

A ROBUST PHONETIC DIGIT RECOGNITION SYSTEM FOR
DIGITS SPOKEN IN AMERICAN ENGLISH
AND ARABIC

By

AHMED MAHMOUD MILYANI

Certificate

Cambridgeshire College of Arts & Technology
Cambridge, England
1969

Master of Science
Oklahoma State University
Stillwater, Oklahoma
1976

Submitted to the Faculty of the Graduate College
of the Oklahoma State University
in partial fulfillment of the requirements
for the Degree of
DOCTOR OF PHILOSOPHY
July, 1981



A ROBUST PHONETIC DIGIT RECOGNITION SYSTEM FOR
DIGITS SPOKEN IN AMERICAN ENGLISH
AND ARABIC

Thesis Approved:

Alio Yahya
Thesis Adviser

Robert J. Mulholland

Richard L. Cummins

Lyle Broemeling

Norman D. Durbin
Dean of the Graduate College

1099214

ACKNOWLEDGMENTS

Praise be to Allah and Peace be upon His messenger, Mohammed. I thank Allah, The Almighty, for His uncountable help and guidance.

I would like foremost to express my sincere gratitude to Dr. Rao Yarlagadda, my thesis adviser and chairman of my doctoral committee, for his inspiration and continued guidance during the entire phase of my graduate study. His dedication and encouragement has contributed significantly to the completion of this dissertation.

My special thanks and appreciation are due to Dr. R. Mullholand, Dr. R. Cummins, and Dr. L. Broemeling, my committee members, for their interest, helpful comments and suggestions throughout this research.

To Dr. R. Schaefer, I would like to express my appreciation for his effort in my behalf in gathering speech data at the OSU Speech Clinic Laboratory. I would like to recognize Dr. L. Ebbesen and Dr. J. Perrault for their aid with the OSU Mechanical Engineering Interdata Computer.

My sincere appreciation are due to the Saudi Arabian Government for their support during the entire phase of my study.

I am grateful to Mrs. Dolores Behrens for her excellent typing of this thesis.

My sincere appreciation also extends to my parents and all the members of my family and my relatives who had their part of sufferings and sacrifices during the entire phase of my study. The great memory of my father who passed to the other world just after the completion of this thesis shall not be forgotten at this moment in life.

TABLE OF CONTENTS

Chapter	Page
I. INTRODUCTION	1
Review of the Literature on Speech Recognition.	2
Digit Recognition	4
Segmentation.	6
Speech Input Rate and Pauses.	7
Human Factors	9
Linear Prediction Analysis.	10
Previous Systems.	10
Thesis Outline.	14
II. SPEECH MECHANISM	16
Introduction.	16
Voice Production.	19
Vocal Pitch and Loudness.	26
Articulation.	28
Classification of Speech Sounds	29
Consonants.	29
Vowels.	38
Speech Acoustics.	39
Sound Propagation	41
Resonance	42
Formant Frequencies	43
Acoustic Characteristics and Perception	44
Effects of Noise.	46
Frequency and Intelligibility	47
Segmental Analysis.	47
Coarticulation.	49
Vowels.	50
III. CONCEPT AND COMPUTATIONAL TOOLS.	52
Introduction.	52
Linear Prediction Analysis.	54
The Autocorrelation Method.	55
Transfer Function Relation.	57
Vocal Tract Division.	57
Reflection Coefficients	57
Parameters in Terms of Vocal Tract Cavity Ratios.	61

Chapter	Page
Short-Time Energy	61
Pre-emphasis.	62
RMS Analysis in Speech Processing	63
Double Smoothing Algorithms	65
Zero Crossings.	75
Cross-Correlation Function.	77
Window Applications	79
Parameters in Dip-Classification.	80
Dip-Classification.	85
Scaling and Normalization	89
Amplitude Normalization	89
Maximum of Frames Normalization	89
Scaling (Energy Normalization).	89
 IV. ACOUSTIC PHONEMIC DIGIT RECOGNITION SCHEME FOR DIGITS SPOKEN IN AMERICAN ENGLISH	91
Introduction.	91
Segmentation.	92
Recognition Scheme.	93
Digit Recognition Flow Chart.	94
Pattern Recognition System.	153
 V. ACOUSTIC PHONEMIC DIGIT RECOGNITION SCHEME FOR DIGITS SPOKEN IN ARABIC	163
Introduction.	163
Arabic Phonemes	163
The Emphatics	165
Velarized Phonemes.	165
Throat and Back-of-Mouth Sounds	166
Minor Differences	169
Segmentation.	172
Digit Recognition Flow Chart.	172
Pattern Recognition Scheme.	230
 VI. SUMMARY AND SUGGESTION FOR FURTHER STUDY	240
Summary	240
Suggestions for Further Research.	242
 REFERENCES.	254
 APPENDIX A - DATA ACQUISITION	260
 APPENDIX B - COMPUTER SUBROUTINES	265

LIST OF TABLES

Table	Page
I. Vowels Used in Digits Spoken in American English, According to the Degree of Constriction and Tongue Hump Position.	34
II. Sequence of Sound Classes of Digits Spoken in American English.	35
III. Vowels Used in Digits Spoken in Arabic, According to the Degree of Constriction and Tongue Hump Position	36
IV. Sequence of Sound Classes of Digits Spoken in Arabic. . .	37
V. Smoothed RMS Energy Peaks and Ratios for Digits Spoken in American English for Fixed Mean and Varying Peaks	117
VI. Smoothed and Quantized RMS Energy Peaks and Their Ratios for Digits Spoken in American English.	118
VII. Sequence of Sound Class Regions of Digits Spoken in American English	151
VIII. RMS Energy Correlation Table for Digits Spoken in American English	154
IX. RMS Energy Correlation Table for Digits in American English.	155
X. BTR Correlation Table for Digits Spoken in American English.	156
XI. BTR Correlation Table for Digits Spoken in American English.	157
XII. RMS Energy Correlation Table for Two Sets of Digits Spoken in American English by Two Different Speakers.	161
XIII. BTR Correlation Table for the Same Two Sets of Digits Used in Table XII.	162

Table	Page
XIV. Transliteration of Arabic Words and Names	164
XV. Comparison Between Beginning and End of English and Arabic Digits	171
XVI. Smoothed RMS Energy Peaks and Their Ratio's for Digits Spoken in Arabic with Varying Amplitude and Fixed Mean.	195
XVII. Smoothed and Quantized RMS Energy Peak Ratio's for Digits Spoken in Arabic	196
XVIII. Sequence of Sound Class Regions of Digits Spoken in Arabic	228
XIX. RMS Energy Correlation Table for Digits Spoken in Arabic	232
XX. RMS Energy Correlation Table for Digits Spoken in Arabic	233
XXI. RMS Energy Correlation Table for Digits Spoken in Arabic	234
XXII. BTR Correlation Table for Digits Spoken in Arabic	235
XXIII. BTR Correlation Table for Digits Spoken in Arabic	236
XXIV. BTR Correlation Table for Digits Spoken in Arabic	237
XXV. RMS Energy Cross-Correlation Table for Two Sets of Digits Spoken in Arabic by Two Different Speakers.	238
XXVI. BTR Cross-Correlation Table for the Same Digits Used in Table XXV.	239

LIST OF FIGURES

Figure	Page
1. The Over-All Isolated Digit Recognition System Block Diagram After Sambur and Rabiner [4]	13
2. X-Ray Tracing of the Speech Mechanism of a Normal Speaker, Showing the Palate at Rest.	17
3. A Schematic View of the Articulators and Places of Articulation Showing the Partitioning Assumed for Classifying the Place and Degree of Constriction of the Vocal Tract During Production of Vowel Sounds . . .	20
4. Normal Articulation Process.	23
5. X-Ray of Cross-Section Area for Different Constrictions. . .	24
6. Schematic Representation of Changes in Glottal Area During a Vibratory Cycle	27
7. Phonemes in American English Digits.	30
8. Vowels Used in Digits 0 - 9 Spoken in American English . . .	31
9. Phonemes of Digits Spoken in Arabic.	32
10. Vowels Used in Digits 0 - 9 Spoken in Arabic	33
11. Distribution of Volume Velocity at the Frequencies of Each of the First Four Resonances of an Ideal Neutral Articulation in Which the Vocal Tract Simulates a Tube of Constant Cross-Sectional Area.	40
12. Schematic Diagram of the Physiological and Acoustical Characteristics of Speech Sound Production	45
13. Concatenation of (N=14) Lossless Tubes of Equal Length . . .	53
14. Equivalent Representation of Predictor Filter F(Z) and Inverse Filter A(Z).	58
15. Equivalent Diagram Representing the Vocal Tract as an Acoustical Tube of Equal Lengths.	58

Figure	Page
16. High and Low Frequency Pre-Emphasis Block Diagram.	64
17. Two Examples of Noisy Signals to be Smoothed	67
18. Median Smoothing of a Sequence with a Discontinuity.	69
19. Effects of Various Median Smoothers on Low-Order Polynomials.	70
20. Simple Smoothing Algorithm	71
21. Double Smoothing Algorithm	73
22. Nonlinear Double Smoothing Algorithm	74
23. Effects of Several Versions of the Smoothing Algorithm on a Speech Intensity Contour.	76
24. Hamming Window	81
25. RMS Plot Used for the Definitions of the Extrema V_i^+ , V_j^- and Their Ratios, R_1 , R_2 and R_{\min}	83
26. Plot of $\text{Log}_{10} R_{\min}$ vs $\text{Log}_{10} V_j^-$ for Estimating the Values of Linear Discriminant Functions LDF1 and LDF2	86
27. Schematic Diagram Showing the Types of Dips, D_1 , D_2 and D_3 , RMS, and BTR, Vowel, Non-Vowel Primary Decision and Combines Decision	87
28. Digit Recognition Flow Chart Based on a Phonemic Feature Detection.	95
29. Smoothed RMS Energy Contour for Digit Zero, i.e. /zIro/ Spoken in American English.	97
30. Smoothed RMS Energy Contour for Digit One, i.e. /wAn/ Spoken in American English	98
31. Smoothed RMS Energy Contour for Digit Two, i.e. /tu/ Spoken in American English.	99
32. Smoothed RMS Energy Contour for Digit Three, i.e. /θri/ Spoken in American English	100
33. Smoothed RMS Energy Contour for Digit Four, i.e. /fɔr/ Spoken in American English	101

Figure	Page
34. Smoothed RMS Energy Contour for Digit Five, i.e. /faIv/ Spoken in American English.	102
35. Smoothed RMS ENergy Contour for Digit Six, i.e. /sIks/ Spoken in American English.	103
36. Smoothed RMS Energy Contour for Digit Seven, i.e. /seven/ Spoken in American English	104
37. Smoothed RMS Energy Contour for Digit Eight, i.e. /eIt/ Spoken in American English	105
38. Smoothed RMS Energy Contour for Digit Nine, i.e. /naIn/ Spoken in American English.	106
39. Smoothed and Quantized RMS Energy Contour for Digit Zero, i.e. /zIro/ Spoken by International Speaker	107
40. Smoothed and Quantized RMS Energy Contour for Digit One, i.e. /wAn/ Spoken by International Speaker.	108
41. Smoothed and Quantized RMS Energy Contour for Digit Two, i.e. /tu/ Spoken by International Speaker	109
42. Smoothed and Quantized RMS Energy Contour for Digit Three, i.e. /θri/ Spoken by International Speaker.	110
43. Smoothed and Quantized RMS Energy Contour for Digit Four, i.e. /fɔr/ Spoken by International Speaker	111
44. Smoothed and Quantized RMS Energy Contour for Digit Five, i.e. /faIv/ Spoken by International Speaker.	112
45. Smoothed and Quantized RMS Energy Contour for Digit Six, i.e. /sIks/ Spoken by International Speaker	113
46. Smoothed and Quantized RMS Energy Contour for Digit Seven, i.e. /seven/ Spoken by International Speaker.	114
47. Smoothed and Quantized RMS Energy Contour for Digit Eight, i.e. /eIt/ Spoken by International Speaker.	115
48. Smoothed and Quantized RMS Energy Contour for Digit Nine, i.e. /naIn/ Spoken by International Speaker.	116
49. Smoothed and Quantized Feature Parameters for Digit Zero Spoken in American English.	120
50. Smoothed and Quantized Feature Parameters for Digit One Spoken in American English	123

Figure	Page
51. Smoothed and Quantized Feature Parameters for Digit Two Spoken in American English	126
52. Smoothed and Quantized Feature Parameters for Digit Three Spoken in American English	129
53. Smoothed and Quantized Feature Parameters for Digit Four Spoken in American English.	132
54. Smoothed and Quantized Feature Parameters for Digit Five Spoken in American English.	135
55. Smoothed and Quantized Feature Parameters for Digit Six Spoken in American English	138
56. Smoothed and Quantized Feature Parameters for Digit Seven Spoken in American English	141
57. Smoothed and Quantized Feature Parameters for Digit Eight Spoken in American English	144
58. Smoothed and Quantized Feature Parameters for Digit Nine Spoken in American English.	147
59. Decision Tree for Digit Identification Scheme for American English Language.	152
60. RMS and BTR Correlation Flow Diagram for Digit Recognition.	159
61. Future RMS and BTR Correlation Flow Diagram for an Efficient Digit Recognition System.	160
62. Tongue Position for the Emphatic-Nonemphatic [s] vs [s] Fricative Consonants in Arabic.	167
63. Smoothed and Quantized RMS Energy Contour for Digit Zero, i.e./ sefr/ Spoken in Arabic, Sample 1	174
64. Smoothed and Quantized RMS Energy Contour for Digit One, i.e. /wâhid/ Spoken in Arabic, Sample 1	175
65. Smoothed and Quantized RMS Energy Contour for Digit Two, i.e. /iθnân/ Spoken in Arabic, Sample 1	176
66. Smoothed and Quantized RMS Energy Contour for Digit Three, i.e. /θalâθh/ Spoken in Arabic Showing Wrong End and Point Detection, Sample 1.	177
67. Smoothed and Quantized RMS Energy Contour for Digit Three, i.e. /θalâθh/ Spoken in Arabic, Sample 1	178

Figure	Page
68. Smoothed and Quantized RMS Energy Contour for Digit Four, i.e. /arbaðh/ Spoken in Arabic, Sample 1.	179
69. Smoothed and Quantized RMS Energy Contour for Digit Five, i.e. /kh^amsðh/ Spoken in Arabic, Sample 1.	180
70. Smoothed and Quantized RMS Energy Contour for Digit Six, i.e. /sIttðh/ Spoken in Arabic, Sample 1.	181
71. Smoothed and Quantized RMS Energy Contour for Digit Seven, i.e. /s^lbpðh/ Spoken in Arabic, Sample 1.	182
72. Smoothed and Quantized RMS Energy Contour for Digit Eight, i.e. /θamānyðh/ Spoken in Arabic, Sample 1.	183
73. Smoothed and Quantized RMS Energy Contour for Digit Nine, i.e. /tIsðh/ Spoken in Arabic, Sample 1.	184
74. Smoothed and Quantized RMS Energy Contour for Digit Zero, i.e. /sefr/ Spoken in Arabic, Sample 2.	185
75. Smoothed and Quantized RMS Energy Contour for Digit One, i.e. /wāhid/ Spoken in Arabic, Sample 2.	186
76. Smoothed and Quantized RMS Energy Contour for Digit Two, i.e. /iθnān/ Spoken in Arabic, Sample 2.	187
77. Smoothed and Quantized RMS Energy Contour for Digit Three, i.e. /θa^lāθðh/ Spoken in Arabic, Sample 2.	188
78. Smoothed and Quantized RMS Energy Contour for Digit Four, i.e. /arbaðh/ Spoken in Arabic, Sample 2.	189
79. Smoothed and Quantized RMS Energy Contour for Digit Five, i.e. /kh^amsðh/ Spoken in Arabic, Sample 2.	190
80. Smoothed and Quantized RMS Energy Contour for Digit Six, i.e. /sIttðh/ Spoken in Arabic, Sample 2.	191
81. Smoothed and Quantized RMS Energy Contour for Digit Seven, i.e. /s^lbpðh/ Spoken in Arabic, Sample 2.	192
82. Smoothed and Quantized RMS Energy Contour for Digit Eight, i.e. /θamānyðh/ Spoken in Arabic, Sample 2.	193
83. Smoothed and Quantized RMS Energy Contour for Digit Nine, i.e. /tIsðh/ Spoken in Arabic, Sample 2.	194

Figure	Page
84. Smoothed and Quantized Feature Parameter for Digit Zero, i.e. /sefr/ Spoken in Arabic	197
85. Smoothed and Quantized Feature Parameter for Digit One, i.e. /wâhid/ Spoken in Arabic	200
86. Smoothed and Quantized Feature Parameter for Digit Two, i.e. /iθnân/ Spoken in Arabic	203
87. Smoothed and Quantized Feature Parameter for Digit Three, i.e. /θaλaθðh/ Spoken in Arabic	206
88. Smoothed and Quantized Feature Parameter for Digit Four, i.e. /arbaðh/ Spoken in Arabic	209
89. Smoothed and Quantized Feature Parameter for Digit Five, i.e. /khλmsðh/ Spoken in Arabic.	212
90. Smoothed and Quantized Feature Parameters for Digit Six, i.e. /sIttðh/ Spoken in Arabic.	215
91. Smoothed and Quantized Feature Parameter for Digit Seven, /sλbðh/ Spoken in Arabic	218
92. Smoothed and Quantized Feature Parameter for Digit Eight, i.e. /θamânyðh/ Spoken in Arabic.	221
93. Smoothed and Quantized Feature Parameter for Digit Nine, i.e. /tIsðh/ Spoken in Arabic	224
94. Decision Tree for Digit Identification for Arabic Language	229
95. Plots of CTR vs BTR for Digit Zero	244
96. Plots of CTR vs BTR for Digit One.	245
97. Plots of CTR vs BTR for Digit Two.	246
98. Plots of CTR vs BTR for Digit Three.	247
99. Plots of CTR vs BTR for Digit Four	248
100. Plots of CTR vs BTR for Digit Five	249
101. Plots of CTR vs BTR for Digit Six.	250
102. Plots of CTR vs BTR for Digit Seven.	251

Figure	Page
103. Plots of CTR vs BTR for Digit Eight.	252
104. Plots of CTR vs BTR for Digit Nine	253
105. Data Acquisition	263
106. Second Order Low-Pass Filter	264
107. Subroutine FRMS.	267
108. Subroutine FBTR.	268
109. Subroutine FDIP.	269
110. Subroutine Class	271

CHAPTER I

INTRODUCTION

Recent advances in computer technology allowed for extensive research and development in the area of speech signal processing [1][10]. One area that has been quite popular is speech recognition [2]. Most of the algorithms used in speech recognition deal with only a limited vocabulary of about 250 words. These algorithms require extensive computer storage and time consuming computation, especially when recognition must be speaker independent. A specialized area of speech recognition is digit recognition, which is the main topic of this thesis. It is clear that the recognition of digits requires the recognition of only 10 words for each language. The limited vocabulary of the spoken digits and the limited number of phonemes it uses, gives hope that an efficient algorithm may be found enabling man to communicate easily with machine. However, the great majority of situations with which we identify the concept of recognition will be found to involve ultimate human perception [11][16]. In fact, the element of human perception [19] is difficult to disassociate, in our thinking, from the concept of recognition. The problem is essentially one of understanding the human anatomy of speech production [19], and acoustics [34][39] in order to find a simple yet more reliable, and adaptive model which aid in the development of a robust phonetic algorithm for recognizing an uttered digit of any presently known language.

This thesis presents several new approaches to digit recognition schemes and emphasizes some newly defined and normalized parameters based on area functions. These parameters are used in automatic phoneme segmentation and feature detection of digits spoken in American English and Arabic. Algorithms for segmenting speech sounds into vowel, vowel-like and non-vowel segments are discussed. In addition, parameters used for identifying vowels and detecting nasal segments, turbulence noise segments, dip-classification, etc. are described. Furthermore, an algorithm, based on the correlation coefficients of the RMS energy together with the back-to-total cavity area ratio, is introduced.

Review of the Literature on Speech Recognition

This section deals with some of the earlier work on speech recognition, and, in particular, digit recognition. With the widespread application and recent growth in the use of digital computers, there has been an increasing need for man to be able to communicate with machines in a manner more naturally suited to humans. The realization of this need has motivated a great deal of research in automatic recognition of speech by computer [1][10]. Although only a moderate degree of success has been obtained in solving the problems associated with machine recognition of continuous speech consisting of a series of spoken digits, a greater degree of success has been obtained in recognition of isolated digits.

Research in the mechanical recognition of spoken connected digits not only furnished the foundation for significant advances in pattern recognition and artificial intelligence, but also gave a better

understanding and deeper insight of what a speech recognition system consists of till now, what will be the new system in the future, what makes speech recognition a difficult problem, and what aspect and prospects of this problem remain unsolved. The aim of this research is to find a better, more economical and viable automatic speech and digit recognition system.

Automatic recognition of speech is that process by which a machine attempts to identify correctly certain speech sounds produced by either a human vocal mechanism or some process simulating the output from the human vocal mechanism. Ideally, a speech recognizer should not be constrained by either limited vocabulary size, vocal differences between the individual speakers providing the stimulus to the machine, or any special conditions imposed upon the manner in which speech may be put into the machine [2][3]. Some special input conditions would include, for example, isolation of words by deliberate and artificial intervals of silence or other than normal signal-to-noise conditions. With suitable constraints on vocabulary size, number of talkers and input conditions, several recognition devices have been built and tested with rather encouraging results [2][5][31].

One of the common constraints on an algorithm is the vocabulary size [5], which requires a forced choice from a group of recognizable words. Word recognition is suitable for certain applications (for example, voice dialing of a telephone or programming a computer where only a limited number of commands is required). However, it also has certain limitations. One inherent limitation in word recognition is the trade-off which exists between the size of the acceptable vocabulary and the time required for recognition. This trade-off is a

result of the fact that word recognition generally involves matching some pattern of a spoken word with a stored library of patterns of recognizable words [3]. As the library of acceptable words is increased the storage limit on the machine must likewise increase. The effective storage limit of the machine may be increased by the use of external storage devices, such as disks, tapes, etc. However, the time required to compare a new word to each of the library words may become so large as to make the method impractical [2].

It is common knowledge that the same word spoken by the same person will have some pattern differences. These pattern differences may make it difficult for recognition. Obviously, speaker independent speech recognition is even more difficult since the design of any successful recognition algorithm demands a quantitative knowledge of the words to be recognized in terms of certain features that may be useful. Therefore, a large amount of data is required from different speakers.

Digit Recognition

In recent years, several schemes for digit recognition have been described [1][3][4][12-15]. Most of these have used decision trees based upon empirically derived rules or parameter measurements such as zero crossings, log energy, LPC coefficients, and LPC error. Also, the use of normalized error and pole frequency, followed by phoneme classification of the beginning and ending of the digit has not fulfilled the requirements for an efficient and accurate digit recognition system.

Digit recognition systems or algorithms can be developed using techniques adopted for general speech recognition, such as word and pattern matching with minor modifications. However it is more practical, efficient and accurate, to adopt a new digit algorithm scheme especially at the final recognition stage. One of the advantages is that digit algorithm requires a much smaller vocabulary than one designed for arbitrary word or phoneme recognition. This may lead to the development of a speaker independent digit recognizer. A speaker-adaptive system that can use comparatively simple pattern-matching algorithms to recognize the input digit is simpler and more accurate than any other known system.

For small vocabulary systems, such as digit recognizers with a large number of potential users, it is not feasible to store training data for every possible user. Also, most systems cannot train themselves on new speakers very rapidly. Thus, the turn-around time of new users is often a major factor limiting the use of speaker dependent systems. Consequently it is worthwhile to build up a speaker independent recognition system. In addition, the variation with time of a speaker's voice characteristics may necessitate frequent updating of his reference patterns. Finally, the design of a speaker dependent digit-recognition algorithm is dependent on the uniqueness of each talker's characteristics; whereas, the design of a speaker independent digit-recognition algorithm requires identification of a set of characteristics based on the uniqueness of the phonemic word features obtained from a large number of speakers of different dialect, accent and nationality. The scheme proposed here is dependent upon some

newly defined parameters for quantifying such common characteristics. Also the uniqueness of these patterns and its special features will contribute to an understanding of the acoustic attributes of speech that reliably distinguish the various sounds.

Segmentation

Speech sounds are usually represented by a finite number of distinguishable mutually exclusive, linguistic elements called phonemes. This raises the possibility of segmentation of the acoustical continuum of speech or digit sounds into discrete parts which can then be associated with specific phonemes.

A discrete representation of the time domain requires a segmentation of the continuous speech or connected-digit waveform into some sort of units or "segments". Using a rectangular window, which is described in Chapter III, the spoken digit or digits data can be segmented into n frames. Customarily, each frame consists of 256 data points, which has a duration of 32 ms [1][7] corresponding to an 8 KHz sampling rate. However, it has been found during this research that a frame of 64 data points gives a more reliable and accurate value for our purpose. In addition, using 128 data points per frame of duration 16 ms maintain almost the same accuracy and is more efficient, because it reduces the computational time.

The most often used definition of the term "segmentation" is phoneme segmentation, in which a given continuous speech sentence is phoneme segmented. The number of frames per segment varies according to the duration of the uttered phoneme.

Connected digit, two or more digits uttered successively at normal speech rate, are characterized by a near-continuous (transition) motion of the vocal apparatus from sound to sound. This motion involves continuous changes in the vocal-tract configuration and its modes of excitation. The set of such transitions is wider and more varied in normal vocabulary than in the limited vocabulary of the ten integers.

Humans are able to perform the function of segmentation in a natural way, although, depending on their culture, they may segment phoneme differently. One individual might segment a certain word into four distinct phonemes, whereas another would insist that there are only three phonemes. This is because some phonemes are indistinguishable variants of the same phoneme depending on their position in a given word or digit. Thus, it is unfair to expect a machine to segment a continuous speech or connected digit waveform into discrete segments so that for every segment there is one and only one phoneme. However, for a general phoneme-digit recognition system, it is desirable to have a procedure which first segments the connected digits into isolated digits, using end point detection used by Rabiner [9][12].

Speech Input Rate and Pauses

All automatic speech or digit recognition systems can be considered as belonging to one or two categories: systems designed for the analysis of continuous (connected) speech and those for analysis of isolated (discrete) speech. The two systems have many features in common, which tend to obscure some of the differences. However, isolated speech systems are defined as those systems that require a short pause before and after utterances that are to be recognized as entities [11][23-27].

The minimum duration of a pause that separates independent utterances is on the order of 100 ms [10]. Anything shorter than 100 ms can be confused with the closure of stop constants in the midst of continuous speech that can produce stop gaps approaching 100 ms in duration. In actuality, a stop gap can exceed a 100 ms duration [2][3][10]. For example, the word "seven" can be spoken with a relatively long silence interval "se-". For a trained speaker, however, a 100 ms minimum separation between digits is a reasonable compromise value.

The speaking rate that can be achieved with isolated speech recognition systems is naturally much less than for connected speech. Speaking rates over 300 words per minute can be achieved quite easily for short intervals of connected speech. The upper bound for an isolated word speaking rate has been measured informally for trained speakers reading digits in random order [6][7][47]. A rate of 120-125 digits/min. was achieved with the best speakers. Each of the digits was classified correctly by a machine capable of recognizing isolated utterances. Measuring these rates is not possible without some objective measure that the words are not connected. The human ear is a fairly good judge of whether a brief pause actually exists between rapidly spoken words [5][19][37]. The speech recognition system may have difficulty locating brief pauses.

Average speaking rates between 30 and 70 isolated words (or phrases) have been achieved in factory environments by individuals using voice input systems during their entire 8-hour working day [6]. These average rates will include peak rates close to 120 words/min. and lower than average rates during light workload requirements.

Another important factor concerns the duration of the utterances to be classified. A limited vocabulary system could be hypothesized in which each of the utterances in the limited vocabulary was of considerable duration. However, the previously used limited vocabulary system accepts isolated words as short phrases 2-4 sec. in duration [6].

Human Factors

Clearly, many factors must be considered in order to choose a suitable, well-defined set of parameters for digit recognition. These parameters must allow for pre-emphasis before they are computed. Consequently, the acoustic and segmental aspects of speech have to be considered too, and so doing, the importance of "transitional patterns" of speech might be discovered. Furthermore, despite individual differences in voice quality or speech acoustics and sex differences on speech frequency, the same speech sounds can be recognized. Moreover, humans have the ability to perceive speech in the presence or absence of noise and among competing messages which result in loss of information [11][37]. Given the same information, machines engineered to recognize digits have limited capabilities as compared to human listeners. Surely other factors must be considered, such as accent, educational background, dialects, pronunciation, colds, any upper respiratory and hearing abnormalities.

Last but not least, contextual clues permit the listener to anticipate what might be or might have been said. When combined with prior knowledge about situations or various conversational subjects, one may be able to follow conversations with less than complete perception of all the acoustic information transmitted. Whereas in a digit

recognition system, the perceptual ability of humans depend not only on the quality of the system to transmit the uttered digits but on the position of the digit with a set of uttered connected digits.

Obviously, knowledge of the spoken languages used is fundamental to speech perception more than digit perception. Speech code must be comprehended or its sound patterns, in order to understand the digit order and expand the possibility of compressed digit transmission and recognition algorithm.

Linear Prediction Analysis

The linear prediction method has become important in the area of speech analysis and synthesis because it gives considerable insight into modeling speech production processes [1]. Also it has been shown previously that the linear prediction digital filter represents a non-uniform acoustic tube model of the vocal tract [11]. The computational aspects of linear prediction analysis and spectral modeling will be discussed in Chapter III. This research is intended to use newly defined parameters based on linear prediction analysis and RMS peak ratio. The new parameters are related to the ratio of front, central and back cavity volume to total cavity volume. The new parameters will be used to develop a pattern recognition scheme for each digit.

