

TN912.3

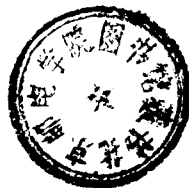
Y 219

0683374

语音信号数字处理

杨行峻 迟惠生等编著

4007/07



21113000911909

电子工业出版社

· 内 容 简 介

语音信号数字处理是在信息高速公路、多媒体技术、办公自动化、现代通信以及智能系统(包括机器人)等新兴领域应用的核心技术之一。近来来,虽然语音低速编码(传输与存储)、合成、识别、增强等已取得了长足的进展,但高新技术的发展又提出了更高的要求。因此,这个领域始终面临着新的挑战。为迎接这一挑战,首先在理论和方法上要有所突破,近年来在语音学、LPC分析、同态处理、矢量量化(VQ)、隐含马尔可夫模型(HMM)和人工神经网络等理论与方法研究方面有了很大发展。本书着眼于这些新发展并总结了编著者的科研成果,着重介绍了有关的基础理论,以及新方法和新技术。全书涉及内容广泛,同时又深入浅出,可读性强。既可作为电子工程、通信、信号与信息处理,以及智能信息处理有关专业大学生和研究生的教材,又可作为相应专业的科技人员自学或研究工作的参考书。

语音信号数字处理
杨行峻 迟惠生等编著
特约编辑 廖寿琪

*

电子工业出版社出版
北京市海淀区万寿路173信箱(100036)
电子工业出版社发行 各地新华书店经销
北京市燕山联营印刷厂印刷

*

开本:787×1092毫米1/16 印张:29 字数:724千字
1995年8月第1版 1995年8月北京第1次印刷
印数:0001—3000册 定价:45.00元
ISBN 7-5053-3147-7/TP-1129

本书编著者名单

清华大学

杨行峻
唐 昆
曹志刚
崔慧娟

北京大学

迟惠生

社科院语言所

杨顺安
李学军

前 言

经过清华大学、北京大学和社科院语言所六位作者近两年时间的共同努力,本书已全部杀青,即将付梓。为了将这几个单位近十年来在语音信号数字处理领域中许多科研和教学成果更好地总结归纳到这本书中,也为了更及时地向读者介绍当今世界语音信号数字处理各个研究领域所达到的新水平,各位作者付出了辛勤的劳动,使得这本篇幅相当大的书能够及时出版,这是令人感到十分欣慰的。回忆约十一年前朱雪龙教授和我以及其他几位同事合作译介 Rabiner 和 Schafer 合著的《语音信号数字处理》一书(科学出版社,1983)时,语音信号数字处理的研究在世界范围内尚处于高速发展的初期,而国内则无论在科研还是教学方面都还刚刚起步。该书的出版对于中国语音信号数字处理科研工作的开展和教学水平的提高无疑是超过作用的。十余年来,语音信号数字处理学科在世界范围内取得了飞速发展,无论在基础研究领域还是在各个特定应用领域(编码、识别、合成、增强等)都出现了许多崭新的算法和高性能的系统,取得了大量突破性的进展。在硬件方面,微机、工作站和 DSP 芯片的迅速更新换代,为各种日益复杂的语音处理算法的实时实现提供了可能性。可以预期,在本世纪的最后几年和下世纪初,这个研究领域的发展速度将更为加快,取得突破性进展的方面将更多,它与高速信息处理、传输和交换(诸如高速信息公路、ISDN、多媒体计算机、人机对话等等)诸方面的关系更加密切。着眼于培养语音信号数字处理领域的新一代科学技术人才,Rabiner 和 Schafer 的那本书已经显得不足了。五年前,作者们已决定写一本满足此要求的书作为大学高年级学生、硕士和博士研究生的教材,同时提供给有关领域科研和工程技术人员作为自学和研究参考用。当时力图将这本书写得尽量系统化,有关的理论基础和典型系统尽量交代得清楚、完备,同时将读者引导到各项研究工作的一些最前沿领域,而对一些因篇幅受限而不能全部写入这本书的内容则尽可能给出详尽的参考文献,以便于读者做进一步研究。

虽然这本书的各位作者做了很大努力,但是,语音信号数字处理毕竟是一个发展迅速、面貌日新月异的领域,在从编写到出版这段时间内,又有很多新发展来不及反映,而且受到编写时间与编著者水平的局限,书中不足和误植之处肯定在

所难免,恳请读者提出宝贵意见。

本书各章的编著分工如下:

一、三、四、五、六、十章由杨行峻撰写

二、十二、十三章由迟惠生撰写

七、八章由唐昆和崔慧娟撰写

十一章由曹志刚撰写

第九章由李学军撰写(参照杨顺安旧稿)

杨行峻负责全书的总体安排和审定。

清华大学电子工程系杨为理教授、社科院语言所吴宗济研究员、林茂灿研究员和北京大学中文系林焘教授,对本书的写作给予了热情的鼓励和多方帮助,在此一并表示衷心的感谢!

