

# 南京航空学院

# 研究生硕士学位论文

研究生姓名 胡磊

专业 计算机应用

研究方向 微机系统结构、应用

指导教师 夏振华高级工程师

一九八五年一月

# VIMS 语言输入微机系统的设计与实现

## 摘要

本文介绍了一种与说话者有关的孤立单词语音识别系统。该系统用带通滤波和线性时间规划的方法提取语音特征，用改进最近邻法进行识别。实验结果表明，该系统已达到了平均 97% 的识别率。文中，介绍了识别系统部分的原理及实现方法；着重分析了最近邻法识别的分类特性及其在实际应用中的限制，在此基础上提出了一种有效的学习方法；此外，还详细地讨论并用多门限法解决了在确定语音首尾时移过界的問題。最后，本文给出了一个应用实例——在 16K 扩展 BASIC 上实现的语音控制 BASIC。

DESIGN AND REALIZATION  
OF  
A VOICE INPUT MICROCOMPUTER SYSTEM (VIMS)

Abstract

In this paper, a speaker-dependent isolated word recognition system VIMS is described. This system extracts the speech feature by bandpass filtering and linear time normalizing, and recognizes the input utterances, using the modified nearest neighbor rules. The experimental result shows that the system has reached an average accuracy of 97%. The principle and method for accomplishing every part of the system are introduced. Especially analysed are the performances of the nearest neighbor rules and some of its limitations in practical use, on the basis of which an effective learning method is proposed. Besides, the problems encountered in determining the endpoints of utterances are discussed in detail and solved with the help of the so-called Multi-Threshold method. Finally, an application example---Voice-controlled BASIC realized on 16K EXTENDED BASIC, is given.

# 目 录

页数

## 摘要

第一章. 引言	1
第二章. 设计原理与子系统结构	7
2.1. 概述	7
2.2. 设计原理及流程	8
2.3. 子系统结构和识别过程	11
第三章. 噪声环境中语音信号的检测	15
3.1. 判首尾时存在的问题	15
3.2. 判别特征的描述	18
3.3. 内存的组织	23
第四章. 数据压缩和参数规定	25
4.1. 参考规定概念	25
4.2. 规定和压缩	27
第五章. 识别	31
5.1. 识别方法简介	31
5.2. NN rule 简介	32
5.3. 识别过程中学习	35
第六章. 应用实例 — VBASIC	42
6.1. VBASIC 的执行过程	42
6.2. 修改方法	44
结束语	46
致谢	47

# 目次(續)

三子

參政文獻	48
附录	50

# 第一章 引言

语音识别技术是人工智能的一个分支。自从 1952 年第一台语音识别装置问世到现在三十多年间，许多研究人员认为在这方面作出了极大的努力，也正是因为他们的辛勤劳动，在语音识别这个领域内，才出现了不少可喜的成果：涌现出了一种识别的新技术。

语音识别系统无论用在什么场合，都显示着它那独特的优点：

(1) 语言，是人们通讯的最自然、最有效的方式，因此，用户勿需接受什么特殊的训练就能熟练地使用语音识别装置。

在我国，电子计算机的应用正在全面发展，形势要求越来越多的人学会使用计算机，——从机关到地方、从大学到中学、从工厂到研究所，都是如此。对计算机专业人员来说，摆在我们面前的任务是如何使计算机的操作更加简便。我国的用户，可以说，大多数不懂得英文，这更给他们操作计算机带来了困难。因此，在我国大力开发语音识别技术就更显得有必要了。

(2) 语音识别装置直接与用户打交道的只有一张嘴，所以它的使用可以不受什么客观条件的限制：既灵活又方便。人们还可以用它，通过无线收发装置，对远距离的设备进行遥控。例如，人们可以用语言对正在从事危险工作的机口人发出指令，让他尚未完成我们想做而

又无法做到的了；又如，在狭小的飞机座舱里，人们无法再设置计算机所需要的键盘和显示器，飞行员也只能用他的手柄与空中设备或起落控制设备进行通讯，在这刻的手段中，语音识别是强有力的一种。

(3) 使用语音识别装置时，只要把麦克风固定在使用者的嘴边(例如使用头式麦克风)，因此，可以把使用者的手、眼睛放起来，同时进行其它的工作。有时，人的手眼忙于一项重要的工作，而同时又想口授一些命令来控制其它机口运行，这时，语音识别更能够发挥它的作用。