Previous Systems

One of the first works on the automatic recognition of spoken digits is based on a circuit analyser and quantizer, followed by a pattern matching network [12]. The circuitry is designed such that the frequency band is divided into two bands, upper and lower, with

the separating frequency of 900 Hz. Note that the first two formants are generally separated by this. The frequency of the maximum syllable rate energy within each band is determined using zero crossing. The plots of formants f_2 versus f_1 for the digits 0, 1, to 9 are stored as reference patterns and the pattern of the input digit is compared statistically with the data in store. Unfortunately the statistical analysis is time consuming. Furthermore, the variability of formants creates problems in regard to the pattern matching. In this early system, accuracy range of 97 to 99 percent may be obtained if the same speaker repeats uttering a random series of digits, with 350 ms pauses between digits. The accuracy may fall to 50 percent when different speakers utter a series of random digits.

In the second method, computer simulations are used to compute the power spectra at 10 ms intervals, and segment the words into vowels and consonants [13]. Vowels are classified into one of 11 categories by a multivariate statistical decision, while an empirically derived decision tree is used to classify consonants into one of three categories. A simulated filter bank is used to transform the waveform of the spoken digits into their power spectrum. Each of the 40 filter banks have a half-power bandwidth of 200 Hz. The convolution interval of the time function and the filter response is evaluated to find the power output filter. The ratio of the power output per filter to the total power is computed. A threshold level is adjusted for each band. Then a recognition procedure based on the power band ratio for vowel-consonant, voiced-unvoiced fricative, etc. is developed. This system has disadvantages that the speaker has to be trained, must speak

clearly and has to pause between digits. Important features are lost due to the use of multiple band-pass filters. Also, large variations in pronunciation cause difficulties in preliminary and final recognition.

The third digit recognition system is based on simple segmentation rules according to articulatory features [14]. The digits are segmented according to manner of articulation. Formants f_1 and f_2 are automatically located and evaluated for vowel-like segments. Each spoken digit is represented by a sequence of segments and the values of the first and second formants at the 1st, 2nd, 3rd, 4th and 5th section of the segment is computed. Also, the maximum value of the first formant within that segment is found. An algorithm based on reference patterns is used to match the input digit against the stored established patterns. Error ratio varies from 1.2 percent to 5 percent if the system is trained, and 20 percent and more if the system is untrained.

The fourth recognition system shown in Figure 1 was implemented by Rabiner and Sambur [1][2] for an isolated speaker independent digit recognition system. The system uses end-point detection, four parameter measurements, segmentation of the utterance into intervals, a preliminary decision tree, and a final class decision digit recognition scheme. End-point detection is based on the algorithm developed by the same authors [1][9][12] and uses self normalized measures of the energy and zero-crossing rate (ZCR) of the speech waveform. The four parameters that are of interest are the rate of zero crossings, the energy, and two pole frequencies obtained from LPC analysis and the LPC residual error. These measurements are made for each frame,

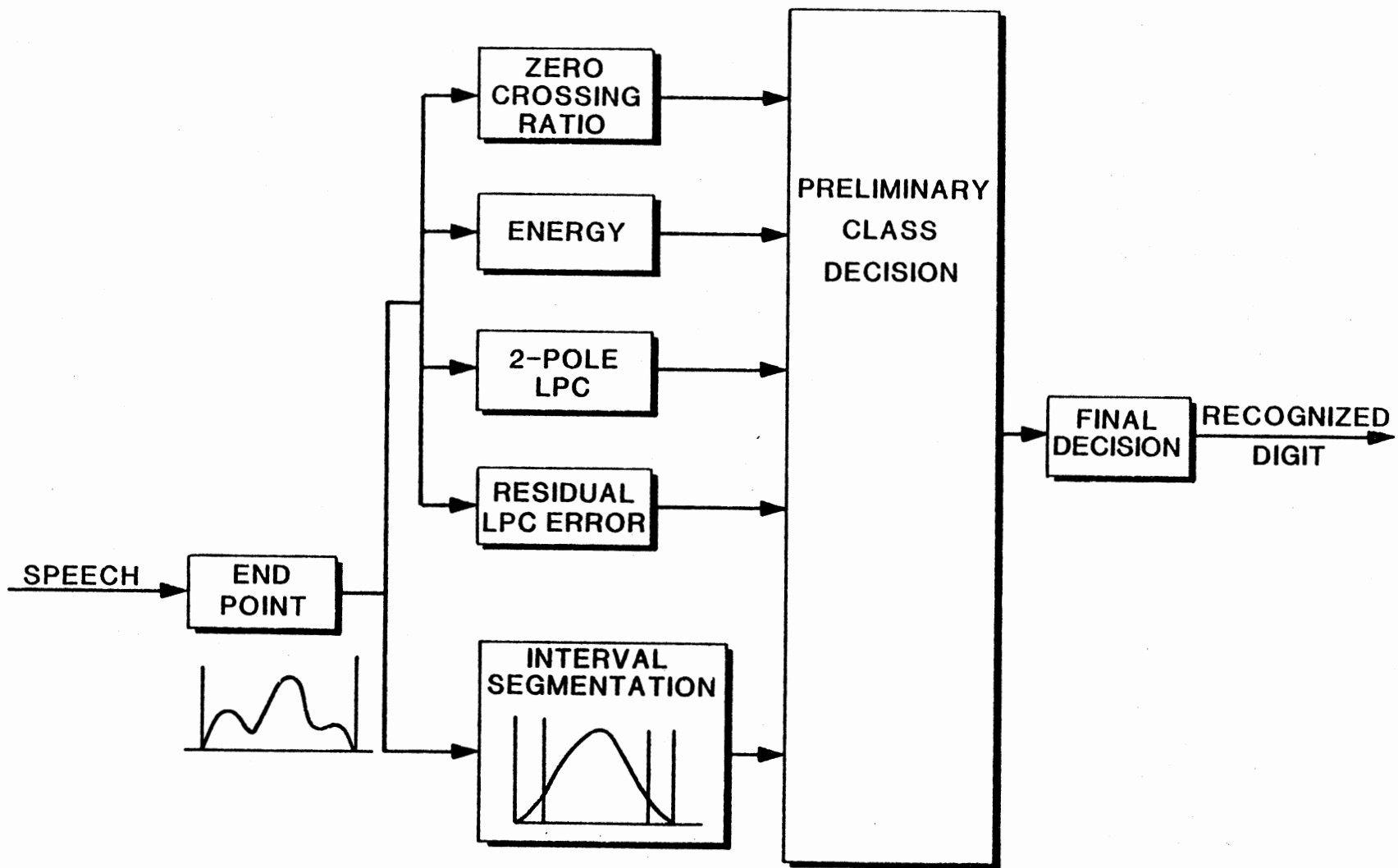


Figure 1. The Over-All Isolated Digit Recognition System Block Diagram After Sambur and Rabiner [4]

where the frame has 256 data points with a sampling rate of 8 KHz. One of the problems with this method is the location of the word boundaries and the dependence of the zero-crossing rate.

The fifth recognition scheme is a modification of the last method [1][4]. A digit segmentation algorithm for continuous digits is added, then a recognition procedure is developed. The accuracy of recognition depends on the preciseness of locating digit boundaries and phoneme segmentation. However, this scheme is less accurate than the one developed previously for isolated digits.

There are other digit recognition systems considered in terms of English, French and other languages. However, the systems given above are the most prominent at the present time.

Thesis Outline

Chapter II deals with anatomy of speech production, and discusses vocal pitch and loudness, articulations, classification of sounds, vowels and consonants. Formant frequencies and its importance in spectrum analysis is discussed briefly. Effect of intensity of sound and noise, frequency and segmental analysis and coarticulation are explained according to their importance to phonemic digit segmentation and recognition.

Chapter III deals with linear prediction analysis (LPA) and its application to the digit recognition scheme. Newly defined parameters are introduced based on area functions derived from LPA. Identification of boundary locations of connected digits using dip-classification is explained, along with phonemic segmentation based on acoustic-phonetic analysis system.

Chapter IV discusses digit recognition scheme for English. The importance of primary segmentation, primary recognition, and final recognition is emphasized. Acoustic-phonetic segmentation and digit recognition flow chart is explained.

Chapter V deals with the Arabic digit phonemes and its relevance to proposed scheme. A comparison table of Arabic and English phonemes is introduced. The modified Arabic digit recognition procedure based on the decision tree is verified. The final digit recognition algorithm, recognition results, accuracy and the correlation matrix are given.

Finally, Chapter VI includes conclusions, suggestions for further research and possibly avenues into speaker verification.

CHAPTER II

SPEECH MECHANISM

Introduction

Acoustical speech waveform results in an acoustic pressure wave which originates from voluntary physiological movements of the human speech mechanism structures as shown in Figure 2. Speech is usually characterized as language that is spoken and heard, and the term is referred to the sounds made by the human vocal apparatus. The generation of sounds of any kind is dependent on the movements in this apparatus. This chapter describes the anatomy and physiology of speech production, its key components as related to the phonetic English and Arabic digits.

The speech waveform is generated due to the variation of pressure above and below the vocal folds. Air is expelled from the lungs into the trachea and then between the vocal folds. The position of the vocal folds across the breath stream allows them to act as a valve to control air flow. Air from the lungs builds up pressure below the vocal folds during glottis closure, i.e. during exhalation or expiration. The vocal folds will be forced open as soon as a sufficient pressure level is reached. The subglottal air pressure will drop slightly, as soon as a puff of air is passed through the glottis,

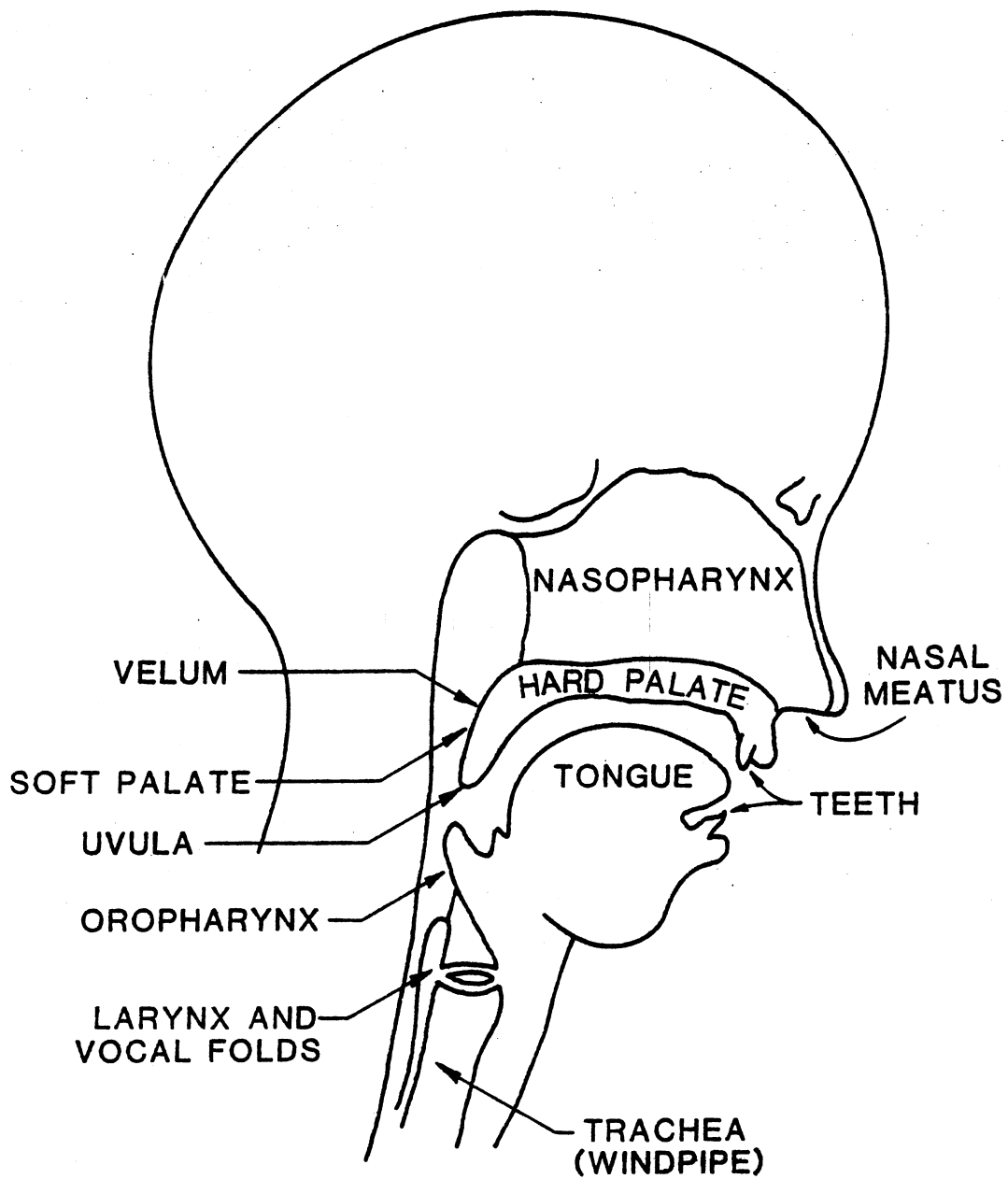


Figure 2. X-Ray Tracing of the Speech Mechanism of a Normal Speaker, Showing the Palate at Rest

then, the vocal folds return to their closed position due to their tension and elasticity. The passage of air through a partially closed glottis causes negative pressure, or suction, which contributes to the complete closure of the vocal folds. This phenomenon is known as the Bernoulli effect [18][19]. The subglottal air pressure will rebuild until it has sufficient force to overcome the forces due to the vocal folds tension and elasticity, and the cycle is repeated. During the generation of voiced sounds the air flowing upward from the lungs causes the vocal folds to open and close at a rate dependent upon the air pressure in the trachea and the physiological adjustment of the vocal folds. This adjustment includes changes in the length, thickness, mucus, and elasticity of the vocal folds. The greater the tension, the higher the perceived pitch of the voice.

The opening between the vocal folds is defined as the glottis. The subglottal air pressure and the time variations in glottal area determine the volume velocity of glottal air flow expelled into the vocal tract. When the movement of the glottis for a complete cycle is repeated about 125 times per second, a tone is generated that has the fundamental frequency of the average adult male voice. The rate at which the glottis opens and closes can be approximately measured acoustically as the inverse of the time interval between observed pitch periods of the acoustic wave. The acoustic energy input to the vocal tract can be determined from the glottal volume velocity wave.

Voice Production

When the vocal cords are relaxed, the air flow is unrestricted through the glottis. When the vocal cords are tensed, their spacing is restricted and the flow of air causes the vocal cords to vibrate in such a way as to modulate the flow of air from the lungs. When the air flow is modulated in this way, the pressure variation of the flow of air into the vocal tract is quasi-periodic and the sound so produced is defined as a voiced sound.

After passing through the larynx, the acoustic pressure and velocity variations of the air are modified by the vocal tract and nasal cavity. The vocal tract is an acoustical tube of non-uniform cross-section, which has its beginning at the vocal cords and ends at the lips [18-20]. The nasal cavity has its beginning at the velum and termination at the nasal meatus. Within the vocal tract and nasal cavity, the pressure and velocity variations of the air are modified by changing the position of the lips, teeth, jaw, tongue, velum and others. These organs are usually referred as the articulators shown in Figure 3. In the production of speech, the articulators are often placed in such a way as to produce a constriction within the vocal tract. This constriction may be made to occur anywhere from the vocal cords (for /h/, as in /wâhid/ for one in Arabic) to the lips (for /f/, as in five) [19][20]. If the constriction is sufficiently narrow, turbulence results, and the vocal tract is said to be frictionally excited and the sound so produced is defined as a fricative. Friction may be produced at a constriction within the vocal tract with or without

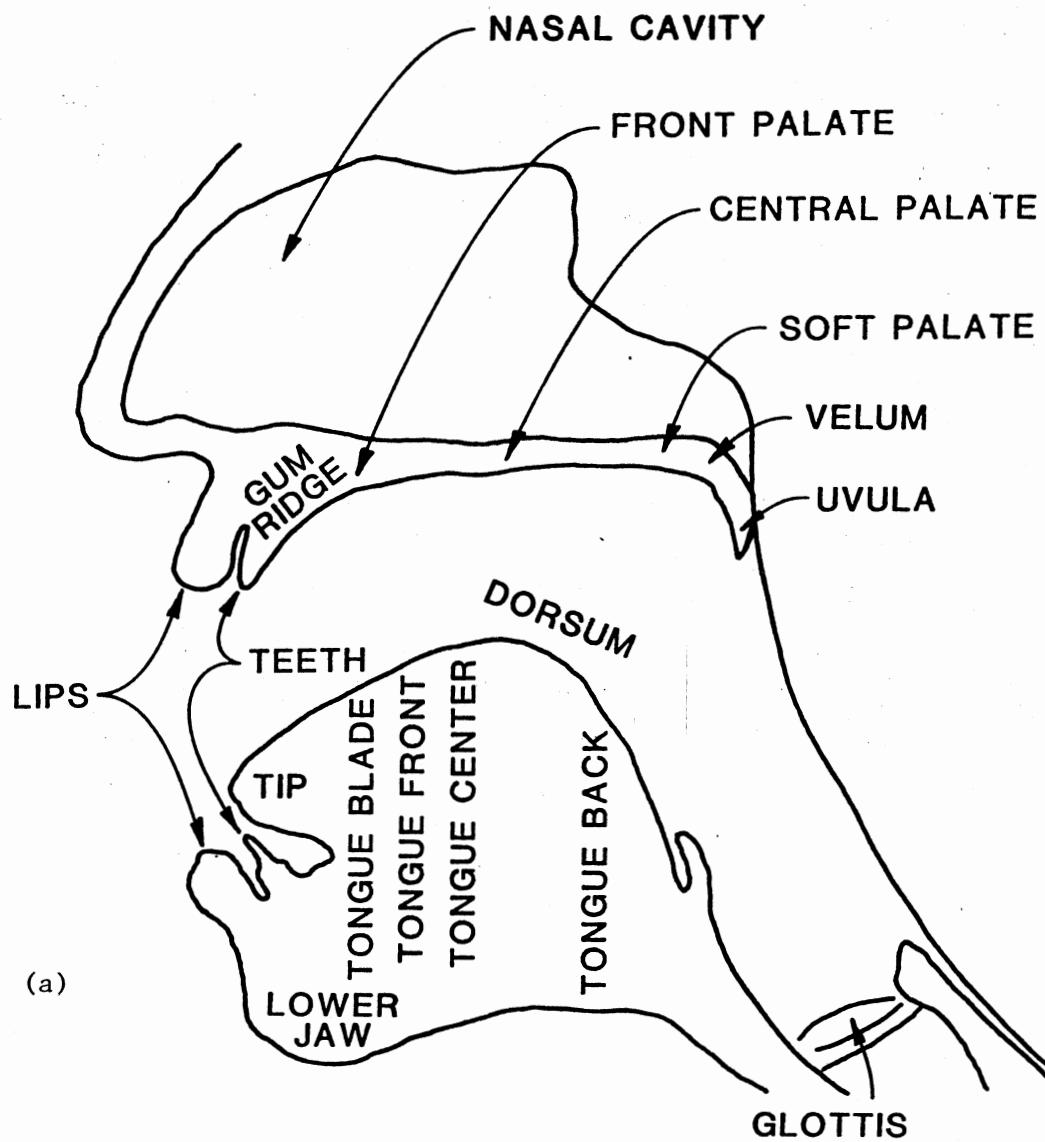
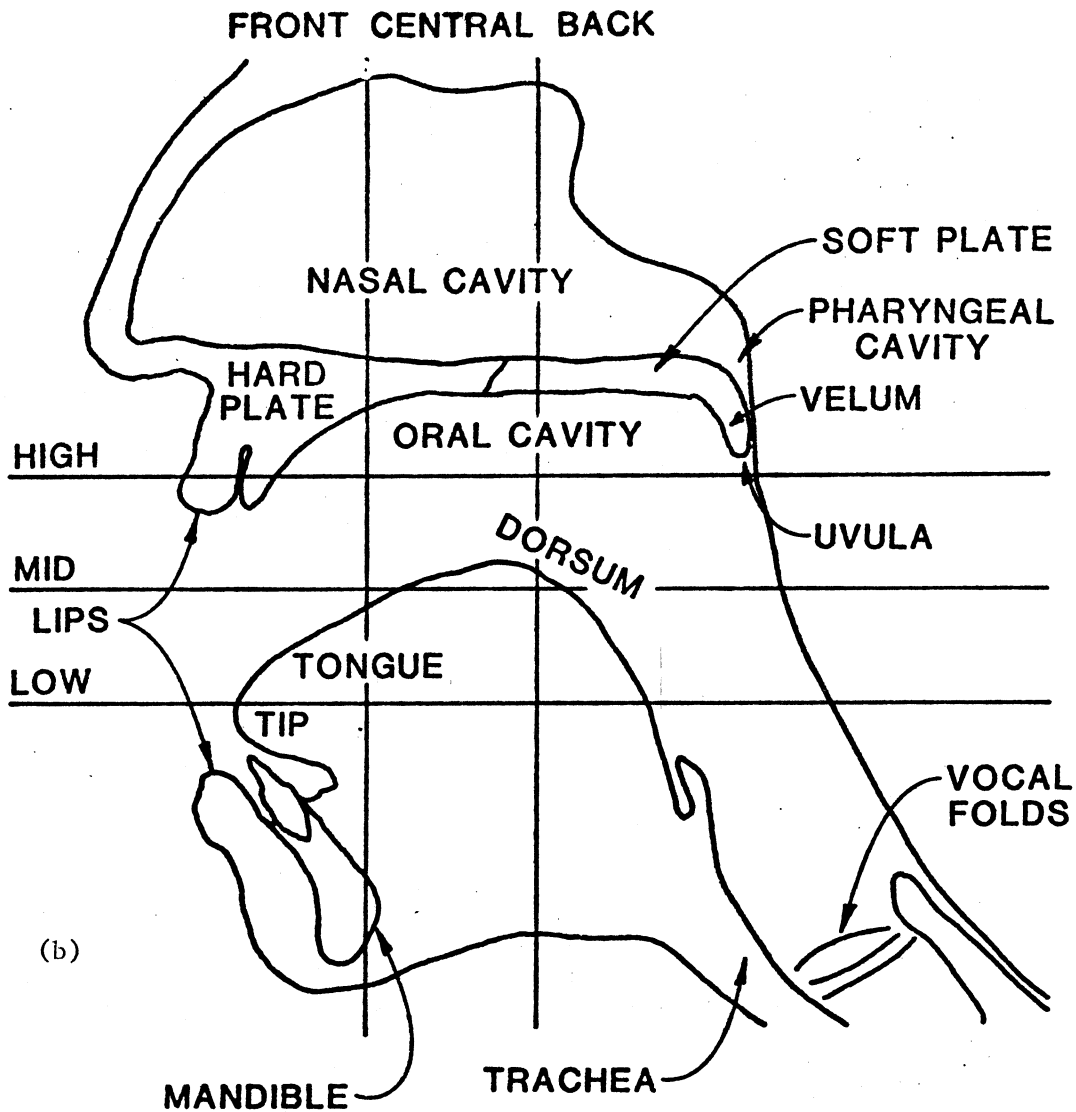


Figure 3. A Schematic View of the Articulators and Places of Articulation Showing the Partitioning Assumed for Classifying the Place and Degree of Constriction of the Vocal Tract During Production of Vowel Sounds

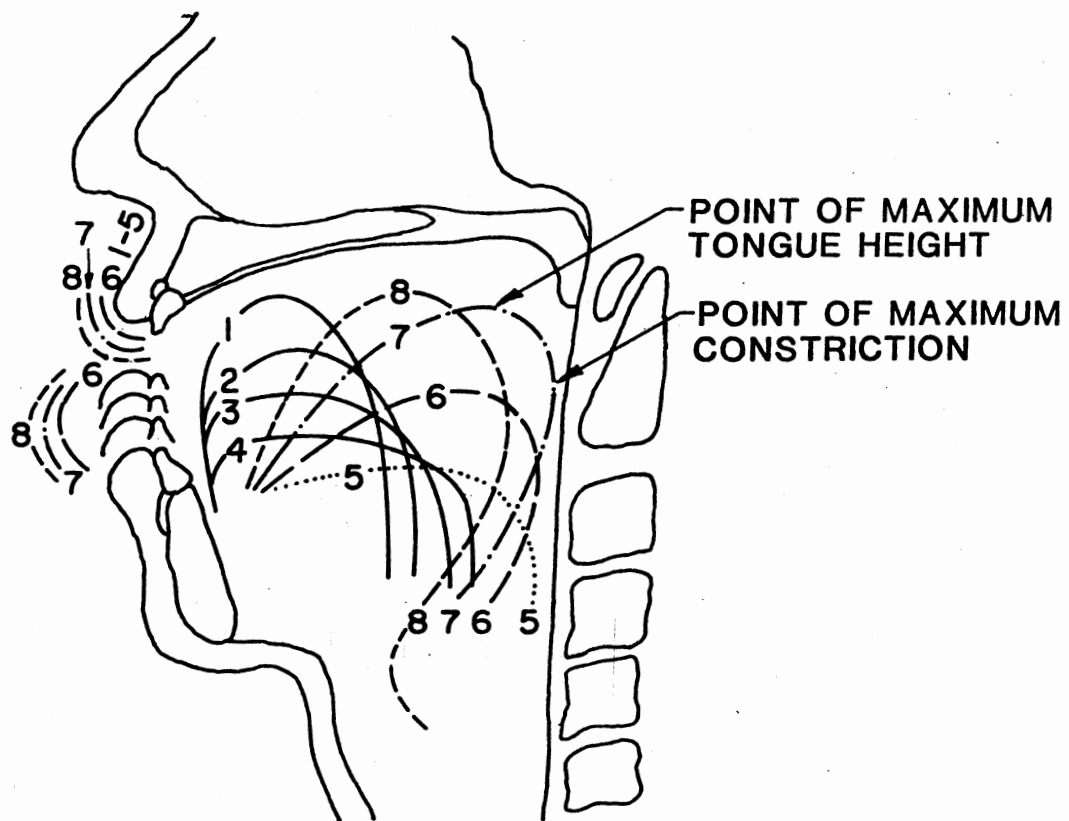


(b)

Figure 3. (Continued)

the presence of voicing. When friction and voicing are both present, the sound produced is called a voiced fricative. When friction is present and the vocal cords are relaxed, the sound produced is called an unvoiced fricative. When the vocal cords are used and when there is no friction in the vocal tract, the sound produced is called a voiced non-fricative.

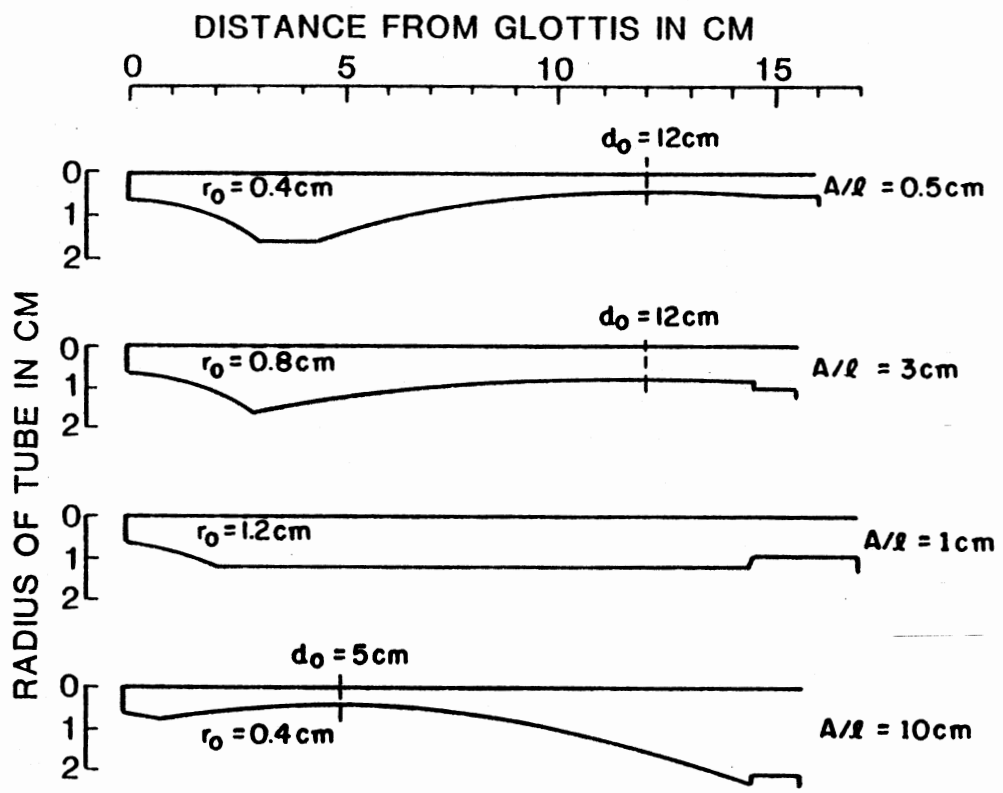
As pointed out earlier, the vocal tract is a non-uniform acoustical tube which is time varying in shape. The major anatomical factors causing this time varying change are the articulators. These articulators cause the cross-sectional area of the lip opening to vary over a range of 0 cm^2 with the lips closed to about 20 cm^2 with the jaw and lips open [10][22][24]. X-ray data show that the cross-sectional area of the vocal tract is controlled primarily by the position and shape of the tongue, as shown in Figures 3-5 [18][24]. The tongue usually forms a constriction or region of minimum cross-sectional area, during the articulation of vowels. The cross-sectional area can vary from 0.3 cm^2 to 10 cm^2 at the lips, with a variation of the restriction radius from $d_0 = .4 \text{ cm}$ to 1.2 cm . The distance from the glottis to the restriction at which the smallest radius is measured varies from $d_0 = 5 \text{ cm}$ to 12 cm , whereas the vocal tract average length is 17 cm , from the glottis to the lips. In general, x-ray results show that during the articulation of vowels, the dimensions of the vocal tract along the length of the tongue are controlled primarily by the position of the tongue constriction and by the degree of tongue constriction. Whereas in the region beyond about 15 cm from the glottis, the mandible and the lips determine the cross-sectional area.



