

MICROCOPY RESOLUTION TEST CHART
NATIONAL BUREAU OF STANDARDS-1963-A



AD-A151 936



IMPLEMENTATION OF A REAL-TIME,
 INTERACTIVE, CONTINUOUS SPEECH
 RECOGNITION SYSTEM

THESIS

Kathy R. Dixon
 First Lieutenant, USAF

AFIT/GE/ENG/84D-26

DTIC FILE COPY

This document has been approved
 for public release and sale; its
 distribution is unlimited.

DTIC
 ELECTE
 MAR 28 1985
 S D
 E

DEPARTMENT OF THE AIR FORCE
 AIR UNIVERSITY

AIR FORCE INSTITUTE OF TECHNOLOGY

Wright-Patterson Air Force Base, Ohio

This document has been approved
 for public release and sale; its
 distribution is unlimited.

85 03 13 103

①

IMPLEMENTATION OF A REAL-TIME,
INTERACTIVE, CONTINUOUS SPEECH
RECOGNITION SYSTEM

THESIS

Kathy R. Dixon
First Lieutenant, USAF

AFIT/GE/ENG/84D-26

DTIC
SELECTED
1985
E

Approved for public release; distribution unlimited.

AFIT/GE/ENG/84D-26

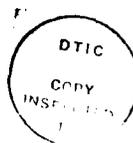
IMPLEMENTATION OF A REAL-TIME, INTERACTIVE,
CONTINUOUS SPEECH RECOGNITION SYSTEM

THESIS

Presented to the Faculty of the School of Engineering
of the Air Force Institute of Technology

Air University

in Partial Fulfillment of the
Requirements for the Degree of
Master of Science



Kathy R. Dixon
First Lieutenant, USAF

Graduate Electrical Engineering

December 1984

Accession For	
NTIS	<input checked="" type="checkbox"/>
DTIC TAB	<input type="checkbox"/>
Unannounced	<input type="checkbox"/>
Justification	<input type="checkbox"/>
By _____	
Distribution/	
Availability Codes	
Dist	Avail and/or Special
A-1	

Approved for public release; distribution unlimited.



Preface

The purpose of this research was to develop a real-time, continuous speech recognition system. Since many computer algorithms existed in the AFIT Signal Processing Laboratory to characterize various aspects of human speech, there was a need to combine these programs into a viable system.

The system described in this paper provides an interactive means of continuous speech recognition by computer. The system is speaker-dependent and requires training with a 70-word vocabulary prior to word recognition.

Completing this research would have been impossible without the assistance of several people. I would like to thank my thesis advisor, Dr. Matthew Kabrisky, Professor of Electrical Engineering at the Air Force Institute of Technology for motivating this research. Also, I would like to thank MAJ Ken Castor for his encouragement and for reviewing this document. In addition, a note of gratitude to my many friends who have put up with me during the highs and lows of this research. I would especially like to thank my friend, CPT Michael Conrad (US Army) for his encouragement and enthusiasm, and for his invaluable typing assistance.

Table of Contents

	Page
Preface	ii
List of Figures	v
List of Tables	vii
Abstract	viii
I. Introduction	1
Background	1
Problem	2
Scope	3
Assumptions	3
Summary of Current Knowledge	4
General Approach	5
Materials and Equipment	6
Sequence of Presentation	6
II. Speech Processing	8
Processor	8
A/D Conversion	12
III. Word Recognition	19
Phoneme Determination	19
Distance Rules	20
Clustering	22
Phoneme Template	23
Word Templates	25
Word Spotting	27
Time Alignment	30
IV. Software Design	31
Training	31
Recognition	34
V. Results	35
A/D Conversion	36
Phoneme Template Creation	37
Distance Matrix Creation	42
Word Phoneme Representations	45

Speech Recognition	49
VI. Conclusions	56
Summary	56
Recommendations	57
Bibliography	60
Appendix A: AFTI F16 Vocabulary	61
Appendix B: Computer Program Structure Charts	63
Appendix C: Program Listings	73
VITA	123

List of Figures

Figure	Page
1. Speech Recognition System	7
2. Acoustic Processor Designed by Hussain	9
3. SAM Configuration File	15
4. Energy in Data Vectors	39
5. Phoneme Template Vector Elements	41
6a. Distance Matrix	43
6b. Distance Matrix	44
7. Phoneme String Before and After Compression	46
8. Word Phoneme Representations from TRAIN	48
9. Phoneme String of Five Seconds of Speech Before and After Compression	50
10. "ONE" Spoken by KRD	51
11. "AIR-TO-AIR" Spoken by MGC	53
12. "AIR-TO-AIR" Spoken by DBS	54
13. AFTI-F16 Tentative Vocabulary	62
14. Legend for Structure Charts	64
15. Structure Chart 0.0: Program TRAIN	65
16. Structure Chart 1.0: Create Phoneme Template	66
17. Structure Chart 1.2: Produce Universal Phoneme Set ..	67
18. Structure Chart 2.0: Create and Print Distance Matrix	68
19. Structure Chart 3.0: Determine and Print Word Representation.....	69
20. Structure Chart 3.2: Create Word Phoneme Representations	70

21.	Structure Chart 4.0: Perform Speech Recognition	71
22.	Structure Chart 4.2: Recognize Input Speech	72

List of Tables

Table	Page
1. ASA-16 Analysis Bands	11
2. Bit Assignments for Control Words.....	12
3. Octal Values for Bit Setting of IBATAL	14
4. SAM Errors	17
5. Minkowski Distance	21
6. Word Phoneme Representations	27

Abstract

This thesis describes

A speech recognition system was designed and implemented to recognize continuous speech in a real time environment (after training). Several techniques were incorporated to characterize phonemes as vectors in space. Through the use of distance rules it was possible to characterize words by a phoneme representation, which could subsequently be used in word recognition. This approach to speech recognition offers several possibilities for future investigations such as varying the Minkowski distances and applying clustering techniques. The algorithm developed was modularized on a hierarchical basis and was user friendly.

Additional keywords: Software design, Computer applications, Computer programs

IMPLEMENTATION OF A REAL-TIME, INTERACTIVE,
CONTINUOUS SPEECH RECOGNITION SYSTEM

I. Introduction

Speech recognition by computer has seen dramatic gains in the past five years. Resulting primarily from increased computer processing power, efficient recognition algorithms can perform near real time. However, recognition systems using these algorithms are constrained to isolated word recognition and have limited vocabularies of between twenty and one hundred words. The major emphasis of present speech research is the development of an algorithm capable of recognizing natural, continuous speech (6:69).

Background

The goal of speech research is to determine a decoding scheme similar to that of the human brain. Many continuous speech recognition algorithms exist, yet none has a significant word recognition rate. There are three primary reasons for this, all of which are a consequence of our lack of understanding of how the human brain deals with the complex operation of speech recognition. First, the speech signal is highly encoded in the brain, and one must acquire an understanding of how such characteristics as intonation, articulation, semantics, and syntax are encoded before one

SAMGEN Rev. 2.10 7/29/82 at 10:49 Filename: SAMCONFIG3.SR

.Answers you gave in the SAMGEN dialog are shown in comment lines.
.Your inputs are immediately preceded by a colon (:) and appear
in the same order as you gave them to SAMGEN.

: Target operating system type :MRD
: Number of DG/DAC 4300 chassis configured: 0
: Fatal error handler name : -1
: Fatal error handler mailbox: -1

DCB.X SAMCO 100 -1 -1

: Number of Analog Subsystem :1

: A/D Con. #1 Device Code :21 Mode :AD Fortran ID = IDS21

: External interrupt handler specified :<NONE>
: Number of pages in Data Channel area :16
: Specifying a starting address for Data Channel area :Y
: Data Channel starting address :IBUFF

DCB.M DBS21 D.IDF+D.INF+D.DCH 21
DCB.I DTS21 SAINI 16 IBUFF
DCB?C -1 -1 DBS21
DCT.M DTS21 000377 INTSA DBS21

DCB.N 921 D.FIF 21 00 AD
DCB.S DBS21 0 AD.IS AD.IN SAIRT
DCB.A

: D/A Con. #1 Device Code :23 Mode :BD Fortran ID = IDS23

: External interrupt handler specified :<NONE>
: Number of pages in Data Channel area :1
: Specifying a starting address for Data Channel area :Y
: Data Channel starting address :IBUFD

DCB.M DBS23 D.IDF+D.INF+D.DCH 23
DCB.I DTS23 SAINI 1 IBUFD
DCB?C -1 -1 DBS23
DCT.M DTS23 000377 INTSA DBS23

DCB.N 923 D.FIF 23 00 BD
DCB.S DBS23 0 BD.IS BD.IN SAIRT
DCB.A

DCB.E

.End of SAMGEN configuration file.

Figure 3. SAM Configuration File

$$\text{IDATA2} = (\text{sample length}) * 800, \quad (2)$$

where sample length is the duration in seconds of each audio input. Variable, IDATA3, specifies the locations where conversion values are placed. In FORTRAN, IDATA3 is an integer array placed in a labeled common block. The block label is the same as that given in the SAMGEN configuration files.

<u>Octal Value</u>	<u>Function</u>
00000K	pulse clock
20000K	DCH clock
40000K	internal clock
60000K	external clock
0-1700K	start channel 0-15
0-17K	final channel 0-15

Table 3. Octal Values for Bit Setting of IDATA1 (1:69)

Configuration files, produced by an interactive dialog with the program SAMGEN, define operating system hardware and operation nodes. For example, SAMCONFIG3 has the source file shown in Figure 3.

The twelve bits in the machine word produce 2^{12} , or 4096 different conversion values. Over a range of 10 volts, this means that each value increment represents .024 volts. Using a one-to-one correspondence between sampled values and integer values, the real value for voltage can be determined by the equation

$$\text{Real Value} = \text{Float}(\text{Integer Value})/32768.*5. \quad (1)$$

Variables (IDATAx words) are passed to the software routines to represent channel use numbers, conversion count, clock source and storage locations for conversion values. IDATA1, occupying one machine word, represents clock source and channel use numbers. There are four clock sources available to the device: pulse, DCH, internal and external. Though all four can be used for A/D conversion, only the external clock is accessed. Channel numbers are given as starting and ending channels. If both are the same only one channel is specified. The bit values of IDATA1 can be set using the the octal values in Table 3. Since an external clock and all sixteen channels are used, IDATA1 equals 61700K.

IDATA2, also occupying one machine word, specifies the total conversion count or the number of data samples. Thus, using an 800 Hz external clock, IDATA2 is computed by

A/D Conversion

A/D conversion is performed by making use of the model 4331 Eclipse analog device and independent software interfaces. Though both A/D and D/A capabilities are present, only the former is addressed.

The Eclipse A/D conversion device operates using device op-codes and conversion data buffers. Organized around a single 12-bit converter and two multiplexers, the device can accept up to 16 different input signals with voltage levels in a $\pm 5V$ range. The conversion values are stored as a machine word with the bit assignments shown in Table 2.

<u>Bit</u>	<u>Assigned Value</u>
0	sign
1-11	storage values
12-15	zero

Table 2. Bit Assignment for Control Words

f ₀ (Hz)	Bandwidth(Hz)	Approximate Band Coverage	
		f ₁	f _h
260	130	203	333
390	130	330	460
520	130	459	589
650	130	588	718
780	130	718	848
910	140	843	983
1060	160	983	1143
1220	180	1133	1313
1400	200	1303	1503
1600	220	1494	1713
1820	250	1699	1949
2070	300	1925	2225
2370	340	2206	2546
3035	1030	2563	3593
4272	1445	3610	5055
5997	2005	5077	7083

Table 1. ASA-16 Analysis Bands (5:12)

Constant voltage level inputs to the ASA-16 chip are maintained by the automatic gain control circuit. Designed to accept 20 mV to 20 V over a 60 dB dynamic range, the AGC outputs voltages of 1 V to 4V, peak to peak.

To define band energies, the ASA-16 chip uses sixteen channels, each composed of a second-order bandpass filter, half-wave rectifier and a low-pass filter. Each channel is sequentially accessed using a sampled-and-held multiplexer. Table 1 shows the center frequencies and bands associated with each filter. A TTL crystal controlled 1 MHz clock is used for timing. Since the processor is built on a single board, both the chip and the board are powered by a common \pm 10 V power supply.

Buffers are used between the chip outputs and the 4331 Eclipse A/D converter inputs, to correct DC offset and satisfy the 200kohm termination resistance that each bandpass filter requires (5:8-41).

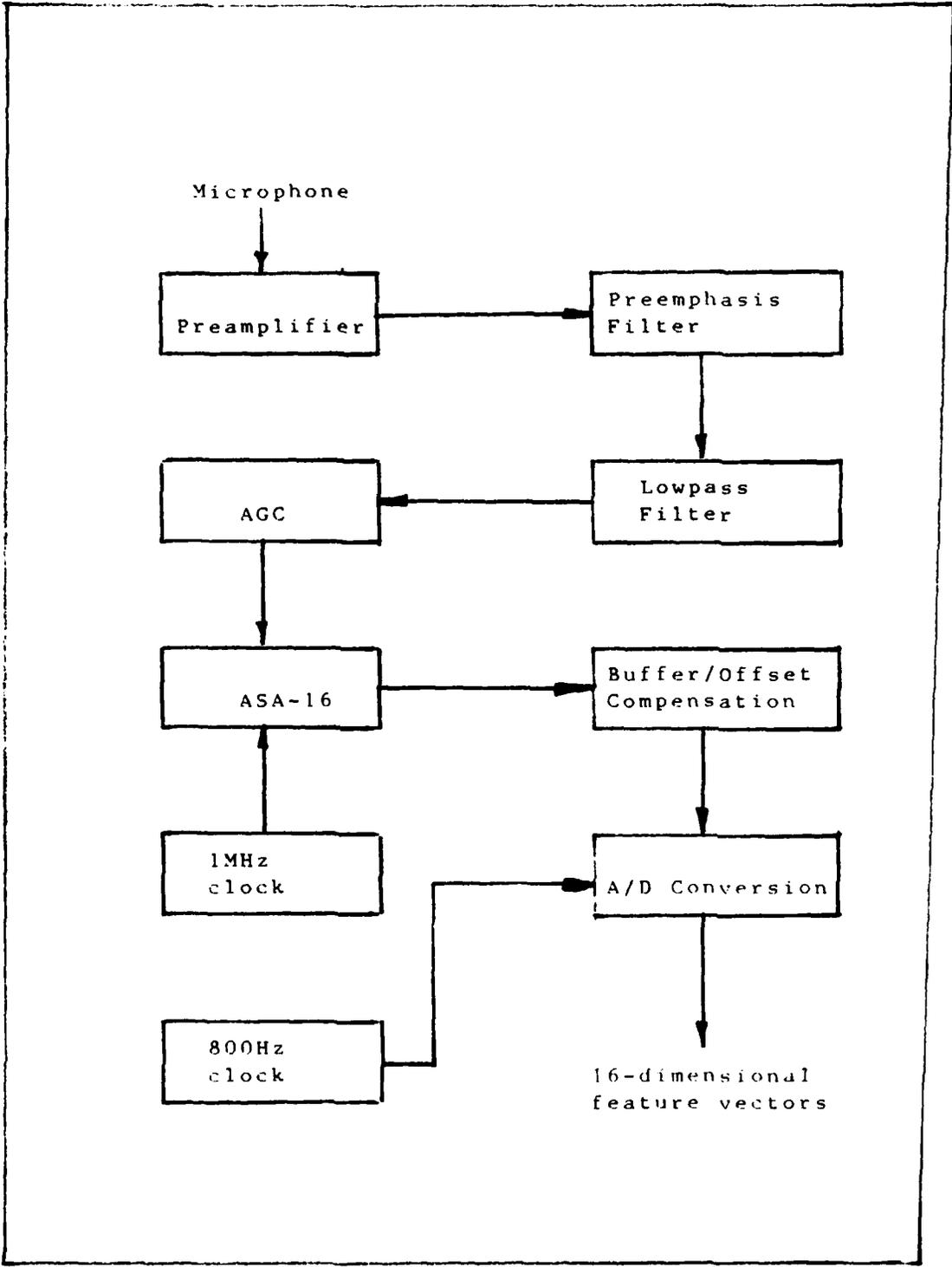


Figure 2. Acoustic Processor Designed by Hussain

II. Speech Processing

The word recognition algorithm developed in this thesis is based on the output of the acoustic processor designed by Ajmal Hussain (5). Built on a single board, the processor conditions the analog speech prior to analog/digital conversion. This chapter discusses the transformation of speech from an analog signal to a set of feature vectors used in pattern recognition.

Processor

Prior to A/D conversion, the audio signal is preconditioned to preserve frequency components and limit signal voltage levels. The primary component of this conditioning process is the ASA-16 spectrum analyzer chip. Figure 2 illustrates what happens to the speech after entering the system via a dynamic microphone and preamplification.

A preemphasis filter and an active low-pass filter band limit the audio signal from 200 Hz to 7000 Hz, as required by the ASA-16 chip. The preemphasis filter, with 6 dB/octave gain above 500 Hz, preserves the signal's high frequency components. The active low-pass filter then defines the signal's upper frequency at 7000Hz (cutoff frequency). A 17 dB passband gain provides proper voltage levels for the automatic gain control circuit which follows.

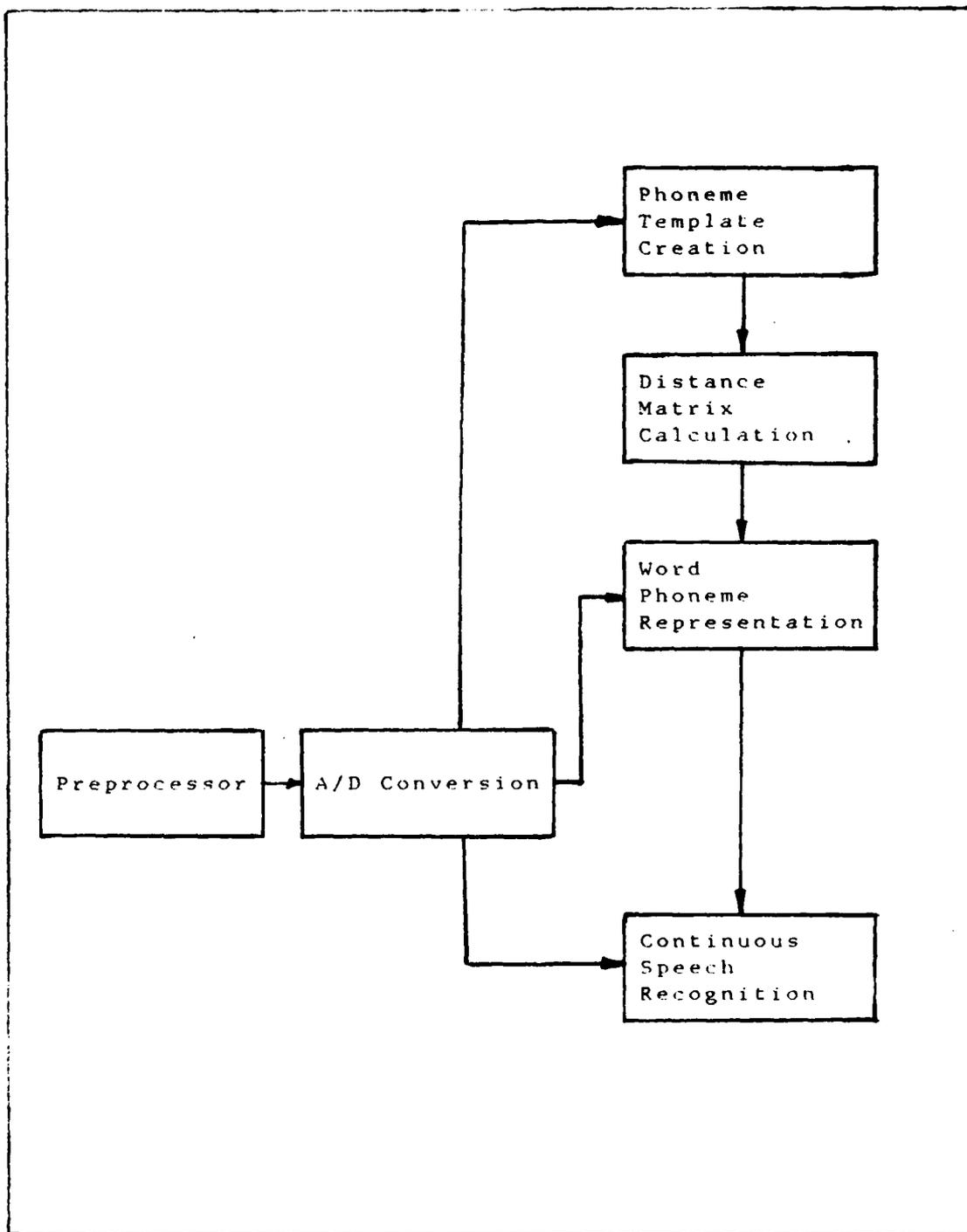


Figure 1. Speech Recognition System

two components of the machine: first, the acoustic processor (developed by Hussain) which provides a vector representation of the speech signal, and the recognition algorithm for determining a phoneme representation for each speaker, establishing templates for each word in the vocabulary (Appendix A), and finding words in a continuous speech phrase. More detailed outlines of these components are described in subsequent chapters.