杨行峻

目 录

前 言

第一章 绪论	(1)
1.1 概述	(1)
1.2 语音信号数字处理的应用	(1)
1.3 语音信号数字处理的基础理论和算法	(2)
1.4 语音信号数字处理的硬件和实用系统	(3)
1.5 全书的组织	(3)
第二章 语音信号的数字表示、基本组成单位、产生模型和短时分析技术	(4)
2.1 概述	(4)
2.2 语音信号的时域波形	(4)
2.3 发声器官	(6)
2.4 音素与音节	(8)
2.5 语音信号的“短时谱”、“语谱图”以及元音和辅音的特点	(9)
2.6 韵母	(13)
2.7 声母	(16)
2.8 基音与四声	(20)
2.9 语音信号产生模型	(21)
2.10 语音信号数字处理中的短时分析技术	(24)
2.11 语音信号的短时自关函数和短时频谱	(27)
2.12 短时基音周期估计	(29)
参考文献	(33)
第三章 听觉系统和语音感知	(34)
3.1 概述	(34)
3.2 听觉系统生理学	(34)
3.3 听觉特性	(40)
参考文献	(47)
第四章 语音信号的同态处理与线性预测编码(LPC)分析	(48)
4.1 概述	(48)
4.2 同态信号处理的基本原理	(49)
4.3 复倒谱和倒谱的特点	(53)
4.4 求复倒谱和倒谱的实用算法	(54)
4.5 同态信号处理在语音信号数字处理中的应用	(59)
4.6 线性预测编码(LPC)分析的基本概念	(62)
4.7 LPC 正则方程组的自关解法和自协方差解法	(66)
4.8 用自关法解 LPC 正则方程的德宾(Durbin)递推算法、格形算法和舒尔(Schur) 递推算法	(69)
4.9 LPC 模型阶数 P 和激励增益 G 的确定,短时分析对于 LPC 参数估计的影响, LPC 分析的频域说明	(80)

4.10	各种 LPC 参数以及它们之间的关系	(83)
	参考文献	(90)
第五章	矢量量化(VQ)	(91)
5.1	概述	(91)
5.2	无记忆 VQ 及其码本形成的 LBG 算法	(92)
5.3	特征矢量和畸变准则的选择	(95)
5.4	无记忆 VQ 系统中减少搜索量、存储量和码本形成中训练量的各种算法	(100)
* 5.5	有记忆的 VQ 系统	(104)
* 5.6	全搜索 VQ 系统中的快速搜索算法	(106)
* 5.7	用随机松弛和模拟退火方法解决 VQ 码本形成算法中平均畸变值陷入局部 最小点的问题	(109)
* 5.8	人工神经网络与 VQ	(116)
	参考文献	(127)
第六章	隐含马尔可夫模型(HMM)	(129)
6.1	概述	(129)
6.2	y 为一维矢量且具有连续正态概率分布函数时 HMM 系统三项问题的解	(130)
6.3	y 为多维矢量且具有连续正态概率分布函数时 HMM 系统三项问题的解	(138)
* 6.4	y 为多维矢量且其概率密度函数为若干正态分布函数线性相加时三项问题 的解	(141)
6.5	离散和半连续 HMM 系统	(144)
* 6.6	语音处理和识别系统中 HMM 系统输出矢量 y 的选择	(149)
6.7	HMM 的各种结构类型	(152)
6.8	求解 HMM 三项问题时的一些具体计算问题	(155)
6.9	两个 HMM 相似度的比较	(160)
	参考文献	(161)
第七章	语音信号的波形编码	(163)
7.1	概述	(163)
7.2	语音编码的性能及其提高的途径	(163)
7.3	语音信号的标量量化	(171)
7.4	自适应差分脉码调制(ADPCM)	(186)
7.5	子带编码(SBC)	(200)
* 7.6	变换域编码(TC)	(211)
	参考文献	(222)
第八章	语音信号的参数编码与混合编码	(225)
8.1	概述	(225)
8.2	线性预测(LPC)声码器	(229)
8.3	多脉冲激励及规则脉冲激励线性预测(MPE-LPC 及 RPE-LPC)声码器	(239)
8.4	码激励线性预测(CELP)声码器	(251)
8.5	多带激励(MBE)声码器	(265)
	参考文献	(283)
第九章	语音合成	(287)
9.1	概述	(287)

9.2	语音产生的声学模型与合成器	(291)
9.3	普通话音节的合成框架及其实现	(301)
9.4	浊音声源的动态特性和普通话的字调模型	(305)
9.5	普通话音节的合成	(308)
9.6	协同发音与普通话词语的合成	(315)
9.7	韵律规则与普通话词语的合成	(319)
9.8	文-语转换系统	(323)
	参考文献	(327)
第十章	语音识别	(330)
10.1	概述	(330)
10.2	孤立词识别系统	(335)
10.3	连接词识别系统	(340)
10.4	采用 HMM 算法框架的连续语音识别	(347)
10.5	HMM 统一框架连续语音识别系统中最佳路径搜索算法	(349)
10.6	HMM 统一框架连续语音识别系统中声学-语音学层的设计	(353)
10.7	HMM 统一框架连续语音识别系统中的句法层设计和语言模型的建立	(360)
* 10.8	说话人自适应(Speaker Adaption)	(368)
* 10.9	关键词确认(Keyword Spotting)	(379)
	参考文献	(382)
第十一章	语音增强	(385)
11.1	语音增强的目的	(385)
11.2	语音和噪声的特性	(386)
11.3	语音增强算法概述	(389)
11.4	基于语音生成模型的增强算法	(391)
11.5	基于短时谱幅度估计的方法	(396)
11.6	短时谱幅度的 MMSE 估计方法	(401)
11.7	语音增强效果测试方法及实验结果	(407)
	参考文献	(410)
第十二章	说话人识别	(412)
12.1	概述	(412)
12.2	说话人识别的基本问题	(412)
12.3	表征说话人特点的基本特征	(415)
12.4	说话人识别的几种方法	(418)
	参考文献	(425)
第十三章	神经网络在语音信号处理中的应用	(426)
13.1	概述	(426)
* 13.2	神经网络初步	(427)
13.3	神经网络语音识别	(435)
13.4	神经网络说话人识别	(442)
13.5	神经网络语音合成	(443)
13.6	神经网络关键词识别	(444)
13.7	神经网络在语音信号处理其它领域中的应用	(447)
	参考文献	(450)