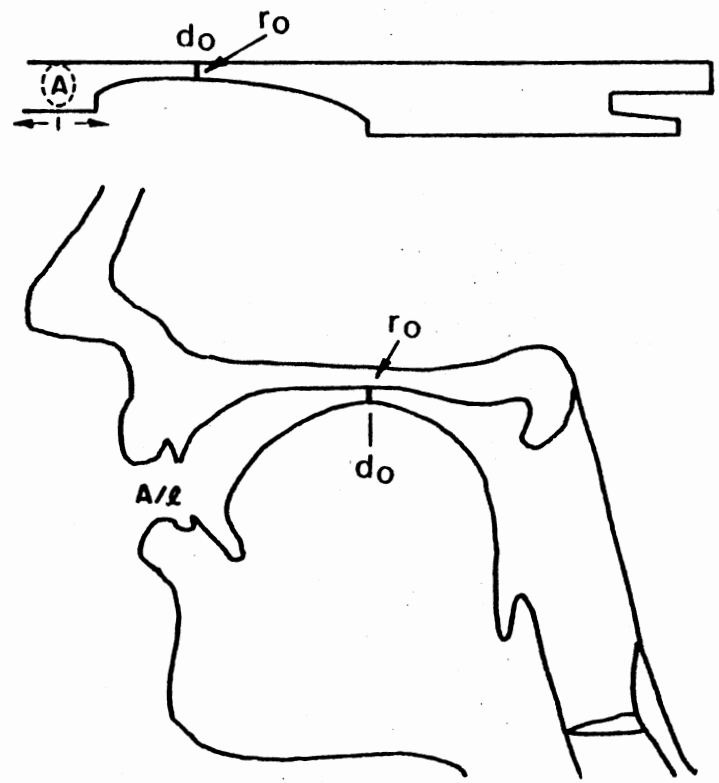
TONGUE AND LIP POSITIONS
FOR THE VOWELS

(1) [i]	(5) [a]
(2) [e]	(6) [ɔ]
(3) [ɛ]	(7) [o]
(4) [æ]	(8) [u]

Figure 4. Normal Articulation Process

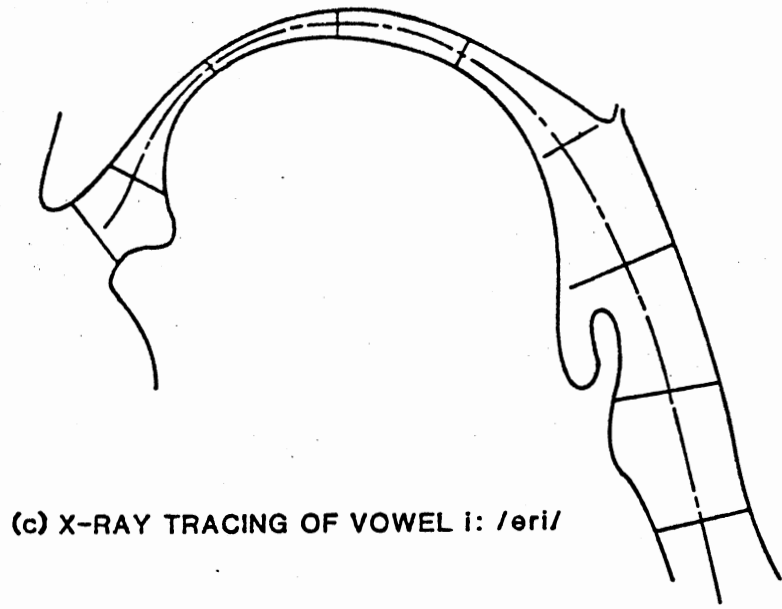


(a)

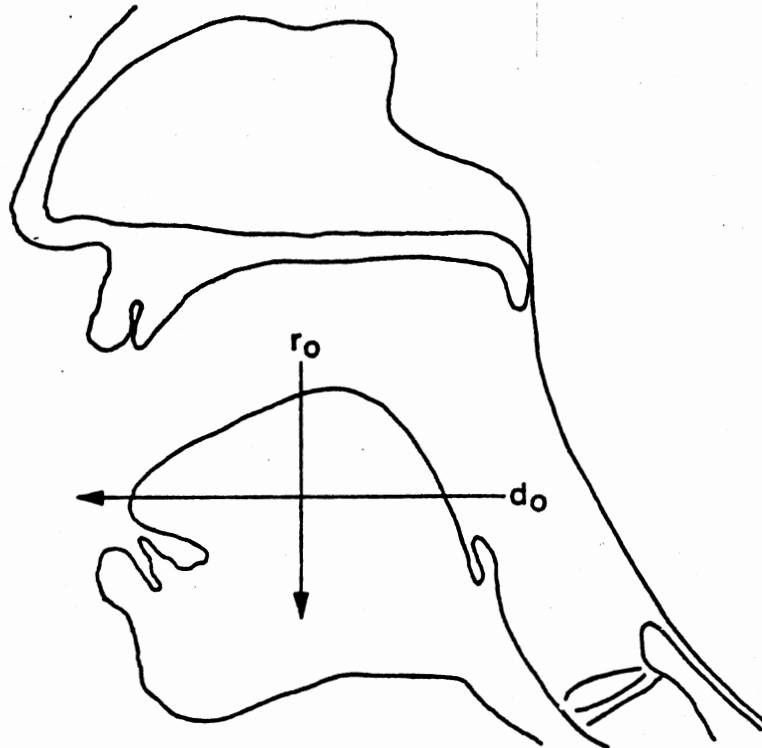


(b)

Figure 5. X-Ray of Cross-Section Area for Different Constrictions



(c) X-RAY TRACING OF VOWEL i: /eri/



(d)

Figure 5. (Continued)

The above discussion is for non-nasal sounds. For nasal sounds, the velum closes the vocal tract from the nasal cavity during the production of these sounds. The nasal sounds /n/ (wʌn/ and /naɪn/ in English digits, for example, /m/ (θamānyðh, for eight in Arabic) and /ŋ/, (/nɔ/ English and Arabic digits do not have this sound) uses the nasal tract.

Vocal Pitch and Loudness

Vocal-fold tension is increased by contraction of the vocal-fold muscle. Releasing more puffs of air in a given period of time increases the frequency and pitch of the tone produced. In other words, the pitch of the human voice changes in accordance with changes in the mass, tension, and length of the vocal folds. Adult male's voices have lower pitch than female's voices because the male larynx is larger and has longer folds than the female larynx [18][19].

Adjustment of subglottal air pressure changes the intensity of the voice; that is, the greater the subglottal air pressure, the more intense the voice. Pitch is primarily a function of the laryngeal system. Loudness is related to intensity, which is primarily a function of the respiratory system. However, these two systems do not work independently of each other, and to maintain good voice control, the speaker must use the larynx and the air stream in a skillful balance. This coordination requires the participation of all the articulator and many components of the nervous system [23] [25].

Figure 6 is a schematic representation of changes in glottal area during a vibratory cycle. The area is determined by movement of the

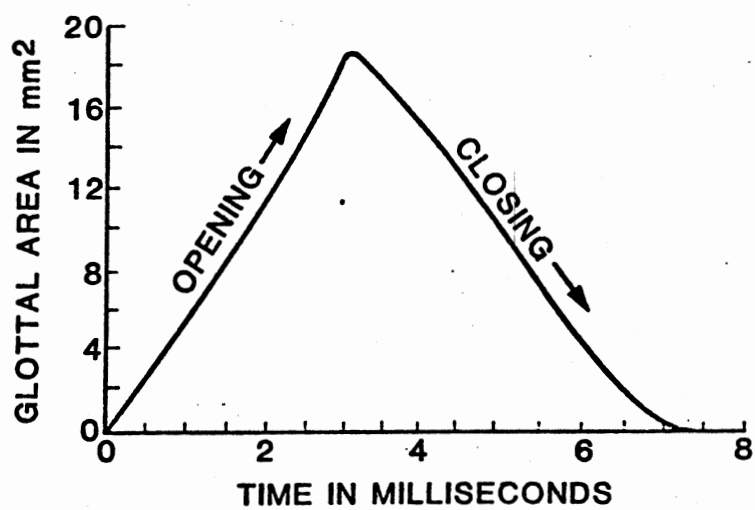


Figure 6. Schematic Representation of Changes in Glottal Area During a Vibratory Cycle

vocal folds, and the steepness of the slope of any segment of the curve represents the velocity with which the vocal folds move [22].

Articulation

X-ray observations of the upper respiratory system show that the larynx opens into a passageway called the pharynx, which in turn opens into two nasal and vocal tracts. The walls of the pharynx, the soft palate (i.e. the velum), tongue, mandible, and lips are somewhat free to move. Their movements change the configuration of the pharyngeal and oral airways or tubes. These structures shown by Figures 3 and 4 are referenced before as articulators because they are involved in the production of speech sounds [22][24]. The soft palate is a muscular continuation of the hard palate. It is very mobile and can move rapidly to close or open the pharyngeal air passage between the oral and nasal cavities. Movement of the pharyngeal walls toward the palate often contributes to the closure of the nasal cavity. The velum is opened during the production of the nasal consonants /m/, /n/ and /ŋ/ (i.e., nG). The velum is open during the production of these three consonants sounds, allowing sound from the larynx to be modified by the nasal cavities as well as by other articulators.

The front teeth contribute to articulation. This can be approximated by the movements and positions of the tongue and the lips. The facial muscles allow the lips to have the proper shape for the production of both vowels and consonants. The lips supplement the tongue fairly well in shaping the vocal tract for the production of vowel and

consonant sounds. The lip sounds are among the first speech sounds that a baby acquires, and these sounds can be found in all oral languages.

Classification of Speech Sounds

It is appropriate to introduce the sounds of speech in terms of the articulators and movements that produce them. The concept of phonemes in continuous speech and connected digit utterances can now be easily appreciated. Therefore sounds belonging to different phonemes is classified according to the movements and positions of the articulators that produce them [1][6][11][24][26].

It is convenient to divide sounds into vowels, diphthongs, semi-vowels, and consonants for spoken digits as shown in Figures 7-10 for digits spoken in English and Arabic, respectively. Tables I-IV present the classification of phonemes for digits spoken in English and Arabic according to the tongue positions and classification of vowel sounds using international phonetic alphabet.

Consonants

Consonants are differentiated by place of constriction, manner of constriction, or the presence or absence of voicing or laryngeal tone. Unvoiced consonants are produced due to the flow of air through the constriction between maximum tongue hump and the palate, while the vocal folds are not vibrating. Any two consonants will differ from each other in terms of one or more of these three features. Place of articulation encompasses the structures from the front to the back of the mouth, and two articulators are necessary to establish a place of constriction [18][22]. For example the sounds produced by closure of

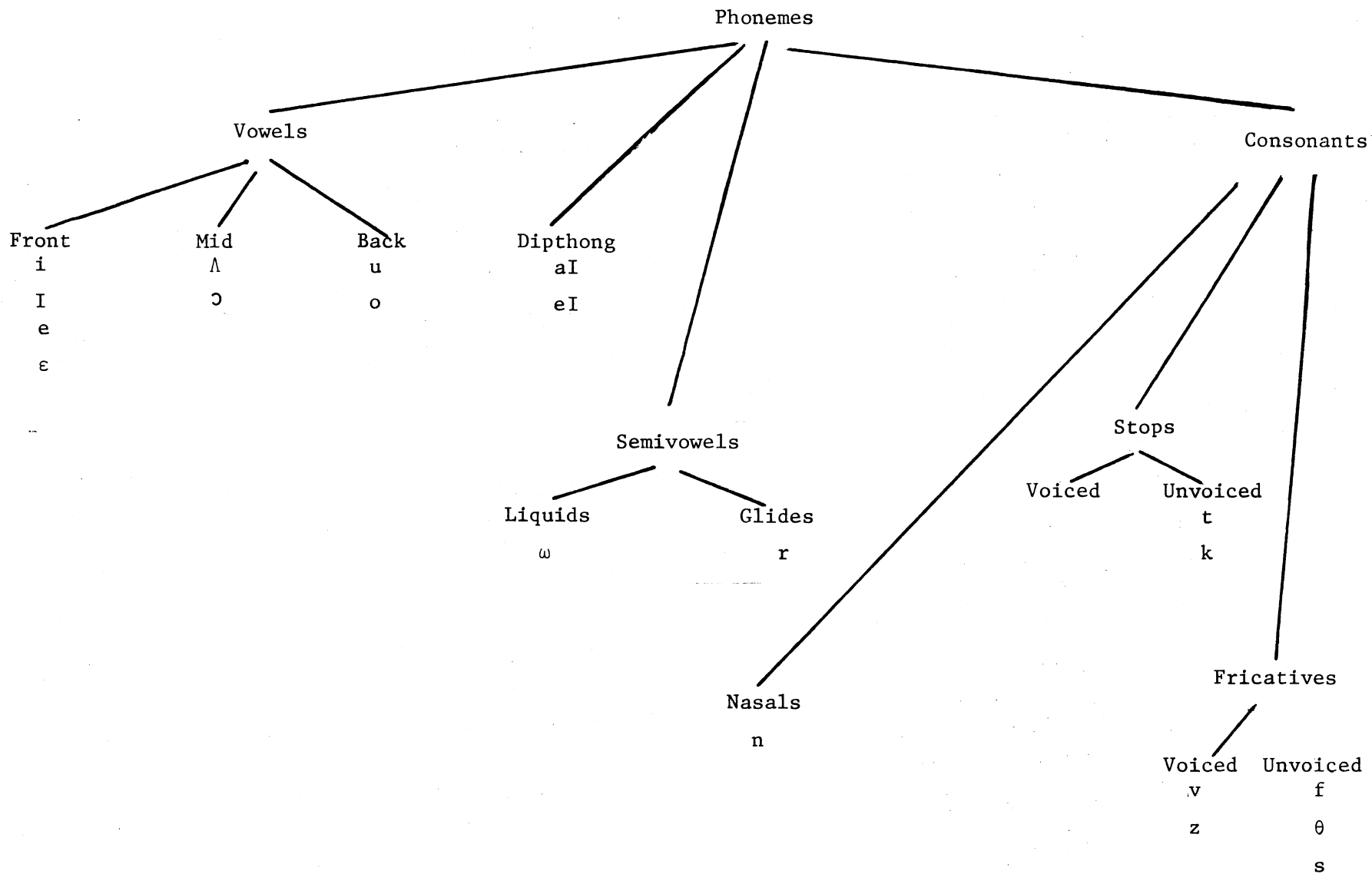
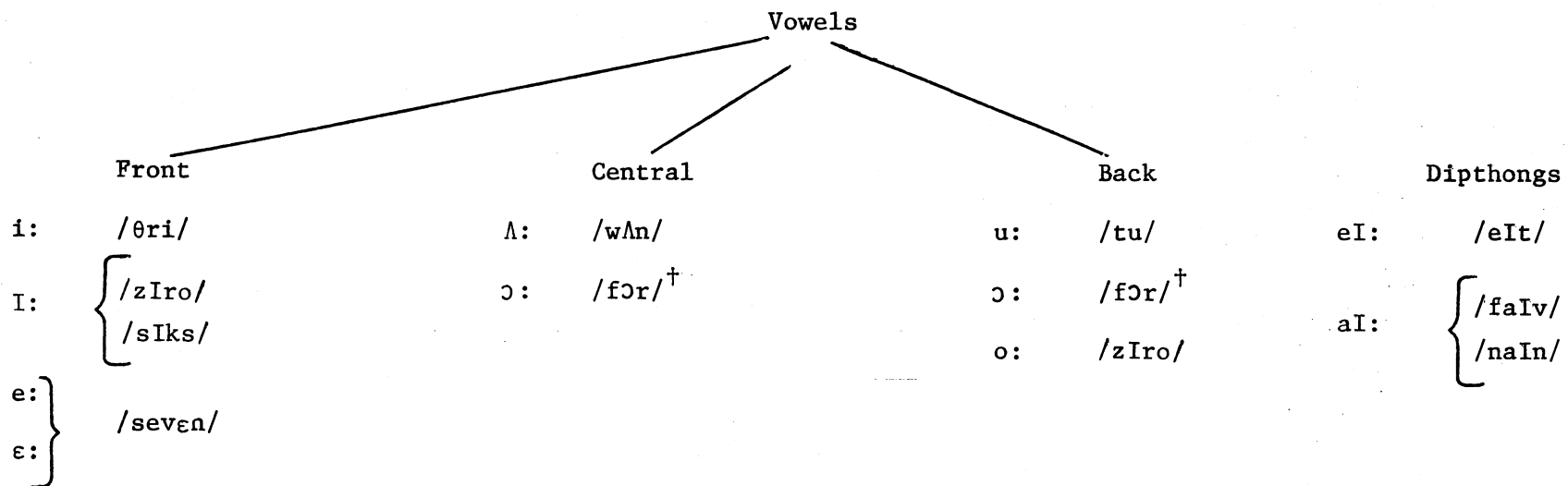


Figure 7. Phonemes in American English Digits



[†]Rabiner considered /ɔ/ as mid-vowel, whereas Flanagan considered /ɔ/ as back-vowel.

Figure 8. Vowels Used in Digits 0 - 9 Spoken in American English

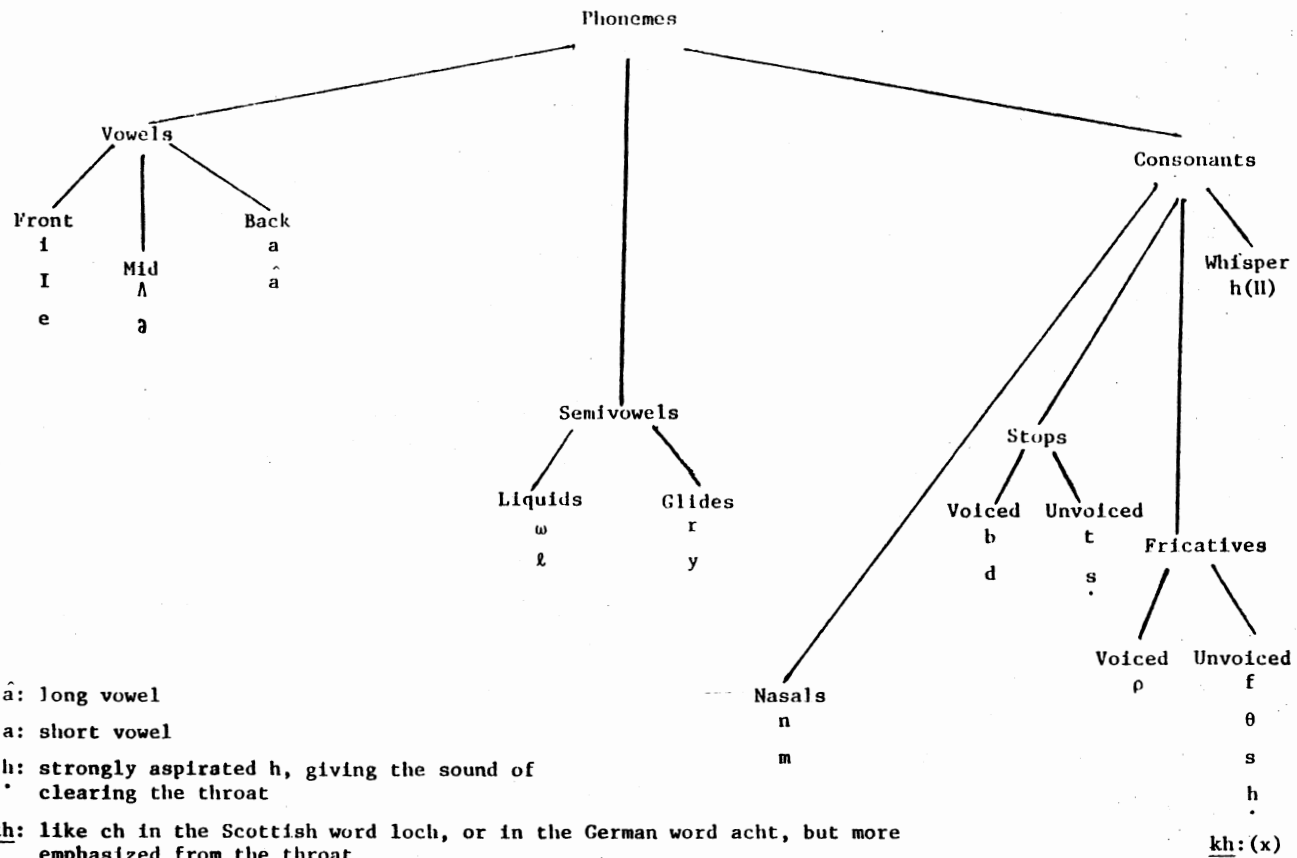


Figure 9. Phonemes of Digits Spoken in Arabic

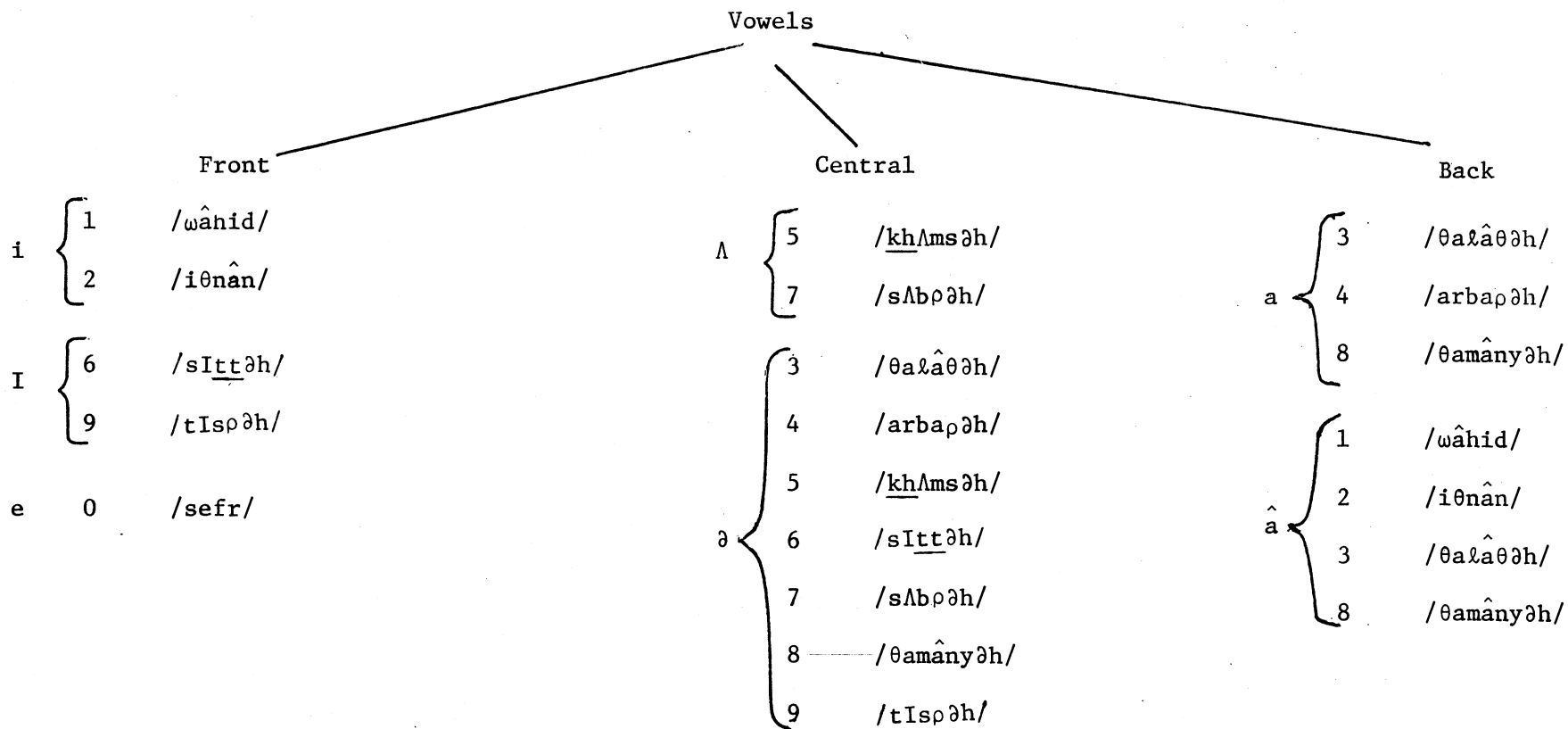


Figure 10. Vowels Used in Digits 0 - 9 Spoken in Arabic

TABLE I

VOWELS USED IN DIGITS SPOKEN IN AMERICAN ENGLISH, ACCORDING TO THE
DEGREE OF CONSTRICTION AND TONGUE HUMP POSITION

Degree of Constriction	Tongue Hump Position		
	Front	Central	Back
High	/i/ : /θri/		/u/ : /tu/
	/I/ : { /sIks/ /zIro/		
Medium	{ /e/ /ɛ/ : } /seven/	/ʌ/ : wʌn	/o/ : /zIro/ /ɔ/ : /fɔr/
Low			
Diphthongs	eI: /eIt/ aI: { /faIv/ /naIn/		

TABLE II
SEQUENCE OF SOUND CLASSES OF DIGITS SPOKEN IN AMERICAN ENGLISH

Digit	Phonemes	Sequences of Sound Classes					Notes
0	/zIro/	Voiced fNLC	FV	VLC	BV		
1	/wAn/	VLC	MV	VLC			
2	/tu/	Unvoiced sNLC	FV	BV			
3	/θri/	Unvoiced fNLC	VLC	FV			
4	/fɔr/	Unvoiced fNLC	BV	MV	(VLC)		
5	/faIv/	Unvoiced fNLC	MV	FV	Voiced fNLC		(Diphthong)
6	/sIks/	Unvoiced fNLC	FV	Unvoiced sNLC	Unvoiced fNLC		
7	/seveŋ/	Unvoiced fNLC	FV	Voiced fNLC	FV	VLC	
8	/eit/	FV	Unvoiced sNLC				(Diphthong)
9	/naIn/	VLC	MV	FV	VLC		(Diphthong)

Voiced, fNLC: voiced fricative noise-like consonant
 Unvoiced, fNLC: unvoiced fricative noise-like consonant
 Unvoiced, sNLC: unvoiced stop noise-like consonant
 VLC: vowel-like consonant
 FV: front vowel
 MV: middle vowel
 BV: back vowel

TABLE III

VOWELS USED IN DIGITS SPOKEN IN ARABIC, ACCORDING TO THE DEGREE
OF CONSTRICTION AND TONGUE HUMP POSITION

Degree of Constriction	Tongue Hump Position		
	Front	Central	Back
High	<i>/i/</i> : { <i>/iθnân/, '2'</i> <i>/wâhid/</i>		
	<i>/I/</i> : { <i>/sIttðh/, '6'</i> <i>/tIspðh/, '9'</i>		
Medium	<i>/e/</i> : <i>/sefr/, '0'</i>	<i>/Λ/</i> : { <i>/khΛmsðh/, '5'</i> <i>/sΛbρðh/, '7'</i>	
		<i>/ð/</i> : { <i>/θaλâθðh/, '3'</i> <i>⋮</i> <i>/tIspðh/, '9'</i>	
Low			<i>/a/</i> : <i>/arbaρðh/, '4'</i> <i>/â/</i> : { <i>/wâhid/, '1'</i> <i>/iθnân/, '2'</i> <i>/θaλâθðh/, '3'</i> <i>/θamânyðh/, '8'</i>

TABLE IV
SEQUENCE OF SOUND CLASSES OF DIGITS SPOKEN IN ARABIC

Digits	Phonemes	Sequences of Sound Classes						
0	/sefr/	Voiced sNLC	FV	Unvoiced fNLC	VLC			
1	/wâhid/	VLC	Long BV	Unvoiced fNLC	FV	Voiced sNLC		
2	/iθnân/	FV	Unvoiced fNLC	VLC	Long BV	VLC		
3	/θalâθðh/	Unvoiced fNLC	MV	VLC	Long BV	Unvoiced fNLC	MV	Whisper
4	/arbaḡðh/	MV	VLC	Voiced sNLC	MV	Voiced fNLC	MV	Whisper
5	/khΛmsðh/	Unvoiced fNLC	MV	VLC	Unvoiced fNLC	MV	Whisper	
6	/sittðh/	Unvoiced fNLC	FV	Unvoiced sNLC	MV	Whisper		
7	/sΛbḡðh/	Unvoiced fNLC	MV	Voiced sNLC	Voiced fNLC	MV	Whisper	
8	/θamânyðh/	Unvoiced fNLC	MV	VLC	Long BV	VLC	MV	Whisper
9	/tIspðh/	Unvoiced sNLC	FV	Unvoiced fNLC	Voiced fNLC	MV	Whisper	

Voiced, sNLC: voiced stop noise-like consonant

Voiced, fNLC: voiced fricative noise-like consonant

Unvoiced, sNLC: unvoiced stop noise-like consonant

Unvoiced, fNLC: unvoiced fricative noise-like consonant

VLC: vowel-like consonant (semi-vowels), Liquids and Glides

FV: front vowel

MV: middle vowel

BV: back vowel

Note: All vowels are short vowel, except /â/, which is a long vowel.

lips (/p, b, m/) as in /arbapðh/ (four in Arabic) and /θamânyðh/ (eight in Arabic). Constrictions produced by lips and teeth results in the production of sounds /f, v/, as in four, five, and seven. During the production of a stop sound, articulators momentarily occlude, or stop, the oral air passage. Air pressure is built behind the occlusion and releasing this results in a stop sound, such as /t/ as in two, and /d/ as in /wâhid/ (one in Arabic). The affricates are much like the stops, except that the pulse of air is sustained a bit longer as in /t/ and /dz/ [24][25]. Fricatives are caused by the approximation of two articulators, thus directing exhaled air through a narrow opening, causing a relatively continuous stream of noise. The nasals are sounds that are made by lowering the velum, thus directing the sound stream through the nose rather than the mouth. The semi-vowel glides /w, r, j/, as in /wAn/, /fɔr/ and /wâhid/ are made with more constriction than vowels, but not enough to cause turbulent air flow. The third type of consonant however is differentiated by the presence or absence of voice, since some consonants are voiceless and some are voiced. The onset of vocal-fold vibration occurs earlier in voiced than in voiceless consonants. That is, vocal-fold vibration is present during portion of the consonants that are perceived as voiceless [18-20]. On studying the phonemes that are used only in spoken digits, the terms vowel-like and non-vowel-like is applied.