Materials and Equipment

The following materials and equipment will be used:

1. Data General Eclipse S/250 Computer,
2. Filter Bank - Acoustic Processor.

Sequence of Presentation

The information contained in this report is presented in the following manner. Chapter 2 provides an explanation of speech processing both the digitization and frequency sampling of the speech signal. Chapter 3 describes the specifics of word recognition. Next, Chapter 4 is a synopsis of algorithm design. Finally, recognition results and appropriate conclusions and recommendations are presented in chapters 5 and 6, respectively.

front ends have been developed using fast Fourier transforms and filter banks. Spectrograms of speech waveforms have also provided information as to the possible waveform characterization using energy components. From these investigations, Seelandt (10) developed a universal feature set representing phonemes found in digitized speech. Then, Seelandt, Montgomery (7), and Hussain (5), succeeded in designing algorithms to recognize both isolated and connected speech by use of distance rules and template matching.

The purpose of this thesis is to combine several techniques developed at the AFIT Signal Processing Laboratory with additional modifications to produce a viable speech recognition system. A complete system, one in which analog speech is input and the recognized words are output, did not previously exist.

Two primary considerations in the design of this system were modularization and user friendliness. Modularity provides easy modification to system components, and user friendliness provides an operational machine which can be demonstrated by someone who may not be a computer expert.

General Approach

The speech recognition system developed from this research is described in Figure 1. Basically, there are

This means that phonemes at word boundaries are transitional and do not reveal much word information.

Summary of Current Knowledge

Research reports show extensive studies in the area of continuous speech recognition (4). Two methods currently being used for continuous speech recognition are segmentation and whole-word template matching. Segmentation is a means by which the acoustic signal is divided into unique sound units. After the word has been segmented, an attempt is made to match the sequence of units to a particular vocabulary word. The problem, though, is that some phonemes are lost in the process of chopping the signal into time slices. The other method, whole-word template matching, is means of matching whole word phoneme representations to the strings of phonemes produced by an acoustic processor. Problems inherent in this technique due to the variability of the word's duration when spoken. In the effort to find a solution some success has been achieved through the characterization of a word by phonemes from the body of the word. Thus, eliminating transition problems (4:574).

The development of a speech recognition machine has been an important area of research at the Air Force Institute of Technology for several years. Several acoustic

features such as formant frequencies, correlation coefficients, or linear predictive coefficients for pattern matching. Once the acoustic signal is encoded, templates are constructed for each word. These templates are subsequently matched to encoded speech for word hypothesis. Although proven effective for isolated speech, this method is difficult to apply to continuous speech recognition which has undefined word boundaries and variable word durations. To remedy this situation, algorithms incorporate partial template matching techniques. Thus, the object of this research was the development of a continuous speech recognition system relying on a universal feature set and partial template matching.

Scope

The computer algorithm developed in this thesis was designed to recognize speaker dependent, continuous speech. Some considerations for acoustic processor error are made by calculating error statistics and incorporating them into the word determination phase of the algorithm.

Assumptions

This research was based on two assumptions. First, the output of the acoustic processor is reliable such that phoneme choices are consistent best guesses. Second, words can be defined by an incomplete phoneme representation.

can hope to perform artificial decoding. Second, speech is a variable signal. Each time a speaker says a particular word, the phoneme representation of the word differs slightly. To the human ear, the variation is almost indistinguishable, but a machine must rely upon mathematical exactitudes. Therefore, any discrepancies are magnified in mechanical speech processing. Finally, although the brain somehow can distinguish the separations between words, these word boundaries are difficult for machines to find because the acoustic signal has no apparent pauses (4:570).

Fast, reliable man-machine communication is becoming a necessity as computers become integrated into today's society. Speech, the most natural form of human communication, seems appropriate in this application (6:64). Already, current testing onboard the Air Force's AFTI F-16, has proven that speech recognition is a definite aid to a pilot taxed to his physical limits by mechanical tasks.

Problem

The major problem in speech recognition lies not in the characterization of the acoustic signal, but in the determination of a decision scheme which uses these characteristics for word recognition. Several decision schemes have been developed which use acoustic signal

In the figure items worthy of particular attention are the device code number for A/D (IDS21), the number of pages in the data channel (16 blocks x 1024 words/block = 16384 words) and the data channel starting address IBUFF (area in common memory) (1:60-70).

To perform A/D conversion, the sequence of software commands is

```
EXTERNAL IDS21
EXTERNAL IDS23
COMMON/IBUFF/IDATA3(16384)
COMMON/IBUFO/IDATAO
DIMENSION IORBA(16)
```

```
CALL DSTRT(IER).
```

```
CALL DOIT[/W](IORBA,device-id,8,IDATA1,IDATA2,IDATA3,
IER).
```

By way of explanation, the first few lines declare both the device identification numbers and the common areas. After that, DSTRT performs initialization, and DOIT requests the conversion operation. Errors are reported by IER values as shown in Table 4. Any additional error is reported by IORBA(14). If IORBA(14) does not equal 40000K after conversion, an external clock interrupt or a clock overrun/underrun has occurred (1:78-85). A complete listing of this procedure can be found in the ATOD.FR source code in Appendix C.

Value	Meaning	Value	Meaning
2179	No - LNK routine in DCB, invalid DCB. Often results from an invalid device-id, so check the device-ids. The first two characters are ID, the third either S, A, or O, and the last two are numbers (e.g., IDS21).	2192	Illegal conversion count -- more than 255 or less than 1 -- for an A/D converter mode in A2; DG/DAC only.
2180	No DCB identifier in IORB, invalid DCB. Same cause as 2179.	2193	Assembly language only. Attempt to move data channel map while IORB is locked. A task tried to change the map while a request was using the window.
2181	Not used. This error should not occur.	2194	Attempt to move data channel map to an address outside the window.
2184	No initializing routine for a device that needs initialization. Same cause as 2179.	2195	Illegal conversion count: less than 1 or more than the device allows.
2185	Output requested to a channel for an illegal device (e.g., output to an A/D converter).	2196	Interrupt occurred from 4222 without a strobe or latch change.
2186	Attempt to set up a locked IORB array. This can happen if a second DSAN/DSOR call uses the same IORB array argument before the original DSAN/DSOR completes.	2197	Assembly language only. Attempt to use data channel map while it is being initialized or moved.
2187	Unable to find free IORB block in IORB array. Can happen if the IORB array was DIMENSIONed too small. A multiple-operation call needs 8 elements + 8 elements per operation.	2198	Assembly language only. Data channel not initialized: use an RMAP call before issuing this mode A2 request.
2188	No DCB exists with specified device-id. Same cause as 2179.	2199	SAM panic code. SAM could not transmit (.IXMT) to the calling task on IORB array completion. SAM aborts the program unless you set up a fatal error handling RECVe task and gave its name to SAMGEN, as described in Chapter 5, "Initial Dialog".
2189	Attempt to use unsupported feature (e.g., mapped call in unmapped system).	2200	External interrupt occurred on a stand-alone analog converter, aborting the request. This error returns from ISA calls only, not from DSAN/DSOR calls.
2190	Attempt to return bad buffer. Will never occur.		
2191	An IDATAx argument gave an illegal clock setting for an A/D or D/A converter.		

Table 4. SAM Errors (1:82-84)

Vectors are sequentially ordered in the data array, IDATA3. Each vector represents a time slice in which all sixteen filters have been sampled. The sampling rate per vector is 25 Hz, or once every 40 msec.

III. Word Recognition

As described in the introduction, word recognition in continuous speech consists of making a phoneme template, constructing word templates from the phoneme template, and matching the word templates to an input phoneme string. This chapter discusses the process in more detail.

Phoneme Representation

Since as early as 1947, phonemes have been used in word recognition. Potter, Kopp, and Green (8) found that sounds had unique spectrograms, and that people who they had trained to visually read these spectrograms could identify sounds in a word spectrogram. In other words, distinct phoneme patterns were discernable in word spectrograms. In 1981, Seelandt (10) investigated digitized speech spectrograms and the concept of phoneme patterns emerged. By combining several time slices or vectors of digitized speech, he produced a set of seventy phonemes. The concept was carried a step further by Hussain (5) who produced single-vector sixteen-dimensional phonemes. The method for phoneme generation developed for this system compares and averages these sixteen-dimensional vectors to produce a set of less than seventy phonemes. The comparison and averaging continue until a phoneme set is defined. This technique is explained more clearly later in this chapter.

Distance Rules

Since they are suited to defining the spatial orientation of n-dimensional vectors, distance rules are key to vector metrics. An understanding of vector relationships allows one to reduce large vector sets and approximate vector strings by known vectors.

Many pattern recognition algorithms incorporate metrics based upon Minkowski distances. Table 5 presents some of these different calculations.

The Minkowski distance between two vectors is computed by the following procedure. The absolute difference between corresponding vector elements is calculated and then raised to the Minkowski power. Next, these values are summed and the resulting sum raised to $(1/\text{Minkowski power})$. The final value is the vector score, or the numerical relationship between the two vectors. This algorithm uses Minkowski 4 which emphasizes the effect of any large discrepancies between compared vectors (5).

$$\text{Minkowski 1 } D_1 = \sum_{i=1}^M |x_{i,m} - x_{i,n}|$$

$$\text{Minkowski 2 } D_2 = \left(\sum_{i=1}^M (x_{i,m} - x_{i,n})^2 \right)^{1/2}$$

$$\text{Minkowski N } D_N = \left(\sum_{i=1}^M (x_{i,m} - x_{i,n})^N \right)^{1/N}$$

where N = Minkowski distance m = vector m
M = Vector length n = vector n
i = vector element

Table 5. Minkowski Distances

Clustering

Rabiner, et al. (9) showed that speech vectors in space tend to cluster. Using this concept, his research group tried clustering analysis on phonemes, words, and speakers. The experiments were successful for isolated word recognition, but clustering became a computational nightmare in continuous speech recognition. As a result clustering was not emphasized in this research. However, some spatial comparisons are performed.

By visualizing the speech vectors as residing in n-dimensional space, one can imagine that some vectors lie closer together than others. Therefore, close vectors can be used to define regions of space, or clusters. As previously mentioned this algorithm redefines nearest vectors by averaging and weighting. Each vector has an initial weight of one, and those vectors which are closest to one another by Minkowski distance are averaged together by the equation

$$X = \frac{(X * \text{weight}(m)) + (X * \text{weight}(n))}{\text{weight}(m) + \text{weight}(n)} \quad (3)$$

The highest-numbered vector is deleted and a new weight is assigned to the new vector by

$$\text{weight}(m) = \text{weight}(m) + \text{weight}(n). \quad (4)$$

The resulting set of vectors represents clusters of similar sounds.

Phoneme Template

By incorporating the above mentioned concepts, a phoneme template is made for each speaker. This template, a set of energy normalized vectors, represents the most unique sounds produced by that speaker. The process is:

1. Vector normalization
2. Vector deletion
3. Vector comparison
4. Vector averaging.

Vector normalization consists of noise removal and energy normalization. Since background noise is present in all environments except an anechoic chamber, its removal from the audio signal allows for better characterization of the signal. The first vector in each speech file is considered to represent the average background noise. Therefore, by subtracting this vector from each of the remaining vectors, the noise is removed. To insure voltage consistency, a noise threshold voltage of 30.5 millivolts (corresponding to the average laboratory background noise) is set. Each element in a vector is checked to insure voltages above the noise threshold. If the value of the

element voltage falls below this threshold, the elements voltage is redefined as zero. Additionally, the vectors are energy normalized, by dividing each vector component by the vector's energy, calculated by

$$E_m = \sum_{i=1}^{16} (X_{i,m})^2. \quad (5)$$

These energy normalized vectors are now considered phonemes.

On the basis of energy thresholding, some vectors are omitted from the speech file. The energy threshold is determined by finding the average energy per vector over the test utterance. The equation is

$$E_{AVG} = \left(\sum_{m=1}^{500} E_m \right) / (\text{number of vectors}). \quad (6)$$

Once normalization and deletion have been performed, the phonemes are compared to one another by Minkowski 4 distance to determine each vectors nearest neighbor. Then, if the total number of vectors is greater than sixty-nine, the two nearest neighbors are averaged together using the

averaging routine described earlier (see Clustering, p.20), and replaced by a single vector which is then energy normalized. After all nearest neighbor averaging have been completed, a check is made to see if the number of remaining vectors is less than seventy (the maximum number of phonemes that the system can handle). If the number of phonemes is greater than that, the entire comparison/averaging process is repeated.

The final outcome of this procedure is a set of vectors, each of which represents a particular phoneme. This vector set is referred to hereafter as the phoneme template. Each phoneme is then assigned a number from one to at most sixty-nine.

Word Templates

After the phoneme set is established, it is used to develop a codebook of phoneme representations for each word in the given vocabulary (Appendix A). The process components are:

1. Normalization
2. Phoneme extraction
3. Word representation by phonemes
4. Vocabulary creation.

As is the case in the creation of phoneme templates, the

vectors in the digitized speech of each word are normalized.

At this point phoneme extraction takes place. Each vector in the phoneme template is compared to each vector in the vector string of the input speech. By calculating the Minkowski \ast distance between a template vector and a string vector, the closest (minimum distance) template vector is determined. Then the phoneme number associated with the template vector is placed in an array, and the process is repeated until all string vectors are represented by phoneme numbers.

Following the process of string representation, the array of phonemes is compressed to at most ten phonemes in the following manner. First, identical adjacent phonemes are represented as one phoneme. Then, the distances between adjacent phonemes are obtained from a distance matrix. Adjacent phonemes with distances between them of less than 11.0 (distance measures are normalized to 100) are compressed by deleting the highest numbered phonemes.

The final outcome of this process is a string of phonemes representative of a vocabulary word. Collectively, these strings represent the entire vocabulary. Table 6 shows the phoneme representations of several words.

Vocabulary WordPhoneme Representation

ADVISE	23	34	66	38	52	61	25	37	14	0
BRAVO	40	9	10	52	41	36	0	31	0	0
CANCEL	17	16	27	0	29	0	37	0	0	0
ECHO	13	28	0	45	44	45	31	0	13	0
FIVE	13	0	61	62	26	25	27	36	0	0
FOUR	59	36	58	31	60	31	0	0	0	0
WEAPON	7	8	66	40	0	45	0	0	0	0
WEST	40	8	65	66	64	20	20	29	14	0
ZERO	20	1	24	6	5	6	40	31	0	0

Table 6. Word Phoneme Representations

Word Spotting

Word spotting is the technique of identifying words in continuous speech. Using the phoneme template and word templates, speech recognition is accomplished by the following procedure:

1. Normalization
2. Phoneme extraction
3. Compression
4. Comparison.

The procedure for normalization, extraction, and compression of input speech is identical to that which is followed in the creation of word templates. Once the input speech is represented as a string of template phonemes, recognition is attempted.

The recognition scheme, based on the scheme described by Hussain (5), is as follows. After having undergone phoneme representation, the input string of phonemes is searched for two or more adjacent zeros. Since zero phonemes represent low energy noise phonemes, adjacent zero phonemes are considered word boundaries. When two word boundaries are found, the next step is to establish the number of phonemes between the boundaries. Less than three phonemes are considered noise, more than eight phonemes means that two words have been spoken.

One word identification is performed by finding the distance (Minkowski 4) between the the first word template phoneme and the first string phoneme. Subsequent word phonemes are then compared to the string phonemes until a match (best fit) has been found for each of the word phonemes. The distances are variations between word

phonemes and string phonemes, and are considered errors. The errors are then summed to obtain the total error by

$$E_{\text{TOTAL}} = \sum_{i=1}^M (D_{m,n})^4 \quad (7)$$

where D is Minkowski 4, m is the string phoneme number, n is the template phoneme number, and M is the number of string phonemes. Error calculations are repeated twice by shifting the word phoneme set one phoneme to the right, and then one phoneme to the left of the original value. Then the three resulting error values are averaged. After all of the vocabulary words have been matched to the input string, the word with the smallest average error is considered to be the correctly recognized word.

For continuous recognition, in which two words exist side by side, a slightly different error calculation is performed. Phonemes are added to a buffer one at a time. Then, after ten additions to the buffer, error statistics are calculated. The algorithm decides at which point the best word match occurs and chooses this word to be the first word. A similar calculation is performed for the second word, however, in this case the last phoneme of the first word is used as the first phoneme of the second. Isolated

Distance Matrix Creation

Seelandt (10) showed that it is more time efficient to look up distances between phonemes than to recalculate them. Therefore, a lower triangular matrix is formed by calculating the distances between the phonemes in the array PHON. These distance values are stored in the array DIS. The number of phonemes in the matrix is stored in position 2432. The code for this subroutine was modified from Hussain's code so that the distance between phonemes are raised to the power 0.25 and then multiplied by 100. The resulting matrix is shown in Figure 6a and 6b. As is the case during phoneme template creation, these values are written to a file called DIST. Program TRAIN is capable of both writing to, and reading from, this file.

vectors could have the same nearest neighbor, it is necessary to find pairs of vectors which are each other's nearest neighbors. This requires a considerable amount of time, but it is important for accurate analysis of vector metrics. After a nearest neighbor pair is found, the vector elements of the lowest-numbered vector are recalculated using Equation 3 and energy normalized. The weight of each vector is recalculated by Equation 4, and the weight of the highest-numbered vector is set to zero. The process is repeated until all nearest neighbor pairs have been found. These vectors are then compressed by the subroutine CMPRS1, and the number of remaining vectors is examined. If this number is greater than sixty-nine, new nearest neighbors are determined and the averaging process is repeated. This entire operation process can take anywhere from five to ten minutes.

After the correct number of vectors has been obtained, the voltage values are stored in the array PHON and passed to TRAIN, which has options for writing them to a file (PHONE) and reading them from a file (PHONE). (Figure 5 depicts the voltage elements for each vector in the phoneme template file. The vectors are horizontally arranged such that phoneme 1 is positioned (1,1), (1,2), (1,3),..., (1,16)).

