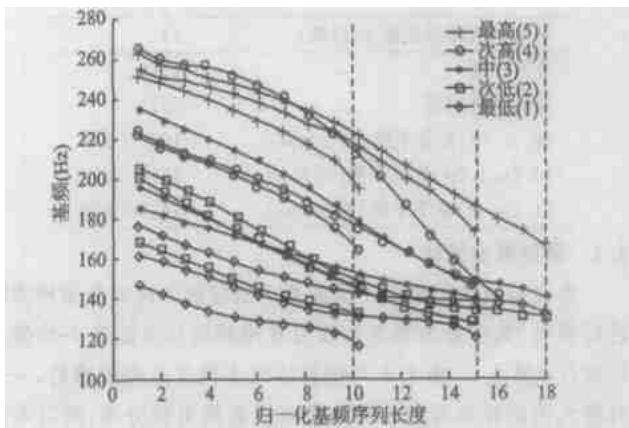


图2 清声母去声音节聚类结果

Fig 2 Cluster result of the syllables in fourth tone with consonant initial

类结果类中心基频曲线图 图中每种类型的曲线各有三条,长

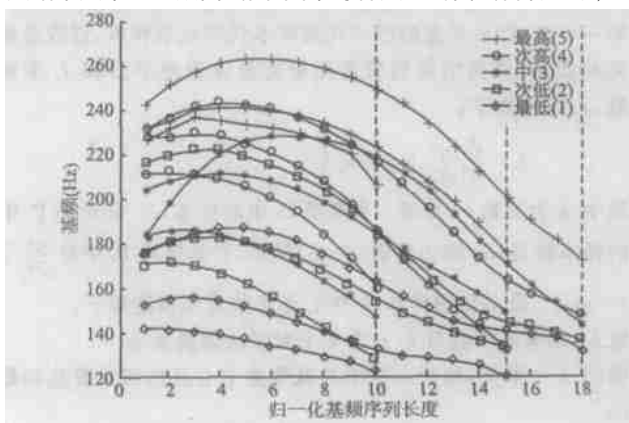


图3 零声母去声音节聚类结果

Fig 3 Cluster result of the syllables in fourth tone with null initial

度不一,分别代表基频序列长度的短、中、长三类 相同长度的基频曲线各有五条,代表了相同长度范围内基频包络各不相同

同的五个聚类 从图中可以看出,聚类中心的基频包络之间有着明显的差异,不同声母类型的去声音节的基频包络之间也存在明显的差异 零声母去声音节基频包络呈抛物线,在曲线头部出现了明显的上凸,这是由于零声母音节没有辅音声母作为与前音节的过渡,而去声音节基频包络又呈从高到低的下降趋势,故在基频曲线上表现出先低后高又下降的形状 而对于清声母音节,由于辅音声母的出现,减弱了前接音节对后

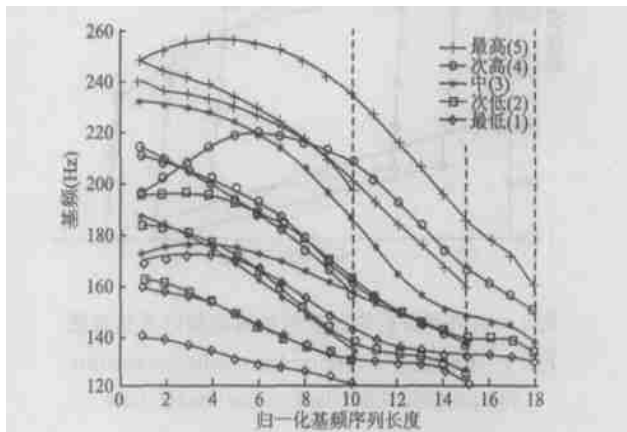


图4 浊声母去声音节聚类结果

Fig 4 Cluster result of the syllables in fourth tone with vowel initial

接韵母的影响,所以没有出现上凸现象 而浊声母去声音节基频包络中部分为平抛曲线,部分为抛物线,也有部分为斜线 这是由于声母浊化现象,声母部分声音信号也呈现出周期性,其发音与零声母类似也受前音节的韵母影响,故基频曲线也会出现上凸现象 下面将以清声母去声音节聚类结果为例,进行声学以及韵律分析