此外，语音识别还允许用户同时控制多台机口；同时与机口和其它的人进行通讯。

语音识别的优点很多，应用范围也很广；特别是随着研制机口人和第三代计算机的需要，它更引起了人们的关注，研究者也越来越多。最近，在日本召开的一次有关第三代计算机的学术会议上，欧洲共同体提出的关于“外接口”的设想，其中就包括语音识别这个研究方向。

语音识别系统有孤立单词识别系统和连续语音识别系统两大类，每一类又可分为说话有关(Speaker-dependent)与说话无关(Speaker-independent)两种。目前，虽然有很多人在研究连续语音识别系统；然而，据我们所知，这些系统还只处在连续数字的阶段，不能进行识别非数字语言，但效果不好。说话有关系统要求每个用户在使用前先对系统进行训练，而说话无关系统

则勿需这样做。但语音有系统而语句内容比较灵活；语音无系统由于首先要对各种人的大易的发言样本进行处理，所以语句内容不易变动。

语音识别是模式识别理论的一个应用，一个语音识别子系统也是一个完整的模式识别子系统，它主要包括三个部分：特征口、特征提取口、分类口。在语音识别这个领域中讨论得最多的是特征提取口，即抽取特征和压缩数据的方法。1952年，Bell电话实验室的Davis等人研制出了第一台孤立单词语音识别口（语音有关），它能识别0~9十个数字。该识别口使用900 Hz 上下的两个通道（带通滤波口）来分析语音信号，立刻提取两个通道的过零率，或过零次数（axis-crossing）作为特征，然后使用所谓“最佳匹配法”（best-match）来作出最后的分类决策[1]。数年之后，Dudley 和 Balashek 开发了一个称为“Audrey”的识别子系统，用十一个频率通道来提取语音的频谱特征，把所得的特征向量与存贮的模式进行比较，最后求得分类结果[2]。应用计算机来进行语音识别大概出现在1959~1960年之间，那时，Denes 和 Mathews 提出了时间规范化（time normalization）的概念，并把它应用于实践，取得了成功[3]。

此后，这种特征提取方法层出不穷——多种子系统，有孤立单词的，也有连续语音的，有语音有关的，也有语音无关的，数目之多令人眼花缭乱。文献[4]对此作了详

细而介绍。这种方法又易程度也各不相同，也较起来也有利弊；但它中间，人们比较重视的方法可能要推萃通滤波口组法和LPC法了(Linear Predictive Coding—线性预测编码)。比较这两种方法，LPC法的特征要也带通滤波法的特征准确，但是，LPC的运算量很大，单机程序实现，在目前的条件下还难以做到实时。

国内，这方面的工作开始得也较晚，但经过努力也取得了不少成果。据我们所知，中科院声学所[5]、[6]、清华电子[7]、沈阳自动化所，使用带通滤波口的方法，先从频谱分析仪上实现了语音有关的孤立单元识别系统。另外，西工大在平板机上实现了这样的系统；沈阳自动化所还对语音无关系统进行了实验。此外，其它的方法也正在被大力开发。

目前，我们正在着手进行人工智能和机器人的研究工作，至于语音识别的重要性，并且结合教研室的科研计划，我们选择了这一课题。考虑到客观条件的限制以及科研任务的要求，设计本课题的宗旨是：在 Cromemco 系统Ⅲ现有的硬件条件下，使用相对简单的电路来进行实验，并且要求识别系统具有较好的实时性。因此，我们选择了通滤波口的方法，不用频谱分析仪，只用 孤立元件设计了四个低频通道来抽取语音的特征，一个高频通道来检测辅音。高频通道的引入使系统检测辅音的能力大大提高了，克服了同类系统难以检测弱摩擦音的缺点；

但音频通道的引入会严重影响子系统的抗干扰性能，为此，本文分析了语音信号和干扰噪音的特点及二者的区别，在此基础上使用了多次判头尾的双门限和自适应门限结合的方法，采用循环取模的存贮四结构，较好地解决了这个问题，提高了子系统的抗干扰性。