1 7197	1 7197	1 7323	1 7324	1 7254	1 7327	1 7299	1 7330	1 7265	1 7293
1 7231	1 7231	1 7303	1 7273	1 7185	1 7285	1 7312	1 7233	1 7235	1 7265
1 7236	1 7262	1 7303	1 7262	1 7169	1 7225	1 7238	1 7244	1 7261	1 7267
1 7236	1 7270	1 7238	1 7272	1 7244	1 7280	1 7296	1 7275	1 7409	1 7321
1 7431	1 7427	1 7568	1 7559	1 7496	1 7491	1 7499	1 7424	1 7584	1 7677
1 7689	1 7638	1 7660	1 7698	1 7663	1 7636	1 7688	1 7485	1 7161	1 7303
1 7334	1 7188	1 7202	1 7319	1 7258	1 7243	1 7285	1 7292	1 7199	1 7278
1 7284	1 7203	1 7231	1 7290	1 7245	1 7329	1 7263	1 7227	1 7340	1 7427
1 7378	1 7268	1 7279	1 7362	1 7318	1 7429	1 7394	1 7324	1 7357	1 7360
1 7298	1 7183	1 7234	1 7330	1 7335	1 7334	1 7379	1 7238	1 7268	1 7228
1 7322	1 7388	1 7318	1 7287	1 7307	1 7295	1 7214	1 7284	1 7283	1 7300
1 7311	1 7237	1 7244	1 7322	1 7391	1 7378	1 7249	1 7252	1 7335	1 7283
1 7282	1 7244	1 7264	1 7264	1 7262	1 7348	1 7317	1 7319	1 7334	1 7317
1 7313	1 7234	1 7274	1 7240	1 7257	1 7277	1 7391	1 7402	1 7387	1 7408
1 7374	1 7336	1 7297	1 7220	1 7294	1 7294	1 7231	1 7418	1 7369	1 7390
1 7520	1 7374	1 7555	1 7472	1 7342	1 7274	1 7241	1 7409	1 7389	1 7480
1 7427	1 7431	1 7349	1 7344	1 7423	1 7352	1 7382	1 7452	1 7364	1 7325
1 7223	1 7201	1 7233	1 7274	1 7321	1 7301	1 7262	1 7202	1 7275	1 7320
1 7193	1 7205	1 7320	1 7301	1 7260	1 7262	1 7318	1 7276	1 7291	1 7284
1 7247	1 7192	1 7280	1 7255	1 7247	1 7252	1 7273	1 7289	1 7257	1 7283
1 7244	1 7281	1 7312	1 7430	1 7442	1 7445	1 7324	1 7210	1 7321	1 7291
1 7250	1 7245	1 7228	1 7307	1 7360	1 7304	1 7345	1 7287	1 7323	1 7222
1 7229	1 7302	1 7301	1 7194	1 7235	1 7320	1 7339	1 7557	1 7557	1 7499
1 7569	1 7361	1 7493	1 7303	1 7194	1 7217	1 7332	1 7314	1 7321	1 7306
1 7253	1 7282	1 7259	1 7288	1 7291	1 7356	1 7431	1 7339	1 7421	1 7575
1 7559	1 7407	1 7472	1 7280	1 7281	1 7271	1 7234	1 7274	1 7333	1 7309
1 7230	1 7278	1 7322	1 7227	1 7286	1 7425	1 7297	1 7264	1 7351	1 7472
1 7468	1 7692	1 7759	1 7710	1 7673	1 7635	1 7568	1 7374	1 7335	1 7323
1 7289	1 7387	1 7387	1 7308	1 7295	1 7402	1 7339	1 7258	1 7210	1 7275
1 7249	1 7236	1 7308	1 7261	1 7258	1 7350	1 7460	1 7308	1 7387	1 7425
1 7373	1 7705	1 7528	1 7527	1 7525	1 7261	1 7279	1 7434	1 7335	1 7173
1 7265	1 7250	1 7265	1 7283	1 7325	1 7364	1 7520	1 7427	1 7449	1 7450
1 732	1 7413	1 7305	1 7308	1 7352	1 7285	1 7288	1 7239	1 7172	1 7282
1 7246	1 7190	1 7208	1 7166	1 7318	1 7263	1 7273	1 7296	1 7256	1 7266
1 7244	1 7218	1 7285	1 7192	1 7284	1 7229	1 7272	1 7240	1 7324	1 7201
1 7218	1 7305	1 7230	1 7274	1 7249	1 7277	1 7218	1 7200	1 7313	1 7481
1 7783	1 7833	1 7993	1 8079	1 7811	1 7544	1 7330	1 7274	1 7223	1 7351
1 7346	1 7237	1 7219	1 7303	1 7248	1 7279	1 7331	1 7458	1 7302	1 7269
1 7316	1 7339	1 7354	1 7348	1 7354	1 7327	1 7353	1 7400	1 7281	1 7301
1 7246	1 7286	1 7277	1 7222	1 7254	1 7222	1 7200	1 7214	1 7285	1 7201
1 7299	1 7239	1 7293	1 7338	1 7313	1 7376	1 7289	1 7287	1 7391	1 7335
1 7357	1 7325	1 7303	1 7271	1 7254	1 7274	1 7338	1 7270	1 7215	1 7270
1 7313	1 7334	1 7368	1 7400	1 7477	1 7434	1 7400	1 7402	1 7509	1 7429
1 7324	1 7399	1 7291	1 7286	1 7274	1 7262	1 7243	1 7230	1 7292	1 7262
1 7264	1 7479	1 7582	1 7582	1 7553	1 7457	1 7468	1 7448	1 7327	1 7292
1 7284	1 7374	1 7315	1 7270	1 7272	1 7296	1 7264	1 7311	1 7289	1 7362
1 7289	1 7310	1 7323	1 7371	1 7324	1 7243	1 7308	1 7325	1 7267	1 7421
1 7304	1 7309	1 7299	1 7200	1 7283	1 7281	1 7278	1 7218	1 7292	1 7285
1 7423	1 7461	1 7390	1 7319	1 7301	1 7322	1 7311	1 7370	1 7422	1 7598
1 7729	1 7661	1 7591	1 7366	1 7255	1 7243	1 7230	1 7266	1 7291	1 7221
ENERGY THREE									
1 7339									

Figure 4. Energy in Data Vectors

In the COMPARE subroutine, the distances between all pairs of vector's are determined by Minkowski 4 distance calculations. Each vector's nearest neighbor is determined by Hussain's algorithm. This concept is carried a step further by the introduction of vector weighting. Since many

into the array. Then the average energy of all the vectors (500 vectors = 10 seconds of speech) is calculated for future use as a threshold. Several alternative energy thresholds were considered before the average energy was considered an optimum threshold.

After energy normalizing each vector (by dividing each vector element by the entire vector's energy), those vectors with energy below the threshold energy are deleted. This is accomplished in the subroutine LOWENERGY which also assigns a weight of one to each vector prior to deletion. Deleted vectors are then weighted zero. (This weighting system dramatically reduces the number of calculations previously used by Hussain to delete vectors). Vectors with weights of zero are bubble sorted to the end of the data array and the number of sixteen-dimensional vectors left is calculated and stored in position IDATA(IDATA2+1) (See CMPS1.FR, Appendix C). If this number is less than seventy the program returns to TRAIN. If not it calls the subroutine COMPARE, which is a comparison and averaging routine.

program can digitize. This is especially important to the various phases of the speech recognition system. A SWAP calls each program for digitizing ten seconds of speech for phoneme template creation (ATOD10), two seconds of speech for word phoneme representations (ATOD2), and five seconds of speech for speech recognition (ATOD5). Each program creates a file (OUT2,OUT5,OUT10, respectively) in which data conversion values or integer voltages are stored. The data are stored sequentially, such that every set of sixteen integer values represents a vector of speech.

An important requirement for the correct processing of speech is that the input speech be between +5 and -5 volts. This level should be checked on the oscilloscope.

Phoneme Template Creation

A prompt from TRAIN requires the user to input ten seconds of speech, from which the set of phonemes will be derived. When the user says a predetermined phrase into a microphone, ATOD10 produces the file OUT10 containing 500 sixteen dimensional vectors. A call is then made to the subroutine TEMPLATE, which directs the creation of the phoneme template. OUT10 is opened and read into a common buffer in which the sampled data points are stored in the array IDATA. Then the energy in each vector is calculated and stored in the array IENERGY. (Figure 4 is a printout of this array) Each vector's energy is entered sequentially

for use in the system described herein required a considerable amount of modifications. Many of Hussain's individual programs appear as subroutines in the current program, making it necessary to provide strict continuity among them, and to develop interface routines so they can become efficient contributors to a unified whole. Furthermore, Hussain's programs require extensive interaction between the user and the machine, often with the human performing tasks which could be better handled automatically (e.g. the creation of phoneme templates which involves manually removing undesirable phonemes). Hussain's results could not be reproduced by other experimenters because his choice of values for some constants used in the algorithms was often completely arbitrary (at least the method for determining them was not revealed). In the interest of a more scientific approach, an attempt has been made in this thesis to derive values for these constants (energy threshold, distance threshold, error constants) from a more logical basis.

A/D Conversion

After the analog speech has gone through the preprocessor, it is converted to a digital signal by one of three programs: ATOD2, ATOD5, and ATOD10. The numbers in the program titles indicate how many seconds of speech each

V. RESULTS

From the very start, this research effort was hampered by time constraints. The thesis report which served as a starting point for this particular investigation was incomplete and contained vague documentation making it necessary to do some preliminary analysis in order to define what had already been accomplished. This coupled with the sheer magnitude of the task of developing a complete, operational speech recognition system where none had existed previously, resulted in the expiration of available time before a thorough quantitative evaluation of the system's performance could be made. Such an evaluation would require the compilation of large amounts of data and systematic tuning of parameters in search of the optimal setting for the control dials. The fine tuning, although important, is not really essential to the validation of this system as "operational", and so the discussion in this chapter will focus on the performance of the different components and the way they are integrated into the system.

The programs developed by Ajmal Hussain (5) provided the foundation for most of this research effort. (The major thrust of the present treatment is to take the various individual parts of the puzzle provided by Hussain(5) and others (10) (7) and combine them into a complete speech recognition system.) However the adaptation of his programs

closeness, value until there are at most ten phonemes left in the string. [An aside: Low energy vectors are considered noise and are represented as zero. Two adjacent zeroes are compressed to one zero, and one zero by itself is ignored, as are all zeroes appearing at the beginning of the string (5).] The array SOUND, containing these phonemes, is then entered into the appropriate location in the matrix VOCAB. After all of the vocabulary words have been processed, PRINTREP prints the words along with their phoneme representations.

Recognition

After completing the training phase, TRAIN enters the recognition mode. It asks the user to input five seconds of speech through the dynamic microphone. Then the program swaps to ATOD5 for A/D conversion and creation of file OUT5. Next, TRAIN passes the phoneme template (PHON), the distance matrix (DIS), the vocabulary words (WORD), and the phoneme word representations (VOCAB) to subroutine SPEECH. The conversion values are read into a buffer, where the vectors are normalized (NORMALIZE) and the phonemes are extracted (EXTRACT). The CMPRS subroutine compresses adjacent phonemes, and RECOG calls Hussain's speech recognition algorithm. The recognized word strings which result are printed on both the screen and the line printer.

equation

$$\text{location} = m(m - 1)/2 + n. \quad (8)$$

A call to PRINTDIS provides a printout of the distances between phonemes, normalized to 100.

After these two steps are completed, TRAIN determines phoneme representations for the words in the vocabulary. Vocabulary words are flashed onto the video screen, one at a time. A swap to ATOD2 enables the A/D conversion of two seconds of speech. Conversion values are subsequently stored in the file, OUT2.

Next, REP is called to read the eight blocks (2048 conversion values) into a buffer. As in TEMP, every 16 values represents a vector of sampled speech. After normalization (NORMALIZE), phonemes are extracted (EXTRACT) from the speech file. This means that the speech file is represented as a string of 250 phonemes. CMPRS goes through the phoneme string comparing each pair of adjacent phonemes. If the two are identical, they are compressed into one, if they are dissimilar the distance between them is examined, and in the case of a distance within the threshold, the phoneme having the higher number is deleted from the string. If the distance lies outside the threshold, no action is taken and the next pair of phonemes is examined. CMPRS continues reducing the phoneme string by increasing the

calculation.

Creating a phoneme template requires digitized speech input and numerous calculations and comparisons. TRAIN requests that the user provide 10 seconds of input speech. A swap with TRAIN calls ATOD10, which performs the A/D conversion and stores the data conversion values into a file named OUT10. Then TEMPLATE reads 16 blocks (8192 conversion values) from the file into a buffer. Every sixteen conversion values represents a sixteen-dimensional vector of speech. The vectors are normalized (NORMALIZE), checked for low energy (LOW ENERGY), and compressed (CMPRS1) by removing the low energy vectors. A comparison and averaging routine called COMPARE finds each vector's nearest neighbor vector. Subsequently, the two nearest neighbors are averaged and, following compression (CMPRS1), a new set of vectors is defined. This process repeats until a maximum of sixty-nine vectors is left. (Seelandt (10) showed 70 phonemes was a viable number for characterizing sounds in speech). These remaining vectors are stored in the array PHON for future use. The value PHON(1121) is the total number of vector phonemes.

Next, TRAIN calls DISTANCE, which creates a lower triangular matrix of the distances between phonemes. The matrix is stored in array DIS, in which the location between two phonemes m and n (where $m > n$) can be found by the

IV. Software Design

The computer algorithms designed for use in this thesis were designed modularly on a hierarchical basis. The motivation for this was the requirement for ease in system modification. The various components can be changed to accommodate various methods of phoneme template creation, word representation, and continuous speech recognition. Therefore, this thesis provides a valuable tool for future research in this area.

The primary driving program of the system is TRAIN. Having two phases (training and recognition), it ties all of the components of a speaker-dependent recognition system together. This chapter provides a discussion of these components.

[NOTE: In addition to the main speech recognition algorithms, several support algorithms were written for A/D conversion, file manipulation, and creation of a vocabulary of words. The source codes for these algorithms are in Appendix C.]

Training

Training is accomplished through subroutine calls from the program TRAIN. In other words, TRAIN has a training mode in which a phoneme template, a distance matrix, and a word phoneme template are made. Each mode is a separate

word recognition is then performed as described above. Finally, error statistics are calculated for the entire string. For continuous recognition the entire process is repeated by choosing the second best guess for the first word, then the third best guess and so on. Following the final iteration, the string having the least total error is accepted as the correctly recognized sequence of words. The entire process repeats for the next phrase indicated by a set of adjacent zeroes.

Time Alignment

Because phonemes are compared on a one-to-one basis after compression, word duration is ignored and time alignment is unnecessary.

Word Phoneme Representation

Once the phoneme template and the distance matrix are formed, the system can start making word phoneme representations. TRAIN reads from a file the list of vocabulary words and stores them as Hollerith strings in the array named (appropriately) WORDS. The vocabulary used in this thesis is that shown in Appendix A. TRAIN's next step is to loop through the vocabulary words to find the phoneme representation of each.

The user is prompted to enter two seconds of speech by saying the word which appears on the computer video screen during the allotted time window. A swap to ATOD2 produces the file OUT2 containing the 125 digitized speech vectors.

TRAIN then calls the subroutine REP, which controls the procedure for finding that word's phoneme representation. As was the case in the creation of phoneme templates, the data vectors are normalized and the low energy vectors are deleted using the previously determined threshold. Each data vector is then compared to the complete set of phoneme vectors by calculating the distance between the data vector and each of the phoneme vectors. After this the phoneme vector whose distance calculation produced the smallest value is used to represent the string vector. This is repeated for the entire string of speech vectors, forming an array of numbers representing the closest phoneme vectors. This string is then compressed (CMPRS) by replacing

After the phoneme representations for all the vocabulary words have been found, they are placed in the array VOCAB. TRAIN then writes them to the file VOCABUL, and reads them from the same. TRAIN also prints the vocabulary words, along with their phoneme representations on the line printer. Figure 8 shows an example of this printout.

ADVISE	23	34	66	38	52	61	25	37	14	0	0
AFFIRMATIVE	58	6	8	52	43	28	31	0	0	0	0
AFT	16	65	25	65	25	0	0	0	0	0	0
AIR-TO-AIR	2	35	33	22	14	16	47	22	45	0	0
AIR-TO-SURFACE	23	33	5	14	22	14	0	0	0	0	0
ALPHA	37	62	61	41	54	12	36	0	36	0	0
ARM	63	41	61	62	43	25	27	0	13	0	0
BACKSPACE	11	61	16	20	37	14	29	14	13	0	0
BEARING	47	34	7	6	22	16	0	0	0	0	0
BRAVO	40	9	10	52	41	36	0	31	0	0	0
CANCEL	17	16	27	0	29	0	37	0	0	0	0
CHAFF	13	14	15	47	33	65	25	27	13	0	0
CHANGE	14	15	14	45	47	34	16	14	13	0	0
CHANNEL	14	15	0	16	22	16	22	0	0	0	0
CHARLIE	15	57	7	26	7	43	44	46	45	0	0
CLEAR	46	64	45	47	48	5	45	0	0	0	0
CONFIRM	21	7	6	44	27	0	0	0	0	0	0
DEGREES	37	21	22	2	56	57	20	0	0	0	0
DELTA	30	38	43	39	12	58	36	63	28	0	0
EAST	13	0	37	0	20	29	0	15	0	0	0
ECHO	16	28	0	45	44	45	31	0	13	0	0
EIGHT	1	17	46	0	0	0	0	0	0	0	0
ENTER	57	17	16	0	30	45	0	0	0	0	0
FAULT	12	41	54	39	12	0	37	0	0	0	0
FIVE	13	0	61	62	26	25	27	36	0	0	0
FLARES	44	38	48	5	3	45	15	14	15	0	0
FORWARD	59	12	40	37	0	37	0	13	0	0	0
FOUR	59	36	58	31	60	31	0	0	0	0	0
FOXTROT	41	26	62	27	29	14	57	13	0	0	0
FREQUENCY	30	48	2	17	37	0	0	0	0	0	0
FUEL	1	23	55	46	31	60	31	13	0	0	0
GUN	28	44	25	43	25	27	0	0	0	0	0
HEADINGS	27	64	34	38	28	16	47	2	20	0	0
HUNDRED	25	28	0	30	0	0	0	0	0	0	0
KNOTS	27	54	52	26	61	0	20	29	20	0	0
LOCK-ON	39	42	53	61	12	44	27	25	27	0	0
MAP	0	0	0	0	0	0	0	0	0	0	0
MARK	12	42	53	42	62	6	36	21	13	0	0
MILES	25	66	12	20	29	20	13	0	0	0	0
MINUS	62	43	26	65	16	63	20	13	0	0	0
MISSEL	63	22	45	46	63	20	29	31	31	0	0
NEGATIVE	22	34	35	47	0	0	0	0	0	0	0
NINE	27	66	43	26	25	22	16	46	0	0	0
NORTH	31	36	12	39	12	44	31	0	13	0	0
NOSE	28	38	7	8	32	37	14	29	20	0	0
ONE	31	50	12	42	39	25	27	36	0	0	0
POINT	13	36	12	11	12	25	27	20	0	0	0
PROFILE	44	8	10	40	31	12	41	25	37	0	0
RADAR	21	22	34	35	34	17	28	12	36	0	0
RANGE	21	33	34	17	18	24	0	14	15	0	0
REPORT	59	36	40	32	31	21	0	0	0	0	0
RHAW	31	12	54	41	25	58	0	0	0	0	0
SEARCH	20	0	46	44	6	30	0	14	15	0	0
SELECT	20	0	31	0	31	25	27	0	0	0	0
SEVEN	20	20	0	45	64	25	28	0	0	0	0
SIX	20	63	22	33	34	47	20	29	20	0	0
SMS	27	31	37	20	29	20	1	20	37	0	0
SOUTH	13	0	37	0	28	25	62	0	13	0	0
STATION	37	20	57	2	14	15	63	13	0	0	0
STRAFE	20	29	0	21	22	33	34	47	2	0	0
TAIL	45	35	64	47	2	19	28	36	0	0	0
TARGET	36	12	43	54	43	7	6	0	2	0	0
THOUSAND	25	54	26	62	61	41	12	31	20	0	0
THREAT	21	22	33	66	65	64	28	0	13	0	0
THREE	21	6	5	33	48	2	0	37	0	0	0
TWO	13	0	29	14	0	1	2	46	0	0	0
WAYPOINT	31	38	34	35	34	17	36	14	13	0	0
WEAPON	7	8	66	40	0	45	0	0	0	0	0
WEST	40	8	65	66	64	20	20	29	14	0	0
ZERO	20	1	24	6	5	6	40	31	0	0	0

Figure 8. Word Phoneme Representations from TRAIN

Speech Recognition

TRAIN enters the recognition mode after the training session is completed, and prompts the user for a five second input of speech. Once again a SWAP call is used for A/D conversion and a file (OUT5) is written to contain sampled data points. The governing subroutine SPEECH then calls the necessary routines to read the data into the common buffer, normalize the data vectors, extract and compress the phoneme number string, and finally, perform speech recognition. All the steps prior to recognition are the same as those in word phoneme representation, with the exception of compression. Here, compression occurs only once, as there is no minimum number of phonemes required. The system is now ready to enter the recognition mode.