第一章 绪 论

1.1 概 述

语音信号数字处理是一门涉及面很广的交叉科学,虽然从事这一领域研究的人员主要来自计算机和通信等学科,但是它与语音学、语言学、数理统计学以及神经生理学等学科,也有非常密切的关系。作为一本为信息处理、通信和计算机科学等领域的高年级大学生、研究生、科研工作者和工程技术人员所写的基础教材,这本书着重从数字信号处理的角度来讨论这个课题。众所周知,语言是人类进行相互通信和交流的最方便快捷的手段。在高度发达的信息社会中用数字化的方法进行语音的传送、储存、识别、合成、增强……是整个数字化通信网中最重要、最基本的组成部分之一。

计算机的高速发展既对语音信号数字处理提出了越来越高的迫切要求(如用语音输入代替键盘输入以实现直接的人机对话),同时也提供了效率不断提高的软、硬件实现手段。另一方面,语音也是人类赖以进行思维的主要工具,因此,这一学科与认知科学和人工智能等研究领域,必然有千丝万缕的联系。近年来,人工神经网络的研究有了飞速发展,语音信号数字处理的各项课题是促使其发展的重要动力之一,同时,它的很多研究成果,也体现在有关语音的各项应用之中。目前,世界科技界正在蓬勃开展的其它一些新研究课题,诸如模糊理论、混沌理论和子波(Wavelet)信号处理等,也都能够在语音信号处理的研究中找到用武之地。

语音信号数字处理涉及一系列前沿科研课题,是目前发展最迅速的信息科学研究诸领域中的一个。正如其它数字信号处理研究课题,语音处理的研究涉及三方面互相密切配合的任务和课题,这就是:应用、算法(包括基础理论和软件)和硬件系统,三者缺一不可。由于这一领域的发展非常迅速,发表的有关文献浩若烟海,进行面面俱到的介绍既不可能也不必要。这里只介绍对当前研究工作有关的且最重要的基础理论和算法,并且迅即将读者引入当前最重要的研究课题(包括应用和系统),而不把精力放在一些支流或已成为历史陈迹的内容上。

1.2 语音信号数字处理的应用

如上所述,几乎语音信号处理的所有研究课题都是受到应用驱动的。以语音编码为例,由于数字化的语音传输和存储,无论在可靠性、抗干扰、速交换、易保密和廉价格等方面都远胜于模拟语音。从 50 代以来,在通信系统中数字化语音所占百分比不断增加。现在已非常清楚,在未来的 ISDN(综合业务数字通信网)、卫星通信、移动通信、微波接力通信和信息高速公路等系统中将无一例外地都采用数字化语音传输和存储。在不到 50 年的时间里,语音编码已有了惊人的发展。最早的标准化语音编码系统是速率为 64kb/s 的 PCM 波形编码器,到 90 年代中期,速率为 4~8kb/s 的波形与参数混合编码器,在语音质量上已逼近前者的水平,且已达到实

用化阶段。

据预测,速率为 2.4kb/s 左右的语音编码器,在未来几年中将在性能和实用化两方面都接近于 64kb/s 的标准 PCM 编码器。语音识别的研究起步较晚,大规模的研究开始于 70 年代初期,近年来已取得了长足的进展。一些中、小词表的孤立或连续语音识别系统已进入市场。目前,研究的重点是实现大词表、非特定人的连续语音识别系统。它可以用于人机直接对话、语音打字机以及两种语言之间的直接通信等一系列重要场合。这是一个难度相当大的高科技课题。在当前,学术界的普遍看法是:在信号处理、计算机、语言学、语音学和人工神经网络等各界学者的通力合作下,这一难题很有可能在本世纪末取得突破性的进展。语音合成是人机对话的另一个重要环节,让机器将文本语言转换成具有人声特点、抑扬顿挫自然流利的口头语言绝非易事,这一研究课题也正日益受到重视。其它一些重要的应用领域还包括语音增强(在强背景噪声或干扰中恢复“干净”的语音)和说话人识别及确认等。以上各个方面都是这本书所要讨论的内容。

