



高等院校信息与通信工程系列教材

SPEECH SIGNAL PROCESSING: PRINCIPLES AND PRACTICE  
(SECOND EDITION)

# 语音信号处理

(第2版)

韩纪庆 张磊 郑铁然◎编著

Han Jiqing Zhang Lei Zheng Tieran

清华大学出版社



· 013031461

TN912.3

18-2



高等院校信息与通信

SPEECH SIGNAL PROCESSING: PRINCIPLES AND PRACTICE  
(SECOND EDITION)

# 语音信号处理

(第2版)

韩纪庆 张磊 郑铁然◎编著

Han Jiqing

Zhang Lei

Zheng Tieran

TN912.3  
18-2

清华大学出版社

北京



北航

C1639963

## 内 容 简 介

本书系统地介绍语音信号处理的基础、概念、原理、方法与应用,以及该学科领域取得的新进展。全书共分9章,其中第1章绪论,介绍语音信号处理及其发展过程。第2章介绍语音产生与人类听觉的机理,传统的线性语音产生模型,以及近年来刚刚兴起的非线性语音产生模型。第3章从语音信号的时域特征入手,引入时频分析的思想,并进一步阐述时频分析中短时傅里叶变换和小波变换在语音信号特征分析中的应用,最后对广泛使用的倒谱特征以及同态解卷积进行介绍。第4章介绍语音信号的线性预测原理、解法、几种推演方法以及线谱对分析法。第5章介绍语音编码的相关知识,包括语音的波形编码、线性预测编码、极低速率语音编码技术,以及相关编码器的性能指标和评测方法。第6章介绍语音识别的基本内容,从基于矢量量化的识别技术到动态时间归正的识别技术,再到隐马尔可夫模型的识别技术,从孤立词识别到连接词识别及连续语音识别技术,再到关键词检出技术,最后还介绍近年来兴起的一些语音识别应用技术,包括语言学模型的自适应、HTK应用以及Lattice结构和混淆网络等。第7章介绍说话人识别的基本原理,主要包括说话人的特征选取、说话人识别的主要方法,以及近年来备受关注的GMM-UBM模型、开集说话人识别的规整技术等。第8章介绍近年来发展迅速的稳健语音识别技术,从影响语音识别性能的环境变化因素分析开始,介绍噪声环境下稳健语音识别技术,以及变异语音识别的技术。第9章介绍语音合成的基本原理、线性预测合成、共振峰合成以及汉语按规则合成,以及最近兴起的基于HMM合成技术等内容。

本书可作为高等院校计算机应用、信号与信息处理、通信与电子系统等专业及学科的高年级本科生、研究生教材,也可供该领域的科研及工程技术人员参考。

本书封面贴有清华大学出版社防伪标签,无标签者不得销售。

版权所有,侵权必究。侵权举报电话:010-62782989 13701121933

### 图书在版编目(CIP)数据

语音信号处理/韩纪庆,张磊,郑铁然编著.--2版.--北京:清华大学出版社,2013.4

高等院校信息与通信工程系列教材

ISBN 978-7-302-30269-8

I. ①语… II. ①韩…②张…③郑… III. ①语音信号处理—高等学校—教材 IV. ①TN912.3

中国版本图书馆CIP数据核字(2012)第234411号

责任编辑:盛东亮

封面设计:李召霞

责任校对:李建庄

责任印制:李红英

出版发行:清华大学出版社

网 址: <http://www.tup.com.cn>, <http://www.wqbook.com>

地 址:北京清华大学学研大厦A座 邮 编:100084

社总机:010-62770175 邮 购:010-62786544

投稿与读者服务:010-62776969, [c-service@tup.tsinghua.edu.cn](mailto:c-service@tup.tsinghua.edu.cn)

质 量 反 馈:010-62772015, [zhiliang@tup.tsinghua.edu.cn](mailto:zhiliang@tup.tsinghua.edu.cn)