The method for recognition has already been discussed (see Word Recognition). Figure 9 shows a string of phoneme numbers representing five seconds of input speech. The second string represents the same speech vectors after compression has been performed. The zeroes represent points of low energy. If a zero appears alone, it is ignored in the recognition phase. If two are side-by-side, the recognition routine recognizes it as a word boundary, and starts counting the number of subsequent phonemes. If less than two phonemes are counted, the routine decides that it has made a mistake and starts counting again at the next boundary. If it counts more than nine adjacent phonemes, it

Figure 11 shows a similar event; however, in this case the word was "AIR-TO-AIR," spoken by a male speaker. The minor inaccuracy of the machine's recognition ("CLEAR CLEAR") in this case can be attributed to the fact that the system had not been retrained to the male speaker's voice.

Another male speaker trained the system and entered the word "AIR-TO-AIR". Both his phoneme string before and after compression are shown in Figure 12, along with the computer algorithm's best guess.

Similar tests were performed for continuous speech recognition. The female speaker who originally trained the system spoke the phrase, "ONE TWO THREE." Then, using her training templates, the computer algorithm recognized "AFT AFT ONE AFT FUEL AFT THREAT AFT AFT". Analyzing this string of words, one could expect the word "AFT" to represent ambient noise. The confusability of "FUEL" for "TWO" and "THREAT" for "THREE" means that not all the sounds present in each word are identified or matched. Those sounds with higher energy are detected, while softer sounds or lower energy phonemes may be deleted as a result of the energy threshold.

In summary, several trials were made using the system trained by both a male and a female speaker. The machine obviously has problems dealing with ambient noise, and usually associates its best guess phoneme string with "AFT" or "EIGHT". Also stops and pauses between syllables such as in "WAYPOINT" appear as zeroes. If the pause duration was too long the system determined that two words were spoken instead of one, and an attempt at best fit was made for two words. The other trend of confusing similar sounding words reveals that vowels are easily spotted, however, low energy fricatives are difficult to detect. This could possibly be remedied by experimenting with voltage thresholding, or the system might be made to be self adjusting. The limitation of phoneme representations to 10 phonemes may be too restrictive, and so can be introducing an excessive number of errors by taxing the systems's inability to properly detect words which are not well suited for such a compression.

VI. Conclusions

Although the system has not been thoroughly tested, several conclusions can still be drawn from the data already accumulated concerning this speech recognition system. In addition, some specific recommendations can be given for system improvement. This chapter summarizes these two areas.

Summary

The speech recognition system which was designed and implemented was speaker dependent and operated near real time (after training). Several techniques were incorporated to characterize phonemes as vectors in space. Through the use of distance rules it was possible to characterize words by a phoneme representation, which could subsequently be used in word recognition. This approach to speech recognition offers several possibilities for future investigation such as varying the Minkowski distances, and the application of clustering techniques.

The system which was designed was user friendly, providing a system of recognition which needed little user interaction. The instructions were kept to a minimum and made easy to understand, thereby taking a lot of the guess work out of trying to understand the programmer's jargon.

The modularity of the system also proved useful by allowing easy modification through the use of comment characters. Any stage of recognition can be changed by simply altering a subroutine. Another advantage to the way this system was designed is that the variables which are passed are kept to a minimum, thus reducing the number of items which must be accounted for.

The objective of designing a speech recognition system capable of operating in real time was met by this research effort. The resulting system uses a phoneme set unique to a particular speaker and partial template matching for continuous word recognition.

Recommendations

Several recommendations can be made for improving system performance. First, the method of energy thresholding should be investigated to determine an optimum range of threshold values, since the present method allows some of the low energy vowels and most consonants to be deleted. Second, word phoneme representations probably should not be limited to a length of ten. Instead, natural breaks should be retained. Third, other methods of word boundary detection could take the arbitrariness out of the present system. Specifically, the use of two adjacent zero energy vectors as an arbitrary word-word boundary could be