Vowels

Classification of vowels require a different analysis from that used for consonants. All vowels are voiced. Although lip configuration is important, differences among vowels are determined primarily by the

position of the tongue tip, tongue edges, and tongue body. Denes and Pinson [21] classify vowels in terms of the position of the highest part of the tongue body. For example, they state that when the tongue body is positioned as high and as far forward as possible without causing turbulence when the lips are spread and when voicing is produced, the vowel /I/ (as in /sIks/ results.

Vowels are classified as front, center, and back and as high, middle, and low relative to the position of the tongue body in the oral cavity as shown in Figures 3 and 4 and Tables I and III. The approximate positions of American English vowels relative to one another can be clearly indicated. Whereas, diphthong sounds are produced whenever a shift in tongue location from that associated with one vowel to that of another vowel results (as in /eIt/ and /naIn/).

Speech Acoustics

An understanding of certain principles of acoustics must be acquired in order to comprehend the production and transmission of speech. Speech is an acoustical phenomenon, and it is a special case of sound production [34]. Although vibration is not sound unless it is heard, but for convenience, it is feasible to accept simply that sound occurs as a result of vibration, as shown in Figure 11.

A source of energy is needed to set a vibrator in motion, for sound generation. Exhalation of air from the lungs may be considered to be the source of energy, and the vocal folds to be the vibrator, for the production of speech. The air in the cavities of the throat and mouth provides the elastic medium for the transmission of these

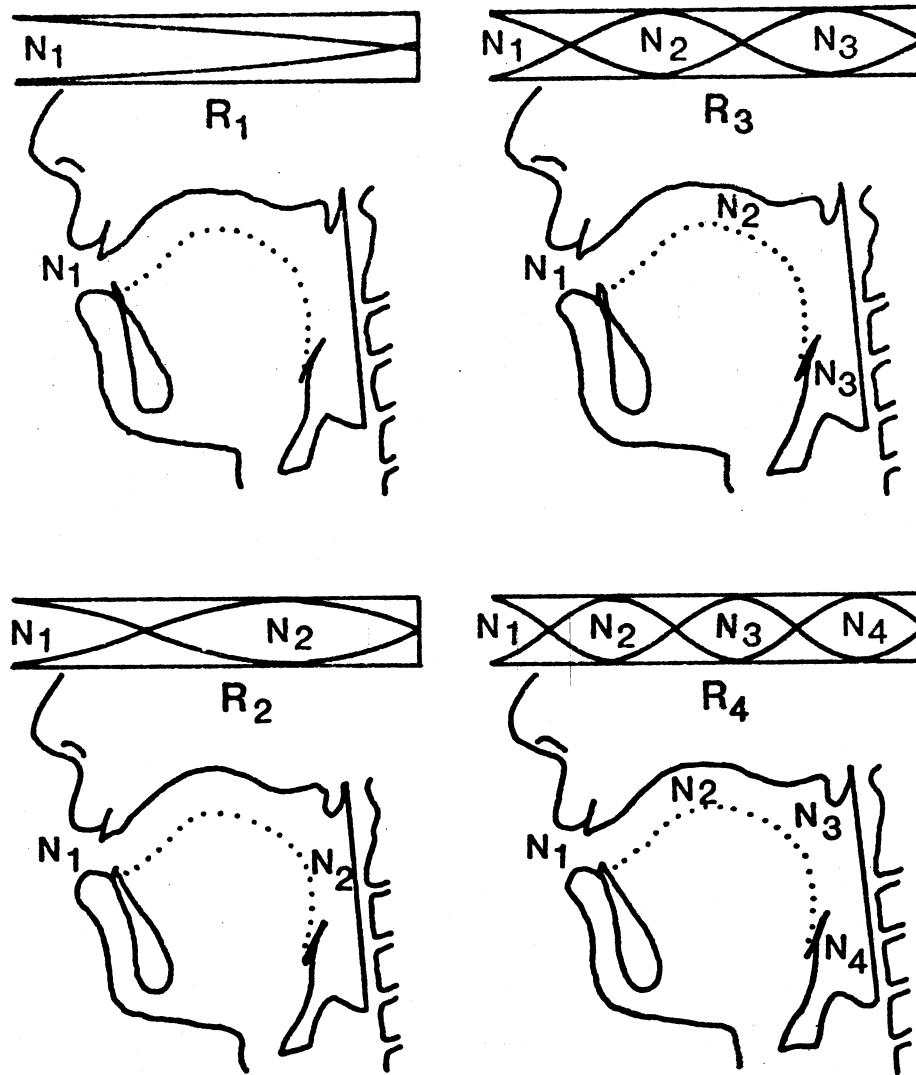


Figure 11. Distribution of Volume Velocity at the Frequencies of Each of the First Four Resonances of an Ideal Neutral Articulation in Which the Vocal Tract Simulates a Tube of Constant Cross-Sectional Area

vibrations or speech sounds. It is feasible to assume that the vibration of the vocal folds causes a series of compression waves or pressure changes in the surrounding air. Since reception and perception of speech by the ear depends on air conduction, special attention must be given to the mode of sound transmission in air.

Sound Propagation

Sound waves in air are described as longitudinal waves, because the molecular movement in air is horizontal, or parallel to the driving force. Sound waves travel as a result of the patterns of molecular displacement in the same direction as the applied force, and the intensity fade away as the distance from the source increases. The density of the material which is defined as "the number of molecules per unit volume", determines the velocity at which sound may be propagated. It is well known that sound travels faster in denser materials. For example, sound travels faster through steel (about 15,000 feet/second) than through air (about 1100 feet/second).

The theory of speech as wave motion and how speech waves are produced and heard is usually included and covered by the acoustic of speech, which is the field of study, that inspired researchers of various specialties during the last decades. Speech sounds defined from their production within the vocal tract have been till now, the most interesting and popular field of research for classical phonetics, well known as articulatory phonetics. The speech wave, which is defined by the sound pressure variations at a point in front of the speaker

has not only been of great concern to speech research of communication engineers, but also of more concern to speech pathology researchers. Complete specifications of the speech wave can be obtained with the aid of modern sound recording and analysis techniques.

Resonance

The vibration of the vocal folds produce only a complex buzzing sound in isolation. However, in normal speech production, the vocal folds cause variation of pressure, which resonate the vocal tract. The quality of voiced sounds produced by the vocal tract resonance, is quite different from that which would result from vocal-fold vibration alone. Hence, it is clear to state the fact that the spectrum of an acoustic signal is influenced considerably by the acoustic environment in which the signal is produced and propagated [35][39].

By examining the concept of resonance, it can be depicted that all objects or volumes of air in open and closed tubes and cavities vibrate more readily at certain frequencies than at others. That is, the cavity will act as a selective filter, passing some frequencies and rejecting other frequencies. In other words, all objects or volumes of air in open or closed tubes, and cavities, have certain natural frequencies of vibration and thus are more responsive to those frequencies. Consequently, when sound with a complex spectrum are produced, the frequencies in that complex waves inherent resonances frequencies, with greater amplitude if these frequencies are the natural frequencies of a given object or volume of air. This is due to the fact that an object will resonate in sympathy of the vibrating

source. Thus, when sounds with a complex spectrum are produced, the cavity resonates more freely and resonance frequencies with greater amplitude are generated if these frequencies are the natural frequencies of a given object, cavity or volume of air.

In the acoustics of speech, the concept of resonance is of great importance and of particular significance. The vocal-fold vibration contains a fundamental frequency with certain harmonics or overtones, since the vibration of the vocal folds is quasi-periodic. The vocal tract will respond to certain frequencies generated by the vocal folds and reject the others since the cavities of the vocal tract act as a resonator with certain natural frequencies. The vocal tract may assume different resonant frequencies, because, the cavities of the vocal tract are modified continuously during the production of speech.

Formant Frequencies

The resonances of the vocal tract are usually called formants, and the frequencies to which they respond more are known as the formant frequencies. The spectrum of vowel sounds have very apparent formant frequencies, and are revealed as resonance peaks, that is, the peaks of maximum energy or amplitude at given frequencies. Different formant frequencies also occur since the configuration of the vocal tract must be changed to produce different speech sounds. The fundamental frequency and resultant harmonics are produced by the vibration of the vocal folds. The resultant resonant frequencies are not the same as the harmonics produced by the vocal folds unless by coincidence

because the formant frequencies are produced from the changes in the time-varying vocal tract [18]. Figure 12 depicts a periodic signal generated by the vocal folds and the resultant harmonic spectrum, which reflects the fundamental frequency of the vocal-fold vibration and the respective harmonics. Also the vocal-tract configuration and the resonant or formant frequencies that occur as a result of vocal-tract resonance due to the vocal-fold output are clearly indicated in Figure 12.

As mentioned earlier, vowels are characterized by quasi-periodic signals and formant frequencies. Consonants are characterized by random signals which may or may not include periodic information as a result of voicing. The voiceless consonant contains only aperiodic vibration, which results from air turbulence produced by changing and constricting the size of the orifice of the mouth and friction or interference in air flow by the articulators such as the tongue, teeth, and lips. Voiced consonants include the quasi-periodic vibration of the vocal folds and aperiodic vibration depicted in voiceless consonants.

Acoustic Characteristics and Perception

The two speech parameters that are widely used in speech analysis are intensity and frequency. The effects of intensity is first considered since intensity changes are the least complicated subject to study. Speech must be sufficiently loud or intense to be clearly understood. Among the various speech sounds of English, the intensity range is 680 to 1, or 28 dB, from the weakest to the strongest speech sound [11][16][18]. The energy or sound pressure of normal speech measured

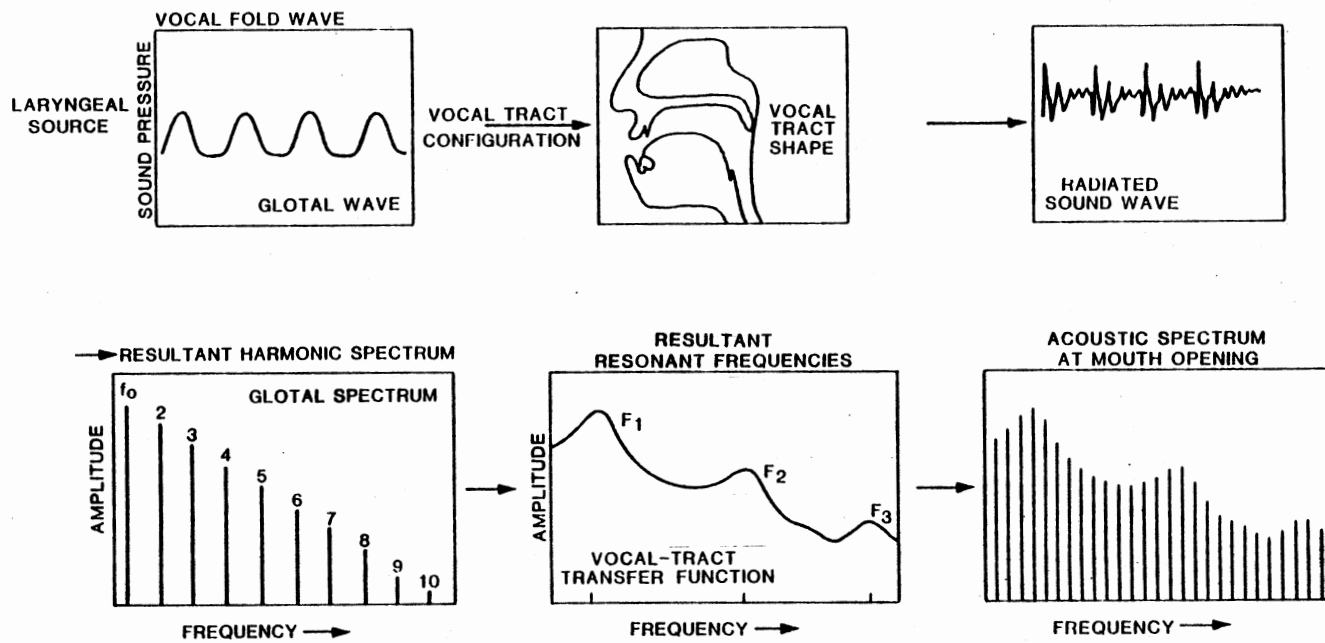


Figure 12. Schematic Diagram of the Physiological and Acoustical Characteristics of Speech Sound Production

at one meter from the lips of the speaker is about 62 dB. Speech signal might reach a sound-pressure level of about 85 dB, if a person talks as loudly as possible, while the softest speech is about 45 dB [16][21].

Vocal-folds vibration supply the internal energy to drive the vowel sounds. The fundamental frequency of the vocal folds is about 125 Hz for males and 250 Hz for females which is an octave higher. The primary energy in vowel sounds is of low frequency, since vowel sounds are quasi-periodic signals composed of the fundamental and related harmonics of the vocal-fold vibration. Also, most of the energy or power of speech are known to be carried (acquired) by vowels and in addition have longer duration than consonants. In fact 60 percent of the energy or power of spoken speech utterance is contained in the frequency range below 500 Hz, which contribute only 5 percent to intelligibility. Also, 60 percent of the intelligibility is contributed by the frequencies above 1000 Hz with only 5 percent of power. Thus most of the power is carried by the vowels, or low-frequency sounds, while the consonants, or higher-frequency speech sounds, which carry lower power, are much more important for speech intelligibility [21].

Effects of Noise

The intelligibility of speech in quiet listening environment is not generally a difficult task. Frequently speech is perceived in the presence of noise, in real-life situations. However, noise does not always interfere seriously in understanding speech unless the noise

and the speech occur in the same frequency range. When this situation occurs, the relationship between the signal-to-noise is important. Speech intelligibility, however, is not affected provided that speech is 100 times more intense than the noise [43]. Speech intelligibility for digits might be reduced by 50 percent if speech and noise are of equal intensities [16]. Speech can be perceived and understood even when it is of lower intensity than noise if the speech and the noise come from two different directions.

Frequency and Intelligibility

Frequency content is obviously important for speech intelligibility. Experiments in which filters are used, depicted important findings, such as the intelligibility of speech for one-syllable words is slightly reduced when frequencies above 1600 Hz are eliminated. But intelligibility is reduced by 25 percent when filtering is extended to 800 Hz. Whereas intelligibility is not significantly affected when frequencies below 1600 Hz are eliminated. In general, the intelligibility of speech is reduced further by about 25 percent when all frequencies below 3200 Hz are filtered out [6][8][9][50].

Segmental Analysis

The perception of vowels depends on the respective formant frequencies of the sound which is revealed by a segmental analysis of vowel sounds. Men, women, and children all use about the same vocal-tract configurations to produce the same vowel sounds, even though different individuals have varying sizes of vocal mechanism and vocal tracts. The relationship among formant frequencies is about the

same even though the formant frequencies may be higher or lower for different people.

Consonant sounds are produced by the constriction of the vocal tract and many are classified as either fricatives (i.e. "s" in /sIks/ or plosive or stops (i.e. "t" as in /tu/). Spectral or frequency differences make it possible to distinguish among fricative sounds. This can be seen from the phoneme in the digit /sIks/, i.e. "s" which has little energy below 4000 KHz, with a resonance peak between 4000 and 7000 Hz [20]. When the turbulent air stream is stopped by closure of the vocal tract and then released, stops or plosives are produced. The initial plosive "t" in the digit /tu/ is more forcefully exploded than the final plosive "t" in the digit /eIt/. Similarly, initial plosive is forcefully exploded in Arabic digits as in /θaλâθθh/ and /ωâhid/. Thus plosives occurring in the initial position of a digit, are exploded more strongly than those occurring at the final position of a digit. Spectral differences are very essential to distinguish plosives among one another and from consonants, such as the voiced plosive "b" in the Arabic digit /arbaρθh/, which has most energy between 500 and 1500 Hz. In contrast, "t" and "d" have a higher spectrum, with energy up to about 400 KHz. Classification of plosives, however, can be influenced significantly by its use in a spoken digit. Phoneme recognition is possible by segmenting speech and using the above analysis.

Methods of segmentation of isolated words into phonemic units have typically utilized information pertaining to rapid changes in the energy concentration in the frequency spectrum of the speech waveform [7][47][48]. The task of performing segmentation, i.e. a subdivision

of the uttered word into discrete consecutive phoneme sections, is a very important basic step towards machine phonemic recognition. In previous research, primary segmentation is used to group together similar acoustic adjacent minimal segments [52]. If the difference between corresponding parameters is less than a minimum, then two parameters should be considered as identical. Transitional segments where the acoustic characteristics vary with time considerably, is very difficult to locate, therefore secondary segmentation procedure is used to correct possible errors of the primary segmentation. This technique involves much computational work, hence it is not efficient.

Several segmentation techniques have been developed for continuous speech, apart from modifying the segmentation techniques used for isolated words, to perform segmentation of connected digits. This leads to the problem of words or digit boundary locations.

The problem of accurately locating the end points of an utterance is actually a specific case of the more general problem of labeling an interval of a signal as silence, unvoiced, or voiced. If one had a perfect technique for this three-level decision, the end point-location problem would be trivially solved. However, such an ideal algorithm does not exist as yet. Therefore, partial solutions to this more specific problem of isolating speech from a noisy background have been examined.

Coarticulation

The recognition of specific sounds, or phonemes, might be influenced in connected rather than in isolated digits. Speech, however,

has been considered to be composed of a sequence of distinctive, separate sounds; by analogy, "beads on a string". This is essentially the case in analysis of sounds in isolated digits, but actually is not the case in an analysis of connected digits or continuing speech. Coarticulation is a term given whenever there is an interaction of associated sounds. Due to the configuration of the vocal tract for any given sound being influenced by the shape required for the previous sound and respiratory movement for production of a following sound, and because the vocal tract is continuously in transition in the production of connected digits or continuous speech, coarticulation occurs [10].

Vowels

The position for a vowel sound can be assumed essentially to mean shaping the resonators in such a way as to produce the desired acoustic effect. It is clear that the vowel sounds are continuants. In other word, the speech mechanism assumes the position for the vowel and holds it with relatively little movement for a measurable fraction of a second while the sound is produced. Although the periods of time involved are small, continuants are characterized by these brief periods of holding. By contrast, the glide sounds are produced while the mechanism is in movement and their identifying characteristics are the result of this movement. In general, a pure vowel is defined as one in which the mechanism is held relatively stable in contrast to the glides in which the movement is the essence of the sound [39][50]. Noting that the term continuants is applied or referred to vowels, fricative consonants and nasals, while stops refers to plosive consonants and glides to inter vowel.

It is the function of the articulatory mechanism to break up and modify the laryngeal tone and to create new sounds within the mechanism itself. In fact, the speech mechanism as a whole not only articulates, i.e. join together, but also separates and molds the sounds delivered to it by the vibrator and resonator mechanisms. In addition it creates new sounds within itself by utilizing the energy supplied by the power mechanism in such a way as to produce within the oral cavity frictional noises that are independent of the laryngeal tone. Because of this, the articulatory mechanism assumes considerable importance to speech researchers.

Speech mechanisms can be divided into four units, distinct functionally but overlapping structurally. These units are: the power mechanism, the vibratory mechanism, the resonator mechanism, and the articulatory mechanism. These various functions and structures are coordinated through the activity of the voluntary nervous system. This coordination is made possible by four types of activity carried on by the nervous system: (1) motor activity that provides the stimuli that causes muscles to contract; (2) sensory reporting that gives information as to how the movements were produced; (3) auditory monitoring that makes possible the setting up of, and conformance to, speech standards; and (4) the associative function that ties up the auditory symbol with its meaning and with the motor pattern necessary to produce it.

CHAPTER III

CONCEPT AND COMPUTATIONAL TOOLS

Introduction

Due to the complex nature of speech process, it is suitable to have a parametric representation of the acoustic waveform which can be used to extract certain desired speech characteristics. Such parameters, used to describe the acoustic waveform over a specified time interval, might include Fourier coefficients, RMS energy, rate of zero crossing or the locations and values of predominant spectral peaks. The Fourier analysis is the most generally used technique for obtaining quantitative information about the speech waveform. Speech is, in general, non-stationary, and can be considered as stationary on a short-time basis. The short-time analysis is discussed in a later section. The segmented speech needs to be described by a set of well-defined parameters, so these can be used in a speech recognition scheme. These parameters must be simple, yet convey qualitative and quantitative information and characteristics of speech signals.

Due to the fact that linear prediction model and the acoustical tube model are equivalent, the reflection coefficients can be obtained from the area functions, and visa versa [17][41]. Figure 13 shows that the vocal tract is considered as a set of interconnected sections of equal length and varying cross-sectional areas. The linear prediction

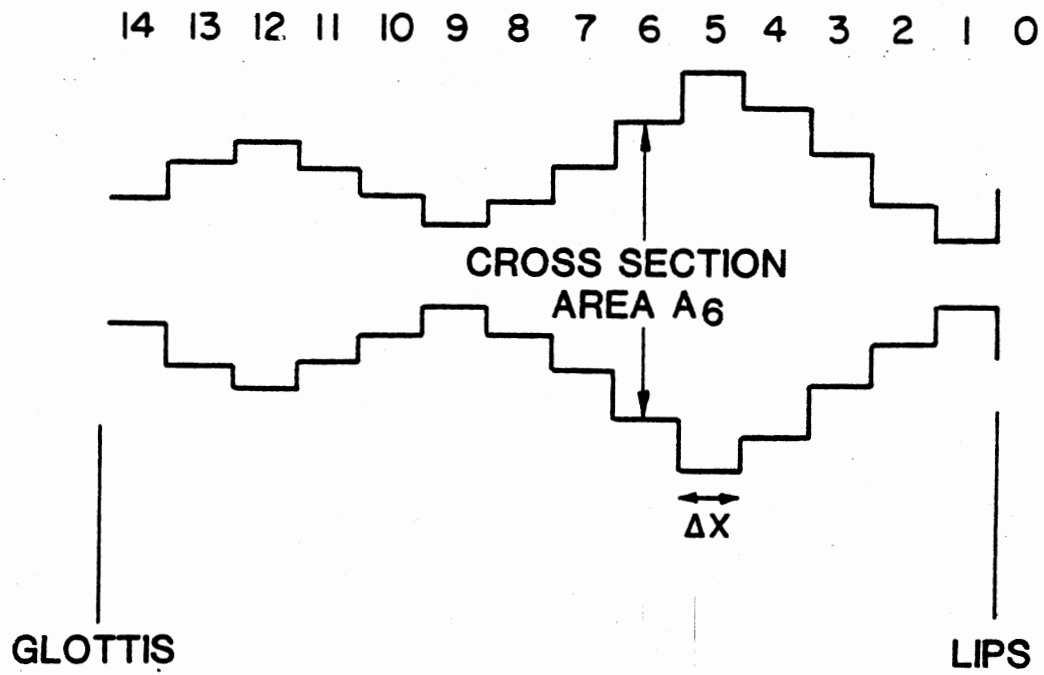


Figure 13. Concatenation of (N=14) Lossless Tubes of Equal Length

analysis model, to be discussed in the next section, is used because it is related to the acoustical tube model. The LPA not only filters out all voiced speech and formants, giving out the pitch period and random noise, but also allows for extraction of some useful parameters needed in the speech algorithm such as area functions. The estimation of the equivalent area function from the reflection coefficients provide for computing either part or the total equivalent area of the vocal tract cavity.

A mathematical discussion of LPA, autocorrelation method, windowing and pre-emphasis is presented. Parameters in terms of the vocal tract cavity ratios are defined, based on the facts obtained from x-ray pictures, that three cavities are observed. The ideal way to compute the RMS energy of the speech waveform is stated. Algorithms used for smoothing the RMS energy is discussed briefly with the aid of a block diagram. Parameters used for dip-classification algorithm are clearly defined. Finally, short-time analysis is briefly introduced. Also, window method is utilized for segmentation scheme and an overlap method is used in the end-point detection algorithm, for detecting digit boundaries accurately.

Linear Prediction Analysis

Linear prediction analysis applies to a class of problems in speech analysis and synthesis in which the present sample is predicted by a linear combination of past samples. The solution is obtained by solving a set of linear simultaneous equations based upon the least-square error criterion of optimality. From the solution, an all-pole digital filter can be derived, which has a discrete frequency response that closely

matches the smoothed spectral characteristics of the analyzed signal [29][44][63]. The significant advantage is that accurate spectral representation is obtained with only a few parameters over a 10 to 32 msec. signal.

An important form of linear prediction has been developed which is referred to as the autocorrelation method [31]. The autocorrelation method is very popular in speech processing, as it allows for a simple implementation when compared to the covariance method [1]. Furthermore, the parameters derived from an acoustic tube model can be related to a set of parameters in the autocorrelation method. Therefore, the autocorrelation method is used in this thesis. A brief review of the autocorrelation method is given below.

The Autocorrelation Method

Let $S(m)$ be the speech signal and $S_n(m)$ be the windowed speech signal.

That is,

$$S_n(m) = S(m+n)\omega(m) \quad (1)$$

where $\omega(m)$ is a window (for example, a Hamming window) of length N . It is clear that $S_n(m)$ is non-zero only for $0 \leq m \leq N-1$. Let the predicted signal $\hat{S}_n(m)$ is expressed by

$$\hat{S}_n(m) = \sum_{k=1}^p a_k S_n(m-k) \quad (2)$$

where a_k 's are some constants that are to be determined and p is the order of the prediction. The prediction error E_n is defined by

$$E_n = \sum_{m=0}^{N-p-1} e_n^2(m) \quad (3)$$

where

$$e_n(m) = S_n(m) - \hat{S}_n(m). \quad (4)$$

The coefficients a_k are obtained by minimizing E_n in (3). This minimization results in a set of normal equations

$$R_n(i) = \sum_{k=1}^p a_k R_n(|i-k|), \quad 1 \leq i \leq p \quad (5)$$

where $R_n(k)$ corresponds to the k th autocorrelation coefficient and is given by

$$R_n(k) = \sum_{m=0}^{N-1-k} S_n(m) S_n(m+k). \quad (6)$$

It is well known that if (5) is expressed in a matrix form, the coefficient matrix is a symmetric Toeplitz. There are several efficient algorithms (Levinson's, Durbin's and Trench's algorithms) [1] available to solve (5).

The Durbin's algorithm (27) is considered to be most efficient and the a_k 's in (5) can be computed using this method. This method is used in this thesis. For completeness, the algorithm is given below in terms of $R_n(k)$. For simplicity, the subscript of R is omitted.

$$E^{(0)} = R^{(0)} \quad (7)$$

$$k_i = \left[R^{(i)} - \sum_{j=1}^{i-1} a_j^{(i-1)} R^{(i-j)} \right] / E^{(i-1)}, \quad 1 \leq i \leq p \quad (8)$$

$$a_i^{(i)} = k_i \quad (9)$$

$$a_j^{(i)} = a_j^{(i-1)} - k_i a_{i-j}^{(i-1)}, \quad 1 \leq j \leq i-1 \quad (10)$$

$$E^{(i)} = (1 - k_i^2) E^{(i-1)}. \quad (11)$$

Equations (7)-(11) are solved recursively for $i = 1, \dots, p$ and the solution for (5) is given by

$$a_j = a_j^{(p)}, \quad 1 \leq j \leq p \quad (12)$$

Transfer Function Relation

The Z-transform of the error sequence $e_n(m)$, $E(Z)$, can be expressed in terms of the Z-transform of the windowed speech signal $S_n(m)$, $S(Z)$, by

$$E(Z) = [1-F(Z)] S(Z) \quad (13)$$

where

$$F(Z) = - \sum_{i=1}^P a_i Z^{-i} . \quad (14)$$

The block diagram representation of Equation (15) is shown in Figure 14.

Vocal Tract Division

In a later section some new parameters are defined based upon fourteen section representation of the vocal tract [8], as shown in Figure 15. In a recent paper, the vocal tract has been considered in terms of two major sections, called the front and back sections as shown by the dotted lines in Figure 15 [8]. X-ray analysis shows that the vocal tract is actually divided into three cavities [18], namely front, central and back, as shown by the solid lines in Figure 15.

Reflection Coefficients

Corresponding to the fourteen sections, a set of reflection

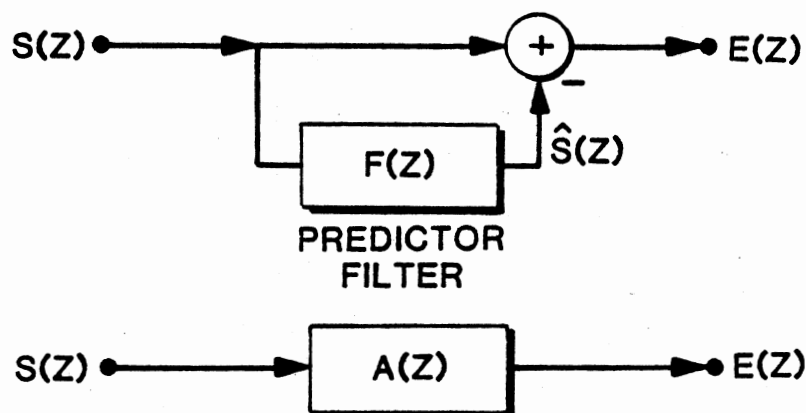


Figure 14. Equivalent Representation of Predictor Filter $F(Z)$ and Inverse Filter $A(Z)$

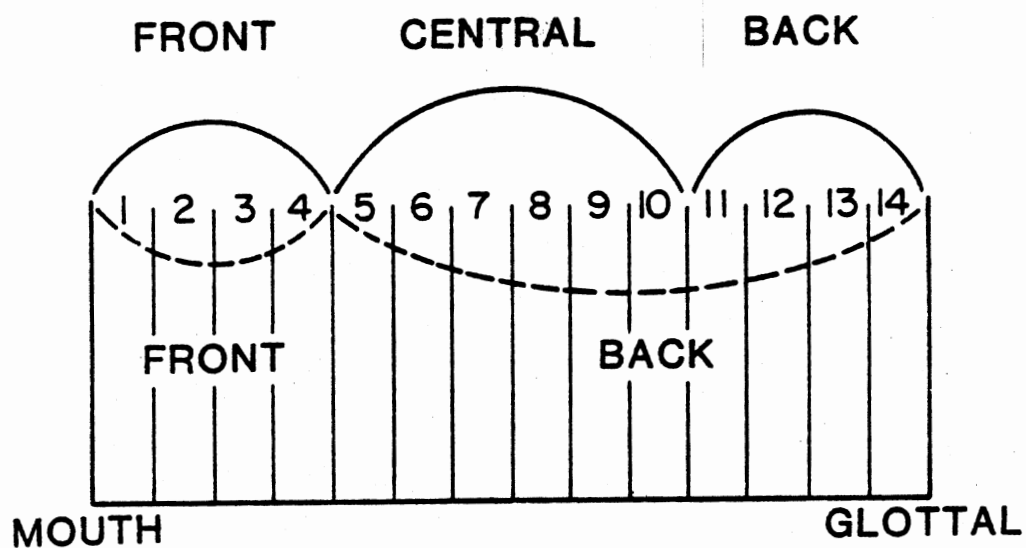


Figure 15. Equivalent Diagram Representing the Vocal Tract as an Acoustical Tube of Equal Lengths

coefficients can be defined. These are

$$k_i = \frac{\frac{A_i}{A_{i+1}} - 1}{\frac{A_i}{A_{i+1}} + 1}, \quad 1 \leq i \leq p = 14 \quad (15)$$

where A_i corresponds to the area of crosssection of the i th section.