### 4.2 聚类的声学分析

聚类结果显示,音节在基频包络上具有一定程度的相似性,聚类内的样本具有可替换性,而类间样本则差异大 图5和图6是从同一个聚类中选出的两个音节样本(“是”)的波形及基频包络图 图7则选自另一个聚类 从图中可以明显看出,似,而与图7音节的差异则很明显;听感试验的结果表明,

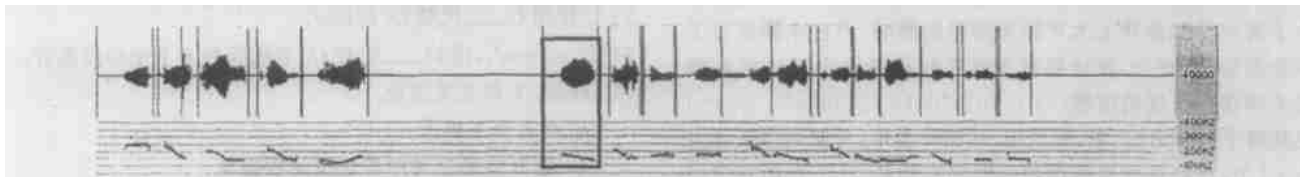


图5 音节1 思必奥转条“是”徽公司系列智力玩具之一

Fig 5 Syllable1:

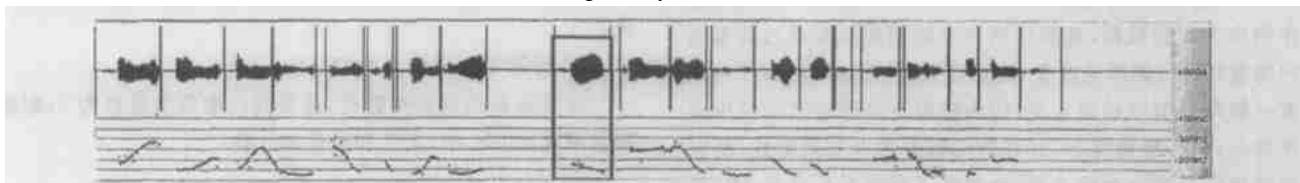


图6 音节2 雄浑宏伟大气磅礴“是”张幼矩画作的基调和风骨

Fig 6 Syllable2:

图 4 音节与图 5 音节不但听感近似, 而且把它们互换后, 两个 句子仍旧比较自然, 而与图 7 音节互换后则在听感上不可接

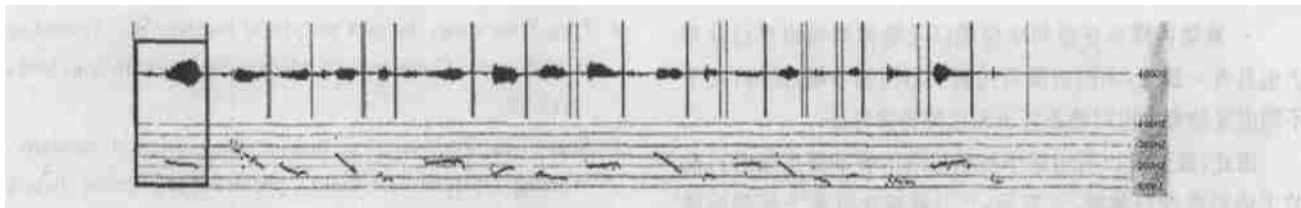


图 7 音节 3 “是”建筑设计师及有关设计工作者的好帮手

Fig 7 Syllable3:

受 从音节所处位置上来看, 前两个“是”都处于句中, 属于韵律短语边界后音节, 而第三个“是”则是句首音节。从这里已经可以看出, 声学特征类似的音节其语境也类似; 而声学特征相异, 则其语境也存在差异。

为了检验预分类以及基于基频包络的聚类结果与听感的关系, 我们设计了听辨实验, 参加实验的共 11 人。实验采用 30 组测试用例, 每组 3 个音节, 其中 2 个音节来自同一个聚类, 另外一个音节选自不相邻的一个聚类。前 15 个用例中的音节是选自时长相同基频不同的聚类, 后 15 个则选自时长不同的聚类。听辨结果为: 1- 15 中正确率为 79%, 16- 30 的正确率为 84%, 整体正确率为 81.5%, 考虑到听辨实验所选的音节都不是相同的有调音节, 且所选聚类差距较近(不相邻但仅间隔一类), 此结果说明预分类以及基于基频包络的聚类使音节在听感上分类效果可以接受。

#### 4.3 聚类的韵律分析

我们从清声母四声音节聚类结果中选出具有代表性的 4 类, 进一步对它们的韵律特征、在句中位置以及韵律边界的分布进行统计分析。

##### 4.3.1 聚类韵律特征分析

下表给出了这 4 类音节的基频、时长特性的统计结果。其中 S-1 即为基频长度在 (0, 0.7] 范围且基频均值最低的一类, S-5 即为基频均值最高的一类, 其余类似。从下表和图 2 中可以看出, 这 4 类样本在韵律特征上具有明显的差异。表中的峰值点为音节基频段波形标注峰值点个数。

聚类	S-1	S-5	L-1	L-5
时长均值	220	188	351	285
基因周期数均值	14	14	33	37
基频均值	131	230	149	220
高音点均值	148	255	177	266
低音点均值	115	194	129	164

表中显示, 基频越高基频周期数目越多, 这与人类的发音习惯相吻合。

##### 4.3.2 聚类音节在句中位置及韵律边界处分布统计

为了研究聚类与句调以及韵律结构的关系, 我们对聚类在句中位置以及韵律边界处的分布进行了统计, 如表 6、7 所示, 其中“S-1”后面括号里的小数为此类音节在句中绝对位置(到句首的音节数/句子的总音节数)的均值。

表 6 体现了音节的韵律参数与位置的关系, 如 S-1 主要分布在句尾短语、非句首词; 表 7 则提醒了韵律参数与韵律结

表 6 聚类音节所在韵律短语以及韵律词在句中位置的分布统计

位置	句首短语	句中短语	句尾短语	句首词	句中词	句尾词
	(%)	(%)	(%)	(%)	(%)	(%)
S-1 (0.85)	19.05	4.76	<b>76.19</b>	0.00	<b>48.81</b>	<b>51.19</b>
S-5 (0.38)	<b>73.33</b>	3.33	23.33	6.67	<b>90.00</b>	3.33
L-1 (0.43)	38.60	29.82	31.58	14.04	<b>85.96</b>	0.00
L-5 (0.25)	<b>75.00</b>	12.93	12.07	<b>46.55</b>	<b>53.45</b>	0.00

构边界的关系, 如 S-1 主要分布在句边界前音节与韵律词边界前音节。可以看出, 聚类音节在句中位置以及韵律边界的分布有着显著差别。我们结合以上统计可以给出这 4 类聚类音节样本集在语境上的描述如表 8 所示, 不同的聚类代表了不