1.3 语音信号数字处理的基础理论和算法

对于语音信号处理的基础理论和各种算法的研究包括紧密结合的两个方面。

一方面是从语音的产生和语音的感知来对其进行研究,前者涉及大脑中枢的言语活动如何转换成人发声器官的运动,从而造成声波的传播;后者涉及耳对声波的搜集并经过初步处理后转换成神经元的活动,然后逐级传递到大脑皮层的语言中枢。这一研究与语音学、语言学、认知学、心理学和神经生理学等密不可分。目前,对于这整个语言链的底层(或称为物理层),其中包括发声器官和耳的功能已经研究得比较透彻,但是对于其上层(即神经元的活动和大脑语言中枢的工作原理)则可以说还很不清楚。

另一方面,是将语音作为一种信号来进行处理。60 年代中期形成的一系列数字信号处理方法和算法:数字滤波器、快速傅里叶变换(FFT)、……与语音信号处理的要求分不开的。嗣后,在 70 年代初期产生了线性预测编码(LPC)和同态信号处理的算法,它们已成为进行语音信号处理最强有力的工具,且广泛应用于语音信号的分析、合成及各个应用领域。80 年代以后,出现了一系列更重要的方法和算法,其中包括语音编码中采用的分析合成方法,简称为 ABS (Analysis By Synthesis)以及各种自适应处理方法和变换方法。在语音识别方面最重要的是与隐含马尔可夫模型 HMM (Hidden Markov Model) 有关的一系列算法以及语言的概率模型。在编码和识别两个方面都非常重要的是与矢量量化(VQ)有关的各种算法。

应该注意,在研究各种算法时科研工作者通常采用两种方法,一种是用概率统计的方法,另一种是用规则的方法(或者说专家系统的方法)。虽然这两种方法互相渗透,不可能截然分开,但是仍能按照其主要观点和方法大致加以划分。从 80 年代至 90 年代中期的发展趋势看,前一种方法略占优势。而后一种方法则逐渐为人工神经网络的方法所取代或与之相结合。可以预期,在本世纪最后几年中,在前一方法继续平稳发展的同时,后一种方法将更加蓬勃地发展。

1.4 语音信号数字处理的硬件和实用系统

绝大多数语音信号数字处理系统需要按照实时方式或称为在线方式工作,这时对于系统的硬件环境要求很高(这里主要指系统的运算速度和内存容量的要求)。

随着语音处理算法的日益复杂,许多语音处理器的运算速度需要达到 10~20MIPS(Million Instructions per Second),在未来几年中这个速度甚至要达到 50MIPS。而在语音识别与合成等领域中对于处理系统的内存容量往往要求达到若干 MB。实用的实时语音信号数字处理系统通常以两种方式实现:第一种是用一台计算机作为主机(微型机、小型机或工作站)插上一块或若干块数字信号处理板来构成,后者由通用或专用的数字信号处理芯片(DSP 芯片)及相应的存储芯片、接口芯片和 A/D、D/A 芯片组成。第二种则由专用或通用的 DSP 芯片及其它辅助芯片构成一个独立工作的系统。前者通常称为非脱机工作系统,用于识别、合成、增强或模拟实验中。后者称为脱机工作系统,用于编码、小词汇表识别与合成等场合。

通用 DSP 芯片的出现及其性能价格比的迅速提高为各种实用化语音信号处理系统的实现铺平了道路。美国 TI 公司在 80 年代中期研制出的第一代 DSP 芯片 TMS32010 和 TMS32020 完成一次乘/累加运算(16 位、定点)需要 200ns,第二代 DSP 芯片 TMS320C25 完成一次乘/累加(16 位、定点)运算需要 100ns,第三代 DSP 芯片 TMS320C30 完成一次乘/累加(32 位、浮点)运算只需要 50ns 且片内的 ROM 和 RAM 和片外可扩充的 RAM 容量都大大增加。此外,美国 AT&T 公司研制出的 DSP-16C 和 DSP-32C,美国 AD 公司研制出的 ADSP21010 和 ADSP21020 等芯片系列与上述 TI 公司的第二代和第三代 DSP 芯片大致处在相似的水平上。第三代 DSP 芯片及更高一代 DSP 芯片的出现将使语音信号数字处理技术的发展和实用化登上一更高的新台阶。

1.5 全书的组织

这本书重点介绍语音信号数字处理的基础理论、算法和应用,其中第二至五章是各种具体应用领域的共同基础部分。为了使做更深入研究的人员得以获得彻底了解,对很多算法都做了详尽的推导。那些仅对这些算法的使用感兴趣的读者,则可以把冗长的推导略去不读。对语音编码感兴趣者,可以进一步阅读第七和第八章。对于语音识别和说话人识别感兴趣的读者,阅读第六、第十和第十二章。对语音合成或增强感兴趣者则可阅读第九或第十一章。为了区分主次,对每一章中的有些节,凡是加上“*”号,表明这是一些内容较深的节次,可以在初学或学习时间不够的情况下将其略去。从每一章所附的文献中,有兴趣的读者还可以找到更多的学习内容。