课 件 下 载: <http://www.tup.com.cn>,010-62795954

印 装 者:北京鑫海金澳胶印有限公司

经 销:全国新华书店

开 本:185mm×260mm 印 张:25.5 字 数:635千字

版 次:2004年9月第1版 2013年4月第2版 印 次:2013年4月第1次印刷

印 数:10001~13000

定 价:44.50元

## 高等院校信息与通信工程系列教材编委会

主 编：陈俊亮

副 主 编：李乐民 张乃通 邬江兴

编 委 （排名不分先后）：

王 京 韦 岗 朱近康 朱世华

邬江兴 李乐民 李建东 张乃通

张中兆 张思东 严国萍 刘兴钊

陈俊亮 郑宝玉 范平志 孟洛明

袁东风 程时昕 雷维礼 谢希仁

责任编辑：盛东亮

# 出版说明

---

信息与通信工程学科是信息科学与技术的重要组成部分。改革开放以来,我国在发展通信系统与信息系统方面取得了长足的进步,形成了巨大的产业与市场,如我国的电话网络规模已位居世界首位,同时该领域的一些分支学科出现了为国际认可的技术创新,得到了迅猛的发展。为满足国家对高层次人才的迫切需求,当前国内大量高等学校设有信息与通信工程学科的院系或专业,培养大量的本科生与研究生。为适应学科知识不断更新的发展态势,他们迫切需要内容新颖又符合教改要求的教材和教学参考书。此外,大量的科研人员与工程技术人员也迫切需要学习、了解、掌握信息与通信工程学科领域的基础理论与较为系统的前沿专业知识。为了满足这些读者对高质量图书的渴求,清华大学出版社组织国内信息与通信工程国家级重点学科的教学与科研骨干以及本领域的一些知名学者、学术带头人编写了这套高等院校信息与通信工程系列教材。

该套教材以本科电子信息工程、通信工程专业的专业必修课程教材为主,同时包含一些反映学科发展前沿的本科选修课程教材和研究生教学用书。为了保证教材的出版质量,清华大学出版社不仅约请国内一流专家参与了丛书的选题规划,而且每本书在出版前都组织全国重点高校的骨干教师对作者的编写大纲和书稿进行了认真审核。

祝愿《高等院校信息与通信工程系列教材》为我国培养与造就信息与通信工程领域的高素质科技人才,推动信息科学的发展与进步做出贡献。

北京邮电大学  
陈俊亮

# 前 言

---

语音信号处理以语音为研究对象,涉及模式识别、心理学、生理学、语言学、人工智能、模式识别等多项研究领域,甚至还涉及说话时表情、手势等人的体态语言信息。由于语音是人们在日常生活中的主要交流手段,因此语音信号处理在现代信息社会中占有重要地位。语音信号处理的研究工作最早可以追溯到 19 世纪 70 年代,在 20 世纪得到了长足的发展,并在 20 世纪 90 年代,随着 IBM、Apple、AT&T、NTT 等著名公司为语音识别的实用化开发投以巨资,语音信号处理技术的应用掀起了热潮。

近年来,随着语音信号处理技术的日益成熟,出现了新的基于语音识别的应用方向,如语音拨号、呼叫中心、移动设备中的嵌入式命令控制、发音学习,以及基于关键词检出的口语会话系统等。随着语音信号处理技术在实际生活中的应用,它已经被广泛地接受和使用。由于语音比其他形式的交互方式具有更多的优势,因此这项技术已经越来越贴近人们的生活。目前,语音信号处理技术处于蓬勃发展时期,不断有新的产品被研制开发,市场需求逐渐增加,具有良好的应用前景。同时,本领域不断出现一些新的技术,如 HTK 的使用、语音识别的 Lattice 结构、语言学模型的自适应技术、基于 HMM 的语音合成技术、开集说话人识别中拒识的处理、最新的基于听觉的音频编码技术等,本书再版的目的是将这些新的技术融合到原来的语音信号处理中,将该学科一些新的发展动态介绍给读者,希望读者在学术思想上受到一些启发。

本书的内容涉及作者承担的国家自然科学基金项目的部分研究成果。可作为高等院校计算机应用、信号与信息处理、通信与电子系统等专业及学科的高年级本科生、研究生教材,也可供该领域的科研及工程技术人员参考。

本书的第 1、2、4 章由韩纪庆编写,第 3、6、9 章由张磊编写,第 4、5、8 章由郑铁然编写。韩纪庆负责全书的总体安排和审定。郑贵滨为书稿中的插图做了大量工作,在此表示感谢。