从这三种不同的识别系统可以看出，人们对特征提取的研究也较多，而对分类口，即识别方法的研究却很少。诚然，特征提取是分类辨别的基础，特征提取口的好坏会直接影响到整个识别系统的性能；然而，我们也不能忽视分类方法的重要性。我们知道，虽然采用相同的特征提取方法，但分类方法不同，系统的性能也不会相同。本文分析了最近邻法的性质，论以了某些子系统在使用最近邻法时的不足，提出了一种识别过程中的学习的方法。最后，作为语音识别系统的一个应用，本文在 Cromemco 16K 扩展 BASIC 的基础上实现了语音控制 BASIC —— VBASIC (Voice BASIC)，并介绍了 VBASIC 编译程序的结构，使用方法，和实现过程中遇到的问题及解决的方法。

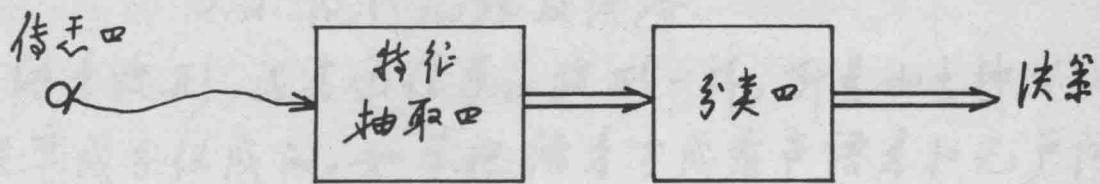
本文的第二章详细介绍了子系统的设计原理、整个识别子系统的结构及子系统的识别过程，另外，还给出了本文所采用的放大、滤波口的原理图。第三章讨论怎样确定语音信号首尾的问题，分析了语音与背景噪音的共同性和二者的区别，最后，根据语音与噪音的各自特点，给出了判别算法。第 9 章给出了参数修正和数据压缩的概念和方法。

第五章详细地分析了最近卸载时的分类属性，以及使用最近卸载时存在的问题，并给出了解决的方法。最后，在本章的第三章，我们给出了一个应用实例——VBASIC，并且介绍了修改 BASIC 编译程序的方法及修改程序。

## 第二章 设计原理与系统结构

### 2.1. 概述

语音识别的过程也是一个模式识别的过程。从理论上说，这个过程可分为两个步骤，一个是特征抽取，另一个是模式分类。但是，为了抽取各种模式的特征，还必须要有一个传感器。这样，我们就可以得到模式识别子系统的一个抽象模型（分类模型），它包括三个部分：传感器、特征抽取口、分类口，它们的连接关系如图[1]所示。



图[1]：模式识别子系统的一个抽象模型

传感器的作用是检测、输入，同时把输入信号转换成适合于机口处理的形式；特征抽取口（也称接收口、特征滤波口、属性检测口、预处理口等）是紧接着传感器的一个环节，它从传感器送来的数据中抽取有关的信息，作为反映该模式的特征。每类模式都有其特征反映该类模式本身共有的特征，它也是一模式区别于它模式的特殊标志。最后，分类口使用这些标志，根据一定的策略作出最后的决策，即把对应的模式分到它应属的类别中去。

语音识别中，传感器是话筒（麦克风），它把声音变为电信号送给后台的两个一部分：特征提取口和分类口。从理论上讲，特征提取口和分类口之间没有以硝的分界线。

一个理想的特征抽取口可以使分类口的工作变得简单，同样，一个全纯的分类口也将无求于特征抽取口；然而，我们将会看到，在实际问题中，这种区分还是有必要的，这是因为特征抽和问题是分类问题更依赖于实际问题本身，区分可以使设计工作更为单一。我们知道，一个适合于各必须识别的特征抽和口会对语音识别毫无用处。

本文根据深思的要求，使用通道滤波口征系数和语音信号的特征。

## 2.2. 设计原理及线路

语音波形，与其它信号的波形一样，都是由多种不同的频率成分组成的。如单起语音分成有声语音和无声语音两大类的话，从频谱图中可以看到（附录[1]中图[1]~[4]），有声语音的低频部分也较丰富，波形具有明显的周期性（这可以从附录[2]的波形图中看到）；从频谱图中还可以看到在频率取某几个位时谱线具有明显的峰位，这几个频率位数是声道的谐振点，通常称之为共振峰。一个元音（其它浊辅音也有相似的性质）有好几个共振峰，它们分布的频率范围大致在  $300 \sim 4000$  Hz 之间。一般说来，前几个——或 3 个——共振峰对区分元音也比较有用。元音不同，共振峰也不同；元音相同，共振峰也不尽相同。例如，附录各 [1]、[2] 和 [3]、[4] 的频谱图共振峰差别较大，而图 [1] 和 [2]、图 [3] 和 [4] 差别相似（注意纵座标的不同）。目前已经公认，共振峰——频谱色谱的一系列峰位，是区别不同元音的主要