第二章 语音信号的数字表示、基本组成单位、产生模型和短时分析技术

2.1 概述

在研究各种语音信号数字处理技术及其应用之前,首先需要了解语音信号的一些重要特点,应知道它是如何由一些最基本的单位组成的,发声器官是如何发出这些音的,在此基础上可以建立一个既实用又便于分析的语音产生模型,这些问题可以归于声学语音学的范畴。通过对于语音信号发声过程的研究以及观察记录的各种语音波形,便可知道语音信号的频谱分量主要集中在 300~3400Hz 的范围内。如果用一个防混叠的带通滤波器将此范围内的语音信号频谱分量取出,然后按 8kHz 采样率对语音信号进行采样,就可以得到离散时域的语音信号。下面将讨论离散时域语音信号或称为数字语音信号。应该注意,为了实现更高质量的语音编解码器或者使语音识别系统得到更高的识别率,某些近代语音系统将此频率范围高端扩展到 7~9kHz,相应的采样率也提高到 15~20kHz。语音信号的另一个重要特点是它的“短时性”。在某些短时段中它呈现出随机噪声的特性,另一些短时段则呈现出周期信号的特征,其它一些是二者的混合。简而言之,语音信号的特征是随时间而变化的。只有在一短段时间间隔中,语音信号才保持相对稳定一致的特征,这短段时间一般可取为 5~50ms。因此,对于语音信号的分析和处理必须建立在“短时”的基础上。最重要的语音信号“短时特征”和“短时参数”包括它的“短时能量”、“短时过零率”、“短时相关函数”、“短时频谱”等。

语音信号的最基本组成单位是音素。音素可分成“浊音”和“清音”两大类。如果将不存在语音而只有背景噪声的情况称为“无声”,那么音素可分成“无声”、“浊音”和“清音”三类。在短时分析的基础上可判断一短段语音属于哪一类。如果是浊语音段,还可测定它的另一些重要参数,如基音和共振峰等。这里将讨论语音信号数字处理的这些基本知识、术语和分析技术。

2.2 语音信号的时域波形

在进行语音信号数字处理时,最先接触到并且也是最直观的是它的时域波形。为了获取一段语音信号的时域波形,首先将语音用话筒转换成电信号,再用 A/D 变换器将其转换为离散的数字化采样信号后存入计算机的内存中,最后将此信号取出,用绘图仪绘成时域波形。图 2-1 所示是一个男青年说的“欢迎你到深圳特区”这段话的语音时域波形。语音是在安静的环境下录取的。采样前经过频带为 0.1~3.4kHz 的带通滤波器进行滤波,采样率为 8kHz。每个采样信号用 12 位进行量化。这段语音的持续时间为 4 秒,图中横轴为时间,纵轴为语音信号的幅度。由于时间轴压缩得很短,从图 2-1 中无法辨别语音波形的细节,但是可以看到语音能量的

起伏,还可以大致分辨出话语中每一个字(音节)在此波形中的位置。为了仔细辨识语音波形,可以把时间轴拉宽。图 2-2(a)和(b)显示了这一段语音的波形细节,其中每一段横线伸展的

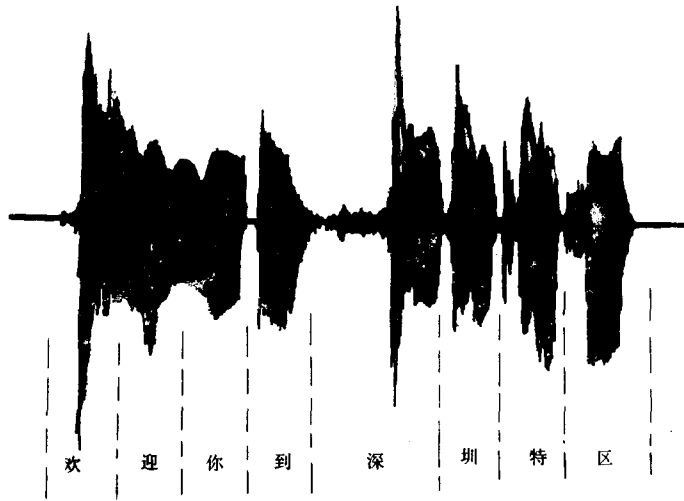
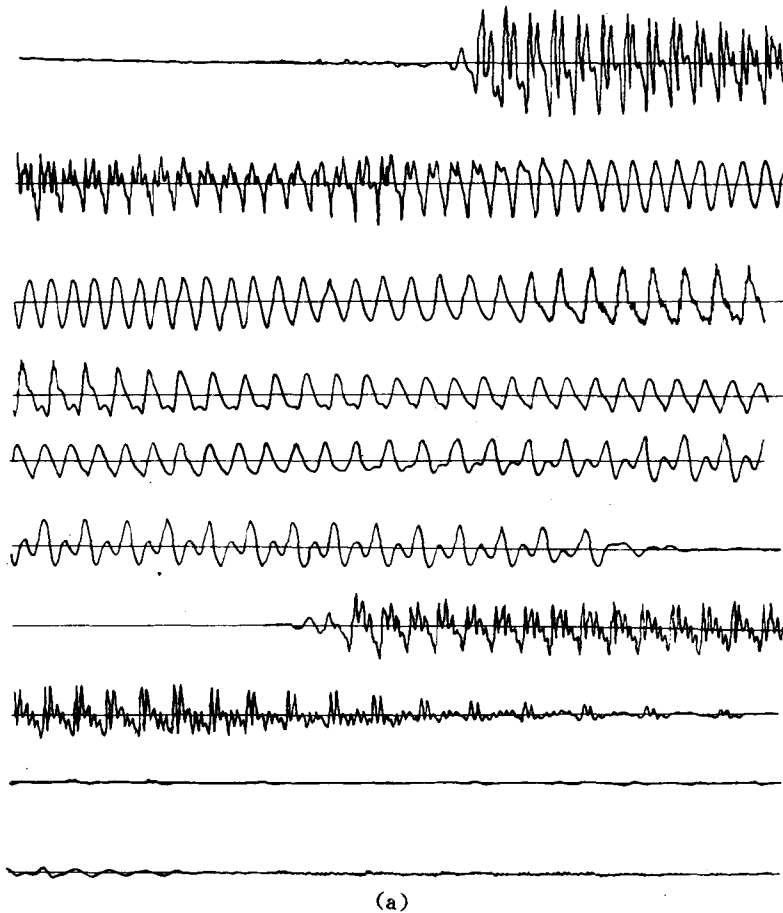