本书虽然是作者从事语音信号处理工作多年的理论与实践的结晶,但因作者水平有限、时间仓促,缺点和错误在所难免,敬请读者批评指正,提出宝贵意见。

作者

于哈尔滨工业大学

2013 年 1 月

# 目 录

<b>第 1 章 绪论</b> .....	1
1.1 语音信号处理的发展 .....	1
1.2 语音信号处理的应用 .....	9
1.3 语音信号处理的总体结构.....	11
参考文献 .....	12
<b>第 2 章 语音信号的声学基础及产生模型</b> .....	13
2.1 语音信号的产生.....	13
2.1.1 语音的发音器官 .....	14
2.1.2 语音的声学特征 .....	16
2.1.3 语音信号在时域和频域的表达 .....	18
2.1.4 汉语中语音的分类 .....	22
2.1.5 汉语语音的韵律特性 .....	24
2.2 语音信号的感知.....	25
2.2.1 听觉系统 .....	25
2.2.2 听觉特性 .....	27
2.2.3 掩蔽效应 .....	28
2.3 语音信号的线性产生模型.....	33
2.3.1 激励模型 .....	33
2.3.2 声道模型 .....	34
2.3.3 辐射模型 .....	35
2.4 语音信号的非线性产生模型.....	35
2.4.1 调频-调幅模型的基本原理.....	36
2.4.2 Teager 能量算子 .....	37
2.4.3 能量分离算法 .....	37
2.4.4 调频-调幅模型的应用 .....	39
参考文献 .....	41
<b>第 3 章 语音信号的特征分析</b> .....	43
3.1 语音信号数字化.....	44
3.1.1 语音信号的采样和量化 .....	44

3.1.2	短时加窗处理	47
3.2	语音信号的时域分析	49
3.2.1	短时能量分析	49
3.2.2	短时平均过零率	50
3.2.3	短时自相关函数和短时平均幅度差函数	52
3.2.4	端点检测和语音分割	55
3.3	语音信号的频域分析	57
3.3.1	滤波器组方法	57
3.3.2	傅里叶频谱分析	57
3.4	传统傅里叶变换缺点及时频分析的思想	61
3.4.1	信号的时频表示	62
3.4.2	不确定原理	64
3.5	Gabor 变换	66
3.6	小波变换在语音信号分析中的应用	68
3.6.1	小波的数学表示及意义	68
3.6.2	小波分析特点	70
3.6.3	小波变换的多分辨分析	72
3.6.4	小波变换在语音处理中应用	74
3.7	语音信号的同态解卷积	76
3.7.1	同态信号处理的基本原理	77
3.7.2	语音信号的复倒谱	78
3.7.3	避免相位卷绕的算法	81
3.7.4	基于听觉特性的 Mel 频率倒谱系数	85
3.8	语音信号特征应用	86
3.8.1	基音周期估计	86
3.8.2	共振峰的估计	92
	参考文献	95
<b>第 4 章</b>	<b>语音信号的线性预测分析</b>	<b>97</b>
4.1	线性预测的基本原理	97
4.2	线性预测方程组的解法	99
4.2.1	自相关法	100
4.2.2	协方差法	102
4.2.3	格型法	103
4.2.4	几种求解线性预测方法的比较	107
4.3	线性预测的几种推演参数	108
4.3.1	归一化自相关函数	108
4.3.2	反射系数	108
4.3.3	预测器多项式的根	109



4.3.4	LPC 倒谱 .....	110
4.3.5	全极点系统的冲激响应及其自相关函数 .....	110
4.3.6	预测误差滤波器的冲激响应及其自相关函数 .....	111
4.3.7	对数面积比系数 .....	111
4.4	线谱对分析法 .....	111
4.4.1	线谱对分析的原理 .....	112
4.4.2	线谱对参数的求解 .....	113
4.5	感知线性预测 PLP 系数 .....	113
	参考文献 .....	114
<b>第 5 章</b>	<b>语音编码</b> .....	<b>115</b>
5.1	波形编码 .....	116
5.1.1	均匀量化 PCM .....	116
5.1.2	非均匀量化 PCM .....	116
5.1.3	自适应量化 PCM .....	117
5.1.4	差分脉冲编码 .....	118
5.1.5	自适应差分脉冲编码 .....	120
5.1.6	增量调制和自适应增量调制 .....	123
5.1.7	子带编码 .....	124
5.1.8	自适应变换域编码 .....	126
5.2	参数编码和混合编码 .....	127
5.2.1	参数编码 .....	127
5.2.2	基于全极点语音产生模型的混合编码 .....	133
5.2.3	基于正弦模型的混合编码 .....	146
5.3	极低速率语音编码技术 .....	151
5.3.1	400bps~1.2Kbps 的声码器 .....	151
5.3.2	识别合成型声码器 .....	152
5.4	语音编码器的性能指标和质量评测方法 .....	153
5.4.1	编码速率 .....	154
5.4.2	顽健性 .....	154
5.4.3	时延 .....	155
5.4.4	计算复杂度和算法的可扩展性 .....	155
5.4.5	语音质量及其评价方法 .....	156
5.5	语音编码国际标准 .....	158
5.6	感知音频编码 .....	158
5.6.1	感知编码的一般框架 .....	159
5.6.2	心理声学模型 .....	160
5.6.3	常用的感知编码标准 .....	162
	参考文献 .....	164