It is clear that k_i is bounded between -1 and $+1$. Also, it follows from (15) that

$$\frac{A_i}{A_{i+1}} = \frac{1+k_i}{1-k_i}, \quad 1 \leq i \leq p - 1 \quad (16)$$

and

$$A_i = \frac{1+k_p}{1-k_p} A_{p+1} \quad (17)$$

Since k_i in (15) is of interest, and k_i is a function of a ratio of area function, it can be assumed that $A_{p+1} = 1$ without losing any generality.

The area's A_i 's are functions of the reflection coefficient as shown in (17). The reflection coefficients are related to the $a_i^{(i)}$ in the Durbin's algorithm. This is shown in Equation (9). The transfer function of a lossless tube model consisting of p sections has the same form as the transfer function derived from the linear prediction analysis. The reflection coefficients r_i 's obtained from the acoustic tube model are related to k_i 's by the equation

$$r_i = -k_i. \quad (18)$$

Using this relationship, Equation (18) can be expressed as

$$A_i = \frac{1-r_i}{1+r_i} A_{i+1} \quad (19)$$

Parameters in Terms of Vocal Tract

Cavity Ratios

Several parameters are defined below for future use in the digit recognition scheme. The front-to-total cavity ratio (FTR) is defined in terms of A_i , the area of cross section of the i th section, by

$$\text{FTR} = \frac{\sum_{i=1}^{p-10} A_i}{\sum_{i=1}^p A_i} \quad (20)$$

where p corresponds to the total number of sections (14 here) in the vocal tract. Note that $A_1, A_2, A_3,$ and A_4 corresponds to the areas of cross-section in the front of the vocal tract (see Figure 15).

Also, $\sum_{i=1}^p A_i$ corresponds to the total sum of the areas of cross-section. The central-to-total cavity ratio (CTR) is defined by

$$\text{CTR} = \frac{\sum_{i=5}^{p-4} A_i}{\sum_{i=1}^p A_i} \quad (21)$$

The back-to-total cavity ratio (BTR) is defined by

$$\text{BTR} = \frac{\sum_{i=p-3}^p A_i}{\sum_{i=1}^p A_i} \quad (22)$$

The front-to-back cavity ratio (FBR) is defined by [8].

$$\text{FBR} = \max (A_1, \dots, A_4) / \max (A_5, \dots, A_p) \quad (23)$$

Finally, the signed front-to-back cavity ratio (SFBR) is defined by

$$\text{SFBR} = \text{Sgn}(k_1) \cdot \text{FBR} \quad (24)$$

where $\text{sgn}(k_1)$ corresponds to the sign of the first reflection coefficient.

The digit recognition uses several other parameters. Before these can be defined, short-time analysis in terms of energy is discussed below.

Short-Time Energy

The amplitude of the speech signal varies appreciably with time, and there is a significant difference between the amplitudes of voiced segments and unvoiced segments. It is convenient to apply the short-time energy of the speech signal to extract the variations in amplitudes. The short-time energy (E_n) is defined [1]

$$E_n = \sum_{m=-\infty}^{\infty} [x(m)\omega(n-m)]^2 \quad (25)$$

where $\omega(n)$ is a window function and $x(n)$ corresponds to the speech signal. Equation (25) can be written as

$$E_n = \sum_{m=-\infty}^{\infty} x^2(m) \cdot h(n-m) \quad (26)$$

where

$$h(n) = \omega^2(n) \quad (27)$$

Equation (28) implies that the signal $x^2(n)$ is filtered by a linear filter with impulse response $h(n)$.

In this thesis the energy is computed using a rectangular window.

$$\begin{aligned} \omega(n) &= 1 & 0 \leq n \leq N - 1 \\ &= 0 & \text{otherwise} \end{aligned} \quad (28)$$

Therefore the short-time energy of a discrete-time signal is defined as

$$E_n = \sum_{m=n-N+1}^n x^2(m) \quad (29)$$

This implies that, the short-time energy at sample n is simply the sum of squares of the N samples $n-N+1$ through n .

It is appropriate to point out that other windows, such as Hamming window, have been used in short-time analysis. For the particular application, rectangular window is used for simplicity. For linear prediction analysis, Hamming window is used. For future use, the Hamming window function is given below.

$$\begin{aligned} h(n) &= 0.54 - .46 \cos(2\pi n/(N-1)), & 0 \leq n \leq N - 1 \\ &= 0, & \text{otherwise.} \end{aligned} \quad (30)$$

In the following the pre-emphasis aspects are discussed.

Pre-emphasis

For the purpose of distinguishing voiced and unvoiced speech segments, the speech signal is passed through a system that emphasizes the

frequency range that is desired. For example, low frequency emphasis can be used if the voiced segment is of concern. Similarly, the high frequency emphasis can be used if the unvoiced sound is of concern. The low and high-frequency pre-emphasis implementations are described by Figures 16 (a and b), respectively.

The low and high frequency pre-emphasized signals can respectively be given by

$$x'_l(n) = x(n) - \mu \cdot x'_l(n-1) \quad (31)$$

$$x'_h(n) = x(n) - \mu \cdot x(n-1) \quad (32)$$

where μ is usually referred to as the pre-emphasis factor.

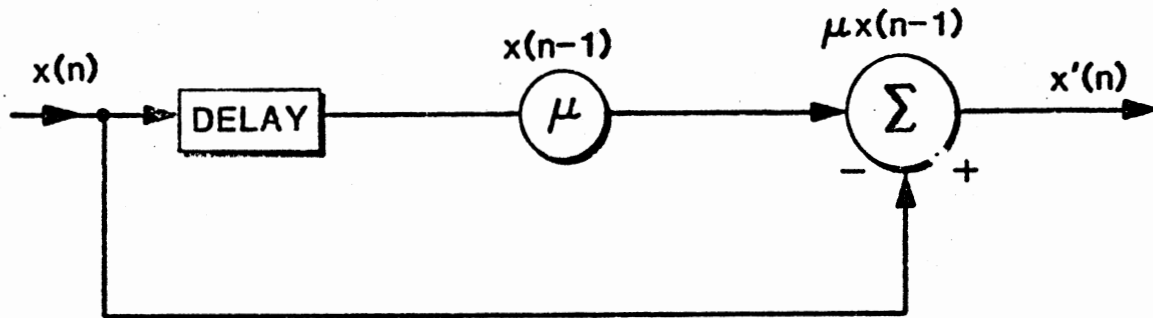
Pre-emphasis values of $\mu = -0.5$, -0.7 and -1.0 have been used, and it has been found that $\mu = -1$ gives decent results and, therefore, $\mu = -1.0$ is used in this thesis.

RMS Analysis in Speech Processing

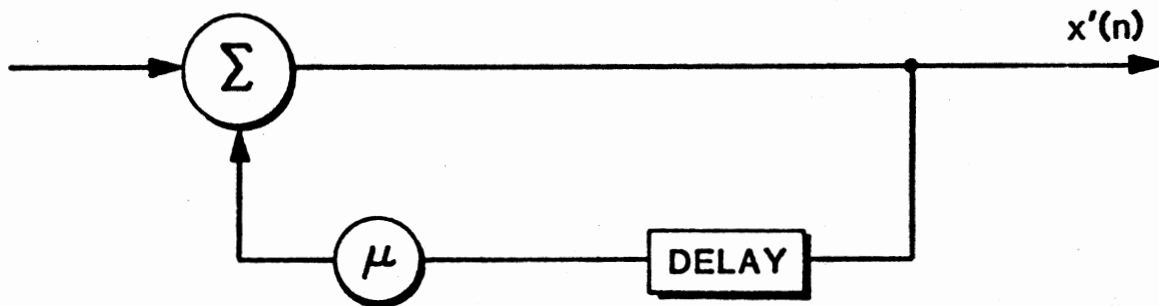
As mentioned before, the short-time energy reflects the amplitude variations of the speech signal. In fact, significant variations of energy is noticed during production of speech sounds especially vowels, semi-vowels, and voiced consonants. Useful phonetic features can be extracted from the speech signal by using the short time energy analysis. Equation (29) will be used here and is redefined here as

$$E_n = \sum_{i=1}^{NPS} x_n^2 \quad (33)$$

where NPS corresponds to the number of data points. In speech processing literature, $(RMS)_n$ is defined by



(a) HIGH-FREQUENCY BLOCK DIAGRAM



(b) LOW-FREQUENCY BLOCK DIAGRAM

Figure 16. High and Low Frequency Pre-Emphasis Block Diagram

$$(\text{RMS})_n = \sqrt{\sum_{i=1}^{\text{NPS}} x_n^2 / \text{NPS}} \quad (34)$$

Note that (34) is strictly not root mean squared value, but it gives the energy per frame, which is a measure of RMS.

In order to extract qualitative useful phonetic features from the RMS energy, smoothing is applied. Smoothing out the undesirable ripples in the RMS energy effectively discriminates voiced segments from unvoiced segments. In order to distinguish significantly between vowels and voiced consonants and detect voiced and unvoiced segments, the RMS speech energy of the utterances must be quantized. Quantization aspects are discussed in a later section.

Double Smoothing Algorithms

In speech processing applications, measurements and processing errors can occur in the data. The data, therefore, will exhibit single- or double- point sharp discontinuities of short duration. The data could also have sharp, isolated discontinuities of very short duration due to imperfect analysis procedures. In order to eliminate all the undesirable roughness and discontinuities in the uttered speech, a suitable and appropriate smoothing algorithm must be utilized. A practical smoothing algorithm has been proposed for speech processing [32], which is a combination of median smoothing and linear filtering. Median smoothing preserves signal discontinuities if the signal has no other discontinuity within $(N/2)$ samples.

The basic concept of a linear smoother is the separation of the signals based on their approximately non-overlapping frequency content. For non-linear smoothers it is more appropriate to consider separating signals based on whether they can be considered smooth or rough (noise-like).

Thus a signal $x(n)$ can be considered as $x(n) = S[x(n)] + R[x(n)]$ [1] where $S(x)$ is the smooth part of the signal x and $R(x)$ is the rough part of the signal x . A non-linearity which is capable of separating $S[x(n)]$ from $R[x(n)]$ is the running median of $x(n)$. The output of the running median smoother, $M_L[x(n)]$, is simply the median of the L numbers, $x(n), \dots, x(n - L + 1)$. Running medians of length L have desirable properties for a good smoother [32].

An ideal compromise is to use a smoothing algorithm based on a combination of running medians and linear smoothings. The running medians provide some smoothing, and the linear smoother can be of a low order system. A 3-point Hanning filter with an impulse response

$$\begin{aligned} h(n) &= 1/4 & n &= 0 \\ &= 1/2 & &= 1 \\ &= 1/4 & &= 2 \end{aligned} \tag{35}$$

is usually adequate, so that delays can be exactly compensated, due to the symmetry of the linear filter.

Linear smoothers are usually used in speech digital signal processing because they obey a superposition principle and they are time or shift invariant. Figure 17 shows two examples of data sequences which are to be smoothed. For case 1, a slowly varying waveform has been corrupted by a high frequency noise component. For this case a

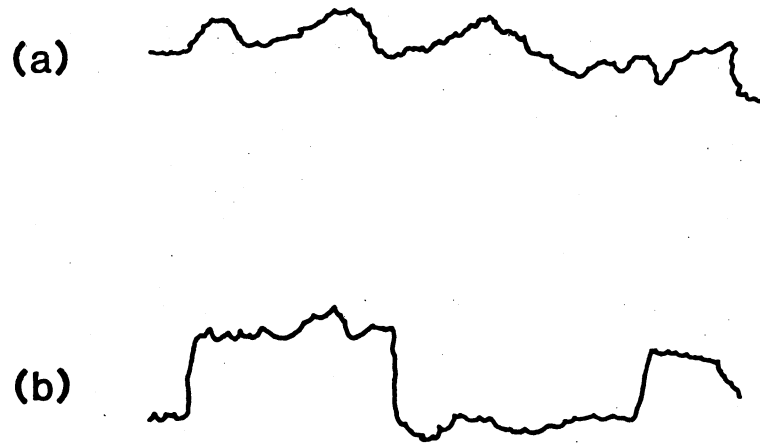


Figure 17. Two Examples of Noisy Signals
to be Smoothed

linear smoother (or low-pass filter) is entirely adequate for filtering out the noise. In case 2, there is noise-like component superimposed on the signal, and the signal displays noticeable sharp discontinuities. The discontinuities here contain much high-frequency energy, and are indistinguishable from the noisy component, as far as their spectral content is concerned. A linear smoother shown in Figure 18 would therefore smear out sharp changes in the data as well as filter out the noise. For cases like the data shown in Figure 17b, a nonlinear smoother is desired, which is capable of preserving sharp discontinuities and filtering out the superimposed noise. The algorithm based on 3-point running medians is illustrated in Figure 18. The input $x(n)$ exhibits sharp discontinuities at $n = 6$ and $n = 11$. The output $y(n)$ is defined as the 3-point-median of $x(n - 1)$, $x(n)$, and $x(n + 1)$, i.e. middle value when these three inputs are ordered in value. If a median greater than 9 is used, the discontinuity would be smoothed out and $y(n)$ would be flat.

An important property of median smoothers is their ability to follow low-order polynomial trends in the data as seen in Figure 19. It is seen in the figure that a 3-point-median follows low polynomial trends, whereas a 7-point-median has smoothed out the quartic polynomial considerably.

As mentioned before, median smoothing preserves sharp discontinuities in the data, but it fails to provide sufficient smoothing of the undesirable noise-like component. An ideal solution is a smoothing algorithm based on a combination of running median and a linear smoother as shown in Figure 20. The output $y(n)$ of the simple smoother

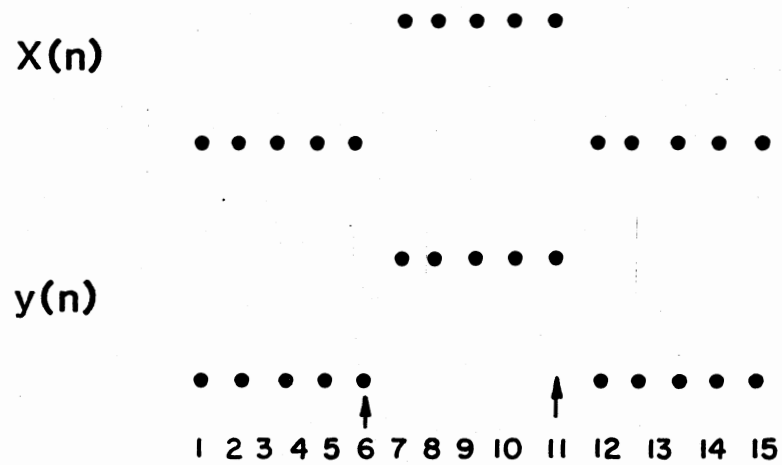
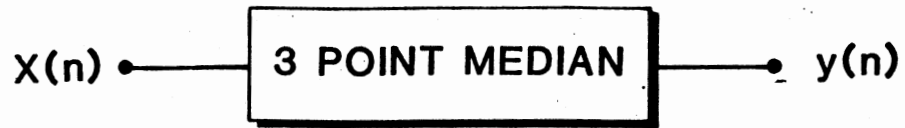


Figure 18. Median Smoothing of a Sequence with a Discontinuity

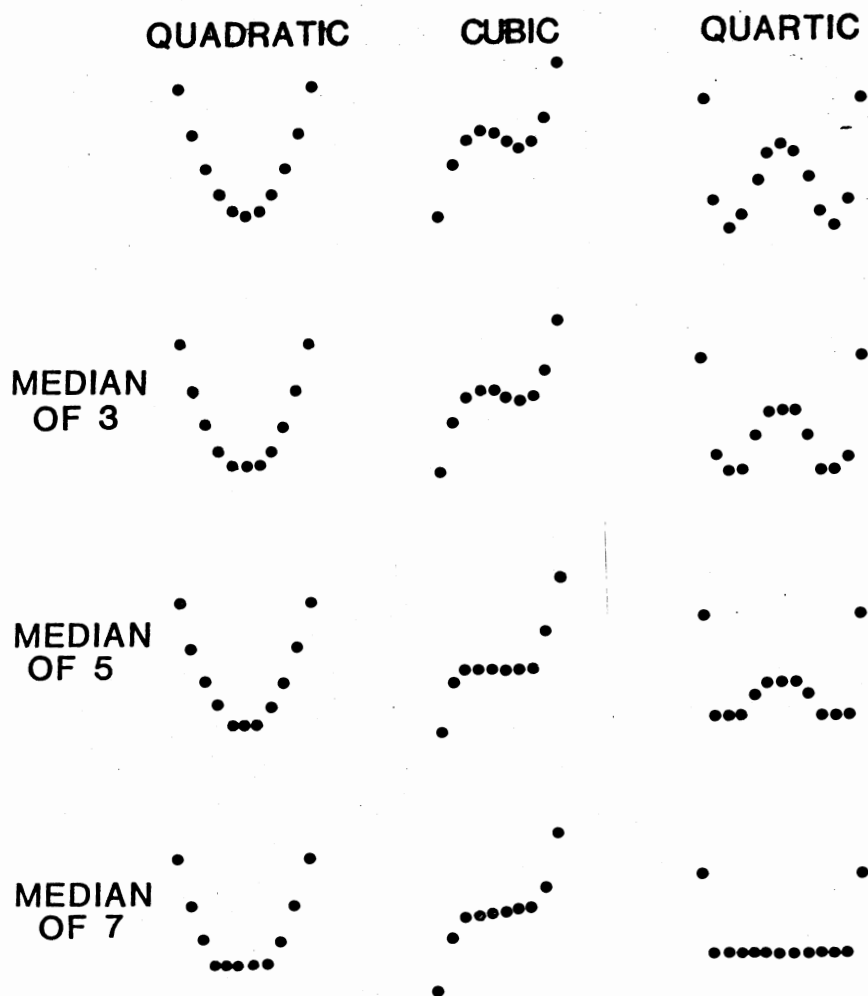


Figure 19. Effects of Various Median Smoothers on Low-Order Polynomials

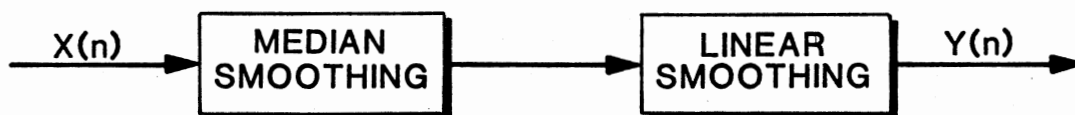


Figure 20. Simple Smoothing Algorithm

in Figure 20 is a smoothed version of $x(n)$, say $S[x(n)]$. If this simple smoother is not adequate to fulfill the requirements, a double smoothing algorithm can be used as shown in Figure 21.

$$\text{since } y(n) \approx S[x(n)] \quad (36)$$

$$\text{then } z(n) = x(n) - y(n) \approx R[x(n)] \quad (37)$$

Smoothing of $z(n)$ yields a correction term, which is added back to $y(n)$ to give $\omega(n)$, the second approximation to $S[x(n)]$. Hence $\omega(n)$ satisfies the relation

$$\omega(n) \approx S[x(n)] + S[R[x(n)]] \quad (38)$$

If $Z(n) = R[x(n)]$, i.e. the smoother is ideal, then $V(n)$, the output of the second smoother, would be identically zero, and the second-order correction would be unnecessary.

In order to implement the system shown in Figure 21, one must account for the delays in each path of the smoother and should be compensated. Median smoother has a delay of $(L - 1)/2$ samples, and each linear smoother has a delay proportional to the number of coefficients in the finite impulse response (FIR) filter.

For the proposed digit recognition algorithm to be discussed in a later chapter, the non-linear smoother shown in Figure 22 is used with a 3-point-median smoother in the front portion and a 5-point-median smoother in the later portion. Different sizes on median smoothers are found to give good results, as discussed below.

For example, 5-point-median has a delay of 2 samples, and a 3-point Hanning window has a delay of 1 sample. Thus the total delay of

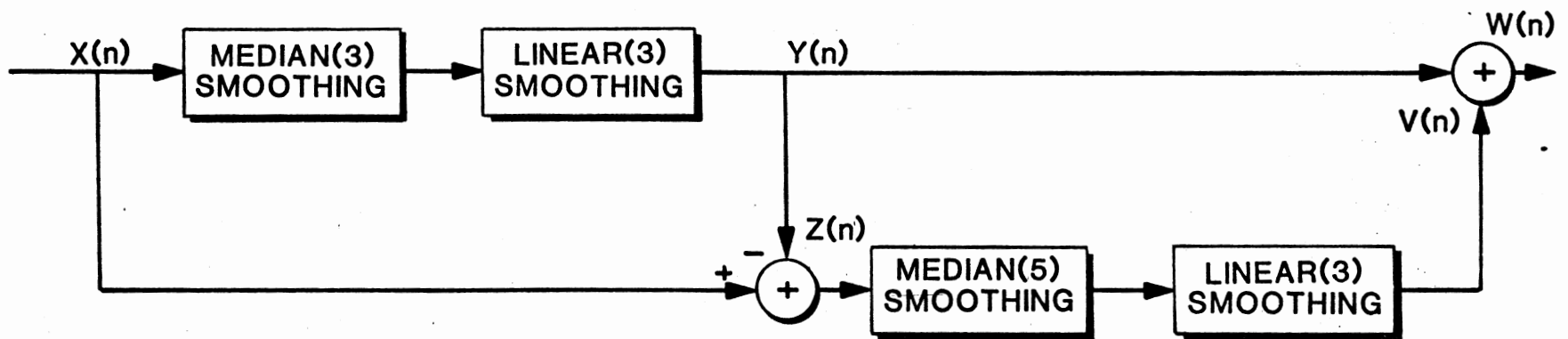


Figure 21. Double Smoothing Algorithm

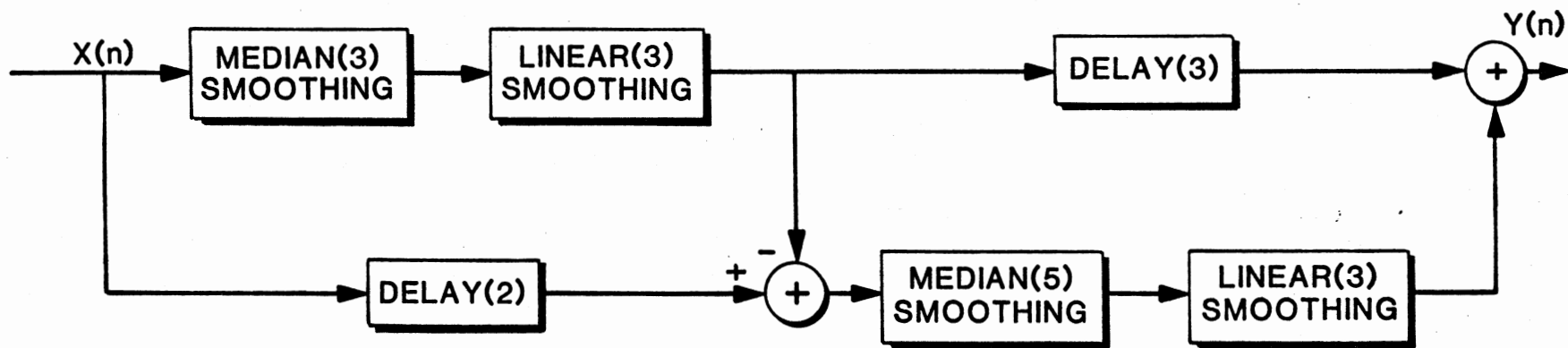


Figure 22. Nonlinear Double Smoothing Algorithm

the first smoother is 3 samples. The remaining important requirement is to implement the system of Figures 20 and 21 so that it provides an algorithm for handling the end points of the data. The effects of several versions of the smoothing algorithm on a speech intensity contour is in Figure 23. The effect of the additional smoothing obtained using higher order medians is clearly seen. Further, the differences between using median smoothing alone and the combination with linear smoothing are significant.

Zero Crossings

The zero-crossing rate is used in the proposed algorithm, and the end point detection algorithm. A brief review of ZCR is discussed below. The rate at which zero crossings occur is a simple measure of the frequency content of the signal. For example, a sine wave signal of frequency F_0 , sampled at a rate F_s , has F_s/F_0 samples per cycle of the sine wave. The long-time average rate of zero crossings is $Z = 2F_0/F_s$ samples, because each cycle has two crossings. But since speech signals are broadband signals then, rough estimates of spectral properties can be obtained using a representation based on the short time average magnitude difference function zero crossing rate. This is [1]

$$Z_n = \sum_{m=-\infty}^{\infty} |\text{sgn}[x(m)] - \text{sgn}[x(m-1)]| \omega(n-m) \quad (39)$$

$$\text{where} \quad \text{sgn}[x(n)] = 1 \quad x(n) \geq 0 \\ = -1 \quad x(n) < 0$$

$$\text{and} \quad \omega(n) = 1/2N \quad 0 \leq n \leq N-1 \\ = 0, \quad \text{otherwise}$$

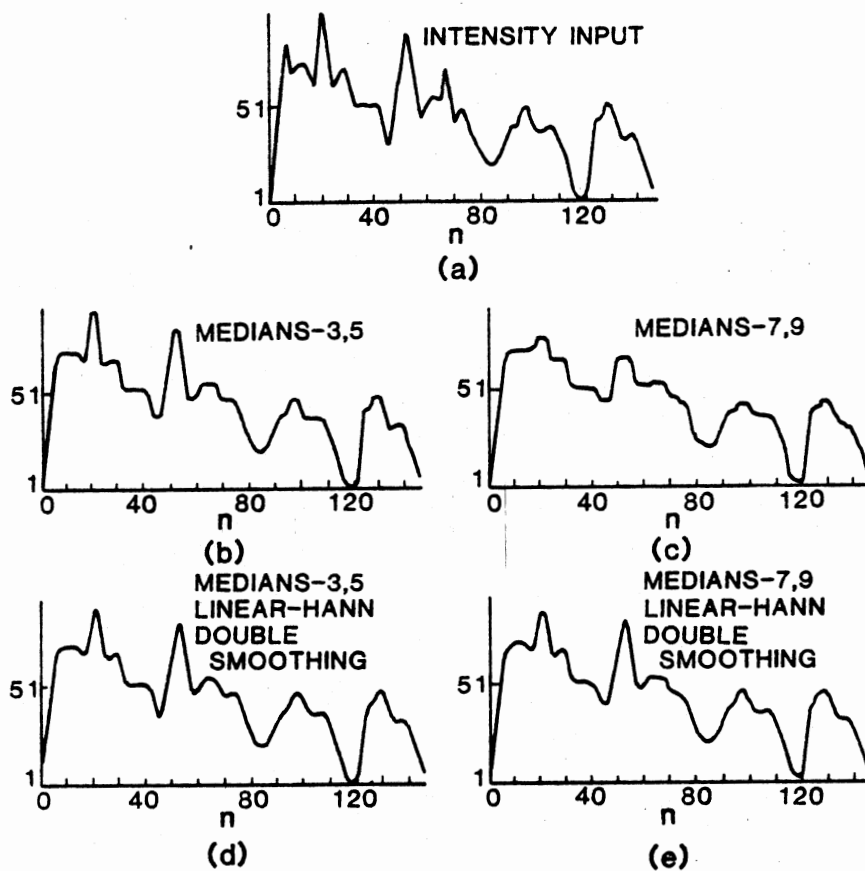


Figure 24. Effects of Several Versions of the Smoothing Algorithm on a Speech Intensity Contour

The computation of Z_n given in (39) appears to be complex. All that is required is to check samples in pairs to determine where the zero crossings occur, and then the average is computed over N consecutive samples. Since a rectangular window of finite length is used, the delay can be exactly compensated. Equation (39) can be simplified as [1].

$$Z_n = \frac{1}{2N} \sum_{m=n-N+1}^n |\text{sgn}[x(m)] - \text{sgn}[x(m-1)]| \quad (40)$$

which can be computed recursively by using

$$Z_n = \left\{ Z_{n-1} + \frac{1}{2N} |\text{sgn}[x(n)] - \text{sgn}[x(n-1)]| - |\text{sgn}[x(n-N)] - \text{sgn}[x(n-N-1)]| \right\} \quad (41)$$

Equation (41) is used in the proposed algorithm.

Cross-Correlation Function

The uniqueness of the parameters used for spoken digits is of importance in the proposed recognition scheme. The BTR, CTR, and RMS values per digit are compared respectively, using the statistical cross-correlation method discussed below [36].