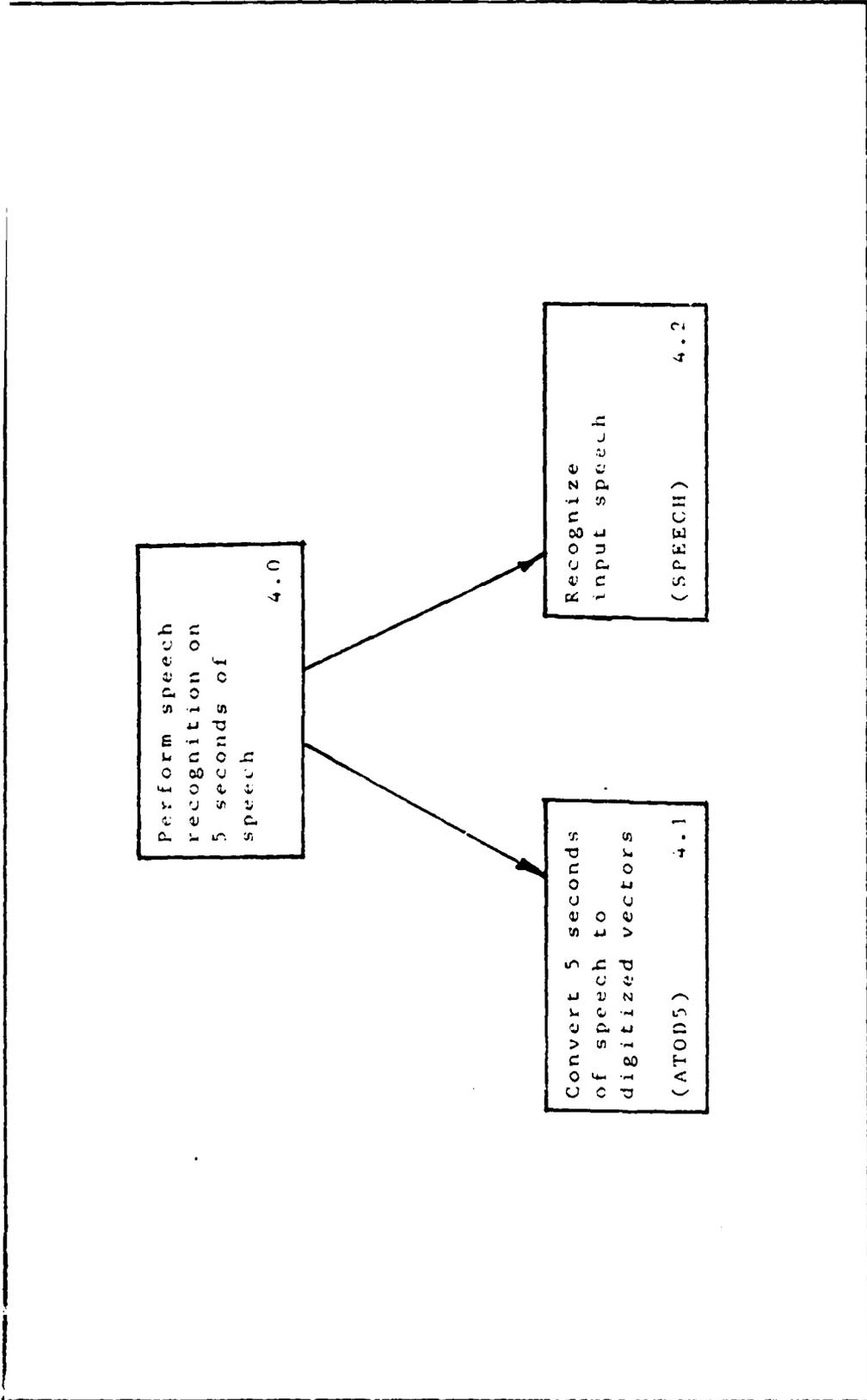


Figure 21. Structure Chart 4.0. Perform Speech Recognition

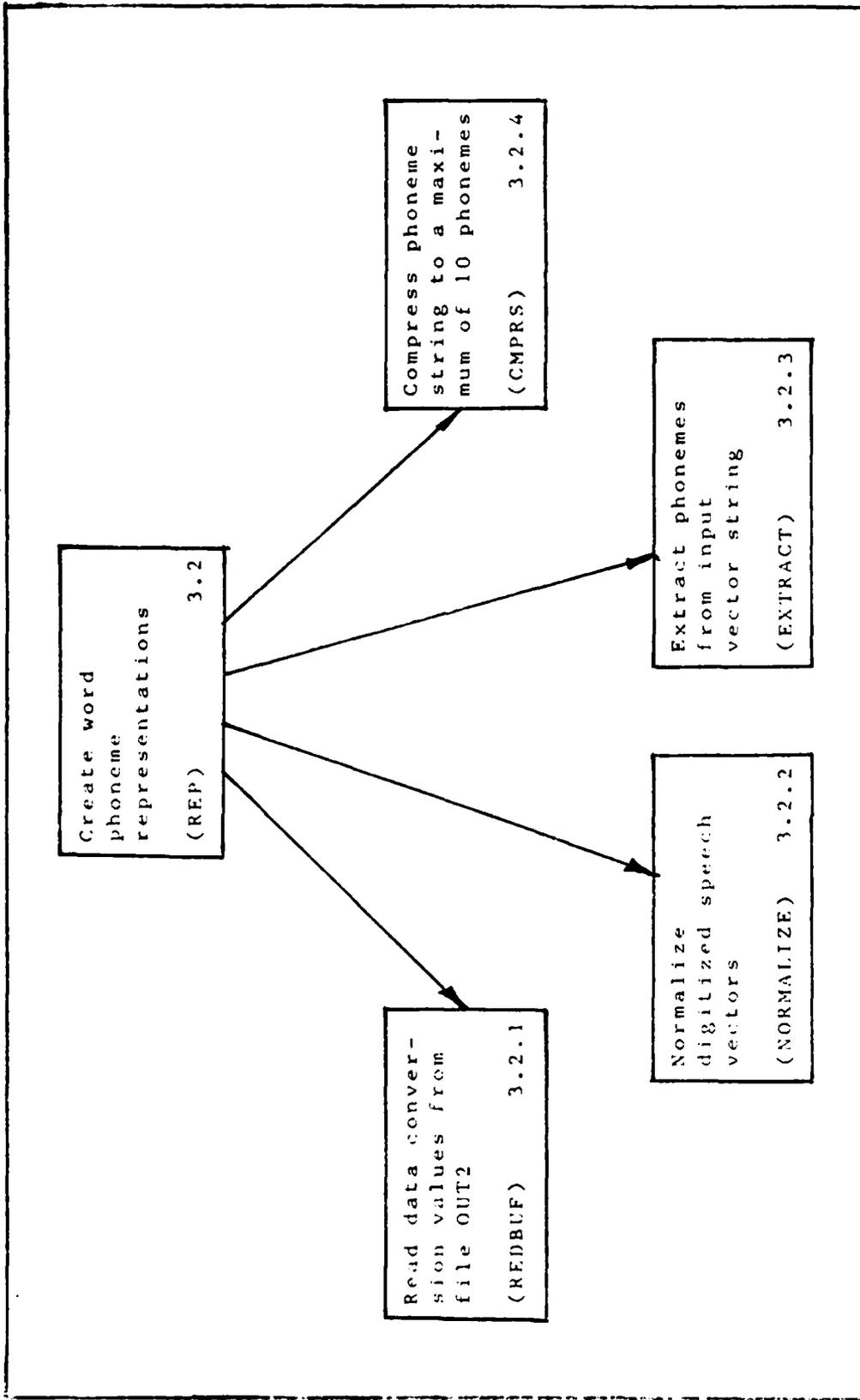


Figure 20. Structure Chart 3.2: Create Word Phoneme Representations

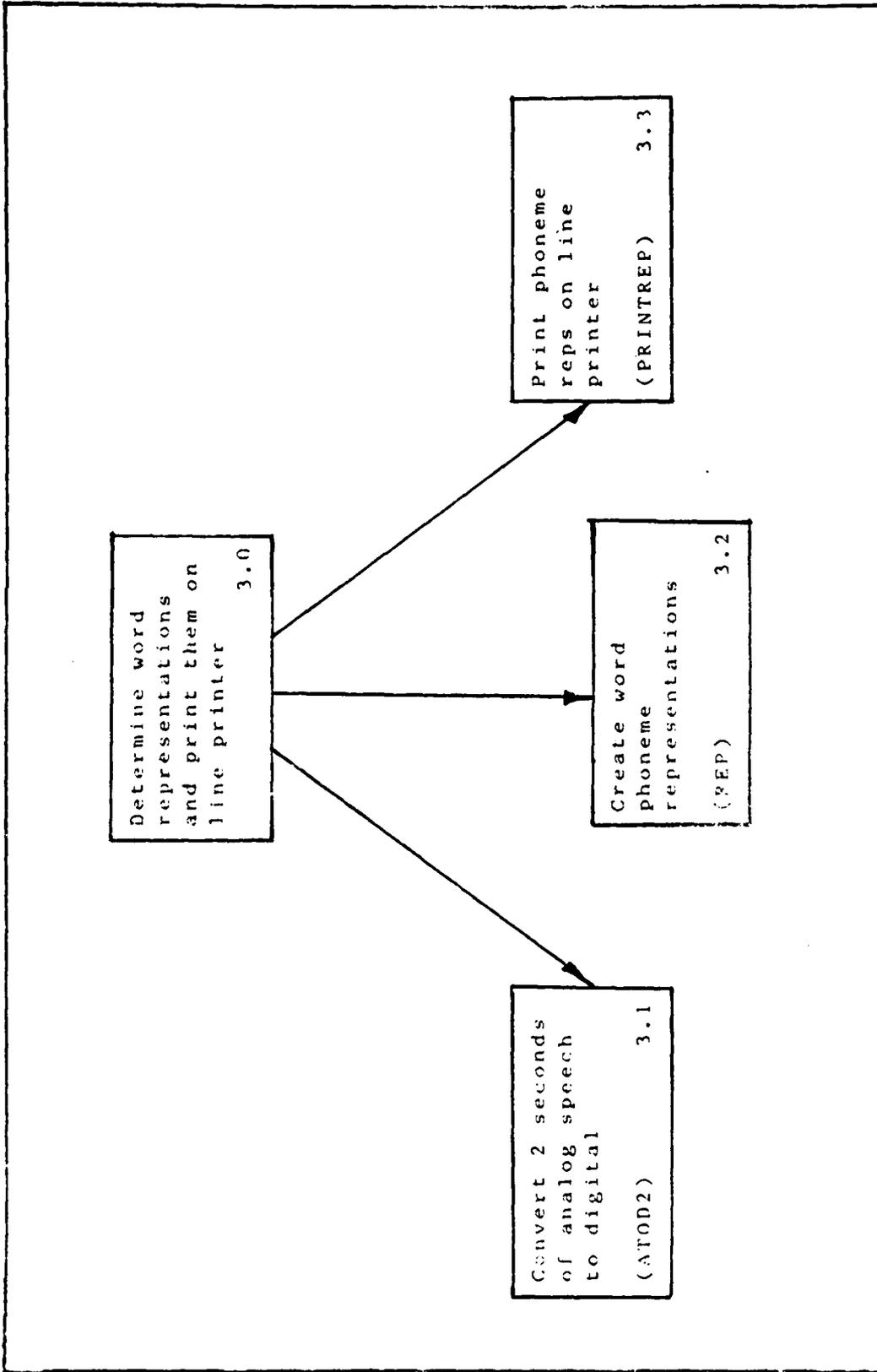


Figure 19. Structure Chart 3.0: Determine and Print Word Representations

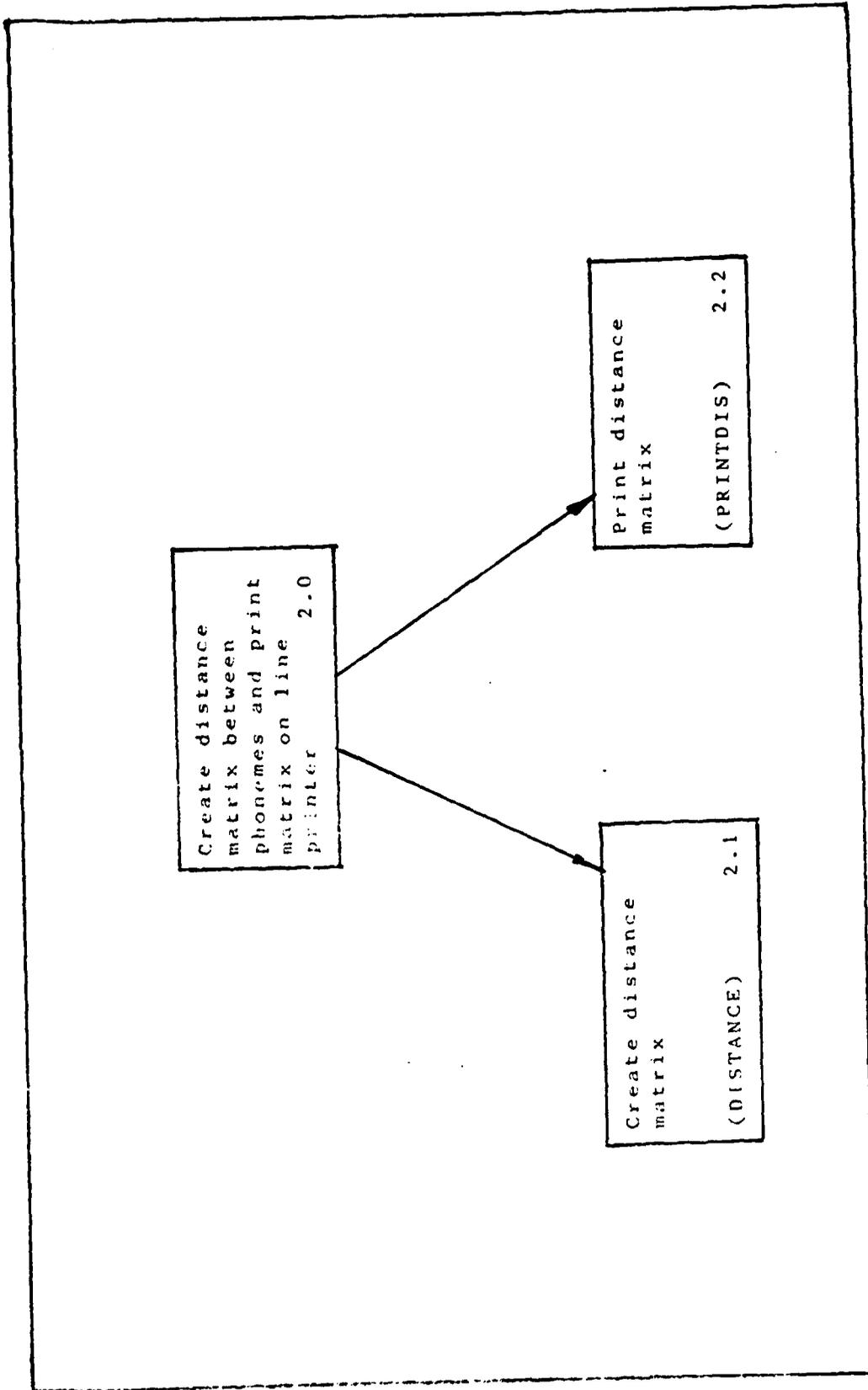


Figure 18. Structure Chart 2.0: Create and Print Distance Matrix

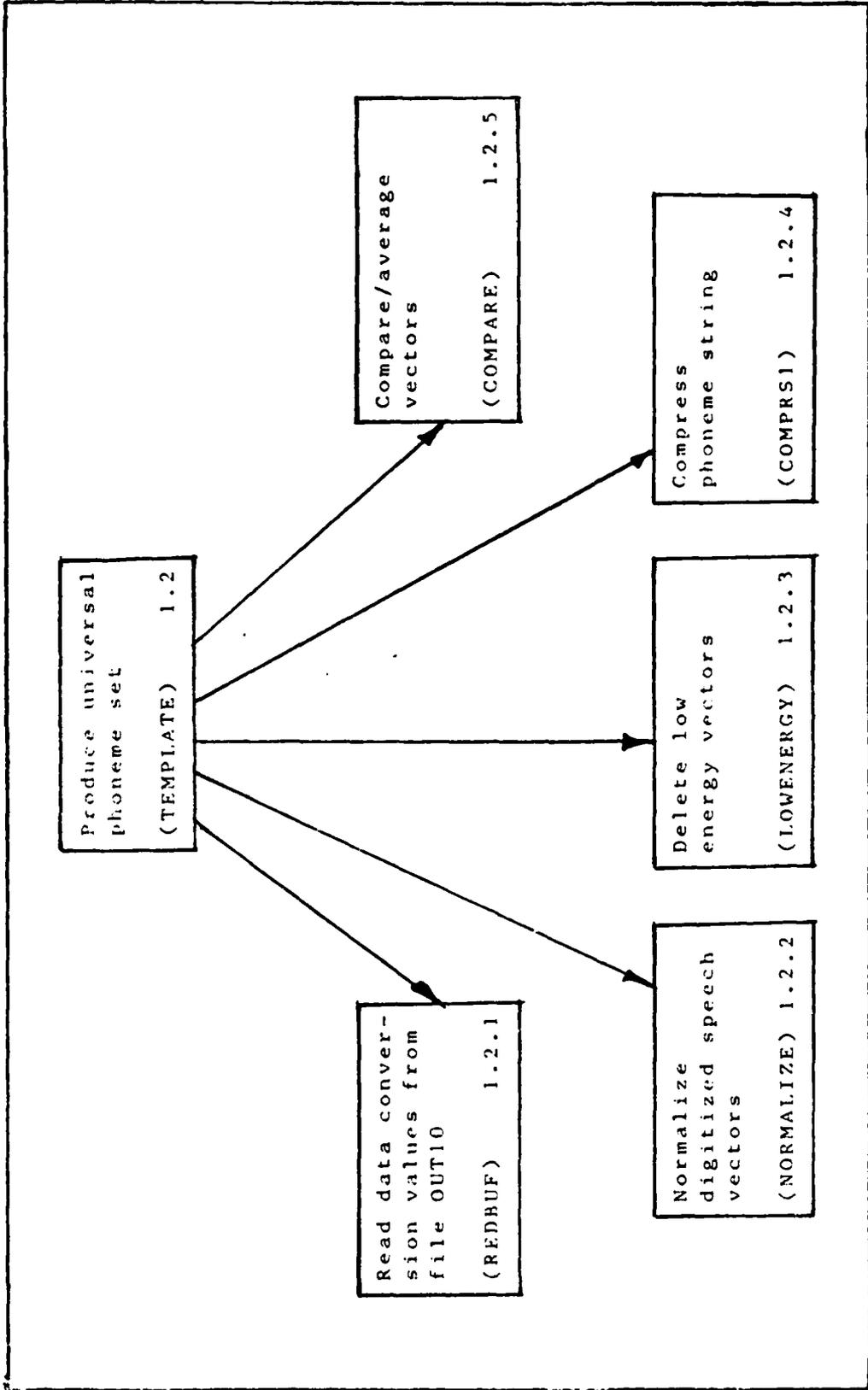


Figure 17. Structure Chart 1.2: Produce Universal Phoneme Set

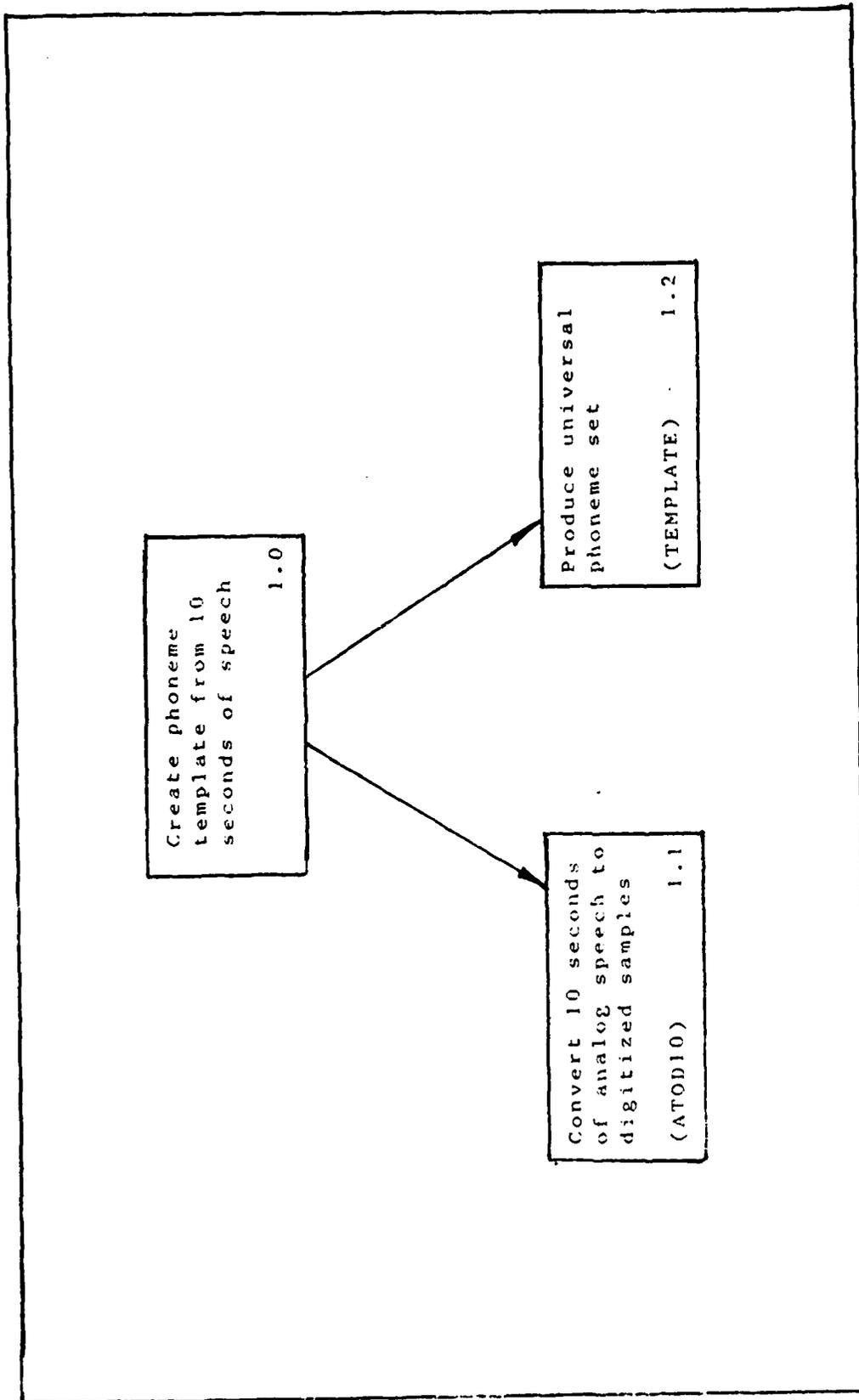


Figure 16. Structure Chart 1.0: Create Phoneme Template

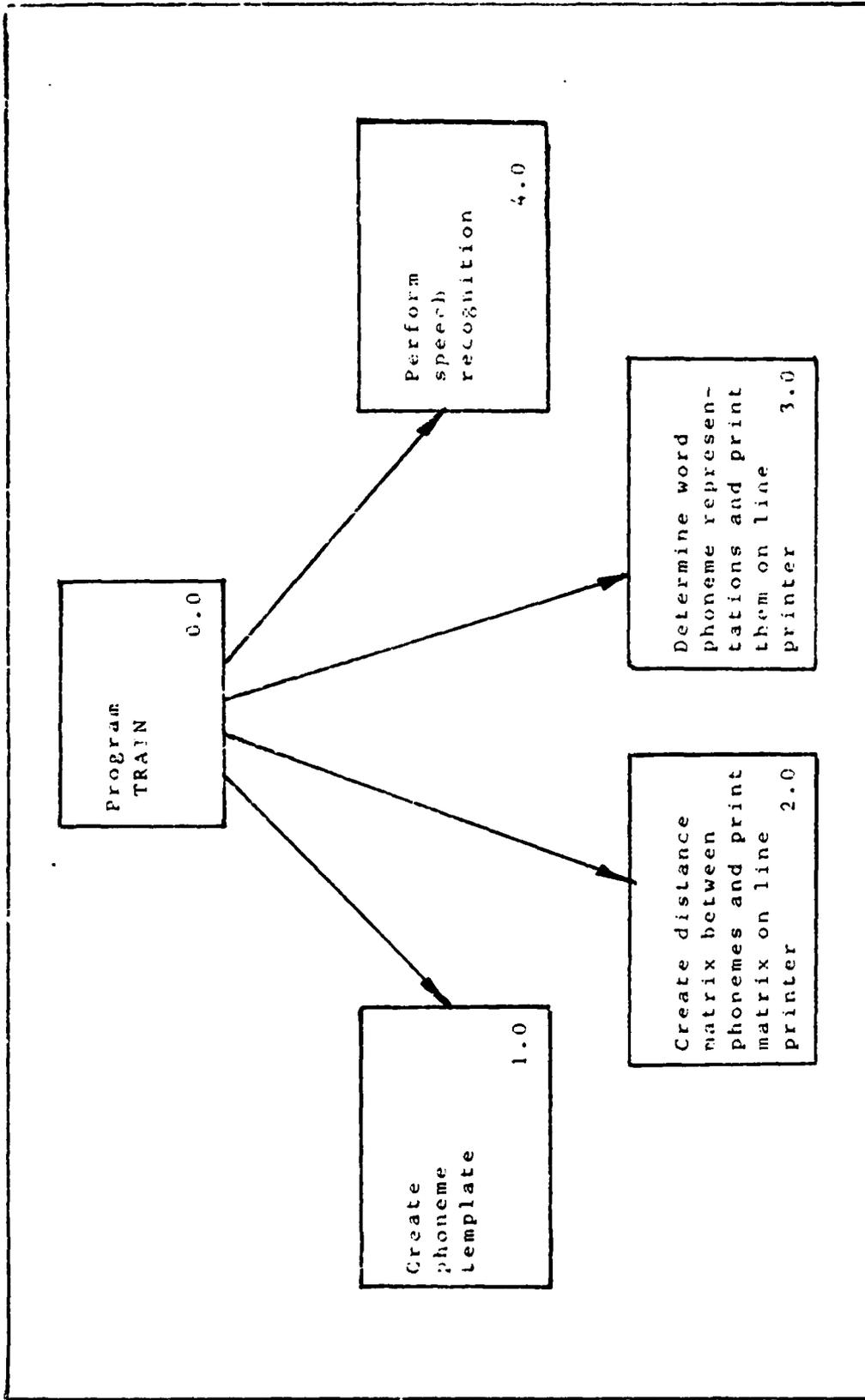


Figure 15. Structure Chart 0.0: Program TRAIN

Appendix B: Computer Program Structure Charts

Structure charts are included to provide a visual description of the proposed model for this speech recognition system. The modularity of the system allows simple modification at any location in the program.

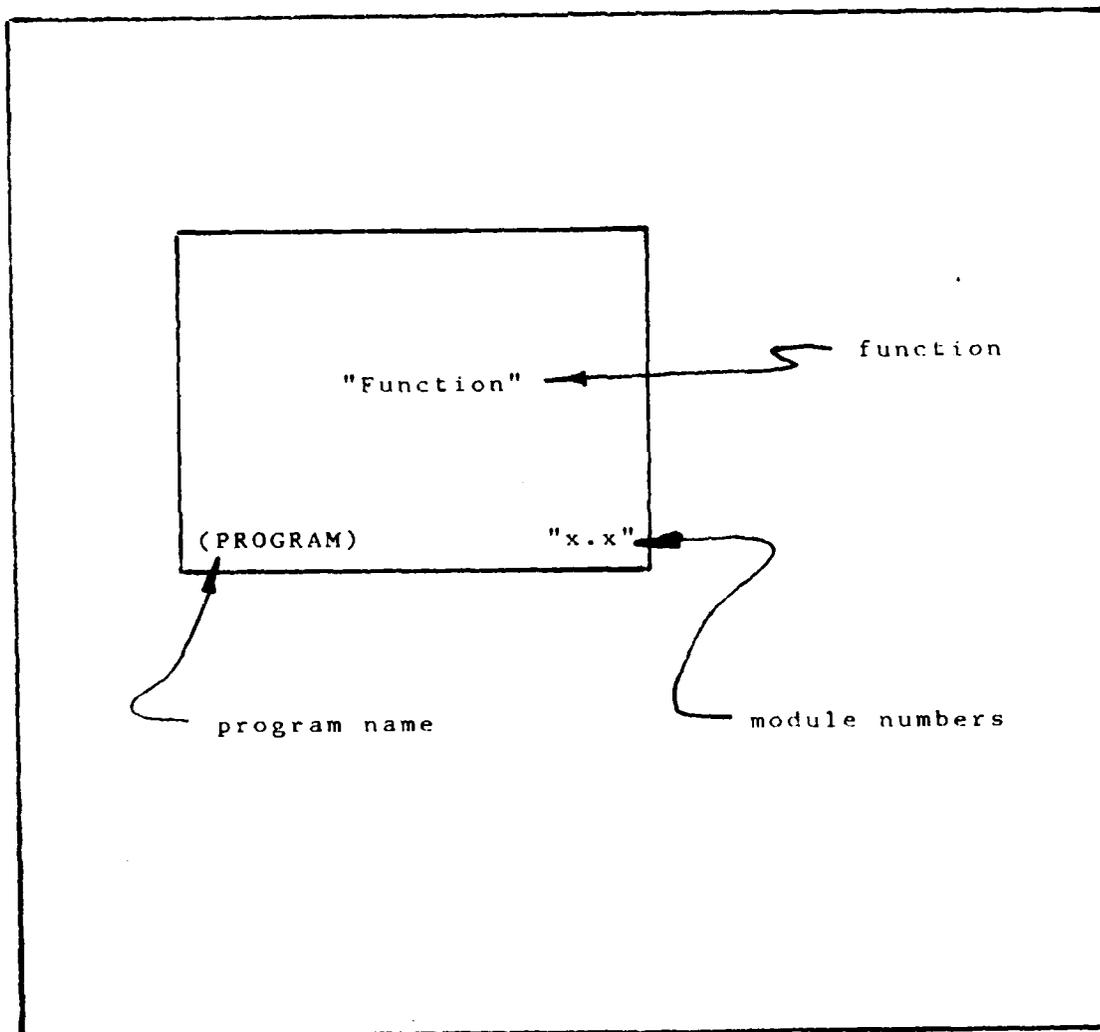


Figure 14. Legend for Structure Charts

APPENDIX B
COMPUTER PROGRAM STRUCTURE CHARTS

Appendix A: AFTI-F16 Vocabulary

Though this speech recognition system can accommodate several different vocabularies, the seventy words below best represent the vocabulary necessary to a military environment.

ADVISE	DELTA	MAP	SEVEN
AFFIRMATIVE	EAST	MARK	SIX
AFT	ECHO	MILES	SMS
AIR-TO-AIR	EIGHT	MINUS	SOUTH
AIR-TO-SURFACE	ENTER	MISSILE	STATION
ALPHA	FAULT	NEGATIVE	STRAFE
ARM	FIVE	NINE	TAIL
BACKSPACE	FLARES	NORTH	TARGET
BEARING	FORWARD	NOSE	THOUSAND
BRAVO	FOUR	ONE	THREAT
CANCEL	FOXTROT	POINT	THREE
CHAFF	FREQUENCY	PROFILE	TWO
CHANGE	FUEL	RADAR	WAYPOINT
CHARLIE	GUN	RANGE	WEAPON
CHANNEL	HEADING	REPORT	WEST
CLEAR	HUNDRED	RHAW	ZERO
CONFIRM	KNOTS	SEARCH	
DEGREES	LOCK-ON	SELECT	

Figure 13. AFTI-F16 Tentative Vocabulary (2)

APPENDIX A
AFTI-F16 VOCABULARY

Bibliography

1. Allen, Gordon R. Expansion of the Eclipse Digital Signal Processing System. MS Thesis GE/EE/82D-16. School of Engineering, Air Force Institute of Technology (AU), Wright-Patterson AFB OH, December 1982 (A124 750).
2. Anderson, Timothy R., Electrical Engineer, personal interview by ILT David B. Stockton. Air Force Aerospace Medical Research Laboratory, Wright-Patterson AFB OH, June 1984.
3. Doddington, George R. and Thomas B. Schalk. "Speech Recognition: Turning Theory to Practice," IEEE Spectrum, 18: 26-32 (September 1981).
4. Haton, Jean-Paul. "Speech Recognition and Understanding," Proceedings of the Sixth International Conference on Pattern Recognition: 570-579. Institute for Electrical and Electronics Engineers, New York, June 1982.
5. Hussain, Ajmal. Limited Continuous Speech Recognition by Phoneme Analysis. MS Thesis GE/EE/83D-31. School of Engineering, Air Force Institute of Technology (AU), Wright-Patterson AFB OH, December 1983 (A138 021).
6. Levinson, Stephen E. and Mark Y. Liberman. "Speech Recognition by Computer," Scientific American, 244: 64-76 (April 1981).
7. Montgomery, Gerard J. Isolated Word Recognition Using Fuzzy Theory. MS Thesis GE/EE/82D-74. School of Engineering, Air Force Institute of Technology (AU), Wright-Patterson AFB OH, December 1982 (A124 851).
8. Potter, R. K., George Kopp, and Harriet Green. Visible Speech. New York: D. Van Nostrand Company, Inc., 1947.
9. Rabiner, Lawrence R., Stephen E. Levinson, Aaron E. Rosenberg, and Jay G. Wilpon. "Speaker-Independent Recognition of Isolated Words Using Clustering Techniques," IEEE Transactions on Acoustics, Speech, and Signal Processing, ASSP-27: 336-349 (August 1979).
10. Seelandt, Karl G. Computer Analysis and Recognition of Phoneme Sounds in Connected Speech. MS Thesis GE/EE/81D-53. School of Engineering, Air Force Institute of Technology (AU), Wright-Patterson AFB OH, December 1981.

This system was intended to present a model for a user friendly continuous speech recognition system. Though it is by no means a perfect system, it should provide a firm foundation for future studies.

changed to a more fuzzy determination. Finally, branching techniques could be used to predict phoneme transitions.