<b>第 6 章 语音识别</b> .....	165
6.1 概述 .....	165
6.2 基于矢量量化的识别技术 .....	167
6.2.1 K-means 矢量量化算法 .....	167
6.2.2 LBG 算法 .....	168
6.3 动态时间归正的识别技术 .....	169
6.3.1 DTW 基本原理 .....	169
6.3.2 模板训练算法 .....	171
6.4 隐马尔可夫模型技术 .....	173
6.4.1 HMM 基本思想 .....	173
6.4.2 HMM 基本算法 .....	176
6.4.3 HMM 算法实现中的问题 .....	180
6.4.4 关于 HMM 训练的几点考虑 .....	186
6.5 连接词语音识别技术 .....	190
6.5.1 连接词识别问题的一般描述 .....	191
6.5.2 二阶动态规划算法 .....	192
6.5.3 分层构筑方法 .....	193
6.6 大词表连续语音识别中的声学模型和语言学模型 .....	197
6.6.1 声学模型 .....	199
6.6.2 统计语言学模型 .....	206
6.6.3 统计语言学模型平滑技术 .....	208
6.6.4 语言学模型自适应技术 .....	212
6.7 大词表连续语音识别中的解码技术 .....	213
6.7.1 图的基本搜索算法 .....	214
6.7.2 面向语音识别的搜索算法 .....	216
6.8 大词表连续语音识别后处理技术 .....	222
6.8.1 语音识别中间结果的表示形式 .....	222
6.8.2 错误处理 .....	224
6.8.3 最小字错误率解码方法 .....	226
6.9 基于 HMM 的自适应技术 .....	231
6.9.1 基于 Bayesian 理论的自适应方法 .....	231
6.9.2 基于变换的自适应方法 .....	232
6.10 关键词检出技术 .....	234
6.10.1 问题描述 .....	235
6.10.2 关键词检出系统的组成 .....	237
6.10.3 垃圾模型建模方法 .....	237
6.10.4 语音解码器的设计 .....	239
6.10.5 关键词确认过程 .....	240

6.10.6	关键词检出系统性能优化	241
6.11	语音识别的应用技术	241
6.11.1	语音信息检索	241
6.11.2	发音学习技术	243
6.11.3	基于语音的情感处理	249
6.11.4	网络环境下的语音识别	252
6.11.5	嵌入式语音识别技术	255
6.12	HTK 工具介绍	256
6.12.1	数据准备阶段	259
6.12.2	模型训练阶段	263
6.12.3	识别阶段	272
	参考文献	273
<b>第 7 章</b>	<b>说话人识别</b>	<b>279</b>
7.1	概述	279
7.2	说话人识别的特征选取	282
7.2.1	特征参数的评价方法	283
7.2.2	说话人识别系统中常用的特征	284
7.3	说话人识别的主要方法	285
7.3.1	与文本有关的识别方法	285
7.3.2	与文本无关的识别方法	286
7.3.3	文本提示型的识别方法	297
7.4	阈值的选取	298
7.5	得分规整	299
7.5.1	零规整(zero normalization)	300
7.5.2	测试规整(test normalization)	300
7.5.3	说话人自适应的测试规整	301
7.5.4	TZ-norm	302
7.5.5	H-norm	302
7.5.6	C-norm	303
7.6	引入区分判别模型的说话人识别	303
7.6.1	SVM	303
7.6.2	基于 SVM 的说话人识别	306
7.6.3	基于 GMM 得分的 SVM 说话人识别	307
7.6.4	基于 GMM 均值超矢量的 SVM 说话人识别	308
7.7	复杂信道下的说话人识别	309
7.7.1	特征映射	310
7.7.2	说话人模型合成	311
7.7.3	扰动属性投影	312

7.7.4	联合因子分析	312
7.8	说话人识别中有待解决的问题	313
	参考文献	315
<b>第8章</b>	<b>顽健语音识别技术</b>	<b>317</b>
8.1	概述	317
8.2	影响语音识别性能的环境变化因素	317
8.3	噪声环境下的顽健语音识别技术	319
8.3.1	基于语音增强的方法	320
8.3.2	通道畸变的抑制方法	325
8.3.3	基于模型的补偿方法	330
8.4	变异语音识别方法	344
8.4.1	变异语音的分析	345
8.4.2	变异语音的分类	346
8.4.3	变异语音的识别	349
	参考文献	354
<b>第9章</b>	<b>语音合成</b>	<b>357</b>
9.1	语音合成的基本原理	358
9.2	参数合成方法	361
9.2.1	线性预测合成方法	362
9.2.2	共振峰合成方法	363
9.3	波形拼接合成技术	369
9.3.1	TD-PSOLA 算法	370
9.3.2	FD-PSOLA 算法	373
9.4	汉语按规则合成	375
9.4.1	韵律规则	376
9.4.2	多音节协同发音规则合成	382
9.4.3	轻声音节规则合成	383
9.4.4	儿化音节的规则合成	384
9.5	基于 HMM 的参数化语音合成技术	385
9.5.1	基于 HMM 参数语音合成系统的训练	385
9.5.2	基于 HMM 参数语音合成系统的合成阶段	390
	参考文献	393