When a set of independent variables are related to or are dependent upon each other, multicollinearity is said to exist among the variables. In the following, the correlation for digits i and j , are considered, where $0 \leq i, j \leq 9$. Furthermore, let t be the frame number and let there be n frames. For a frame t , x_{ti} and x_{tj} represent the RMS value for digit i and j respectively. The same analysis can be

used for the parameters BTR and CTR. Also, features of several sets of digits can be obtained and stored.

The correlation coefficients are denoted by the symbol $R_{x_{ti}, x_{tj}}$ and is given by

$$R_{x_{ti}, x_{tj}} = \frac{\sum_{t=1}^n (x_{ti} - \bar{x}_i)(x_{tj} - \bar{x}_j)}{\left[\sum_{t=1}^n (x_{ti} - \bar{x}_i)^2 \quad \sum_{t=1}^n (x_{tj} - \bar{x}_j)^2 \right]^{1/2}} \quad (42)$$

where $n = \max$ (number of frames for digit i, number of frames for digit j), and where the mean for digit i is:

$$\bar{x}_i = \frac{\sum_{t=1}^n x_{ti}}{n} \quad (43)$$

and the mean for digit j is:

$$\bar{x}_j = \frac{\sum_{t=1}^n x_{tj}}{n} \quad (44)$$

Note that the parameters will be padded by zero for the digit that has fewer number of frames.

If $n = \text{minimum}$ (number of frames for digit i, number of frames for digit j) in (43), then the padding will not be necessary. The results will not be as good for other cases of n . These aspects will be discussed in the next chapter.

It can be shown that $R_{x_{ti}, x_{tj}}$ is always between -1 and 1. A value of $R_{x_{ti}, x_{tj}}$ close to 1 indicates that the independent variables x_{ti} and x_{tj} are highly related or correlated. In other words a value of $R_{x_{ti}, x_{tj}}$ close to 1 denotes that x_{ti} and x_{tj} have similar patterns; that is, their first derivatives are almost the same in RMS, BTR, or CTR. Similarly, a value of $R_{x_{ti}, x_{tj}}$ close to -1 indicates that x_{ti} and x_{tj} have opposite patterns; that is, their first derivatives are reversely related. A value of $R_{x_{ti}, x_{tj}}$ close to 0 indicates that x_{ti} and x_{tj} are not correlated; that is, the independent variables x_{ti} and x_{tj} have no similarity in their RMS, BTR or CTR patterns.

Window Applications

Speech is a continuously time varying process as mentioned before. There must be a finite number of points that be used at any time, which requires segmentation. Since speech signals are stationary on a short-time basis, it is appropriate to consider Equation (25), which can be redefined as

$$E(n) = \sum_{m=0}^{N-1} [\omega(m)x(n-m)]^2 \quad (45)$$

where $\omega(m)$ is a weighting sequence or window which selects a segment of $x(n)$, and N is the number of samples in the window. For the simple case of $\omega(m) = 1$, $E(n)$ is the sum of the squares of the N most recent values of $x(n)$.

It is to be expected that the function $E(n)$ would display the time varying amplitude properties of the speech signal. Equation (45)

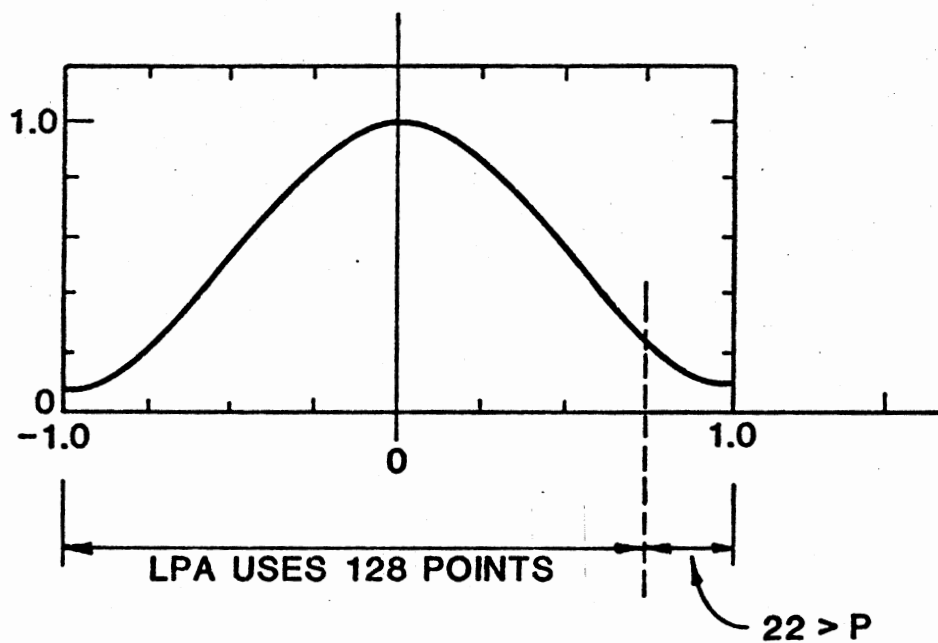
requires careful interpretation. First there is a choice of a window.

The effect of the window on the time dependent energy representation is explained by the properties of the rectangular window given by (28), and the Hamming window by (30). The rectangular window described by (29), corresponds to applying equal weight to samples in the interval $(n-N+1)$ to n . It has been found [1], that the bandwidth of Hamming window is about twice the bandwidth of a rectangular window for the same window length. The Hamming window has lower side lobes than the comparable rectangular window.

It has already been discussed previously that the LPA using autocorrelation method computes the short-time autocorrelation function $R_n(0), \dots, R_n(p)$ where $R_n(k)$ is given by Equation (6) and p is the order of the filter with the limit $0 \leq k \leq p$. If LPA is applied for m data points where $0 \leq m \leq N - 1 - k$, then N becomes greater than $m + 1 + k$ which indicates that for m points, LPA window be applied for $m + 1 + k$ points in order to avoid taper effect. For example, 128 points LPA with 14th order needs about 150 points windowing. Also the sketch of Figure 24 given below, justifies the application of the autocorrelation method every 128 points, with a 150 points Hamming window. In other words, the window length N should be greater than the number of data points plus the order of the filter, for example $N \geq m + 1 + p$ and $N \geq 128 + 1 + 14 \geq 143$.

Parameters in Dip-Classification

For segmentation and end point detection, the dip-classification of an RMS energy contour is utilized. A brief review along with some



$$R_N(K) = \sum_{m=0}^{N-1-K} S_N(m) S_N(m+K)$$

Where $R_N(K)$ Corresponds to the K th Autocorrelation Coefficients

Figure 24. Hamming Window

of the parameters used in the proposed algorithm are discussed below. The RMS energy is first smoothed and then normalized to maximum level of 100, to emphasize weakly pronounced voiced sounds. A plot of RMS energy versus time is obtained for given uttered digits as shown in Figure 25.

Let V_i^+ represent the first positive, V_j^- denote the 1st negative peak, or dip, and V_{i+1}^+ be the second positive peak, etc. It can be seen that there is always a minimum or dip between successive peaks.

Let V_j^- be the RMS dip value and

$$\text{Let } R_1 = V_j^- / V_i^+ \quad (46)$$

$$R_2 = V_j^- / V_{i+1}^+ \quad (47)$$

Where

$$R_{\min} = \text{Min} \{R_1, R_2\} \quad (48)$$

$$\text{Also Let } X_1 = \log_{10} V_j^- \quad (49)$$

$$\text{and } X_2 = \log_{10} R_{\min} \quad (50)$$

Then the functions used for the segmentation algorithm were obtained as follows [8].

$$Z_1 = a_1 \log_{10} V_j^- + a_2 \log_{10} R_{\min} \quad (51)$$

$$Z_2 = a_3 \log_{10} V_j^- + a_4 \log_{10} R_{\min} \quad (52)$$

The coefficients a_1 , a_2 , a_3 , and a_4 are computed from the design samples on the basis of Fisher's method [38]. Z_1 and Z_2 are computed

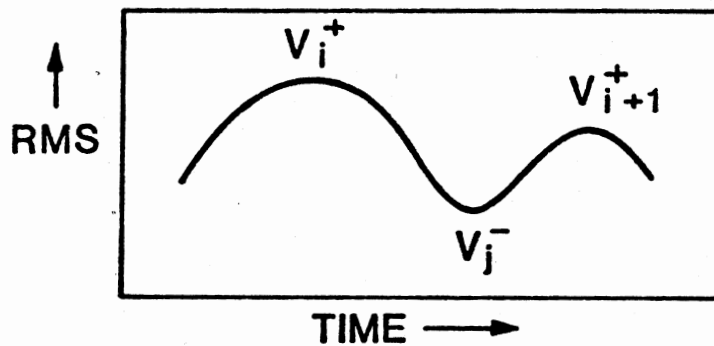


Figure 25. RMS Plot Used for the Definitions of the Extrema V_i^+ , V_j^- and Their Ratios, R_1 , R_2 and R_{\min}

from a set of measurements of X_1 and X_2 . Then,

$$Z_1 = 0.63 X_1 + 0.77 X_2 \quad (53)$$

$$Z_2 = 0.87 X_1 + 0.50 X_2 \quad (54)$$

A plot of the two parameters $\log_{10} V_j^-$ versus $\log_{10} R_{\min}$ can be obtained in order to determine the values of two linearly discriminant functions, which are used to minimize the misclassifications of vowels and voiced consonants, as shown below.

$$Z_1 = 0.63 X_1 + 0.77 X_2 - 1.03 \quad (\text{LDF1}) \quad (55)$$

$$Z_2 = 0.87 X_1 + 0.50 X_2 - 0.83 \quad (\text{LDF2}) \quad (56)$$

Two other important functions are SL1 and SL2, referred to as the slicing functions, which are used to determine the level of the dip, by intersecting the slope at that point, as shown below.

$$\text{SL1} = V_j^- + C_1 (V_i^+ - V_j^-) \quad (57)$$

$$\text{SL2} = V_j^- + C_2 (V_{i+1}^+ - V_j^-) \quad (58)$$

The values of the constants C_1 and C_2 are the means of the statistical distributions of $(V_i^+ - V_j^-)$ and $(V_{i+1}^+ - V_j^-)$ as shown below.

$$\text{SL1} = \text{Sum of } (V_i^+ - V_j^-) / N_{\max} \quad (59)$$

$$\text{SL2} = \text{Sum of } (V_{i+1}^+ - V_j^-) / (N_{\max} - 1) \quad (60)$$

$$\text{Then } SL1 = v_j^- + 0.3 (v_i^+ - v_j^-) \quad (61)$$

$$SL2 = v_j^- + 0.17 (v_{i+1}^+ - v_j^-). \quad (62)$$

The two constants (-1.03) and (-0.83) in (55) and (56) are the values of the two linear discriminant functions LDF1 and LDF2 respectively as shown in Figure 26 [8]. The LDF1 is defined by taking the vowels as one class and the sonorants (nasals, liquids, and semi-vowels) as the other class. The LDF2 is defined by taking the vowels as one class and the obstruents (all the consonants except the sonorants) as the other class.

The constants (or threshold) of the LDF1 of -1.03 was determined for the design samples so as to minimize the misclassification of vowels as sonorants and to eliminate the misclassification of sonorants as vowels [8]. Likewise, the constant of the LDF2 of -0.86 was determined so as to minimize the misclassification of obstruents as vowels and to eliminate the misclassification of vowels as obstruents.

The type of dips D_1 , D_2 , D_3 shown in Figure 27, are classified by the sign of Z_1 and Z_2 given by Equations (55) and (56) and the slicing functions $SL1$ and $SL2$, given by (61) and (62). The identification for vowel and non-vowel procedure will be discussed in the next paragraph.

Dip-Classification

For segmentation and digit boundary detection, the dip-classification of the smoothed RMS contour is utilized. As a first step, the RMS dips and peaks are extracted from the smoothed RMS function and stored. If the logarithm of a peak value ($\log_{10} v_i^+$) is smaller than

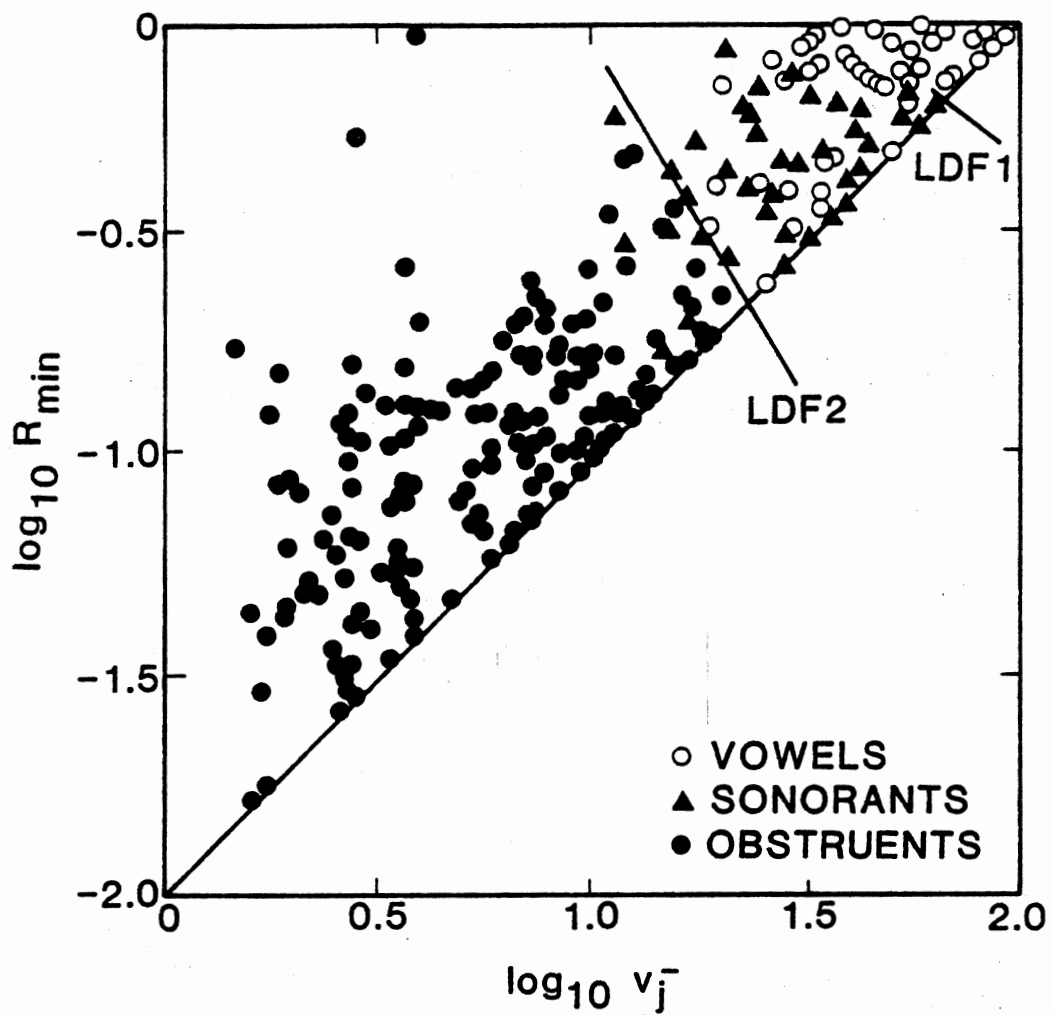


Figure 26. Plot of $\log_{10} R_{\min}$ vs $\log_{10} v_j$ for Estimating the Values of Linear Discriminant Functions LDF1 and LDF2

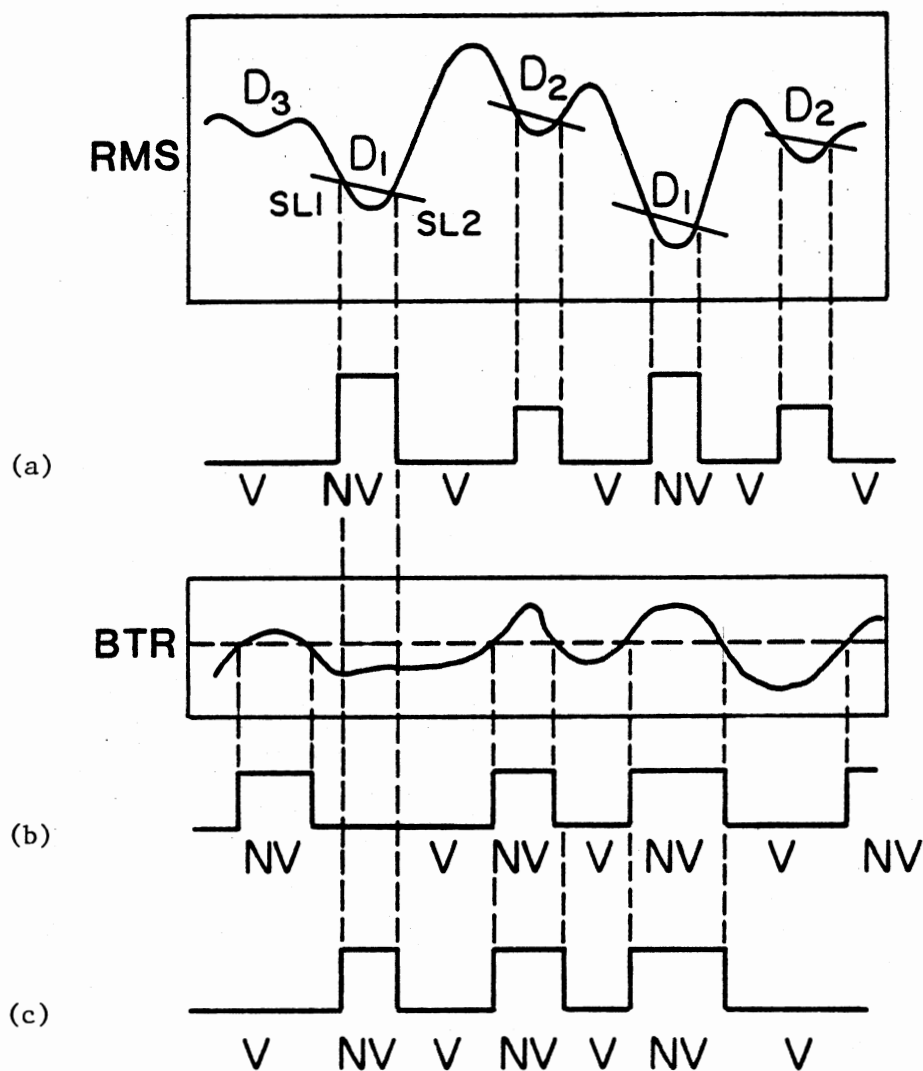


Figure 27. Schematic Diagram Showing the Types of Dips, D_1 , D_2 and D_3 , RMS, and BTR, Vowel, Non-Vowel Primary Decision and Combines Decision

1.0 or if the SFBR value is larger than a threshold of 4.0 around the peak, that peak and the following dip are not considered, since logarithmic RMS peak values less than 1.0 rarely occur for vowels, and the two functions Z_1 and Z_2 are not computed. Accordingly if the important constraint is met, i.e. if $\log_{10} V_i^+$ is greater than 1.0 and the SFBR value is smaller than a threshold of 4.0, then the functions Z_1 and Z_2 are computed. Hence whenever a dip is detected as a significant dip, that dip is classified into one of the three types of dips, dip 1 (D_1), dip 2 (D_2), or dip 3 (D_3) according to the sign of Z_1 and Z_2 . Therefore the dip is an element of dip 1 (D_1) if both Z_1 and $Z_2 \leq 0$. The dip is an element of dip 2 (D_2) if $Z_1 \leq 0$ and $Z_2 > 0$. Finally the dip is an element of dip 3 (D_3) if both Z_1 and Z_2 are greater than zero. In summary,

$$\text{dip} \in D_1 \quad \text{if} \quad Z_1 \leq 0 \text{ and } Z_2 \leq 0$$

$$\text{dip} \in D_2 \quad \text{if} \quad Z_1 \leq 0 \text{ and } Z_2 > 0$$

$$\text{dip} \in D_3 \quad \text{if} \quad Z_1 > 0 \text{ and } Z_2 > 0$$

Figure 27 illustrates the three types of dips, D_1 , D_2 and D_3 . Dip 1 indicates non-vowel-like dip, dip 2 is an ambiguous dip, and dip 3 indicates a vowel-like dip.

A vowel, vowel-like and non-vowel decision is made from the RMS energy dips as shown in Figure 27(a). Another decision of vowel and non-vowel only is based on the BTR plot threshold, as shown in Figure 27(b). A final "OR" decision is obtained based on the decision obtained from the RMS and BTR contour, as shown in Figure 27(c).

Scaling and Normalization

The parameters of digit signal segments can be more robust and less variable by normalization or scaling techniques. Applying normalization techniques requires certain rules to represent a digit signal sound faithfully. Improved signal representation through appropriate scaling may speed up the overall recognition process. These techniques include, energy, amplitude and time normalization.

Amplitude Normalization

A practical form of amplitude normalization is achieved by linearly quantizing the RMS energy into 100 levels, assigning 100 to the maximum value within each utterance.

Maximum of Frames Normalization

The aim of this method is to align the time of occurrence of the unknown digit utterance events to see how they match. In order to find best similarities between the unknown pattern and the referenced pattern, the maximum number of frames of a given set of digits is used. Digits having lower number of frames are padded with zero.

Scaling (Energy Normalization)

Scaling is a form of energy normalization in which each point of the incoming data is multiplied by a factor to fix the mean. This factor is derived from the inverse of the mean of the signal as shown below.

Let M_α denote the mean value of the signal, then the new value of the signal is given by

$$\text{New value} = (\text{signal level}) \left(\frac{\text{constant}}{M_\alpha} \right)$$

$$M_\alpha = \frac{\text{sum of square values}}{\text{no. of points}}$$

This technique tends to increase the peak value considerably whenever wrong end point detection results in a greater number of points than the number of points within the correct digit boundaries. Therefore, amplitude normalization is adopted.

CHAPTER IV

ACOUSTIC PHONEMIC DIGIT RECOGNITION SCHEME

FOR DIGITS SPOKEN IN AMERICAN

ENGLISH

Introduction

Phonemic digit recognition is one of the steps towards simplifying communications between man and machine. It is the process whereby an operator can use spoken digit commands that can be recognized by a phonemic digit recognition system. Generally man's communication with machines has been structured according to the operational requirements of the machine. Learning the "language" of the machine and manipulation of special dials or keys in the proper sequence and format is required to communicate with machines. Any deviation from this unnatural machine language can produce errors which are not easily detectable because of the complexities of the rules for proper communication between man and machine.

The main objective is to develop a phonemic digit recognition system which makes it simple for humans to "talk" directly to a machine, without any intermediate keying or handwritten steps. The operator would provide instructions in his natural language, using digit commands to control mechanical systems such as entering and

leaving restricted areas, postal zip codes, banking, inventory, etc. If connected digit is to be recognized, a suitable method of segmentation into recognizable units is required. Since phonemic digit recognition is the goal of this study, a highly accurate method of segmentation, using acoustical and phonetic feature parameters is required. The machine should be able to correlate and recognize the discrete acoustical waveform by segmenting the waveform into a sequence of elements so that the spoken phonemes can be localized and the end points are detected and located. In the following, a brief discussion on the segmentation is given.

Segmentation

Methods for segmentation of isolated and connected digits into phonemic units have typically utilized information pertaining to rapid changes in the energy contour of the given utterance. Rapid changes in speech parameters, such as energy and pole frequency derived from two-pole LPA model is used for detecting voiced segments, because these parameters have high value for vowels and develop a dip for vowel-like consonants. Since the ZCR rate and the normalized error show high values for unvoiced phonemes, and low values for voiced phonemes, they are utilized together to detect unvoiced segments. The ZCR rate and the energy signal are used for end point detection of the spoken digit. In fact the problem of accurately locating the beginning and end of an utterance is actually a special case of the more general problem of labeling an interval of a signal as silence, unvoiced, or voiced. If there is a perfect technique for this three level decision, the end

point-location problem would be solved. However, such an ideal algorithm does not exist yet. Therefore, it is very important to develop an appropriate algorithm for end-point detection of connected digits spoken in any environment. Consequently, it is considered worthwhile to consider dip-classification algorithm for locating digit boundaries. The dip-classification scheme was discussed in the previous chapter.

Recognition Scheme

Most digit recognition systems, known up till now, lack the usage of the natural phonetic features of the spoken digit. As discussed in Chapter I, some recognition procedures are based on finding the spectrum energy for both vowels and voiced consonants. The application of more than one band-pass filter in the above mentioned system causes permanent loss of useful phonetic and acoustic features and some useful information relating to the transition regions and formant peak locations. The above system is not efficient as the computational requirements are severe. Rabiner and Sambur [4][12] improved the efficiency of this system by using silence, voiced and unvoiced segment detection based on energy and pole frequency, ZCR and normalized error measurements. But this scheme lacks the accuracy of detecting digit boundaries and has the disadvantage of training the system. In addition the system must know the number of spoken digits in order to estimate the number of boundaries to be located. Also, the speaker must be trained, and the best uttered digit is used for the recognition scheme. Furthermore, a 2-pole LPA model can't extract the normalized error that eliminates the higher order formants, and a model is needed to account for f_2 , f_3 , f_4 , etc. Another disadvantage is the

artificial introduction of silence at the beginning and end of the connected digits, which limits the usefulness of the system in cases where the machine cannot locate digit boundaries, due to the miscalculation of the ZCR threshold. Hence misclassification and wrong recognition may result.

In order to solve some of the previous mentioned digit recognition system problems, and minimize the computational time, a suitable method of segmenting the spoken digits into recognizable phonetic units based on area functions and RMS energy contour is utilized. Since the RMS energy depends on the amplitude of the signal, an appropriate scaling or normalization is needed prior to parameter measurements. Scaling and normalization is discussed in Chapter II.

Digit Recognition Flow Chart

The digit recognition scheme for digits in English is shown in Figure 28. The analog data is first low-pass filtered with a cut-off frequency of 4 KHz and sampled at 8000 Hz/sec. The first step in the digit recognition scheme is the end-point detection [9][12]. Following the flow chart, one can see that two different ideas are used. The first one is based on the RMS energy peak ratio and threshold to detect the uttered digit. The second one depends on the BTR and SFBR parameters derived from area function of the vocal tract, via LPA. Before implementing the first idea, the incoming digit signal is segmented into 128 points per frame using a rectangular window, because it is feasible to assume that the vocal tract shape will remain unchanged within this short frame or segment.

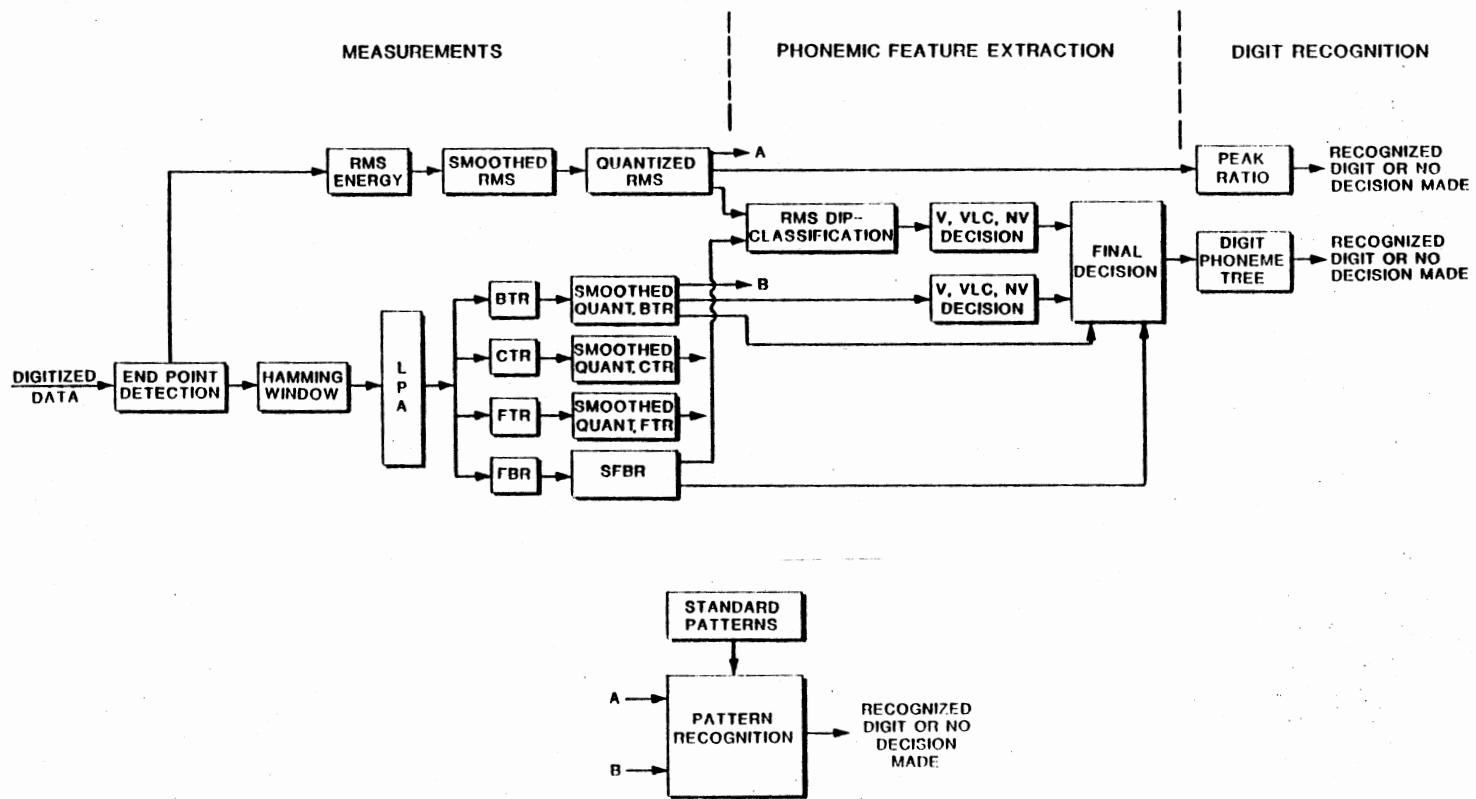


Figure 28. Digit Recognition Flow Chart Based on a Phonemic Feature Detection

From Figure 28, it can be seen that there are three main parts in the overall scheme of digit recognition. These are measurements, phonetic feature detection and, finally, the digit recognition itself. In the first part corresponding to the first idea, the RMS energy, smoothed RMS energy and quantized RMS energy are computed per frame. With a frame length of 128 points, rectangular window is used here for frame segmentation. Two sets of representative smoothed and quantized RMS plots for digits zero through nine in English are given in Figures 29-48 respectively. From these plots, the two largest peaks are obtained and are given in Tables V and VI for ten digits. Also, peak-to-peak ratios of two largest peaks are computed. It is assumed that this ratio is always less than one. From Table V, it can be seen that for digits four, six and eight, there is only one peak and therefore, the ratios are not given. However, from Table VI, there are two peaks for the digits four, indicating an inconsistency.