The continuation of this research should not be difficult. Investigations into energy thresholding should be made so that the system might become self-adjusting. Additionally, other methods (besides Hussain's choice of Minkowski 4) of distance calculations should be studied to which will yield the best results. Also, the compression technique for extracted phonemes could be investigated to determine an optimum threshold for adjacent phoneme distance and to see if deleting the highest-numbered vector is really necessary. Montgomery(8) applied techniques involving average branching factors for phoneme transitions which could easily be incorporated into this system. Though ambient noise was not removed in the manner originally described, the effect of background noise on energy thresholded vectors is another subject worthy of further investigation. The choice of vocabulary for this system may not have been the best suited for it. Perhaps the phonetic alphabet (ALPHA, BRAVO, CHARLIE,...) would have more validity as these characters are easily distinguishable, having been designed for low confusability in the presence of high levels of background noise. Finally, the entire system would probably work more quickly if implemented on an array processor, since most of the calculations incorporate vector combinatorics.

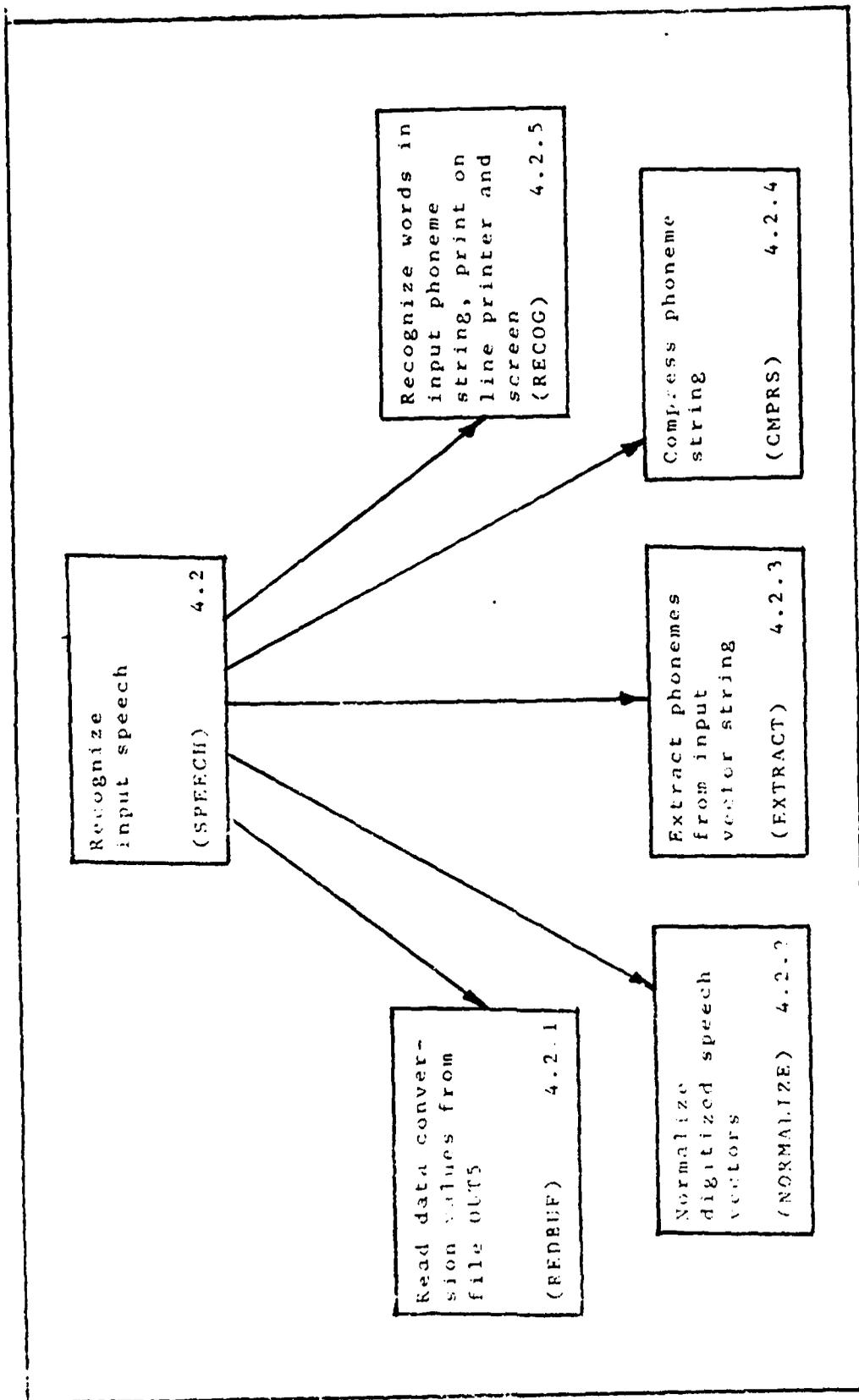


Figure 22. Structure Chart 4.2: Recognize Input Speech

APPENDIX C
COMPUTER PROGRAM LISTINGS

Appendix C: Computer Program Listings

The source code for all computer programs and support routines, written in FORTRAN V, is included in this Appendix. Two macrofiles ATODMC and TRAINMC are used to load all routines necessary for speech recognition. These programs appear in the following sequence:

PROGRAM	Page
ATODMC.MC	76
ATOD2.FR	77
ATOD5.FR	79
ATOD10.FR	81
DISS.FR	83
PHONS.FR	84
REPS.FR	85
TRAINMC.MC	86
CMPRS.FR	87
CMPRS1.FR	89
COMPARE.FR	90
DISTANCE.FR	92
EXTRACT.FR	93
FINDWORD.FR	95
LOWENERGY.FR	97
NEWSCR.FR	98
NORMALIZE.FR	99

PRINTDIS.FR	100
PRINTREP.FR	102
RECOG.FR	103
REDBUF.FR	106
REDDIS.FR	107
REDPHON.FR	108
REDREP.FR	109
REDWRDS.FR	110
REP.FR	111
SPEECH.FR	112
TEMPLATE.FR	113
TRAIN.FR	114
VTYPE.FR	118
WRTDIS.FR	119
WRTPHON.FR	120
WRTREP.FR	121
VOCAB.FR	122

*****ATODMC.MC*****

Function: load ATOD programs required by TRAIN.

RLDR/P 2/K ATOD10 SAMCONFIG3 @SAMLIB@

RLDR/P 2/K ATOD5 SAMCONFIG3 @SAMLIB@

RLDR/P 2/K ATOD2 SAMCONFIG3 @SAMLIB@

```

C*****
C
C Title: ATOD2.FR
C Author: 1Lt Kathy R. Dixon
C Date: Nov 84
C
C Function: performs A/D on Eclipse for 5 sec
C
C Command Line:
C
C RLD R/P 2/K ATOD2 SAMCONFIG3 @SAMLIB@
C*****

```

```

EXTERNAL IDS21 ;external input device
EXTERNAL IDS23 ;external output device
COMMON/IBUFF/IDATA3(16384) ;input data buffer
COMMON/IBUFO/IWAST ;output data buffer

INTEGER IORBA(16),DEVICE

DEVICE=21 ;input device
IDATA1=61700K ;external clock
IDATA2=1600 ;conversion count

TYPE"<CR>
* start<BEL><CR>"

CALL DSTRT(IER) ;initialize A/D device
IF(IER.NE.1)CALL ERROR("DSTRT ERROR")

CALL DOITW(IORBA,IDS21,8,IDATA1,IDATA2,IDATA3,IER)
IF(IER.NE.1)TYPE"DOIT ERROR",IER
IF(IORBA(14).NE.40000K)TYPE"IORBA(14) RETURN",
* IORBA(14)

TYPE"<CR>
* stop<BEL><CR>"

CALL DFILW("OUT2",IER)
IF((IER.NE.1).AND.(IER.NE.13))TYPE"DFILW ERROR",IER

CALL CFILW("OUT2",2,IER)
IF(IER.NE.1)TYPE"CFILW ERROR",IER

CALL OPEN(1,"OUT2",2,IER)
IF(IER.NE.1)TYPE"OPEN FILE ERROR",IER

CALL WRBLK(1,0,IDATA3,8,IER)
IF(IER.NE.1)TYPE"WRBLK ERROR",IER

```

```
CALL CLOSE(1, IER)  
IF(IER.NE.1)TYPE"CLOSE ERROR", IER
```

```
CALL EXIT  
END
```

```

C *****
C
C Title: ATOD5.FR
C Author: 1Lt Kathy R. Dixon
C Date: Nov 84
C
C Function: performs A/D on Eclipse for 5 sec
C
C Command Line:
C
C RLD R/P 2/K ATOD5 SAMCONFIG3 @SAMLIB@
C
C *****

```

```

EXTERNAL IDS21 ;external input device
EXTERNAL IDS23 ;external output device
COMMON/IBUFF/IDATA3(16384) ;input data buffer
COMMON/IBUFO/IWAST ;output data buffer

```

```

INTEGER IORBA(16),DEVICE

```

```

DEVICE=21 ;input device
IDATA1=61700K ;external clock
IDATA2=4000 ;conversion count

```

```

TYPE"<CR>"

```

```

* start<BEL><CR>"

```

```

CALL DSTRT(IER) ;initialize A/D device
IF(IER.NE.1)CALL ERROR("DSTRT ERROR")

```

```

CALL DOITW(IORBA,IDS21,8,IDATA1,IDATA2,IDATA3,IER)
IF(IER.NE.1)TYPE"DOIT ERROR",IER
IF(IORBA(14).NE.4000K)TYPE"IORBA(14) RETURN",

```

```

* IORBA(14)

```

```

TYPE"<CR>"

```

```

* stop<BEL><CR>"

```

```

CALL DFILW("OUT5",IER)
IF((IER.NE.1).AND.(IER.NE.13))TYPE"DFILW ERROR",IER

```

```

CALL CFILW("OUT5",2,IER)
IF(IER.NE.1)TYPE"CFILW ERROR",IER

```

```

CALL OPEN(1,"OUT5",2,IER)
IF(IER.NE.1)TYPE"OPEN FILE ERROR",IER

```

```

CALL WRBLK(1,0,IDATA3,16,IER)
IF(IER.NE.1)TYPE"WRBLK ERROR",IER

```

```
CALL CLOSE(1, IER)  
IF(IER.NE.1)TYPE"CLOSE ERROR", IER
```

```
CALL EXIT  
END
```

```

C*****
C
C Title: ATOD10.FR
C Author: 1Lt Kathy R. Dixon
C Date: Nov 84
C
C Function: performs A/D on Eclipse for 10 sec.
C
C Command Line:
C
C RLD R/P 2/K ATOD10 SAMCONFIG3 @SAMLIB@
C*****
EXTERNAL IDS21 ;external input device
EXTERNAL IDS23 ;external output device
COMMON/IBUFF/IDATA3(16384) ;input data buffer
COMMON/IBUFO/IWAST ;output data buffer

INTEGER IORBA(16),DEVICE

DEVICE=21 ;input device
IDATA1=61700K ;external clock
IDATA2=8000 ;conversion count

TYPE"<CR>
* start<BEL><CR>"

CALL DSTRT( IER) ;initialize A/D device
IF( IER.NE.1)CALL ERROR("DSTRT ERROR")

CALL DOITW( IORBA,IDS21,8, IDATA1, IDATA2, IDATA3, IER)
IF( IER.NE.1)TYPE"DOIT ERROR", IER
IF( IORBA(14).NE.40000K)TYPE" IORBA(14) RETURN",
* IORBA(14)

TYPE"<CR>
* stop<BEL><CR>"

CALL DFILW("OUT10", IER)
IF( IER.NE.1.AND. IER.NE.13)TYPE"DFILW ERROR", IER

CALL CFILW("OUT10",2, IER)
IF( IER.NE.1)TYPE"CFILW ERROR", IER

CALL OPEN(1,"OUT10",2, IER)
IF( IER.NE.1)TYPE"OPEN FILE ERROR", IER

CALL WRBLK(1,0, IDATA3,32, IER)
IF( IER.NE.1)TYPE"WRBLK ERROR", IER

```

```
CALL CLOSE(1, IER)  
IF(IER.NE.1)TYPE"CLOSE ERROR", IER
```

```
CALL EXIT  
END
```

```
C*****
C
C Title: DISS.FR
C Author: 1Lt Kathy Dixon
C Date: Nov 84
C
C Function: Prints distance matrix from a disk file.
C
C*****
```

```
REAL DIS(2432)
```

```
CALL OPEN(1,"DIST",2,IER)
IF(IER.NE.1)TYPE"OPEN ERROR",IER
READ(1,100)(DIS(I),I=1,2432)
WRITE(12,101)(DIS(I),I=1,2432)
```

```
100
```

```
FORMAT(G11.5)
```

```
101
```

```
FORMAT(8G11.5)
```

```
CALL CLOSE(1,IER)
```

```
IF(IER.NE.1)TYPE"CLOSE ERROR",IER
```

```
RETURN
```

```
END
```

```
C*****
C
C   Title:  PHONS.FR
C   Author: 1Lt Kathy Dixon
C   Date:   Nov 84
C
C   Function: Reads up to 70 16-dimensional vectors
C             from a file, PHONE and prints them on
C             the line printer.
C
C*****
```

```
INTEGER PHON(1130)
```

```
CALL OPEN(1,"PHONE",2,IER)
IF(IER.NE.1)TYPE"OPEN ERROR",IER
READ(1,100)(PHON(I),I=1,1130)
WRITE(12,101)(PHON(I),I=1,1130)
```

```
100
```

```
101
```

```
FORMAT(I6)
FORMAT(2X,16I6)
CALL CLOSE(1,IER)
IF(IER.NE.1)TYPE"CLOSE ERROR",IER

RETURN
END
```

```
C*****
C
C Title: REPS.FR
C Author: 1Lt Kathy Dixon
C Date: Nov 84
C
C Function: Reads word phonemes from a file
C VOCABUL.
C
C*****
```

```
INTEGER VOCAB(10,70),NOWORDS
NOWORDS=70
```

```
CALL OPEN(1,"VOCABUL",2,IER)
IF(IER.NE.1)TYPE"OPEN ERROR",IER
READ(1,100)((VOCAB(I,J),I=1,10),J=1,NOWORDS)
WRITE(12,101)((VOCAB(I,J),I=1,10),J=1,NOWORDS)
100 FORMAT(I3)
101 FORMAT(10I3)
CALL CLOSE(1,IER)
IF(IER.NE.1)TYPE"CLOSE ERROR",IER

RETURN
END
```

AO-R151 936

IMPLEMENTATION OF A REAL-TIME INTERACTIVE CONTINUOUS
SPEECH RECOGNITION SYSTEM(U) AIR FORCE INST OF TECH
WRIGHT-PATTERSON AFB OH SCHOOL OF ENGI... K R DIXON
DEC 84 AFIT/GE/ENG/84D-26 F/G 9/2

2/2

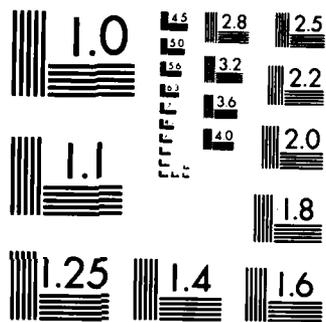
UNCLASSIFIED

NL

END

FILED

ERIC



MICROCOPY RESOLUTION TEST CHART
NATIONAL BUREAU OF STANDARDS 1963-A

*****TRAINMC.MC*****

Function: Loads TRAIN and required subroutines

RLDR TRAIN NEWSCR TEMPLATE DISTANCE PRINTDIS REDWRDS REP^
PRINTREP SPEECH REDBUF NORMALIZE LOWENERGY CMPSI^
COMPARE EXTRACT CMPS RECOG FINDWORD VTYPE^
REDDIS REDPHON REDREP WRTDIS WRTREP WRTPHON @FLIB@

```

C*****
C
C Title:  CMPRS.FR
C Author: Capt. Ajmal Hussain
C Modified by 1Lt Kathy Dixon
C Date:   Nov 84
C
C Function:
C Compresses input speech vectors after phoneme
C extraction.
C*****

```

```

SUBROUTINE CMPRS(IDATA2,PHON,DIS,LIB,ISWITCH,SOUND,J)
INTEGER SOUND(250),LIB(250),PHON(1130),IDATA2,ISWITCH
REAL DIS(2432),LDIS

NOS = (1+ISWITCH)*125           ; no. of vectors for
                                ; processing
LDIS=10.0                       ;set distance threshold

DO 76 I=1,250
76 SOUND(I)=0

C WRITE(12,35)(LIB(I),I=1,NOS)   ;option to write
35 FORMAT(25I3)                 ;phoneme string to line printer

1 DO 806 K=1,5
DO 809 I=1,((IDATA2/16)-1)
IF(LIB(I).EQ.0)GO TO 807
IF(LIB(I+1).EQ.0)GO TO 807
IF(LIB(I).EQ.LIB(I+1))GO TO 807
N=LIB(I)
P=LIB(I+1)
IF(N.GT.P)GO TO 810
Q=N
N=P
P=Q
810 IF(DIS(((N*(N-1))/2)+P).GE.LDIS)GO TO 807
IF(LIB(I).LT.LIB(I+1))LIB(I+1)=LIB(I)
IF(LIB(I).GT.LIB(I+1))LIB(I)=LIB(I+1)
807 CONTINUE
809 CONTINUE
806 CONTINUE

J=1
DO 805 I=1,(IDATA2/16)
IF((LIB(I).EQ.0).AND.(J.EQ.1))GO TO 805
IF((LIB(I).EQ.0).AND.(LIB(I-1).NE.0))GO TO 805

```

```

IF(LIB(I).EQ.LIB(I+1))GO TO 805
LIB(J)=LIB(I)
J=J+1
805 CONTINUE
DO 808 L=J,(IDATA2/16)
808 LIB(L)=0

IF(ISWITCH.EQ.1)GO TO 1001 ;check to see if
LDIS=LDIS+0.5 ;need to compress
IF(J.LE.10)GO TO 1001 ;to 10 phonemes
GO TO 1 ;threshold

1001 DO 900 I=1,250 ;store word phonemes
900 SOUND(I)=LIB(I) ;into SOUND

1000 CONTINUE
C WRITE(12,35)(LIB(I),I=1,NOS) ;option to print
;phonemes on line printer
;after compression

RETURN
END

```

```

C*****
C
C   Title:  CMPRS1.FR
C   Author: 1LT Kathy Dixon
C           Based on subroutine by CAPT Ajmal Hussain
C   Date:   Nov 84
C
C   Function:
C           Deletes vectors with weight 0.
C
C*****

```

```

SUBROUTINE CMPRS1(IDATA2,IWEIGHT,IWGHTT)

COMMON/IBUFF/IDATA(8192)

INTEGER IDATA2 ,IWEIGHT(500),IWGHTT(500)

J=1
M=1

DO 1 I=1,500
1  IWGHTT(I)=0

DO 2 I=1,IDATA2,16 ;check vector components
IF(IWEIGHT(INT((I+16)/16)).EQ.0)GO TO 2

DO 3 K=0,15
IDATA(J)=IDATA(I+K) ;replace deleting vector
J=J+1 ;with next vector

IWEIGHT(M)=IWEIGHT(INT((I+16)/16))
M=M+1

2  CONTINUE

DO 4 I=M,500

DO 5 K=1,16
5  IDATA(((I-1)*16)+K)=0

4  IWEIGHT(I)=0

DO 6 I=1,500
6  IWGHTT(I)=IWEIGHT(I)

RETURN
END

```

```

C*****
C
C Title: COMPARE.FR
C Author: 1Lt Kathy Dixon
C Nearest phoneme code derived from
C CAPT Ajmal Hussain
C Date: Nov 84
C
C Function:
C Compares normalized data vectors. Finds each
C vector's nearest vector. Averages two vectors,
C replaces lowest numbered vector with new vector,
C sets vector components of second vector to 32000.
C
C*****