表 7 聚类音节在韵律结构边界的分布统计

边界	句边界		韵律短语边界		韵律词边界		音节 (%)
	后音节	前音节	后音节	前音节	后音节	前音节	
	(%)	(%)	(%)	(%)	(%)	(%)	
S-1	0.00	<b>35.71</b>	0.00	9.52	16.67	<b>30.95</b>	7.14
S-5	0.00	0.00	3.33	0.00	<b>70.00</b>	3.33	<b>23.33</b>
L-1	0.00	0.00	<b>28.07</b>	19.30	3.51	<b>49.12</b>	0.00
L-5	<b>36.21</b>	0.00	8.62	0.00	<b>32.76</b>	7.76	14.66

同的句中位置以及韵律结构边界, 同时数据也表明, 不同的位置以及韵律结构边界也可能会表现出近似的韵律特征。例如句边界后音节与句首短语词边界后音节同属 L-5 聚类。

表 8 聚类的语境描述

聚类	分布描述
S-1	句边界前音节, 句尾短语词边界前音节 (主要在句尾短语中)
S-5	句首短语词边界后音节 (主要在句首短语的非句首词中)
L-1	句首短语词边界后音节, 句中、尾短语词边界前音节
L-5	句边界后音节, 句首短语词边界后音节 (主要在句首短语中)

## 5 结 论

本文以语音合成语料库中的去声音节为对象, 以基频包络为特征, 进行了聚类实验, 实验结果表明:

· 聚类内样本听感近似, 具有互换性, 而聚类间则差异明显

· 聚类内样本在语句中位置以及韵律结构边界的分布上也具有一致性, 不同的聚类代表了不同的语境, 同时, 处于不同语境的音节也可能表现出相似的韵律特征

因此, 我们可以利用基于基频包络的音节聚类来对目前的大语料库进行裁减。一方面, 它只裁减在听感上近似的样本, 从而保证了裁减后的音库在声学特性上损失较小; 另一方面, 聚类对应的语境同样具有代表性, 从而保证了压缩后的音库在语境上比较完备。这是实现高质量小型 TTS 系统的基础。此外, 在不同位置或韵律层次边界处的音节, 可能表现出相似的韵律特征, 那么我们在进行录音文本设计时, 就可以对这些位置或韵律层次边界进行合并以简化文本语料; 另一方面, 聚类的结果在前后音联上并没有表现出较强的相关性, 那么在小型 TTS 系统进行基元选取时可以简化处理, 以位置以及韵律层次边界作为选音的重点

#### References

- 1 Chen Fang Natural sounding embedded text-to-speech systems [C]. Proceeding of 5<sup>th</sup> national Conference On Modern Phonetics, Beijing, 2001, 10, 302-306
- 2 Sun Jin-cheng, Yi Li-fu Desing and analysis of layered corpus ofr speech synthesis [C]. Proceeding of National Conference on Acoustics, 2002, 377-378
- 3 Feng Yong-qiang Duration analysis of mandarin [C]. Proceeding of 5<sup>th</sup> National Conference On Modern Phonetics, Beijing, 2001, 10, 66-69.
- 4 Shen Jiong Pitch range of tones and intonation of mandarin, Working papers in experimental phonetics [M]. Beijing: Peking University Press, 1985, 96-101.
- 5 Feng Long Duration of the Initial and Final of Mandarin, Working papers in experimental phonetics [M]. Beijing: Peking University Press, 1985, 171-179.
- 6 Jiawei Han, Micheline Kamber, Data mining: concepts and techniques [M]. Beijing: Higher Education Press, 2001.

#### 附中文参考文献:

- 2 孙金城, 易立夫 分层语音合成数据库设计与分析 [C]. 全国声学学术会议, 2002, 377-378
- 3 冯勇强 等, 汉语语音音节时长统计分析 [C]. 第五届全国现代语音学学术会议论文集, 北京: 2001. 10, 66-69
- 4 沈炯 北京话声调的音域和语调 [A]. 北京语音实验录 [M]. 北京: 北京大学出版社, 1985, 96-101.
- 5 冯隆 北京话语流中声韵调的时长 [A]. 北京语音实验录 [M]. 北京: 北京大学出版社, 1985, 171-179