特征号[8]。下表给出的汉语(普通话)十元音的前两个共振峰可帮助说明这一性质。

表：汉语元音的共振峰(引自[8])

元音	ɑ	ɔ	e	i	u	ü	ɛ	ər	ɪ	ɿ
$F_1(H_z)$ 男	905	631	598	337	377	341	587	662	448	426
	1061	772	853	394	491	415	670	797	511	522
$F_2(H_z)$ 男	1236	1113	1134	2380	641	2132	2122	1578	1993	1498
	1464	1131	1389	3098	836	2660	2428	1763	2381	1874

那些无声语音，它们没有以上的共振峰，谱线较平，而上覆盖的频带很宽，高频分量(5K以上)极其丰富，见附录各[5]。

以上的分析是设计带通滤波器组——频谱分析器的依据，工作原理可大致表述如下：

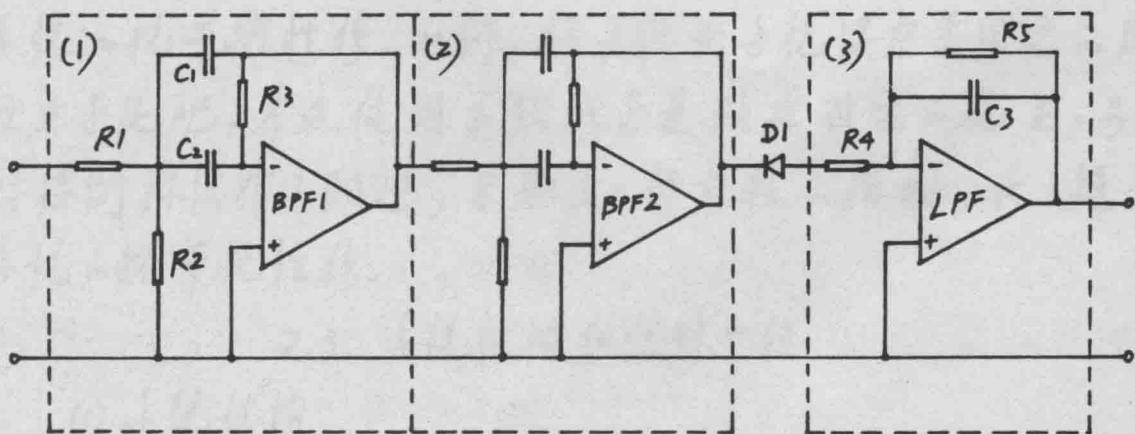
语言信号输入后进入由一维带通滤波电路组成的频谱分析器，滤波器的数目和频率分割均可有较大不同，一般可由10~16个滤波器覆盖各种不同语音信号的频率范围。通过带通滤波器对语音信号进行频率分析以后，再经过一个整流和低通滤波电路(截止频率 $f_c = 25 Hz$ )加以平均，至分析器的输出端便得到了语音信号的频谱色谱。

A/D转换器对输入的频谱色谱进行采样，将得到的数字量送入计算机进行处理(数据压缩和分类)。

清华大学用一线数字式动态频谱分析仪完成了上述功能[7]，在200~4000Hz这个频率范围内使用了14个通道

未抽取谐振位，取得了较好的效果。在选择方案时，考虑到客观的条件，我们设计了一组带通滤波器和整流、低通滤波电路半代替动态频谱分析仪。

通道数目的选择没有唯一的标准，从理论上讲，通道越多越好，越多，抽取的特征越准确；但是，在实际的应用中，考虑到谐振峰造价、复杂性等因素，通道的数目不可能做到很大。文献[7]由于使用了现成的分析仪，所以选择了14个通道，但在本文中，一共设计了5个通道，4个低频通道和一个高频通道。设计时，主要考虑让四个低频通道覆盖元件向着两个共振峰（主要的共振峰），高频通道检测辅助。每个通道都由两个二阶带通滤波器和一个低通滤波器组成，见各[2]。各中