The range for the threshold values for the largest peak (P_1, P_2 in Tables V and VI) and the ratio for the two largest peaks have been established from previous measurements. These ranges are given in Tables V and VI. First, the measured ratio is compared with the established range. If it does not match with any given range, it will compare with the largest peak range. If it cannot match this range either, then the first idea based on RMS did not work. From the measurements, it has been found that RMS method gives the actual digit about half the time. The other half of the time, it indicates that no decision can be made using RMS plots. This completes the first idea in the digit recognition scheme.

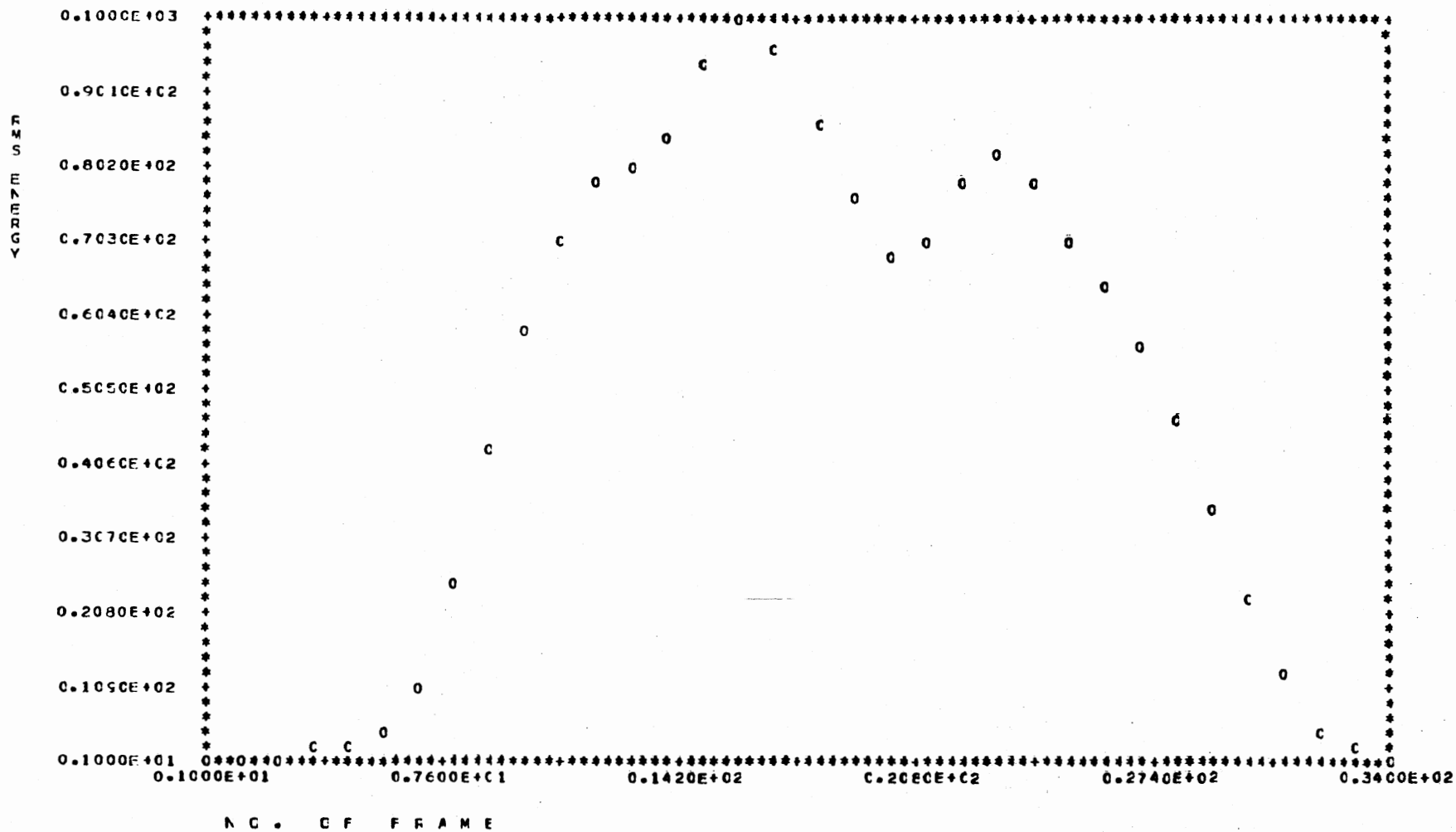


Figure 29. Smoothed RMS Energy Contour for Digit Zero, i.e. /zIro/ Spoken in American English

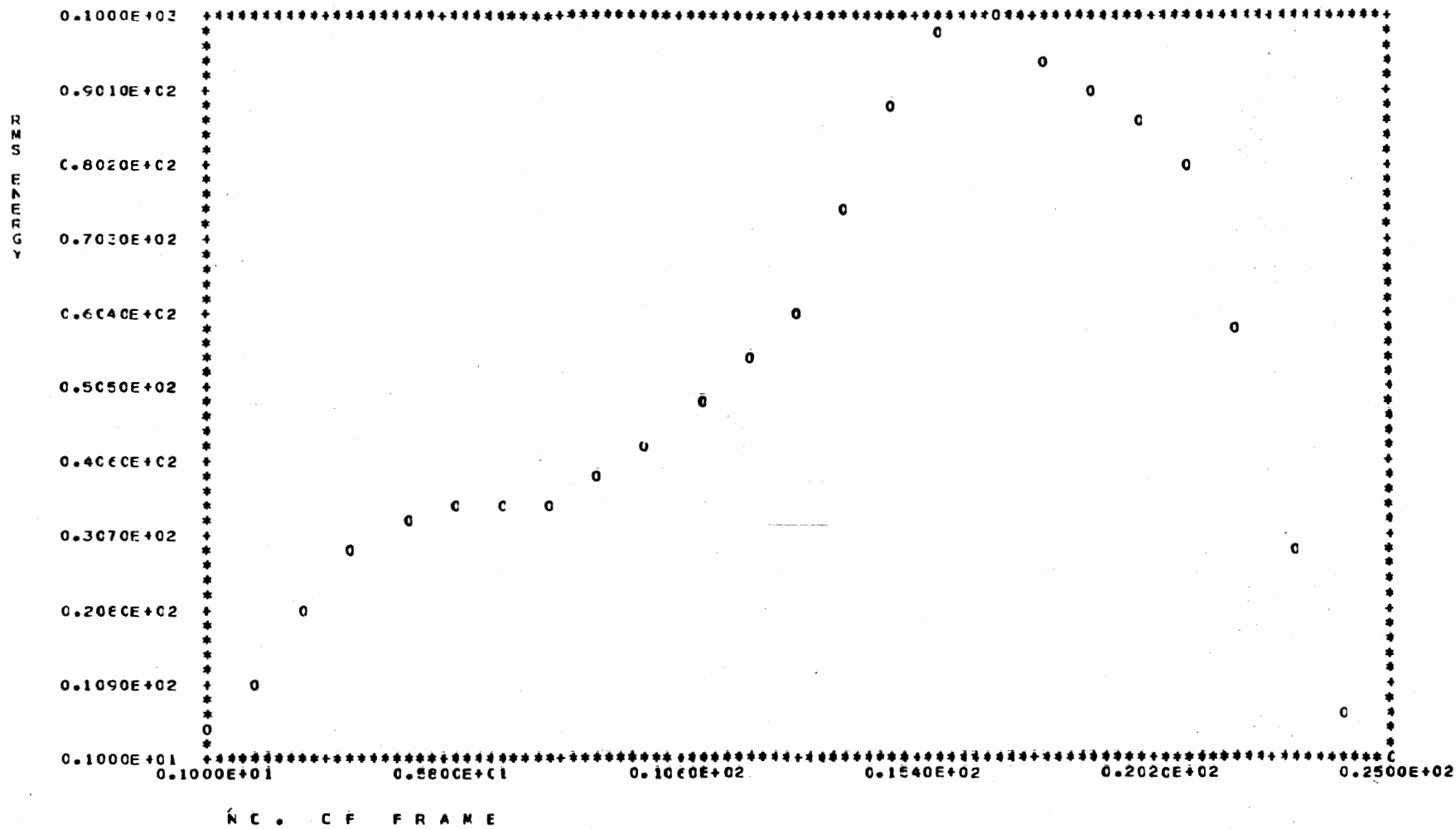


Figure 30. Smoothed RMS Energy Contour for Digit One, i.e. /wAn/ Spoken in American English

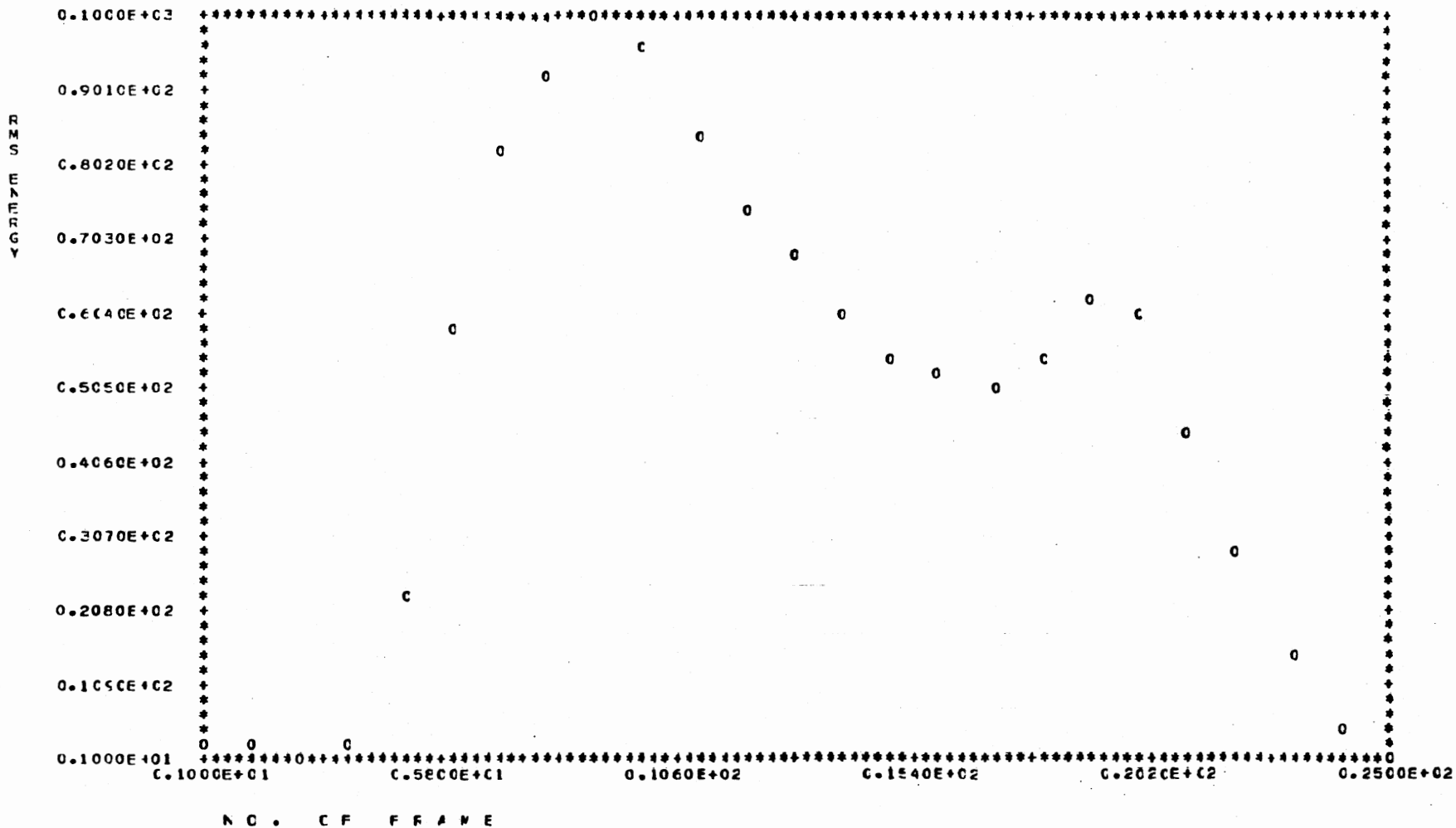


Figure 31. Smoothed RMS Energy Contour for Digit Two, i.e. /tu/ Spoken in American English

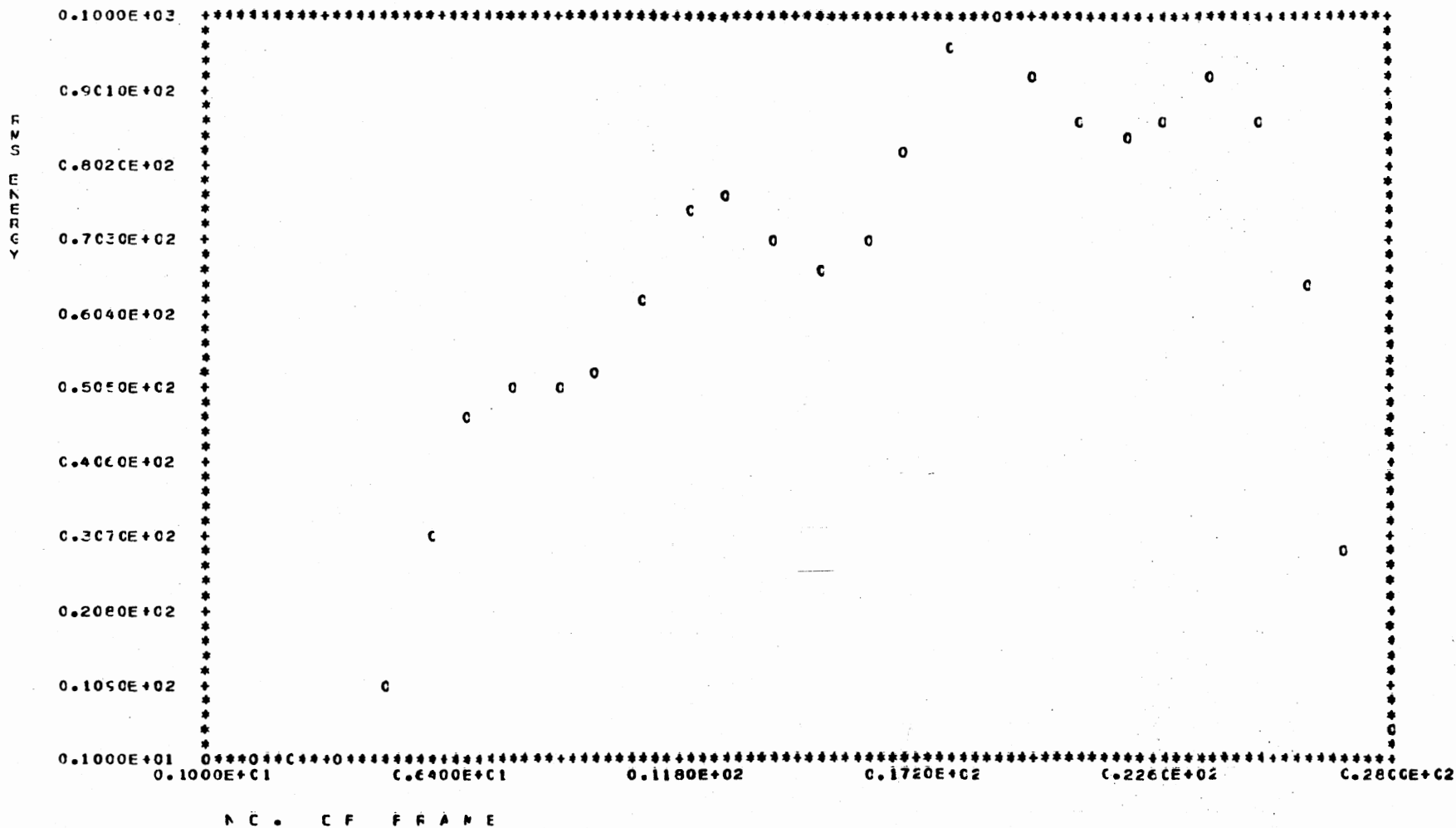


Figure 32. Smoothed RMS Energy Contour for Digit Three, i.e. /θri/ Spoken in American English

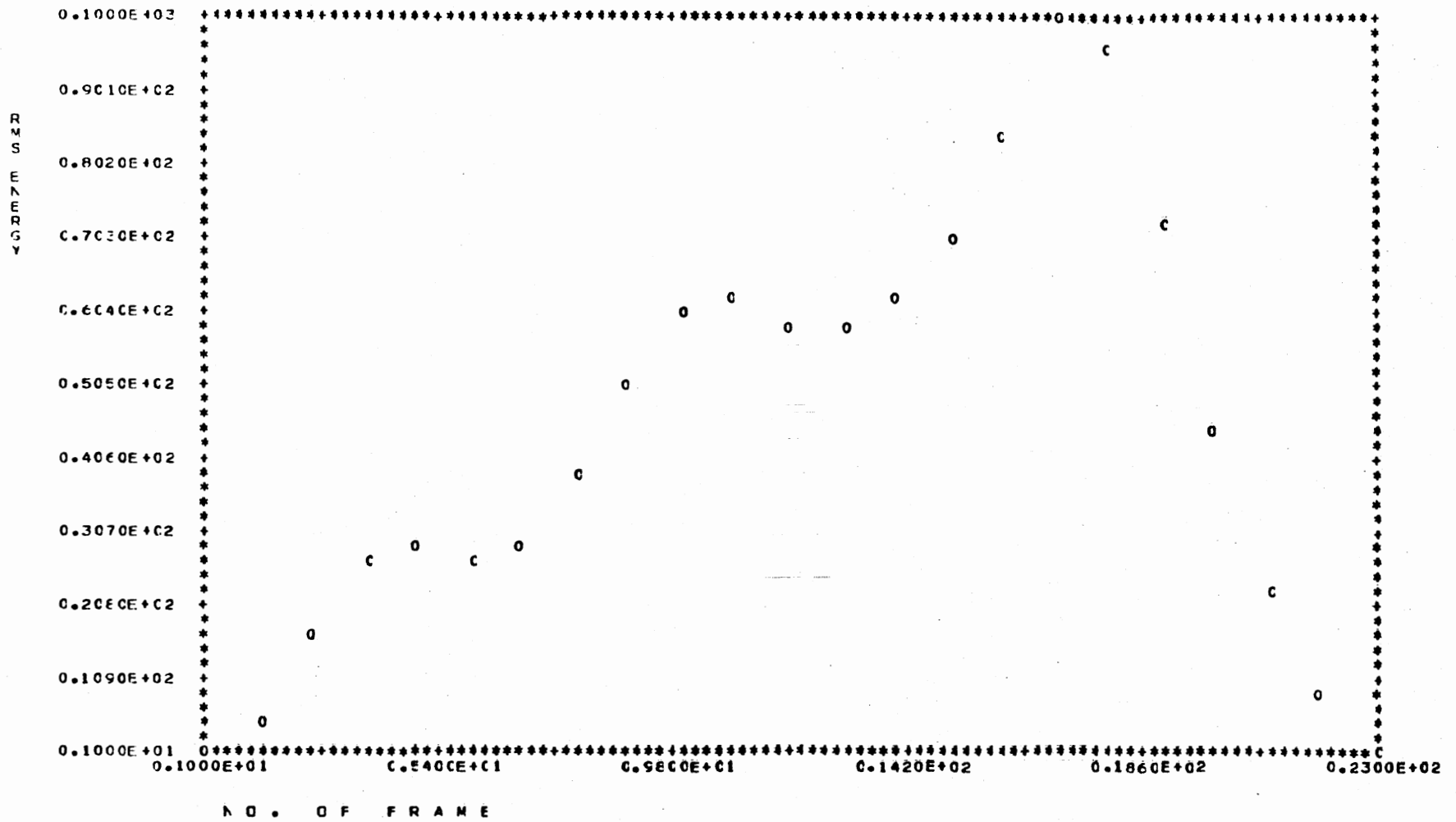


Figure 33. Smoothed RMS Energy Contour for Digit Four, i.e. /fɔr/ Spoken in American English

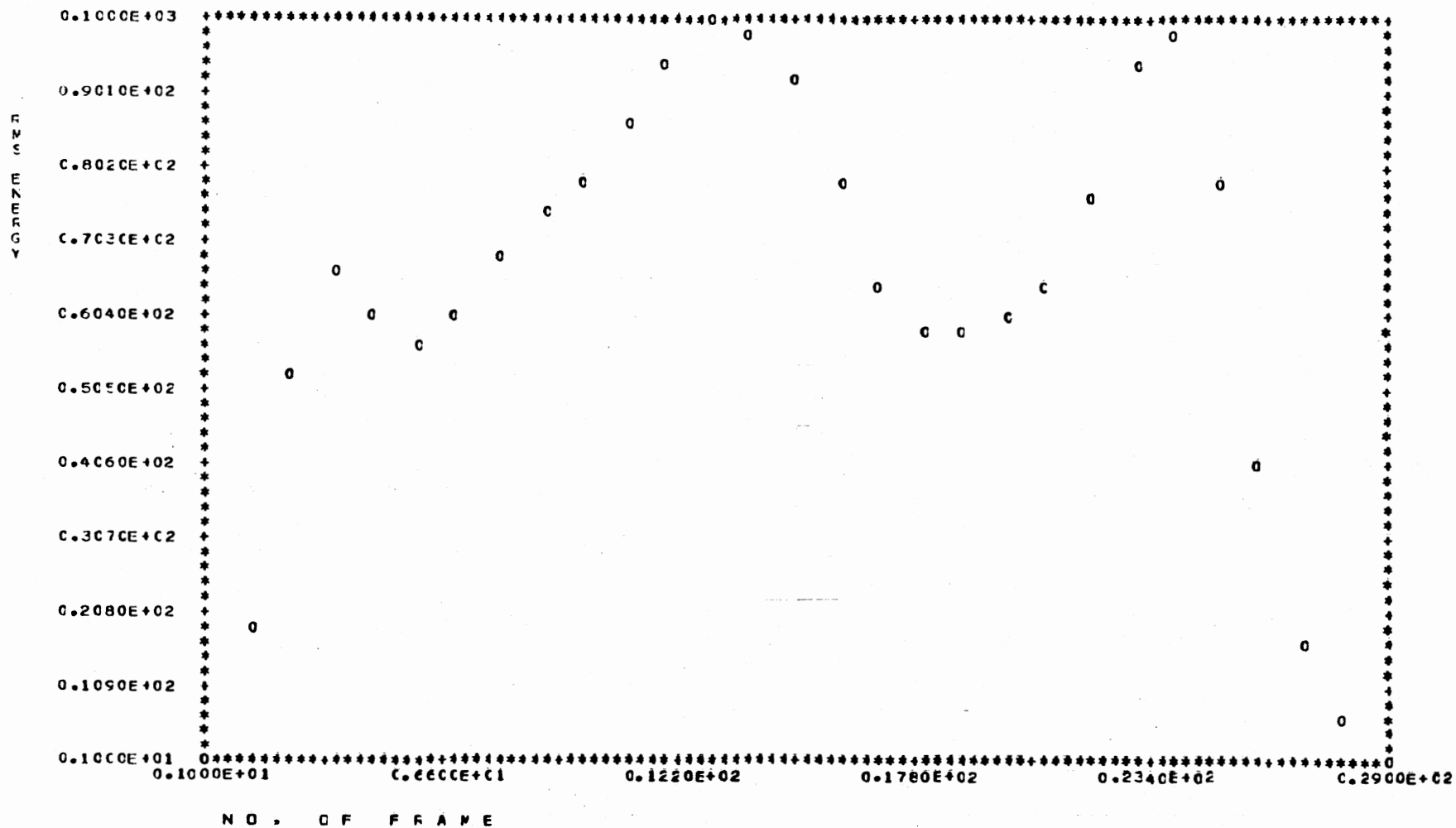


Figure 34. Smoothed RMS Energy Contour for Digit Five, i.e. /faɪv/ Spoken in American English

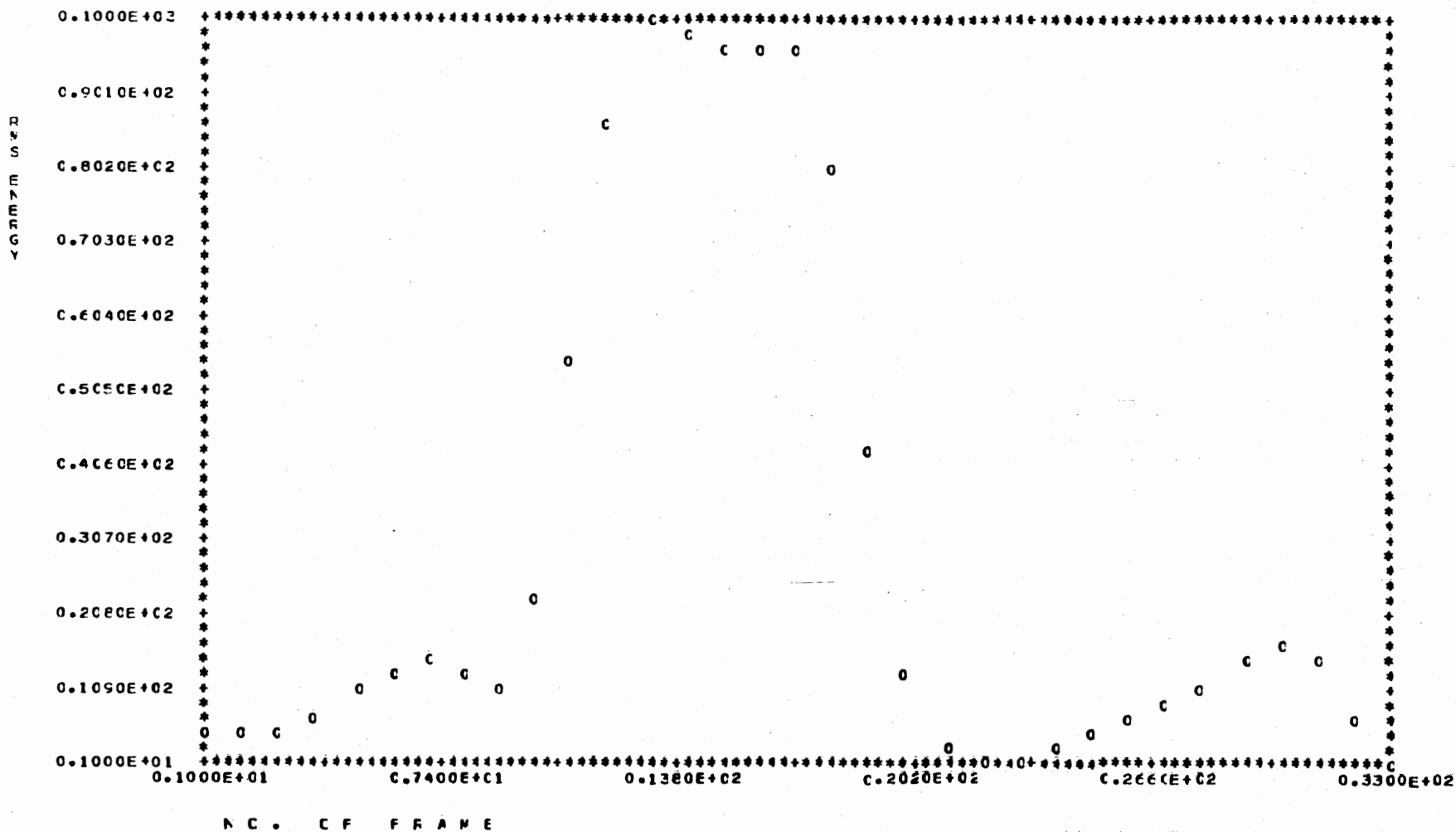


Figure 35. Smoothed RMS Energy Contour for Digit Six, i.e. /siks/ Spoken in American English