# 第 1 章 绪 论

语言是人类最重要的交流工具,它自然方便、准确高效。随着社会的不断发展,各种各样的机器参与了人类的生产活动和社会活动,因此改善人和机器之间的关系,方便人对机器的操纵就显得越来越重要。随着电子计算机和人工智能机器的广泛应用,人们发现,人和机器之间最好的通信方式是语言通信。而语音是语言的声学表现形式;要使机器听懂人的语言并能使用人类的语言进行表达,需要做很多工作,这就是研究了几十年的语音识别和语音合成技术。而随着移动通信的迅猛发展,人们可以随时随地通过电话进行交流,其中语音压缩编码技术发挥着重要的作用。上述这些应用领域构成了语音信号处理技术的主要研究内容。

语音信号处理是语音学与数字信号处理技术相结合的交叉学科,它和认知科学、心理学、语言学、计算机科学、模式识别和人工智能等学科联系紧密;语音信号处理技术的发展依赖于这些学科的发展,而语音信号处理技术的进步也会促进这些学科的进步。

## 1.1 语音信号处理的发展

语音信号处理的研究工作最早可以追溯到 1876 年贝尔发明的电话,它首次完成了用声电—电声转换来实现远距离传输语音的技术。1939 年 Dudley 研制成功了第一个声码器,从此奠定了语音产生模型的基础,这一工作在语音信号处理领域具有划时代的意义。1947 年贝尔实验室发明了语谱图仪,将语音信号的时变频谱用图形表示出来,为语音信号的分析提供了一个有力的工具。1948 年美国 Haskins 实验室研制成功“语图回放机”,它把手工绘制在薄膜片上的语谱图自动转换为语音,可以进行语音合成。共振峰合成方法就是源于这一思想。

对语音识别而言,它的研究相对较晚,起源于 20 世纪 50 年代。语音识别技术的根本目的是研究出一种具有听觉功能的机器,能接受人类的语音,理解人的意图。由于语音识别本身所固有的难度,人们提出了各种限制条件下的研究任务,并由此产生了不同的研究领域。这些领域包括:按说话人的限制,可分为特定说话人语音识别和非特定说话人语音识别;按词汇量的限制,可划分为小词汇量、中词汇量和大词汇量的识别;按说话方式的限制,可分为孤立词识别和连续语音识别等。最简单的研究领域是特定说话人小词汇量孤立词的识别,而最难的则是非特定说话人大词汇量连续语音的识别。

1952 年贝尔实验室的 Davis 等人研制了特定说话人孤立数字识别系统。该系统利用每个数字元音部分的频谱特征进行识别。1956 年 RCA 实验室的 Olson 等人也独立地研制出 10 个单音节词的识别系统,系统采用从带通滤波器组获得的频谱参数作为语音的特征。1959 年 Fry 和 Denes 等人尝试构建音素识别器来识别 4 个元音和 9 个辅音,采用频谱分析

和模式匹配来进行识别决策,其突出贡献在于,使用了英语音素序列中的统计信息来改进词中音素的精度。1959年MIT林肯实验室的Forgie等人,采用了声道的时变估计技术对10个元音进行识别。

20世纪60年代初期,日本的很多研究者开发了相关的特殊硬件来进行语音识别,如东京无线电研究实验室Suzuki等人研制的通过硬件来进行元音识别的系统。在此期间,开展的很多研究工作对后来近二十年的语音识别研究产生了很大的影响。RCA实验室的Martin等人在20世纪60年代末开始研究语音信号时间尺度不统一的解决办法,开发了一系列的时间归正方法,明显地改善了识别性能。与此同时,当时苏联的Vintsyuk提出了采用动态规划方法来解决两个语音的时间对准问题。尽管这是动态时间弯折算法(dynamic time warping,DTW)的基础,也是连接词识别算法的初级版,但Vintsyuk的工作并不为学术界的广大研究者所知,直到20世纪80年代大家才知道Vintsyuk的工作,而这时DTW方法已广为人知。