```

```

SUBROUTINE COMPARE(IDATA2)

COMMON/IBUFF/IDATA(8192)

INTEGER IDATA5(500),NOVECT, IDATA2,TEMP4, IDATA6
REAL DIFF(500)
DOUBLE PRECISION REAL TEMP,TEMP1

TEMP=0
TEMP1=9.0E60

IDATA6=IDATA(IDATA2+1)*16

DO 104 J=1, IDATA6,16
DO 102 K=1, IDATA6,16

IF(J.EQ.K) GO TO 103

DO 101 L=0,15
101 TEMP=TEMP+(FLOAT(IDATA(J+L))-FLOAT(IDATA(K+L)))*.4

IF(TEMP.GE.TEMP1) GO TO 103
TEMP1=TEMP

IDATA5(INT((J+15)/16))=INT((K+15)/16)
DIFF(INT((J+15)/16))=TEMP
103 TEMP=0
102 CONTINUE

TEMP1=9.0E60
104 CONTINUE

```

```

TEMP=0

KKK=0

NOVECT=IDATA(IDATA2+1)

DO 107 I=1,NOVECT
107 IF(DIFF(I).GT.TEMP)TEMP=DIFF(I)

DO 108 I=1,NOVECT
108 DIFF(I)=(DIFF(I)/TEMP)

DO 111 I=1,NOVECT
J=IDATA5(I)
TEMP=DIFF(I)

DO 2 JJ=1,NOVECT
K=0
IF(IDATA5(JJ).EQ.I)TEMP1=DIFF(JJ)
IF((TEMP1.EQ.0).OR.(TEMP1.EQ.100))GO TO 2
IF(TEMP1.LT.TEMP)K=JJ
IF(J.EQ.K)GO TO 10
2 CONTINUE

GO TO 111

10 DO 15 KL=1,16
IDATA(((I-1)*16)+KL)=(IDATA(((I-1)*16)+KL)
* +IDATA(((J-1)*16)+KL))/2
15 IDATA(((J-1)*16)+KL)=32000

DIFF(J)=100

DO 20 JJJ=1,NOVECT
IF((IDATA5(JJJ).EQ.I).AND.(DIFF(JJJ).NE.100))
* DIFF(JJJ)=0
20 CONTINUE

111 CONTINUE

DO 222 I=1,NOVECT
222 IF(DIFF(I).EQ.100)KKK=KKK+1

IDATA(IDATA2+1)=IDATA(IDATA2+1)-KKK

RETURN
END

```

```

C*****
C
C   Title:  DISTANCE.FR
C   Author: Capt. Ajmal Hussain
C           Modified by 1 Lt Kathy Dixon
C   Date:   Nov 84
C
C   Function:
C           Finds Minkowski 4 distance between phonemes in
C           phoneme template.
C*****

      SUBROUTINE DISTANCE(PHON)

      REAL DIS(2432)
      INTEGER PHON(1130)
      DOUBLE PRECISION REAL TEMP, TEMP1

10      DO 10 I=1,2432                ;zero distance matrix
      DIS(I)=0.

      TEMP=0
      TEMP1=0
      I=1
      DO 31 J=1,(PHON(1121)*16),16
      DO 32 K=1,(PHON(1121)*16),16
      IF(K.GT.J) GO TO 35
      DO 33 L=0,15
33      TEMP=TEMP+(FLOAT(PHON(J+L))-FLOAT(PHON(K+L)))**4
                                   ;M-4 calculation

      IF(TEMP.GT.TEMP1)TEMP1=TEMP    ;find largest distance
      DIS(I)=TEMP                   ;store distance
      I=I+1                          ;increment DIS
35      TEMP=0
32      CONTINUE
31      CONTINUE
      DO 34 I=1,(((PHON(1121)*(PHON(1121)-1))/2)
* +PHON(1121))
34      DIS(I)=((DIS(I)/TEMP1)**0.25)*100 ;normalize DIS to
                                   ;TEMP1
      DIS(2416)=PHON(1121)          ;store no. phonemes

      CALL WRDIS(DIS)

      RETURN
      END

```

```

*****
C
C Title: EXTRACT.FR
C Author: Capt. Ajmal Hussain
C Modified by 1Lt Kathy Dixon
C Date: Nov 84
C
C Function:
C Extracts phonemes from input speech file.
C
*****

```

```

SUBROUTINE EXTRACT(IDATA2,PHON,IENERGY,ENRGY,LIB)

COMMON/IBUFF/IDATA(4096)

INTEGER LIB(250),PHON(1130)
REAL IENERGY(500),ENRGY
DOUBLE PRECISION REAL TEMP,TEMP1

DO 1 I=1,250
1 LIB(I)=0 ;zero array to hold
;extracted phonemes

TEMP=0
I=1
M=0
TEMP1=9.0E60

DO 87 L=1, IDATA2, 16

IF(IENERGY(I).LT.ENRGY) GO TO 801 ;check energy
;threshold

DO 86 K=1,(PHON(1121)*16),16

DO 85 J=L,(L+15)
TEMP=TEMP+(FLOAT(IDATA(J)-PHON(K+M)))**4 ;m-4
;distance
85 M=M+1 ;between template

IF(TEMP.GE.TEMP1)GO TO 82 ;vector

TEMP1=TEMP
LIB(I)=(K+15)/16

82 TEMP=0
M=0

86 CONTINUE
801 TEMP1=9.0E60
87 I=I+1

```

RETURN
END

```

C*****
C
C Title: FINDWORD.FR
C Author: Capt. Ajmal Hussain
C Modified by 1Lt Kathy Dixon
C Date: Nov 84
C
C Function:
C This routine compares a phoneme string with word
C strings in a library based upon a distance matrix
C to give the word in the library which is the best
C match.
C
C*****

```

```

SUBROUTINE FINDWORD(PHON,IPHON,VOCAB,MAT,TEMP3,L)

```

```

INTEGER I,J,K,L,M,IPHON(250),LIB(700),NOWORD
INTEGER VOCAB(10,70),PHON(1130)
DOUBLE PRECISION REAL TEMP1,TEMP,TEMP3
REAL MAT(2432),PEN

```

```

70 TEMP3 = 9.0E 60 ;initialize variables
TEMP = 0
COUNT = 0
NOWORD=70
J=1
L=0

```

```

DO 1 I=1,70
DO 2 K=1,10
LIB(J)=VOCAB(K,I)
2 J=J+1
1 CONTINUE

```

```

C start comparison

```

```

DO 71 M = 1,(NOWORD*10),10 ;library, each word a
;maximum of 10 phonemes

```

```

DO 72 K = -1,1 ;shift phoneme string one phoneme
;left, none and one phoneme right to
;account for error in first phoneme
;string

```

```

DO 73 I = 1,10 ;compare phoneme at a time for each
;word in library

```

```

IF ((I+K).EQ.0) GO TO 73 ;skip first phoneme when
;string shifted left one

```

```

                                ;phoneme

C      if both phonemes zero error value unchanged

      IF ((LIB(M+I-1).EQ.0).AND.(IPHON(I+K).EQ.0)) GO TO 73

C      if both phonemes not zero add distance between
C      phonemes to error value

      IF ((LIB(M+I-1).NE.0).AND.(IPHON(I+K).NE.0)) GO TO 74

C      if one phoneme zero only add penalty to error value

      P=LIB(M+I-1)
      PEN=0

      DO 135 JJ=0,15
135    PEN = PEN+ (FLOAT(PHON((P*16)-JJ))**4)
      PEN =((PEN)**0.25)

      TEMP1 = TEMP+PEN
      GO TO 75

74    N = IPHON(I+K)                ;find distance between
                                   ;phonemes from
                                   ;distance matrix
      P = LIB(M+I-1)
      IF(N.GE.P) GO TO 76
      Q = N
      N = P
      P = Q
76    TEMP1 = MAT(((N*(N-1))/2)+P)

75    TEMP = TEMP + TEMP1           ;add distance to
      COUNT = COUNT + 1
73    CONTINUE                     ;average error value
                                   ;and find word
      TEMP = TEMP/COUNT            ;match with minimum
                                   ;error value

      IF(TEMP.GT.TEMP3) GO TO 77
      TEMP3 = TEMP
      L = M

77    TEMP = 0                     ;initialize variables
                                   ;for next word

      COUNT = 0
72    CONTINUE
71    CONTINUE

      RETURN
      END

```

```

C*****
C
C Title: LOWENERGY.FR
C Author: Capt. Ajmal Hussain
C Modified by Kathy Dixon
C Date: Nov 84
C
C Function:
C Checks energy normalized vectors for energy less
C than a particular threshold. Sets vector weights
C to 0.
C*****

```

```

SUBROUTINE LOWENERGY(IDATA2, IENERGY, ENRGY, IWG)

INTEGER IDATA2, IWG(500)
REAL IENERGY(500), ENRGY

COMMON/IBUFF/IDATA(8192)

K=0

DO 1 I=1,500
1 IWG(I)=1

DO 2 L=1,(IDATA2/16)
IF(IENERGY(L).GE.ENRGY) GO TO 2 ;check phoneme
IWG(L)=0 ;energy
K=K+1
2 CONTINUE

IDATA(IDATA2+1)=(IDATA2/16)-K ;store no. of
;REMAINING VECTORS

RETURN
END

```

```

C*****
C
C Title: REP.FR
C Author: 1Lt Kathy R. Dixon
C Date: Nov 84
C
C Function:
C Creates phoneme representation of a word.
C Stores representation in array, SOUND.
C
C*****

```

```

SUBROUTINE REP(ENRGY, PHON, DIS, SOUND)

COMMON/IBUFF/IDATA(4096)

INTEGER PHON(1130), SOUND(250), LIB(250), IDATA2
INTEGER ISWITCH, ISTOP, IFILE
REAL DIS(2432), IENERGY(500), ENRGY

IDATA2=1600 ;conversion count
ISTOP=8 ;last block to read
IFILE=2 ;file to read
ISWITCH=0 ;switch to reduce input
;phoneme string to max 10

CALL REDBUF(IFILE, ISTOP) ;read data from file

CALL NORMALIZE(IDATA2, IENERGY) ;normalize vectors

CALL EXTRACT(IDATA2, PHON, IENERGY, ENRGY, LIB) ;extract
;phonemes

CALL CMPRS(IDATA2, PHON, DIS, LIB, ISWITCH, SOUND, J)
;compress phoneme string

RETURN
END

```

```
C*****
C
C Title: REDWRDS.FR
C Author: 1Lt Kathy R. Dixon
C Date: Nov 84
C
C Function:
C Reads vocabulary: words from a file, WORDS.
C*****
```

```

SUBROUTINE REDWRDS(NOWORDS,WORD)

INTEGER WORD(7,70),NOWORDS

CALL OPEN(1,"WORDS",2,IER)
IF(IER.NE.1)TYPE"OPEN ERROR",IER
100 READ(1,100)(WORD(1,I),I=1,NOWORDS)
FORMAT(S14)
CALL CLOSE(1,IER)
IF(IER.NE.1)TYPE"CLOSE ERROR",IER

RETURN
END
```

```
C*****
C
C   Title:  REDREP.FR
C   Author: Kathy Dixon
C   Date:   Nov 84
C
C   Function:  Reads word phonemes form a file.
C
C*****
```

```
      SUBROUTINE REDREP(NOWORDS,VOCAB)
```

```
      INTEGER VOCAB(10,70),NOWORDS
```

```
      CALL OPEN(1,"VOCABUL",2,IER)
```

```
      IF(IER.NE.1)TYPE"OPEN ERROR",IER
```

```
      READ(1,100)((VOCAB(I,J),I=1,10),J=1,NOWORDS)
```

```
100
```

```
      FORMAT(I3)
```

```
      CALL CLOSE(1,IER)
```

```
      IF(IER.NE.1)TYPE"CLOSE ERROR",IER
```

```
      RETURN
```

```
      END
```

```
C*****
C
C Title: REDPHON.FR
C Author: Kathy Dixon
C Date: Nov 84
C
C Function: Reads 70 16-dimensional vectors
C           from a file PHONE.
C
C*****
```

```
      SUBROUTINE REDPHON(PHON)
```

```
      INTEGER PHON(1130)
```

```
      CALL OPEN(1,"PHONE",2,IER)
      IF(IER.NE.1)TYPE"OPEN ERROR",IER
      READ(1,100)(PHON(I),I=1,1130)
      FORMAT(I6)
      CALL CLOSE(1,IER)
      IF(IER.NE.1)TYPE"CLOSE ERROR",IER
```

100

```
      RETURN
      END
```

```
C*****
C
C   Title:  REDDIS.FR
C   Author: 1Lt Kathy Dixon
C   Date:   Nov 84
C
C   Function: Reads distance matrix from a disk file.
C
C*****
```

```
      SUBROUTINE REDDIS(DIS)
```

```
      REAL DIS(2432)
```

```
      CALL OPEN(1,"DIST",2,IER)
```

```
      IF(IER.NE.1)TYPE"OPEN ERROR",IER
```

```
      READ(1,100)(DIS(I),I=1,2432)
```

```
100
```

```
      FORMAT(G11.5)
```

```
      CALL CLOSE(1,IER)
```

```
      IF(IER.NE.1)TYPE"CLOSE ERROR",IER
```

```
      RETURN
```

```
      END
```

```

C*****
C
C   Title:  REDBUF.FR
C   Author: 1Lt Kathy R. Dixon
C   Date:   Nov 84
C
C   Function:
C           Reads A/D data from a file.
C
C*****

```

```

SUBROUTINE REDBUF(IFILE,ISTOP)

INTEGER IFILE,ISTOP

COMMON/IBUFF/IDATA(8192)

IF(IFILE.EQ.10)GO TO 8           ;read OUT10
IF(IFILE.EQ.5)GO TO 7          ;READ OUT 5

CALL OPEN(1,"OUT2",2,IER)
IF(IER.NE.1)TYPE"OPEN ERROR",IER

GO TO 9

8  CALL OPEN(1,"OUT10",2,IER)
   IF(IER.NE.1)TYPE"OPEN ERROR",IER

   GO TO 9

7  CALL OPEN(1,"OUT5",2,IER)
   IF(IER.NE.1)TYPE"OPEN ERROR",IER

9  CALL RDBLK(1,0,IDATA,ISTOP,IER)
   IF(IER.NE.1)TYPE"RDBLK ERROR",IER

CALL CLOSE(1,IER)
IF(IER.NE.1)TYPE"CLOSE ERROR",IER

RETURN
END

```

```
DO 182 I=1,10
182 CALL VTYPE(REJ(Z,I),WORD)
GO TO 79
174 W=W+U
V=W-1
GO TO 171

179 TOT(Y)=TOT(Y)/X
Y=Y+1
X=1
FLAG=1
V=1
W=1
GO TO 171

1000 RETURN
END
```

```

        S=1
        GO TO 172

C      CONTINUOUS RECOGNITION

177    V=1
        W=1
        X=1
        Y=1
        FLAG=0

171    DO 176 I=1,10
176    TWORD(I)=0

        TEMP3=9.0E60

        DO 173 I=1,10
        TWORD(I)=WRD(V)
        V=V+1
        IF(I.LE.2)GO TO 173

        CALL FINDWORD(PHON,TWORD,VOCAB,DIS,TEMP,L)
        IF(Y.LE.1)GO TO 180

        IF((L.EQ.REJ((Y-1),(Y-1))).AND.(FLAG.EQ.1))GO TO 186

180    IF(TEMP.GE.TEMP3)GO TO 173

        TEMP3=TEMP
        T=L
        U=I
        GO TO 173

186    FLAG=0
173    CONTINUE

        REJ(Y,X)=T
        TOT(Y)=TOT(Y)+TEMP3
        X=X+1
        IF((W+U).LT.S)GO TO 174
        IF(Y.LE.4)GO TO 179

        TOT(Y)=TOT(Y)/X
        S=1
        TEMP3=9.0E60

        DO 181 I=1,Y
        IF(TOT(I).GE.TEMP3)GO TO 181
        TEMP3=TOT(I)
        Z=I

181    CONTINUE

```

```

C*****
C
C      Title:  RECOG.FR
C      Author: Capt Ajmal Hussain
C              Modified by 1Lt Kathy Dixon
C      Date:   Nov 84
C
C      Function:
C              Performs continuous speech recognition. Prints
C              recognized string on video screen and line printer.
C
C*****

```

```

      SUBROUTINE RECOG(PHON,VOCAB,WORD,DIS,LIB,J)

      INTEGER LIB(250),VOCAB(10,70),J,K,I,L,M,N,P,Q,R,S,T
      INTEGER U,V,W,X,Y,Z,TWORD(10),WRD(250),REJ(10,10)
      INTEGER FLAG,WORD(7,70),PHON(1130)
      REAL DIS(2432),TOT(5)
      DOUBLE PRECISION REAL TEMP,TEMP1,TEMP3

      LEN1=2
      LEN2=9

      DO 188 I=1,70
      DO 188 K=1,10
188      REJ(K,I)=0

      DO 187 I=1,5
187      TOT(I)=0

      R=1
      S=1
172      DO 175 I=1,250
175      WRD(I)=0

      79      IF(R.GE.(J+1))GO TO 1000
      IF(LIB(R).EQ.0)GO TO 170
      WRD(S)=LIB(R)
      R=R+1
      S=S+1
      GO TO 79
170      IF(S.GT.LEN1)GO TO 178
      R=R+1
      GO TO 79

178      IF(S.GT.LEN2)GO TO 177

      CALL FINDWORD(PHON,WRD,VOCAB,DIS,TEMP,L)

      CALL VTYPE(L,WORD)

```

```
C*****
C
C Title: PRINTREP.FR
C Author: 1Lt Kathy R. Dixon
C Date: Nov 84
C
C Function:
C Prints table of vocabulary words and associated
C phoneme representaition
C
C*****
```

```

SUBROUTINE PRINTREP(WORD,VOCAB,NOWORDS)
INTEGER WORD(7,70),VOCAB(10,70),NOWORDS
WRITE(12,100)(WORD(1,I),(VOCAB(J,I),J=1,10),
* I=1,NOWORDS)
100 FORMAT(2X,S14,10X,10I4)
RETURN
END
```

```

K=820
DO 502 I=41,DIS(2416)
WRITE(12,58)
WRITE(12,54)I
DO 503 J=1,40
WRITE(12,54)(INT((DIS(((I*(I-1))/2)+J))))
503 CONTINUE
502 CONTINUE
WRITE(12,59)
WRITE(12,51)IOP
WRITE(12,58)
WRITE(12,58)
WRITE(12,55)
DO 505 I=41,DIS(2416)
WRITE(12,54)I
505 CONTINUE
DO 506 I=41,DIS(2416)
WRITE(12,58)
WRITE(12,54)I
DO 507 J=41,DIS(2416)
IF(J.GT.I)GO TO 507
WRITE(12,54)(INT((DIS(((I*(I-1))/2)+J))))
507 CONTINUE
506 CONTINUE

500 RETURN
END

```

```

C*****
C
C   Title:  PRINTDIS.FR
C   Author: Capt. Ajmal Hussain
C           Modified by 1Lt Kathy Dixon
C   Date:   Nov 84
C
C   Function:
C           Prints lower triangular distance matrix.
C
C*****