图 2-1 一段语音信号的时域波形



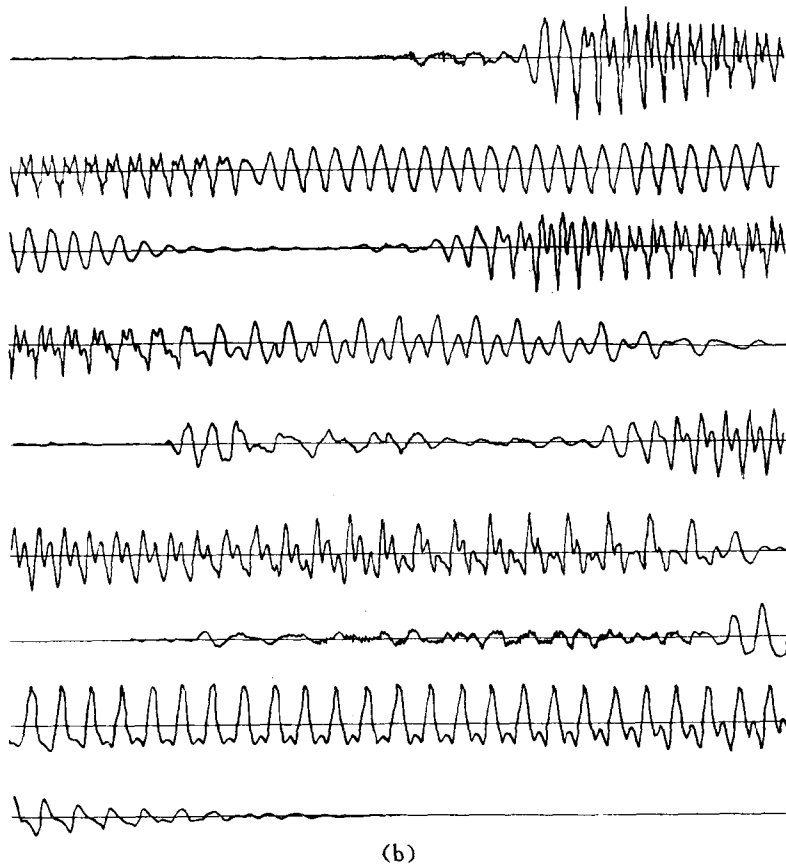


图 2-2 图 2-1 时域波形的展开图

范围是 200ms, 上段的末尾与下一段的起点相连接, (a) 和 (b) 相连接。由此图看出, 语音信号具有很强的“时变特性”。在有些段落中它具有很强的周期性, 有些段落中又具有噪声特性, 而且周期性语音和噪声语音的特征也在不断变化之中, 只有在较短的时间间隔中 (例如 20~200ms) 才可认为语音信号的特征基本保持不变。这一特点是语音信号数字处理的一个重要出发点。

2.3 发声器官

发声器官由三部分组成: 喉、声道和嘴。下面分别介绍它们的结构和功能。

2.3.1 喉

喉位于气管的上端, 其顶视解剖结构如图 2-3 所示。实际上它是气管末端的一圈软骨构成的一个框架, 前方稍高处的软骨称为甲状软骨, 前后方环成一圈的称为喉部环形软骨。喉中有两片肌肉, 称为声带, 它们的一侧由甲状软骨支撑, 另一侧则由两块杓状软骨支撑和控制。后者又与环形软骨连接。当它们分开时声带是张开的, 空气可自由地流过喉和气管 (见图 2-3(a)),

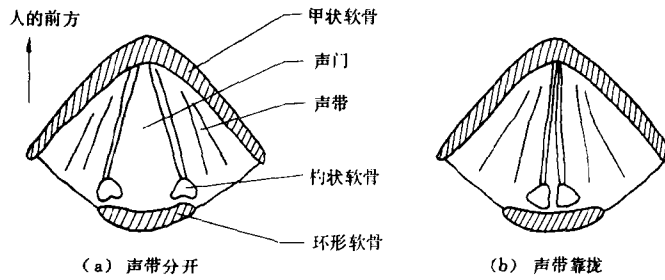


图 2-3 喉的解剖结构

正常呼吸时就处于这种情况。当它们合拢时，声带闭合将喉封住，在吃东西时食物就不会落入气管。两片声带之间的空隙称为声门。说话时两片声带在杓状软骨的作用下相互靠近但不完全封闭，这样声门变成一条窄缝(见图 2-3(b))。当气流通过这个窄缝隙时其间的压力减小，从而两片声带完全合拢使气流不能通过。在气流阻断时压力恢复正常，因此声带间的空隙再次形成，气流再次通过。这一过程周而复始的进行，就形成了一串周期性的脉冲气流送入声道。它的典型波形如图 2-4 所示。这一周期气流脉冲串的周期称为“基音周期”，用 T_p 表示，其倒数称为