值得一提的是20世纪60年代中期,斯坦福大学的Reddy就开始尝试用动态跟踪音素的方法来进行连续语音的识别。后来Reddy加入到卡内基梅隆大学,多年来在连续语音识别上开展了卓有成效的工作,直至现在仍然在此方面居于领先地位。

20世纪70年代之前,语音识别的研究特点是以孤立词的识别为主。20世纪70年代语音识别研究在多方面取得了诸多的成就,在孤立词识别方面,日本学者Sakoe给出了使用动态规划方法进行语音识别的途径——DTW算法,它是把时间归正和距离测度计算结合起来的一种非线性归正技术。这是语音识别中一种非常成功的匹配算法,当时在小词汇量的研究中获得了成功,从而掀起了语音识别的研究热潮。Itakura利用语音编码中广泛使用的线性预测编码(linear predictive coding,LPC)技术,通过定义基于LPC频谱参数的合适的距离测度,成功地将其扩展到语音识别中。以IBM为首的一些研究单位还着手开展了连续语音识别的研究,AT&T的贝尔实验室也开展了一系列非特定说话人语音识别方面的研究工作。

应该指出的是,20世纪70年代起人工智能技术开始被引入到语音识别中来。美国国防部的高级研究规划局(Advanced Research Projects Agency,ARPA)组织了有卡内基梅隆大学等五个单位参加的一项大规模语音识别和理解的研究计划,当时专家们认为:要使语音识别研究获得突破性进展,必须让计算机像人那样具有理解语言的智能,而不必过多地在孤立词识别上下功夫。在这个历时五年的庞大的研究计划中,最终在语言理解、语言的统计模型等方面积累了经验,其中卡内基梅隆大学完成的Hearsay-II和Harpy两个系统效果最好。在这两个系统中,引用了“黑板模型”来完成底层和顶层之间不同层次的信息交换和规则调用,成为以后其他专家系统研究工作中的一种规范。但从整体上看,这个计划并没有取得突破性的进展。

20世纪70年代末80年代初,Linda、Buzo、Gray等提出了矢量量化(vector quantization)码本生成的方法,并将矢量量化技术成功地应用到语音编码中,从此矢量量化技术不仅在语音识别、语音编码和说话人识别等方面发挥了重要的作用,而且很快推广应用到其他领域。这一时代,语音识别的研究重点之一是连接词识别,典型的工作是进行数字串的识别。研究者提出了各种连接词语音识别算法,大多数工作是基于对独立的词模板进行拼接来进行匹配的方法,如两级动态规划识别算法、分层构筑(level building)、帧同步(frame synchronous)

分层构筑方法等。这些方法都有各自的特点,广泛用于连接词识别当中。

20世纪80年代开始,语音识别研究的一个重要进展,就是识别算法从模式匹配技术转向基于统计模型的技术,更多地追求从整体统计的角度来建立最佳的语音识别系统。隐马尔可夫模型(hidden markov model, HMM)技术就是其中的一个典型;尽管开始的时候仅有较少的单位采用这种模型,但由于该模型能很好地描述语音信号的时变性和平稳性,具有把从声学—语言学到句法等统计知识全部集成在一个统一框架中的优点,因此从20世纪80年代起,它被广泛地应用到语音识别研究中。直到目前为止, HMM方法仍然是语音识别研究中的主流方法。HMM的研究使大词汇量连续语音识别系统的开发成为可能。20世纪80年代末,美国卡内基梅隆大学用VQ/HMM实现了997词的非特定人连续语音识别系统SPHINX,这是世界上第一个高性能的非特定人、大词汇量、连续语音识别系统。此外, BBN的BYBLOS系统, 林肯实验室的识别系统等也都具有很好的性能。这些研究工作开创了语音识别的新时代。

从20世纪80年代后期和90年代初开始, 神经网络(artificial neural network, ANN)的研究异常活跃, 并且被应用到语音识别的研究中。进入20世纪90年代后, 相应的研究工作在模型设计的细化、参数提取和优化, 以及系统的自适应技术等方面取得了一些关键性的进展, 使语音识别技术进一步成熟, 并且出现一些很好的产品。许多发达国家, 如美国、日本、韩国, 以及IBM、Microsoft、Apple、AT&T、NTT等著名公司都为语音识别系统的实用化开发研究投以巨资。