各[2]: 通道线路各

各中，虚线框①、②内的电路及元件参数完全相同，是两个二阶无陷增益多路反馈带通滤波器[9]。两个滤波器的级联是为了使幅度响应曲线更加陡峭。二极管具有单向导电性，作整流用。框③的电路是低通滤波电路，它和二极

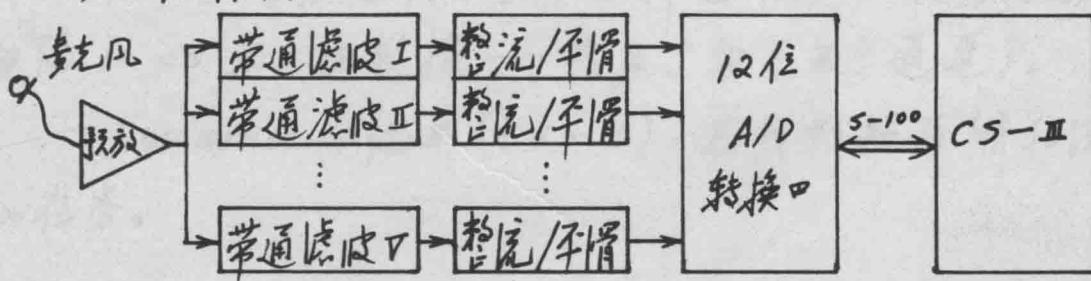
管一起作抽取频谱色络之用。

五个通道的中心频率分别为 $225$ 、 $450$ 、 $900$ 、 $1800$ 、 $7200$ Hz，品质因数(Q值)为3。五个通道中，每个的放大倍数各不相同，滤波器的中心频率愈高，放大倍数也愈大，这是为了补偿人的声道特征，因为它在音频范围内产生的信号要比低频的十倍以上，所以在确定放大倍数时应尽量使多种不同的频率成分有同样较大的幅度。

值得一提的是，音频通道的使用使子系统能够极大地控制辅音：声，尤其是摩擦音：声，这使本子系统控制摩擦音的能力大大超过了所有不使用音频通道的子系统；但是，这也正是由于音频通道的使用，使子系统对外部噪音更为敏感，因为外部噪声大都具有较高的频率分量，这严重地影响了子系统的抗干扰性能。也许，许多识别子系统没有采用它的原因之一就在于此吧。本文使用音频转换单元充分利用它的优点；另外，对检测特征作了改进，采取了一些有效的措施，大大提高了子系统的抗干扰性能。

### 2.3. 子系统结构和识别过程

#### (1) 子系统结构

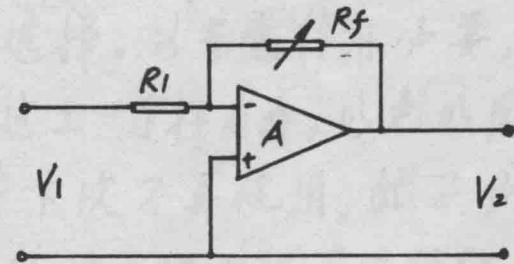


图[3]: 语音识别子系统的硬件框图

答[3] 该市是整个语音识别系统的硬件框架，它由以下几部分组成：

• 麦克风：它的作用是把语音信号变为电信号。这里采用的是灵敏度较低的驻极体话筒。

• 预放大器：麦克风输出的信号很微弱，以至采用灵敏度较低的动圈式话筒，信号就更为弱小。带通滤波器虽然有一定放大倍数，但不能满足要求；因此，话筒的输出信号在进入带通滤波器之前，还必须经过预放大。答[4] 是它的原理图。 $R_f$  是负反馈电阻，阻值大小可调， $R_f$  改变，放大倍数也随之改变，这样可以使系统适应不同灵敏度的麦克风。



答[4]：预放原理图

• 滤波、整流、平滑：得到语音信号。详见前节。

• A/D 转换器：使用的模/数转换器是 Cromemco 的 AIM-12 (12 BIT ANALOG INPUT MODULE)，它通过 5-100 端线与 Cromemco 系统Ⅲ相联，把输入的模拟电压转换成 12 位的数字量送给主机。采样速率为 100 Hz，也就是说每隔 10 ms 采一组数据——共五十分（五个通道）。

• Cromemco 系统Ⅲ (CS-III)：完成数据压缩和决策的任务。