```

```

SUBROUTINE PRINTDIS(DIS)

REAL DIS(2432)

IOP=1
L=DIS(2416)
WRITE(12,51)IOP
51  FORMAT(50X,"MAT",I2)
WRITE(12,58)
IF(DIS(2416).LE.40)GO TO 52
L=40
WRITE(12,55)
55  FORMAT("+",3X,Z)
52  DO 53 I=1,L
WRITE(12,54) I
53  CONTINUE
54  FORMAT("+",I3,Z)
K=1
DO 56 I=1,L
WRITE(12,58)
WRITE(12,54) I
DO 57 J=1,L
IF(J.GT.I)GO TO 57
WRITE(12,54)(INT(DIS(K)))
K=K+1
57  CONTINUE
56  CONTINUE
58  FORMAT(1X)
IF(DIS(2416).LE.40)GO TO 500

WRITE(12,59)
59  FORMAT("1")
WRITE(12,51)IOP
WRITE(12,58)
WRITE(12,58)
WRITE(12,55)
DO 501 I=1,40
WRITE(12,54)I
501 CONTINUE

```

```

C*****
C
C   Title:  NORMALIZE.FR
C   Author: Capt. Ajmal Hussain
C           Modified by Kathy Dixon
C   Date:   Nov 84
C
C   Function:
C           Normalizes sixteen dimensional vectors
C           of digitized speech.
C*****

```

```

SUBROUTINE NORMALIZE(IDATA2,IENERGY)

COMMON/IBUFF/IDATA(8192)

INTEGER IDATA2
REAL IENERGY(500)
DOUBLE PRECISION REAL TEMP

TEMP=0
K=1
J=1
L=1

DO 5 I=1,(IDATA2/16)
TEMP=0

DO 2 J=K,(K+15)
TEMP=TEMP+FLOAT(IDATA(J))**2      ;energy calculation
CONTINUE

TEMP=(SQRT(TEMP)/32000)
IENERGY(I)=ABS(TEMP)             ;store energy values
DO 4 J=K,(K+15)
IDATA(J)=FLOAT(IDATA(J))/TEMP    ;energy normalize
                                   ;vectors
5   K=K+16

RETURN
END

```

```
C*****
C
C   Title: NewScr
C   Author: Lt Allen
C   Date: Dec 82
C
C   Function:
C       This routine erases the screen by typing 24
C       blank lines.
C
C   Compile command:
C       FORTRAN NEWSR
C
C*****
```

```
      SUBROUTINE NEWSR
```

```
      DO 10 I=1,24
      TYPE
10    CONTINUE
```

```
      RETURN
      END
```

```
C*****
```

```

C*****
C
C   Title:  SPEECH.FR
C   Author: 1Lt Kathy R. Dixon
C   Date:   Nov 84
C
C   Function:
C           Calls routines to recognize continuous speech.
C
C*****

```

```

SUBROUTINE SPEECH(ENRGY,PHON,DIS,VOCAB,WORD)

```

```

COMMON/IBUFF/IDATA(4096)

```

```

INTEGER PHON(1130),VOCAB(10,70),LIB(250),SOUND(250)
INTEGER ISTOP,IFILE,IDATA2,WORD(7,70),ISWITCH,J
REAL DIS(2432),IENERGY(500),ENRGY

```

```

IDATA2=4000                ;conversion value
ISTOP=16                   ;last block to read
IFILE=5                    ;file to read
ISWITCH=1                  ;compress input string once

```

```

CALL REDBUF(IFILE,ISTOP)   ;read data

```

```

CALL NORMALIZE(IDATA2,IENERGY) ;normalize vectors

```

```

CALL EXTRACT(IDATA2,PHON,IENERGY,ENRGY,LIB)
                                ;extract phonemes

```

```

CALL CMPRS(IDATA2,PHON,DIS,LIB,ISWITCH,SOUND,J)
                                ;compress phonemes

```

```

CALL RECOG(PHON,VOCAB,WORD,DIS,SOUND,J)
                                ;determine words in
                                ;input string
RETURN                          ;print to screen
END

```

```

C*****
C
C   Title:  TEMPLATE.FR
C   Author: 1Lt Kathy Dixon
C   Date:   Nov 84
C
C   Function:
C       Performs necessary calls to support routines
C       to produce a phoneme template from 500 vectors
C       of input speech. The template is stored in the
C       sequentially in the array, PHON.
C
C*****

```

```

SUBROUTINE TEMPLATE(ENRGY)

```

```

INTEGER PHON(1130),NOVECT,IDATA2
INTEGER WEIGHT(500),IWGHT(500),IWGHTT(500)
REAL IENERGY(500),ENRGY

```

```

COMMON/IBUFF/IDATA(8192)

```

```

NOVECT=500
IDATA2=8000 ;conversion count
ISTOP=32 ;no. blocks to read
IFILE=10 ;file to read
ENRGY=0

```

```

CALL REDBUF(IFILE,ISTOP) ;read data into buffer

```

```

CALL NORMALIZE(IDATA2,IENERGY) ;normalize vectors

```

```

C   WRITE(12,347)(IENERGY(I),I=1,500) ; option to write
347   FORMAT(10G11.5) ; energy values to line
;printer

```

```

DO 88 I=1,500
88   ENRGY=ENRGY+IENERGY(I)

```

```

ENRGY=(ENRGY/500)

```

```

CALL LOWENERGY(IDATA2,IENERGY,ENRGY,WEIGHT)
;low energy vectors

```

```

TYPE"Calculating template"

```

```

1   CALL CMPRS1(IDATA2,WEIGHT,IWGHTT) ;delete vectors

```

```

DO 99 I=1,500
99   WEIGHT(I)=IWGHTT(I)
NOVECT=IDATA(IDATA2+1)

```

```
TYPE"Number of remaining vectors ",NOVECT
IF(NOVECT.LT.70)GO TO 2
```

```
CALL COMPARE(IDATA2,WEIGHT,IWGHT)
```

```
;compare and average
```

```
44 DO 44 I=1,500
WEIGHT(I)=IWGHT(I)
GO TO 1
```

```
;nearest vectors
```

```
2 DO 5 I=1,1130
5 PHON(I)=0
```

```
;zero phoneme array
```

```
3 DO 3 I=1,1120
PHON(I)=IDATA(I)
```

```
;store phoneme template
```

```
PHON(1121)=NOVECT
```

```
;store no. of phonemes
```

```
CALL WRTPHON(PHON)
```

```
RETURN
END
```

```

C*****
C
C      Title:      TRAIN.FR
C      Author:     1Lt Kathy R. Dixon
C      Date:       Jul 84
C
C      Function:
C          This routine drives the system for training
C          and speech recognition for a specified user.
C
C*****

```

```

      INTEGER PHON(1130),SOUND(250),VOCAB(10,70)
      INTEGER WORD(7,70)
      REAL DIS(2432),ENRGY

```

```

1      CALL NEWSCR                                ;erase screen

```

```

      TYPE"<CR>
*      Welcome to the AFIT Speech Recognition Project.<CR>
*      <CR>
*      Using AFIT studies in speech research, a machine<CR>
*      was developed to recognize continuous speech.<CR>
*      Though training is required, the ultimate goal <CR>
*      of this work is to develop a machine which is<CR>
*      speaker independent.<CR>
*      <CR>
*      Continue? [Y]"

```

```

      CALL GCHAR(ICCHAR,IER)
      IF(ICCHAR.EQ.78)GO TO 111

```

```

      CALL NEWSCR

```

```

      TYPE"<CR>
*      Press CR, then say the following<CR>
*      phrases into the microphone after<CR>
*      prompted by the word, start.<CR>
*      Only 10 seconds of speech will<CR>
*      be accepted.<CR>
*      <CR>
*      <CR>
*      CHANGE FREQUENCY TO THREE FIVE SEVEN<CR>
*      LOCK-ON TARGET AT TWO THOUSAND MARK<CR>
*      ARM STATION ALFA BRAVO CHARLIE FOXTROT<CR>
*      MAP AIR-TO-SURFACE MISSLE THREAT"

```

```

      CALL GCHAR(ICCHAR,IER)

```

```

      CALL SWAP("ATOD10.SV",IER)                ;swap for A/D

```

```

IF(IER.NE.1)TYPE"SWAP error ",IER      ;10 sec speech
CALL NEWSCR
CALL TEMPLATE(ENRGY)      ;create phoneme template
CALL NEWSCR
TYPE"Calculating distance matrix"
CALL GCHAR(ICCHAR,IER)
CALL REDPHON(PHON)
CALL DISTANCE(PHON)      ;create distance matrix
CALL REDDIS(DIS)
CALL NEWSCR
TYPE"Printing distance matrix"
CALL GCHAR(ICCHAR,IER)
CALL PRINTDIS(DIS)      ;print distance matrix
34 CALL NEWSCR
DO 15 I=1,70
DO 15 J=1,10
15 VOCAB(J,I)=0      ;zero VOCAB matrix
NOWORDS=70      ;no. vocabulary words
CALL REDWRDS(NOWORDS,WORD)      ;read vocabulary
;from file
TYPE"<CR>"
* Press CR, a word will appear<CR>
* on the screen. Press CR again<CR>
* then say the word. Repeat. "
CALL GCHAR(ICCHAR,IER)
DO 99 II=1,NOWORDS
CALL NEWSCR
100 WRITE(10,100) WORD(1,II)      ;write word to screen
FORMAT(2X,S14)
CALL SWAP("ATOD2.SV",IER)      ;swap for A/D

```

```

IF(IER.NE.1)TYPE "SWAP error",IER      ;5 sec speech

CALL REP(ENRGY,PHON,DIS,SOUND)

DO 63 K=1,10
63  VOCAB(K,II)=SOUND(K)                ;store word rep
99  CONTINUE

CALL NEWSCR

TYPE"Printing word phoneme representations"

CALL WRTREP(VOCAB,NOWORDS)

CALL PRINTREP(WORD,VOCAB,NOWORDS)      ;print word rep

CALL REDREP(NOWORDS,VOCAB)

CALL NEWSCR

C*****

TYPE"<CR>"
*   Training is complete.  Testing <CR>
*   can now be done.  An asterisk will<CR>
*   appear on the screen.  Press CR, <CR>
*   then say the test word.  The machine's<CR>
*   guess will then be typed on the screen.<CR>
*   <CR>
*   Continue? [Y]"

CALL GCHAR(ICHAR,IER)
IF(ICHAR.EQ.78)GO TO 111

77  CALL NEWSCR

TYPE"*"                                ;request speech input

CALL GCHAR(ICHAR,IER)

CALL SWAP("ATOD5.SV",IER)              ;swap for A/D
IF(IER.NE.1)TYPE"SWAP ERROR",IER      ;5 sec speech

CALL SPEECH(ENRGY,PHON,DIS,VOCAB,WORD) ;speech
                                       ;recognition

```

```
TYPE"<CR>
* Continue? [Y] "

CALL GCHAR(ICCHAR,IER)
IF(ICCHAR.EQ.78) GO TO 73

GO TO 77

73 CALL NEWSR

TYPE"<CR>
* This concludes the session. <CR>
* End session? [Y] "

CALL GCHAR(ICCHAR,IER)
IF(ICCHAR.EQ.78) GO TO 111
GO TO 112

111 TYPE"<CR>
* Start again? [N] "

CALL GCHAR(ICCHAR,IER)
IF(ICCHAR.EQ.89) GO TO 1

112 CALL EXIT
END
```

```
C*****
C
C Title: VTYPE.FR
C Author: 1Lt Kathy Dixon
C Modelled after Hussain's WTYPE
C Date: Nov 84
C
C Function:
C This routine displays the word specified on the
C H19 terminal in video on a single line.
C
C*****
```

```
      SUBROUTINE VTYPE(L,WORD)
```

```
      INTEGER L,WORD(7,70)
```

```
      IF(L.EQ.0) GO TO 10          ;skip zero numbered words
```

```
      L=(L-1)/10+1
```

```
      WRITE(12,15)L
```

```
      WRITE(12,14)WORD(1,L)
```

```
      WRITE(10,11)WORD(1,L)      ;display word
```

```
11      FORMAT(S14)
```

```
14      FORMAT("+",S14,Z)
```

```
15      FORMAT(2X,I4)
```

```
10      RETURN
```

```
      END
```

```
C*****
```

```
C*****
C
C      Title:  WRTDIS.FR
C      Author: Kathy Dixon
C      Date:   Nov 84
C
C      Function:  Writes Distance matrix to a disk file.
C
C*****
```

```
      SUBROUTINE WRTDIS(DIS)
```

```
      REAL DIS(2432)
```

```
      CALL CFILW("DIST",2,IER)
      IF(IER.NE.1)TYPE"CFILW ERROR",IER
      CALL OPEN(1,"DIST",2,IER)
      IF(IER.NE.1)TYPE"OPEN ERROR",IER
      WRITE(1,100)(DIS(I),I=1,2432)
      FORMAT(G11.5)
      CALL CLOSE(1,IER)
      IF(IER.NE.1)TYPE"CLOSE ERROR",IER
```

100

```
      RETURN
      END
```

```
C*****
C
C   Title:  WRTPHON.FR
C   Author: Kathy Dixon
C   Date:   Nov 84
C
C   Function:  Writes 70 16-dimensional vectors to
C              a file, PHONE
C*****
```

```
      SUBROUTINE WRTPHON(PHON)
```

```
      INTEGER PHON(1130)
```

```
      CALL CFILW("PHONE",2,IER)
```

```
      IF(IER.NE.1)TYPE"WRTPHON CFILW ERROR",IER
```

```
      CALL OPEN(1,"PHONE",2,IER)
```

```
      IF(IER.NE.1)TYPE"WRTPHON ERROR",IER
```

```
100
```

```
      WRITE(1,100)(PHON(I),I=1,1130)
```

```
      FORMAT(I6)
```

```
      CALL CLOSE(1,IER)
```

```
      IF(IER.NE.1)TYPE"CLOSE ERROR",IER
```

```
      RETURN
```

```
      END
```

```
C*****
C
C Title: WRTREP.FR
C Author: Kathy Dixon
C Date: Nov 84
C
C Function: Writes word phonemes to a file, VOCABUL.
C
C*****
```

```
SUBROUTINE WRTREP(VOCAB,NOWORDS)
```

```
INTEGER VOCAB(10,70),NOWORDS
```

```
CALL CFILW("VOCABUL",2,IER)
```

```
IF(IER.NE.1)TYPE"CFILW ERROR",IER
```

```
CALL OPEN(1,"VOCABUL",2,IER)
```

```
IF(IER.NE.1)TYPE"OPEN ERROR",IER
```

```
100
```

```
WRITE(1,100)((VOCAB(I,J),I=1,10),J=1,NOWORDS)
```

```
FORMAT(I3)
```

```
CALL CLOSE(1,IER)
```

```
IF(IER.NE.1)TYPE"CLOSE ERROR",IER
```

```
RETURN
```

```
END
```

```

C*****
C
C Title: VOCAB.FR
C Author: 1Lt Kathy Dixon
C Date: Nov 84
C
C Function:
C Creates a file of English words.
C
C Command Line:
C RLD R VOCAB NEWS CR @FLIB@
C*****

```

```

INTEGER W(7,70),L

ACCEPT"NUMBER OF VOCABULARY WORDS ",L
TYPE" * * * *"
DO 200 I=1,L
ACCEPT"WORD "
100 READ(11,100)W(1,I)
FORMAT(S14)

200 CONTINUE

ACCEPT"PRINT VOCABULARY [Y]"
CALL GCHAR(ICAR,IER)
IF(ICAR.EQ.78)GO TO 111

CALL NEWS CR

WRITE(10,150)(W(1,I),I=1,L)
WRITE(12,150)(W(1,I),I=1,L)
150 FORMAT(2X,5S14)

CALL CFILW("WORDS",2,IER)
IF(IER.NE.1)TYPE"CFILW ERROR",IER

CALL OPEN(1,"WORDS",2,IER)
IF(IER.NE.1)TYPE"OPEN ERROR",IER
WRITE(1,100)(W(1,I),I=1,70)
CALL CLOSE(1,IER)
IF(IER.NE.1)TYPE"CLOSE ERROR",IER

111 CALL EXIT
END

```

VITA

Kathy Renee Dixon was born on 27 January 1959 in Roswell, New Mexico. She graduated from high school in Orlando, Florida in 1977 and attended the University of Central Florida, Orlando, Florida from which she received the degree of Bachelor of Science in Engineering in April 1982. She entered the Air Force on active duty in May 1982 and received her commission from Officer Training School in August 1982. She served as a Project Engineer in the Signal Processing Laboratory, Air Force Institute of Technology until entering the School of Engineering in June 1983.

Permanent address: 520 Sheppard Road
Orlando, Florida 32820

AD-A151936

CLASSIFICATION OF THIS PAGE

REPORT DOCUMENTATION PAGE

1. SECURITY CLASSIFICATION UNCLASSIFIED		1b. RESTRICTIVE MARKINGS	
2. CLASSIFICATION AUTHORITY		3. DISTRIBUTION/AVAILABILITY OF REPORT Approved for public release, distribution unlimited	
4. DOWNGRADING SCHEDULE		5. MONITORING ORGANIZATION REPORT NUMBER(S)	
6. ORGANIZATION REPORT NUMBER(S) AFIT/ENG-26		7a. NAME OF MONITORING ORGANIZATION	
8. PERFORMING ORGANIZATION Department of Engineering (City, State and ZIP Code) Air Force Institute of Technology Patterson AFB, Ohio 45433	8b. OFFICE SYMBOL (If applicable) AFIT/ENG	7b. ADDRESS (City, State and ZIP Code)	
9. FUNDING/SPONSORING AGENCY (City, State and ZIP Code)	8b. OFFICE SYMBOL (If applicable)	9. PROCUREMENT INSTRUMENT IDENTIFICATION NUMBER	
10. SOURCE OF FUNDING NOS.		10. SOURCE OF FUNDING NOS.	
10. SOURCE OF FUNDING NOS.		PROGRAM ELEMENT NO.	PROJECT NO.
10. SOURCE OF FUNDING NOS.		TASK NO.	WORK UNIT NO.
10. SOURCE OF FUNDING NOS.		10. SOURCE OF FUNDING NOS.	
11. AUTHOR(S) Dixon, B.S.E., 1st Lt, USAF			
12. REPORT NUMBER AFIT/ENG-26	13b. TIME COVERED FROM _____ TO _____	14. DATE OF REPORT (Yr., Mo., Day) 1984 December	15. PAGE COUNT 135
16. ABSTRACT			
17. KEY WORDS			
18. COSATI CODES		18. SUBJECT TERMS (Continue on reverse if necessary and identify by block number)	
GROUP	SUB. GR.	Recognition, Speech, Phonemes, Speech Recognition, Continuous Speech Recognition	
02			
(Continue on reverse if necessary and identify by block number)			
IMPLEMENTATION OF A REAL-TIME, INTERACTIVE, CONTINUOUS SPEECH RECOGNITION SYSTEM			
Chairman: Dr. Matthew Kabrisky, Professor of Electrical Engineering, Air Force Institute of Technology			
19. DISTRIBUTION/AVAILABILITY OF ABSTRACT UNLIMITED <input checked="" type="checkbox"/> SAME AS RPT. <input type="checkbox"/> OTIC USERS <input type="checkbox"/>		21. ABSTRACT SECURITY CLASSIFICATION UNCLASSIFIED	
20. RESPONSIBLE INDIVIDUAL Matthew Kabrisky, Professor, AFIT		22b. TELEPHONE NUMBER (Include Area Code) 513-255-5276	22c. OFFICE SYMBOL AFIT/ENG

73, 83 APR

EDITION OF 1 JAN 73 IS OBSOLETE.

UNCLASSIFIED SECURITY INFORMATION THIS PAGE

UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE

A speech recognition system was designed and implemented to recognize continuous speech in a real-time environment (after training). Several techniques were incorporated to characterize phonemes as vectors in space. Through the use of distance rules, it was possible to characterize words by a phoneme representation. This approach to speech recognition offers several possibilities for future investigations such as varying Minkowski distances and applying clustering techniques. The algorithm which was developed was modularized on a hierarchical basis and was user friendly.

UNCLASSIFIED

END

FILMED

4-85

DTIC