近年来, 语音识别研究工作更趋于解决在真实环境应用时所面临的实际问题, 这可从作为国际语音识别研究热点风向标的NIST(national institute of standards and technology)评测情况反映出来: 其评测的语音类型已从最初的朗读语音到广播语音, 再到后来的交谈式电话语音(conversational telephone speech), 发展到目前真实场景的会议语音。相对于广播语音, 交谈式电话语音增加了相应的难度, 具体表现在: 发音多为自发的口语语音, 存在着大量的不流利(如犹豫词、重复、更正等)现象, 同时, 语音内容和词汇的随机性明显增加。此外, 针对实际的电话线路, 噪声的影响较大。2002年, 美国国防部先进研究项目局(Defense Advanced Research Projects Agency, DARPA)提出了一个“EARS-Effective, Affordable and Reusable Speech-to-text(高效低耗可重用语音文字转化)”的项目, 把NIST的语音评测推到了又一个新的时代——丰富的语音文本(rich transcription, RT)转写, 其要求不仅将语音所对应的文字显示出来, 而且要将语音中的其他丰富信息, 如文字之间的标点符号、句词之间的停顿、说话人等也能同时识别出来。从2004年的评测结果看, 对广播语音和电话语音的词错误率(word error rates, WERs)已分别下降到10%和15%以下。从2005年起, NIST评测的语音类型转变为英语会议语音, 包括磋商式会议(conference meeting)和演讲式会议(lecture meeting), 其特点是研究真实会议场景中多人多方对话时的口语语音识别。相对于交谈式电话语音, 会议语音又增加了相应的难度, 表现在: 必须解决会议场景中处于不同位置上说话人语音数据的有效采集问题, 以及在多人交谈相互语音有少部分重叠时各自语音的分离问题。为此, NIST评测中开始提供采用远离用户, 且处于空间上多个位置、摆放形式多样的多麦克风或麦克风阵列采集来的现场数据作为评测的语料。从2007年进行的评测结果看, 会议语音的词错误率在40%~50%之间。2009年的评测内容基本与2007年相同, 所不同的是仅进行磋商式会议语音的评测, 同时为各个测试任务定

义了视频和音视频的输入条件。

目前无论从 NIST 评测的内容看,还是欧美发达国家的关注点看,研究真实场景中多人多方对话时的口语语音识别是当前语音识别的研究热点之一。从处理口语语音与朗读语音的方法看,其不同之处在于声学模型的自适应(acoustic adaptation)和发音词典自适应(lexicon adaptation)方面。声学模型自适应常采用基于最大似然线性回归(maximum likelihood linear regression, MLLR)和最大后验概率(maximum a posteriori, MAP)的方法。这两种方法是当前最为有效的自适应方法,许多新的自适应方法都是从二者中派生出来的。发音词典自适应常采用发音变化建模(pronunciation variation modeling)相关技术,主要研究由说话方式、语速、口音等带来的影响。

口语语音识别的另一个挑战是缺乏建立在大量口语文本语料之上良好的语言模型。朗读语音识别器所使用的统计语言模型,实际上都要依赖于大规模的训练语料,但是同样量级的口语语言的文字脚本还难以实现。口语语音中的不连贯进一步增加了语言模型估计的难度。目前研究者正致力于多种口语语言模型的建模方法研究。

当前语音识别研究的另一个趋势是,不再只单纯地关注大词表连续语音识别的精度,而是从实际的应用角度出发,积极探索机器对人类的语音进行感知与理解的途径和方法。而从整个计算领域的发展趋势看,近年的研究热点之一是普适计算,计算的模式与物理位置也正从传统的桌面方式逐步向以嵌入式处理为特征的无处不在的方式发展,比较典型的是移动计算方式。因此对语音处理而言,探讨在典型的移动方式下的语音感知与理解机制,实现能根据用户的语音内容及所处的音频场景,并借助其他辅助信息(如地理位置、时间等)自主地感知和理解用户的意图及情感倾向,从而提供更智能化、人性化的人机交互手段,具有重要的理论意义与现实意义。同时,随着网络技术和移动计算技术的迅速发展,出现了网络环境下的语音识别技术、嵌入式和计算资源有限时的语音识别技术、语种识别技术、基于语音的情感处理技术等一些新的研究方向。

在国内,20世纪50年代末就有人尝试用电子管电路进行元音识别,而到了70年代才由中科院声学所开始了计算机语音识别的研究。在此之后,有关专家也开始撰文介绍这方面的工作。从20世纪80年代开始,很多单位陆续参加到这一行列中来,它们纷纷采用不同的方法,开展了从最初的特定说话人中、小词汇量孤立词识别,到非特定说话人大词汇量连续语音识别的研究工作。20世纪80年代末,以汉语全音节识别作为主攻方向的研究已经取得了相当大的进展,一些汉语语音输入系统已向实用化迈进。四达技术开发中心、星河公司等相继推出了相应的实际产品。清华大学、中科院声学所在无限词汇的汉语听写机的研制上获得成功。20世纪90年代初,四达技术开发中心又与哈尔滨工业大学合作推出了具有自然语言理解能力的新产品。在国家“863计划”支持下,清华大学和中科院自动化所等单位在汉语听写机原理样机的研制方面开展了卓有成效的工作。北京大学在说话人识别方面也做了大量的工作。