图 2-4 典型的声门脉冲串波形

为“基音频率”，用 f_p 表示。 f_p 值取决于声带的尺寸和特性，也决定于它所受的张力。男性说话者的 f_p 值大致分布在 60~200Hz 范围内，女性说话者和小孩的 f_p 值在 200~450Hz 之间。用上面所述的方式发出的语音是“浊音”(Voice)。

2.3.2 声道

气流从喉向上经过口腔或鼻腔后从嘴或鼻孔向外辐射，其间的传输通道称为声道。声道的解剖结构(纵剖面)如图 2-5 所示。口腔的上顶分成两部分。前部是一块称为硬腭的骨头，它的作用是将口腔和鼻腔分开，并且支撑上排牙齿。后部由肌肉和连接组织构成，称为软腭，软腭的终端是小舌。当软腭在肌肉的作用下卷起贴在鼻腔的后壁上时，鼻腔和口腔相互隔开；反之，二者连通在一起。硬腭前部的骨头较厚，其中固定着牙齿，沿

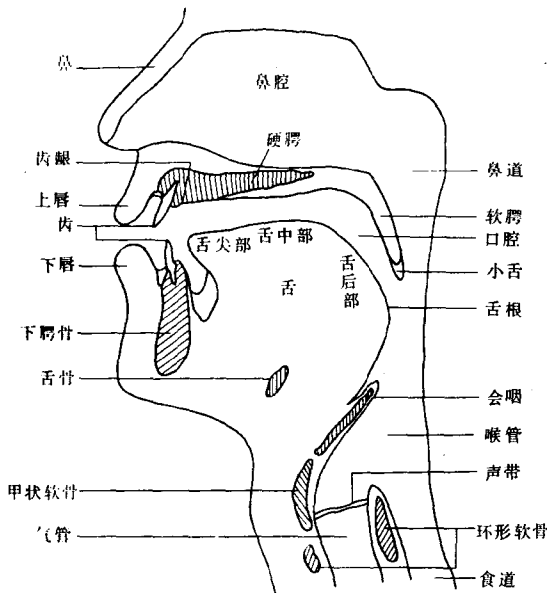


图 2-5 头的纵剖面，其中显示了各主要发声器官

着牙齿的一圈边缘称为齿龈。口腔下部是由肌肉构成的舌头,它的前部与下腭相连,后部和喉部的骨头及头部其它骨头相连。

气流流过声道时犹如通过一个具有某种谐振特性的腔体。输出气流的频率特性既取决于声门脉冲串的特性。又取决于声道的特性。为了便于分析,可以把声道当作一段无损声管如图

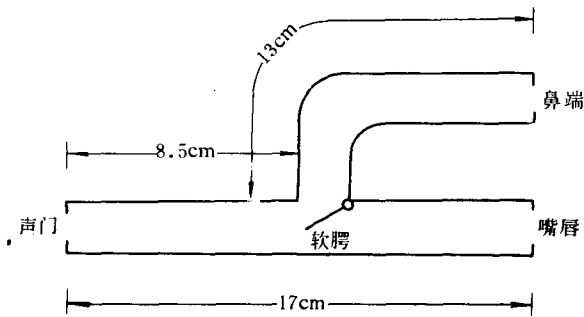


图 2-6 声道的无损声管模型

2-6 所示,其中鼻腔和口腔是否耦合,取决于软腭的位置。有耦合时发出的语音称为“鼻音”,否则为非鼻音。对成年男性而言,声道的口腔段的长度为 17cm 左右,而鼻腔段的长度约为 13cm。实际上声道的横截面积并非常数,所以声道模型中的声管应是一变截面积声管,而声道的频率特性主要取决于声道截面的最小值(一般称之为收紧点)出现的

位置。这一收紧点的位置又主要由舌的位置来控制。

语音的另一种产生方式是声门完全封闭,这时声道不是受声门周期脉冲气流的激励而是利用口腔内存有的空气释放出来而发声。由于该气流通过一个狭通道时在口腔中形成湍流,因而明显地具有随机噪声的特点。相应的语音称为“清音”(Unvoice)。汉语发音中的韵母如[a],[i],[u],[o]等均为浊音,某些声母如[s],[sh],[h],[x],[f]等为清音,另一些声母如[z],[zh],[j]等兼具二者的特点。[n]和[ng]是鼻音韵母。[m],[n],[l]是鼻音声母。

2.3.3 嘴

嘴的作用是完成声道的气流向外辐射。嘴的张开形状会影响语音频谱的形状,但是其作用较之声道而言是次要的。粗略而言,可以根据发音时嘴唇张开的圆形程度将一个音划归“圆唇音”或“非圆唇音”。