近年来,随着改革开放的不断进行,我国的国际地位与日俱增,汉语语音识别越来越受到重视,国外很多著名的公司,如 Microsoft、IBM、Motorola、Intel 等都在国内设立了研发机构,并且都将汉语语音识别作为主攻方向之一。IBM公司于1997年推出了汉语连续语音识别系统 ViaVoice,输入速度平均每分钟可达150字,平均最高识别率达到95%,并具有“自我”学习的功能。2000年发布的 ViaVoice 千禧版,用户可以通过语音导航到电脑桌面



及浏览网页。1998 年微软(Microsoft)投资 8000 万美元在中国筹建微软中国研究院(2000 年更名为微软亚洲研究院),开发的重点方向之一就是语音识别。1998 年 Intel 提出了基于 Intel 架构发展语音技术的构想,向软件开发厂商提供包括信号处理库、识别库、图像处理库在内的高性能语音函数库支持,1999 年 Intel 又和 L&H 公司合作,推出语音识别软件开发包 Spark3.0,其中包括 Spark 语音识别引擎和软件开发工具箱。微软也推出了基于 .net 的语音识别引擎。国内一些著名企业也投入大量资金开始资助语音识别方面的研究,如盛大创新院、比亚迪公司等。

尽管语音识别技术研究已经取得了很大的成绩,但到目前为止离广泛的应用尚存在距离。很多因素影响着语音识别系统的性能,如实际环境中的背景噪声、传输通道的频率特性、说话人生理或心理情况的变化,以及应用领域的变化等都会导致语音识别系统性能的下降,甚至不能工作。研究语音识别系统顽健性(robustness)问题受到了研究者的广泛重视,国内外很多单位都开展了大量的工作。但到目前为止,所做的工作大都是针对某一种或两种影响因素进行补偿的研究,综合考虑各种影响因素补偿方法的研究还很少。

语音识别通常是指能识别出相应的语音内容,除此之外,它还有一种特殊的形式——说话人识别。说话人识别不必识别出语音信号的具体内容,而只要鉴别出该语音是哪个说话人发出的即可。从实现的技术手段上看,说话人识别和语音识别一样,都是通过提取语音信号的特征,并建立相应的参考模板来进行分类判断。说话人识别问题,最初是在第二次世界大战期间,美国国防部向贝尔实验室提出的课题。目的是根据窃听到的电话语音来判断说话人是哪一位德军高级将领,这对分析当时的德军战略部署具有重要的意义。该项目持续进行了三年,但并未达到预期的目的。

说话人识别研究的早期工作,主要集中在人耳听辨实验和探讨听音识别的可能性方面。随着语音识别研究的不断深入,说话人识别研究也获得了突飞猛进的发展。语音识别中很多成功的技术,如矢量量化(vector quantization, VQ)、隐马尔科夫模型等都被应用到说话人识别中。

20 世纪 90 年代, Rose 等提出了单状态的 HMM,即后来的高斯混合模型(gaussian mixture model, GMM),它是一个顽健的参数化模型。Matsui 等比较了基于连续 HMM 的说话人识别方法,发现识别率是状态和混合数的函数。同时,识别率与总的混合数有很强的关联性,但与状态数无关。这意味着不同状态间的转移信息对文本无关的说话人系统而言是没有作用的,因此,高斯混合模型 GMM 得到了与多状态 HMM 几乎相同的识别性能。正是上述工作,使得 GMM 建模方法在说话人识别研究中得到了越来越多的重视。特别是 Reynolds 等对高斯混合模型 GMM 以及通用背景模型(universal background model, UBM)的详尽介绍后,由于 GMM-UBM 具有简单有效,以及具有较好的顽健性等特点,迅速成为当今与文本无关的说话人识别中的主流技术,并由此将说话人识别技术带入了一个新的阶段。20 世纪 90 年代另一项重要的研究工作是,针对说话人确认中,说话人自身的似然度的得分变异的规整技术,出现了很多关于得分规整的算法,比较典型的如基于似然比(likelihood ratio)和后验概率(a posteriori probability)的技术。为了降低计算规整算法的计算复杂性,相继出现了群组说话人(cohort speakers)等方法。与此同时,说话人识别技术与其他语音研究方向的结合更加密切,比如针对对话/会议中包含多人的说话人分割与聚类技术,音频元数据(metadata)的检索研究等也得到了很多研究人员的关注。