2.4 音素与音节

语音流由音素结合而成的最小单位,同时也是发声的最小单位是“音节”(Syllable),音节可以结合成更大的单位——“词”。词进一步可结合成“节奏群”、“句子”等等。音素的英语对应词是 phoneme,可以认为它是语音的最基本组成单位。事实上,同一音素与不同音素结合时,发音是有差异的。例如,[sh]这个音素在发“诗”([shi])这个音与发“书”([shu])这个音时,发音方式不完全一致,前者是非圆唇音,而后者是圆唇音。对于同一音素,它的各种不同发音方式称为“音素变体”(Allophone)。一个音节由元音(Vowel)和辅音(Consonant)构成。元音构成一个音节的主干,无论从长度看还是能量看,元音在音节中都占主要部分。所有元音都是浊音。辅音则出现在音节的前端或后端或前后两端。在汉语普通话中,每个音节都是由“辅音-元音”构成的(其中包括只有元音而没有辅音的纯元音音节,例如“啊”,这种情况称为“零辅音”),这种结构称为“C-V 结构”。在其它语系中还可以出现“V-C 结构”或“C-V-C”结构。在汉语中辅音也称为声母,元音也称为韵母。

单独发声的一个音节或是语音流中的任何一个音节都可能由 9 个部分组成,如图 2-7 所

示。其中1~4段属于声母(辅音)段,6~9段属于韵母(元音)段,第5段是二者的过渡段。对一个具体指定的音节而言,有可能只包括其中的某几段,但是第7段(主要元音段)是每一个音节都具有的。各段的特点及其发音机制将结合各个声母和韵母进行解释。

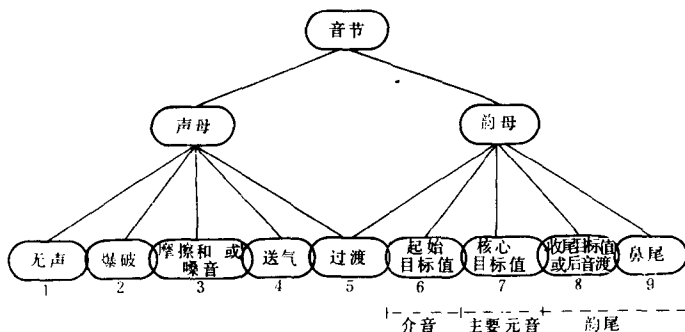


图 2-7 汉语普通话的音节结构框架(引自[1])

2.5 语音信号的“短时谱”、“语谱图”以及元音和辅音的特点

语音信号的最重要特征表现在它的“短时频谱”(简称为“短时谱”)上。如果从语音流中利用加窗的方法取出其中的一个短段,再对其进行傅里叶变换,就可以得到该段语音的短时谱。采用数字信号处理(DSP)的手段,可以在数字计算机上非常方便、快捷地完成这一任务。图 2-8 所示是一段浊音和一段清音的时域波形及其短时谱,语音的采样率是 10kHz,窗长为 50ms(相应的样点数为 500),窗形为哈明窗(在 2.10 中将较详细地讨论与此有关的短时分析问题)。浊音的短时谱有两个特点:第一,有明显的周期性起伏结构,这是因为浊音的激励源为周期脉冲气流。第二,频谱中明显地具有几个凸起点,它们的出现频率与声道的谐振频率相对应。这些凸起点称为“共振峰”(Formant),其频率称为共振峰频率。共振峰按频率由低到高排列为第一共振峰、第二共振峰,……,相应的频率用 F_1 、 F_2 、……来表示。一般浊音中可以辨别的共振峰有 5 个,其中前 3 个(尤其是前 2 个)对于区别不同语音是至关重要的。清音的短时谱则没有这两个特点,它十分类似于一段随机噪声的频谱。

在 DSP 技术发展起来以前很久,人们早就用一种特殊仪器——语谱仪来分析和记录语音信号的短时谱。它将语音信号(经话筒变成了电信号)送进一排频率依次相接的窄带滤波器,各窄带滤波器的输出记录在一卷按一定速度旋转的记录纸上(各滤波器的由低到高按频率排列),信号强则记录得浓黑一些,反之则浅淡一些。由此得到的即是语音信号的语谱图,此图的水平方向是时间轴,垂直方向是频率轴,图上或深或浅的黑色条纹表征各个时刻的短时谱。图 2-9 给出了 [i], [æ], [ə], [ɔ], [a], [u] 这六个美国英语元音单独发声时的时域波形和语谱图,其中与时间轴平行的几条深黑色带纹称为“横杠”(Bar),它们相应于短时谱中的几个凸出点,也就是共振峰。由横杠的频率及宽度可以确定相应共振的频率和带宽。在一个语音段的语谱图中,有没有横杠存在是判断它是否为浊音的重要标志。

图 2-10 给出了若干辅音配以元音 [a] 发音时产生的时域波形图和语谱图,它们的花纹比较复杂,其中比较典型的花纹是横杠、乱纹和冲直条(语谱图中出现与时间轴垂直的一条窄黑条)。每一种辅音包括上面几种典型花纹中的一种或几种,它们与该辅音发音的特点有密切关系,这将在 2.7 中进行详细讨论。