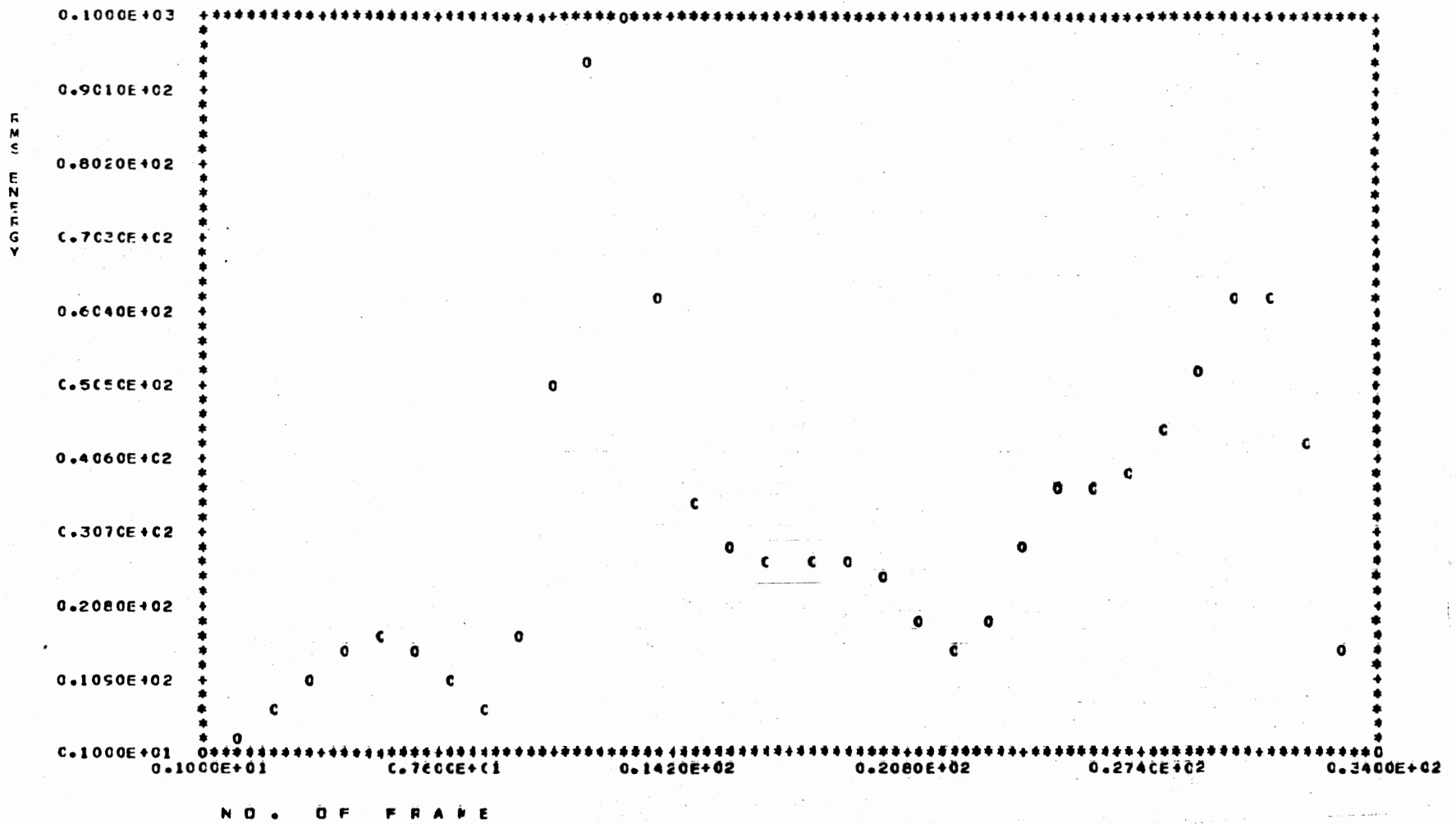


Figure 36. Smoothed RMS Energy Contour for Digit Seven, i.e. /seven/ Spoken in American English

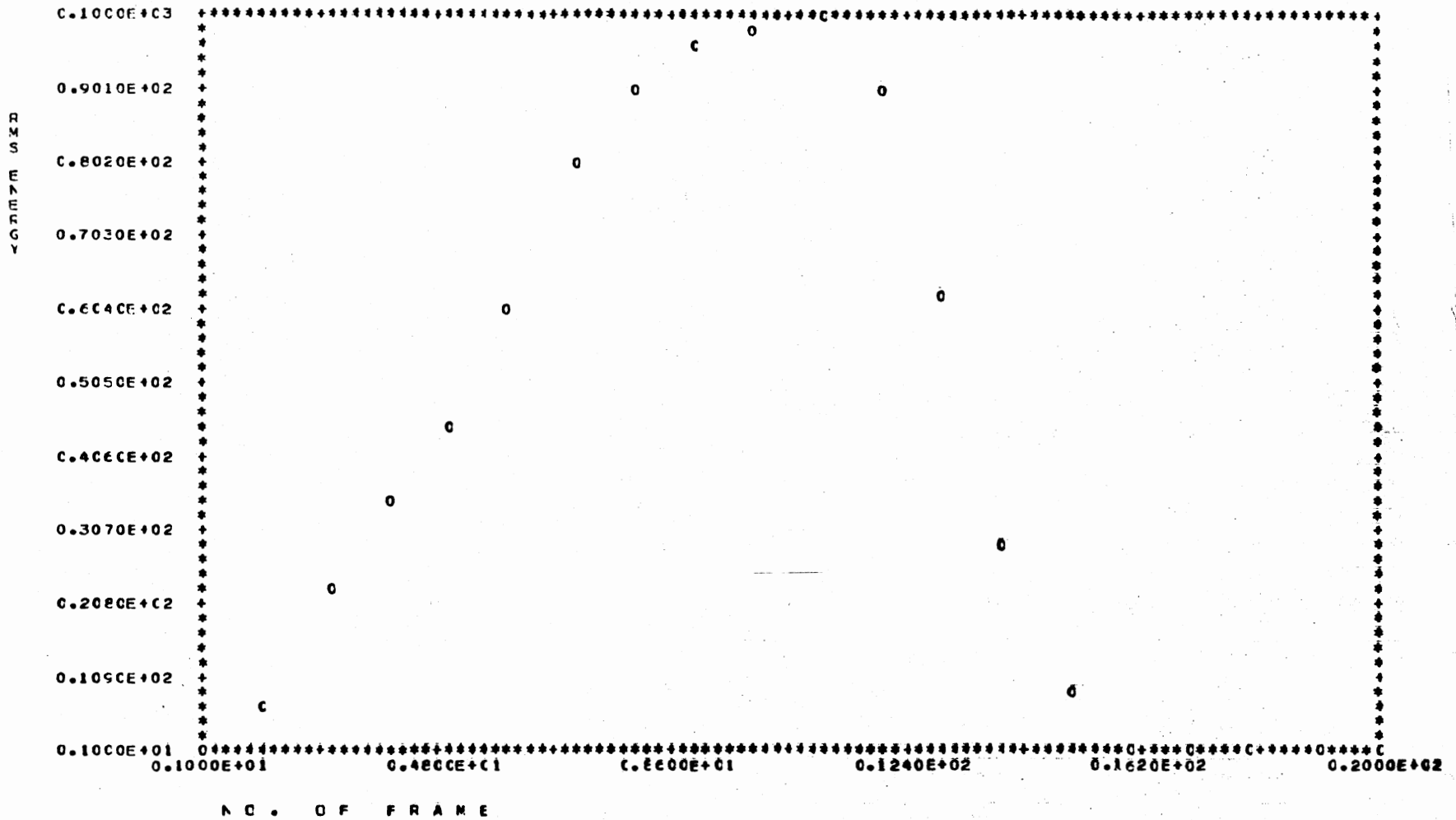


Figure 37. Smoothed RMS Energy Contour for Digit Eight, i.e. /eIt/ Spoken in American English

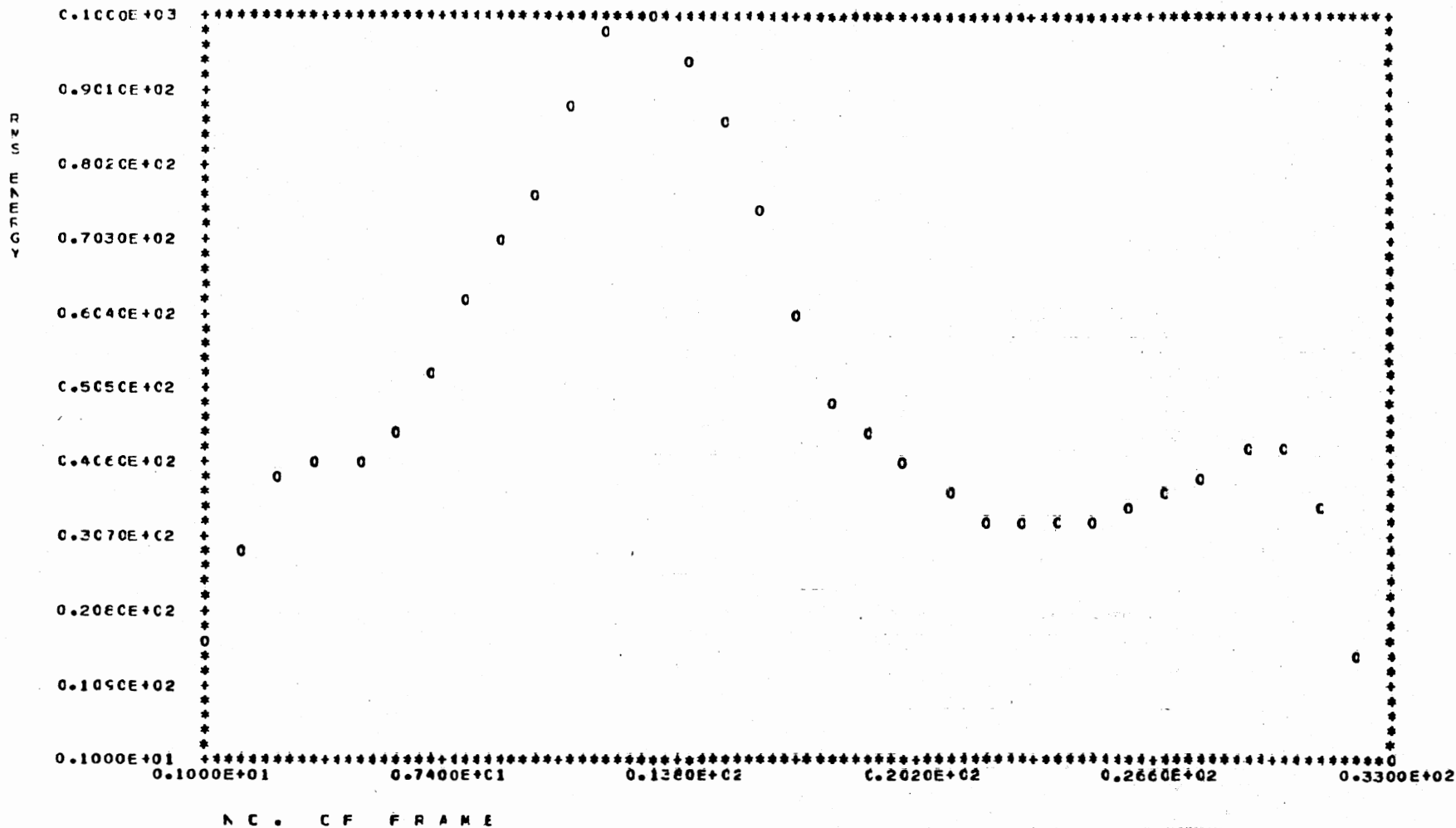


Figure 38. Smoothed RMS Energy Contour for Digit Nine, i.e. /naIn/ Spoken in American English

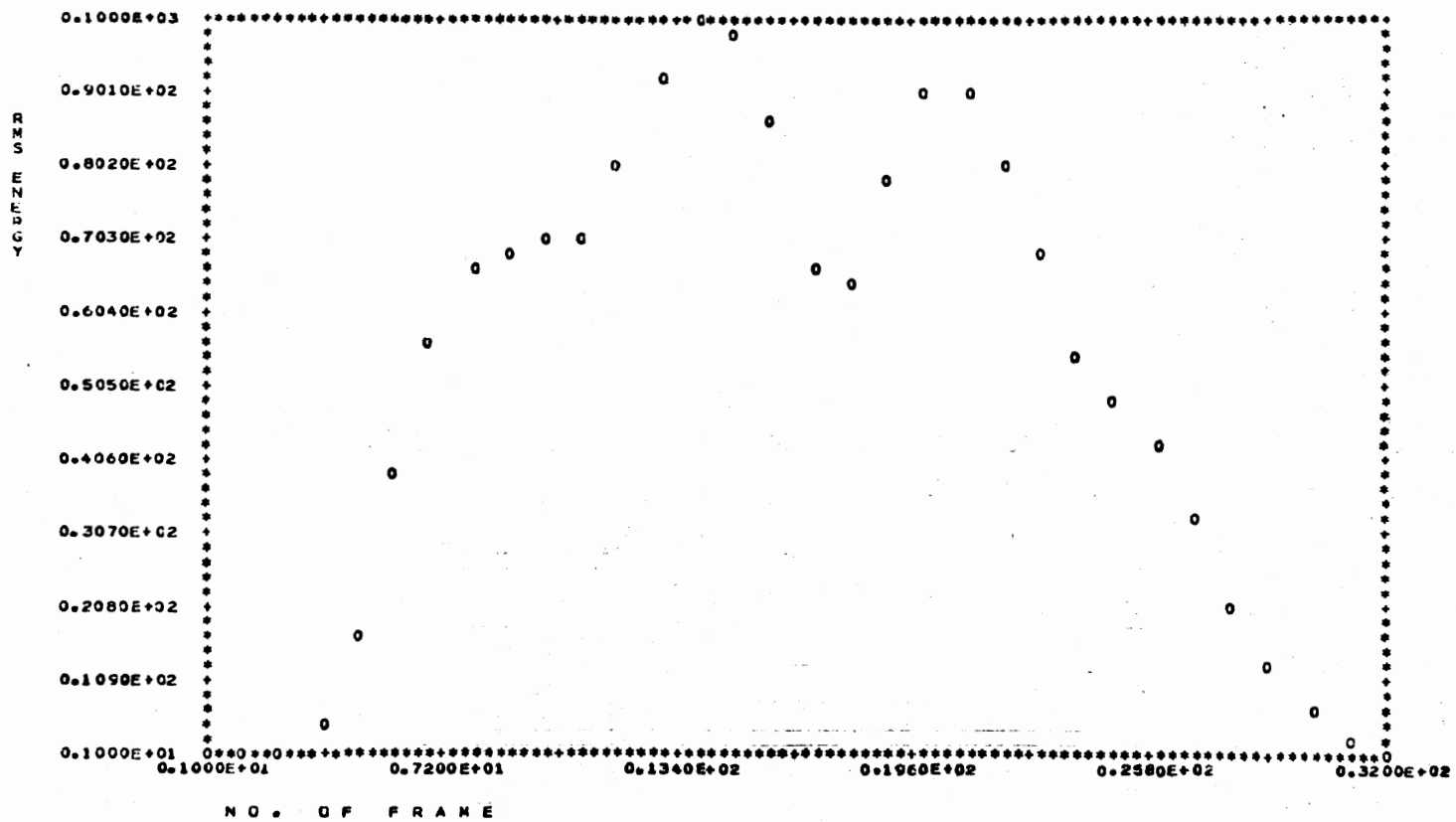


Figure 39. Smoothed and Quantized RMS Energy Contour for Digit Zero, i.e. /zIro/ Spoken by International Speaker

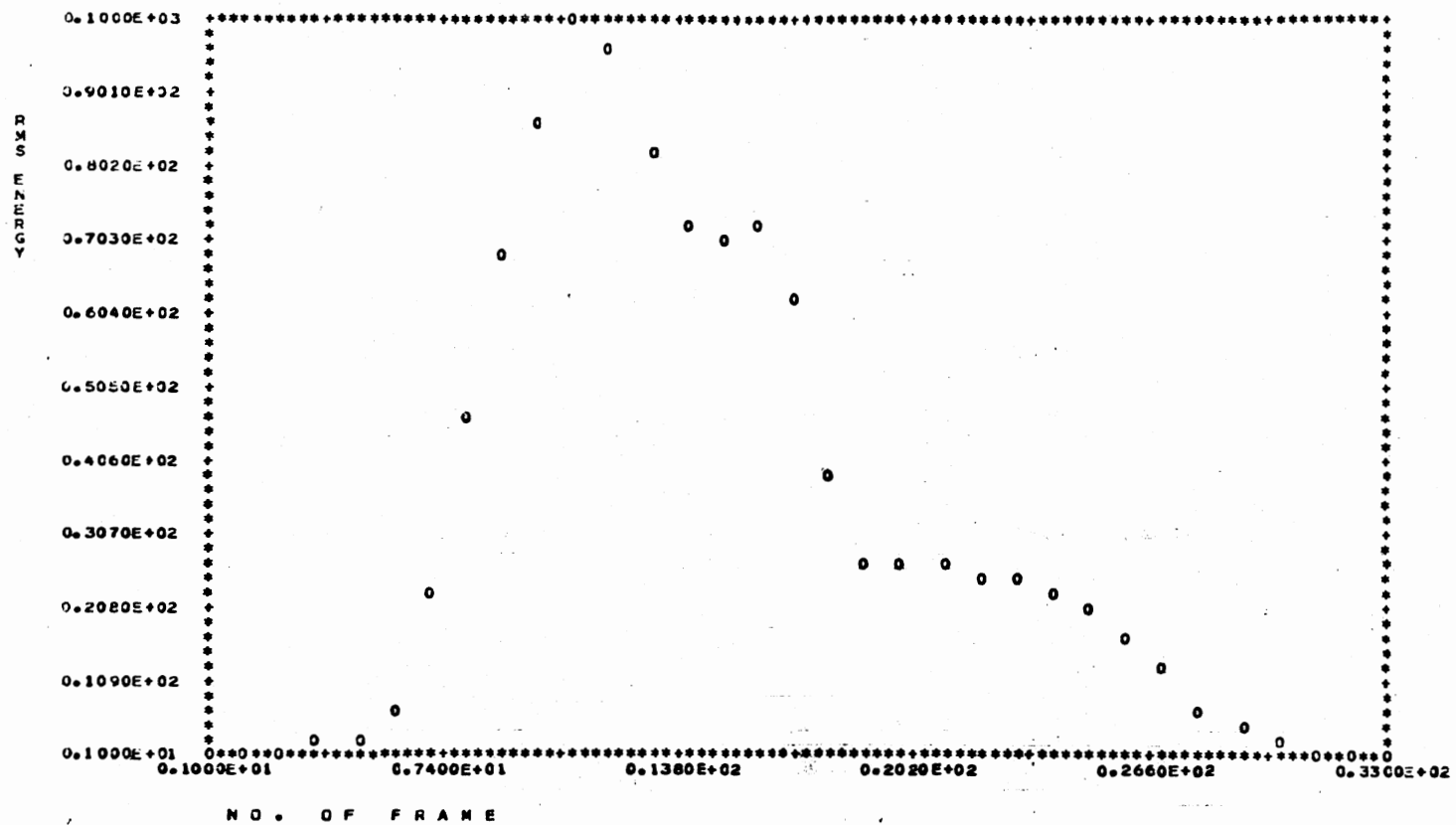


Figure 40. Smoothed and Quantized RMS Energy Contour for Digit One, i.e. /wAn/ Spoken by International Speaker

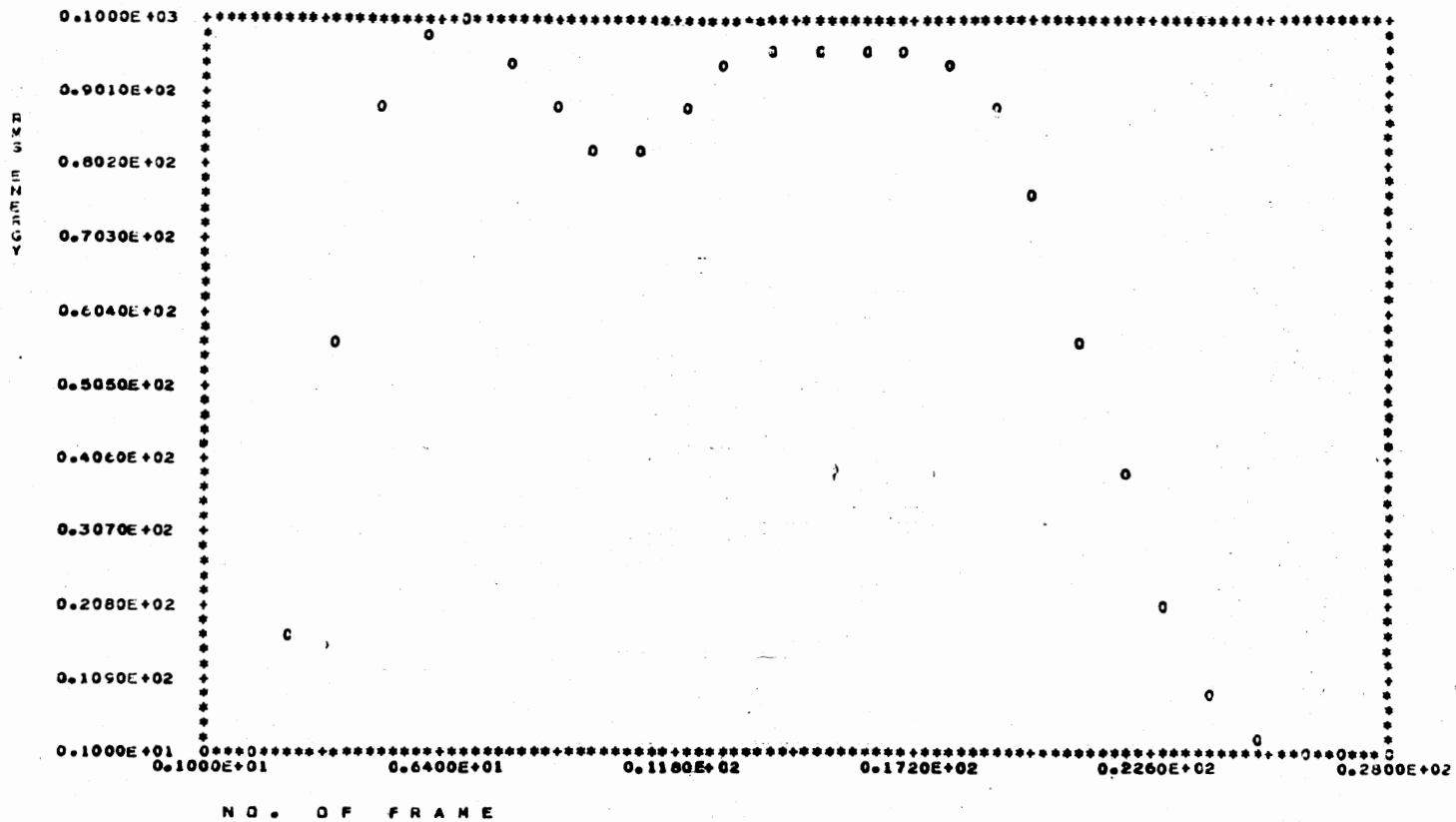


Figure 41. Smoothed and Quantized RMS Energy Contour for Digit Two, i.e. /tu/ Spoken by International Speaker

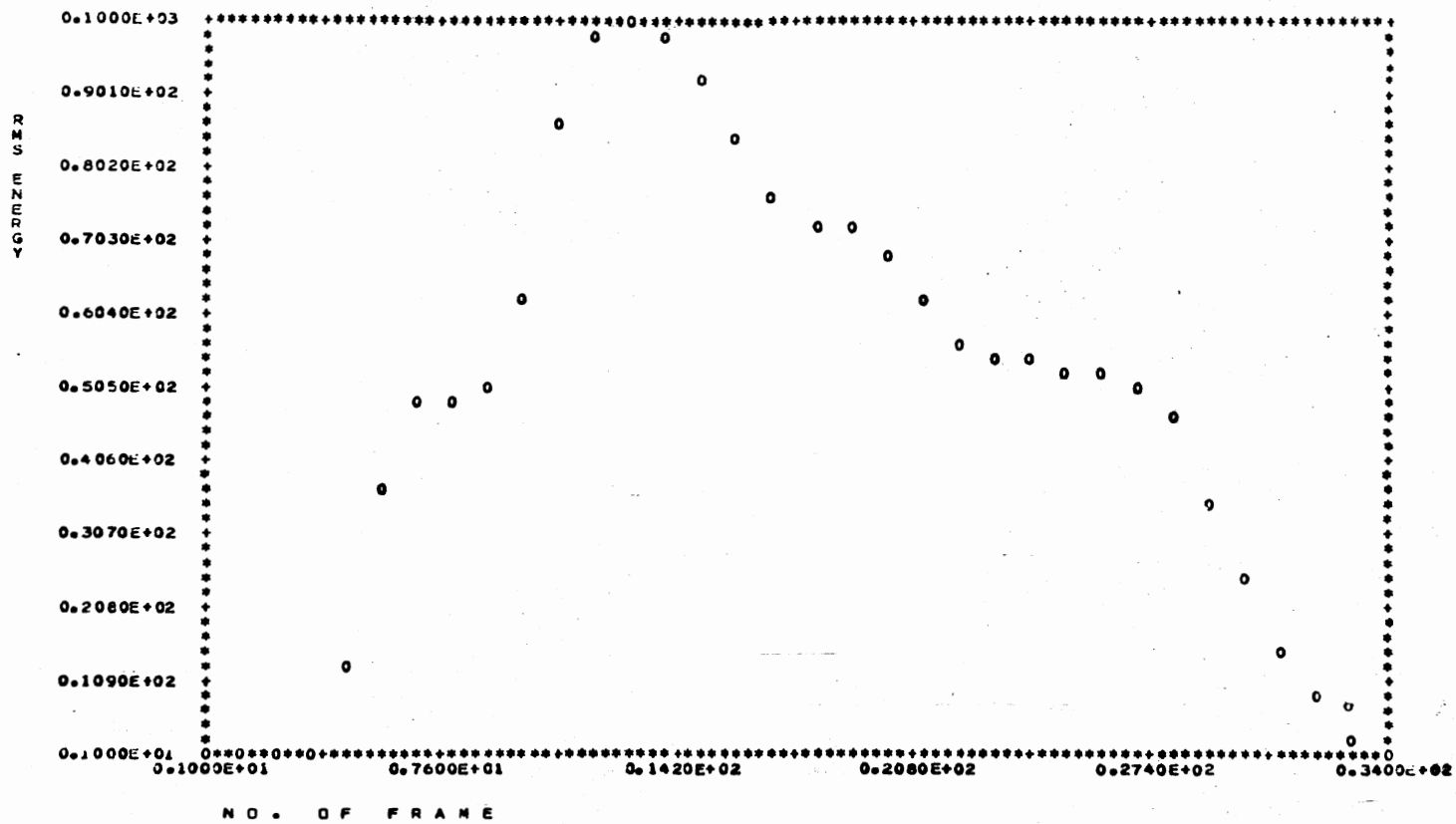


Figure 42. Smoothed and Quantized RMS Energy Contour for Digit Three, i.e. /θri/ Spoken by International Speaker

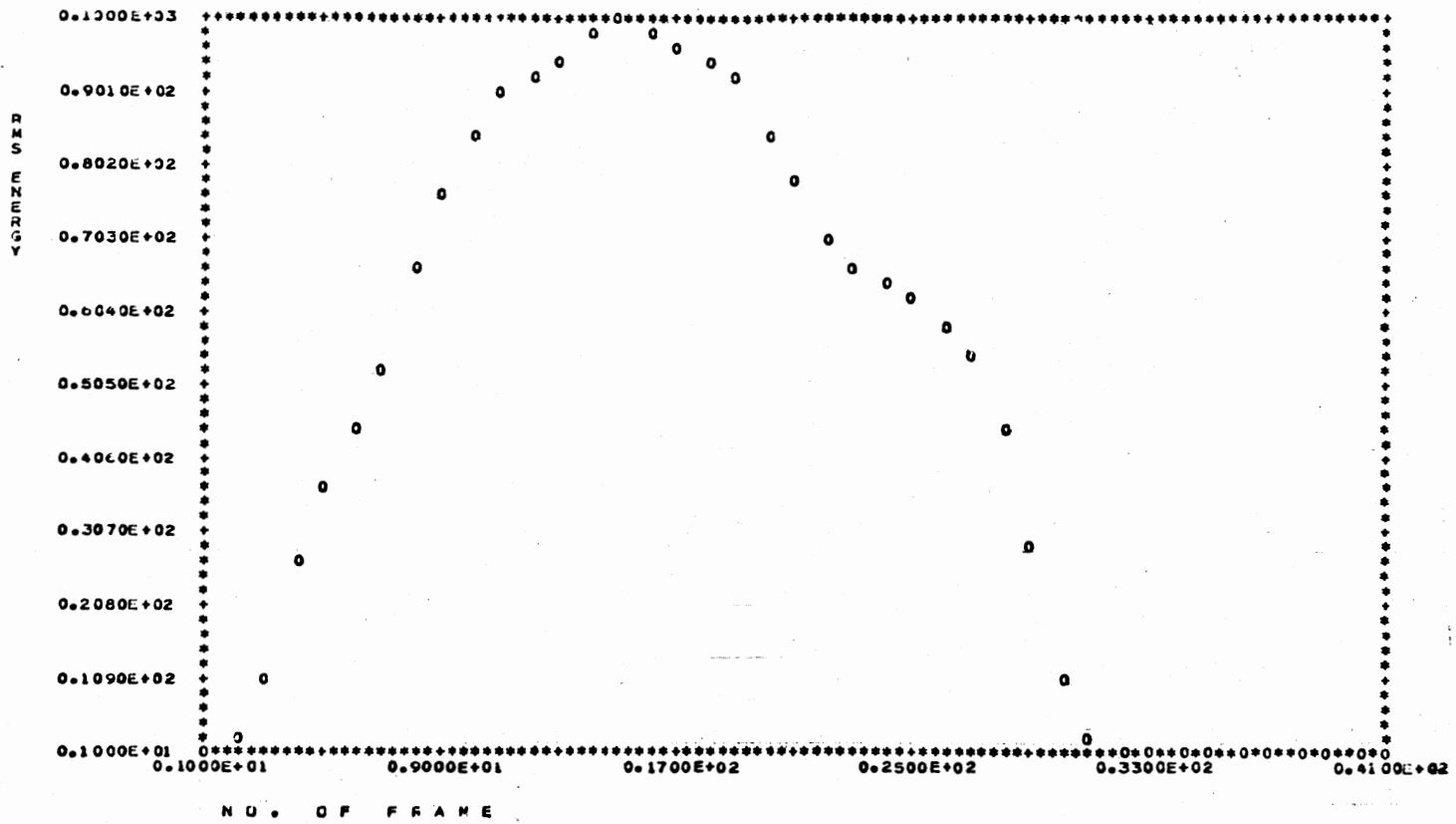


Figure 43. Smoothed and Quantized RMS Energy Contour for Digit Four, i.e. /fɔr/ Spoken by International Speaker

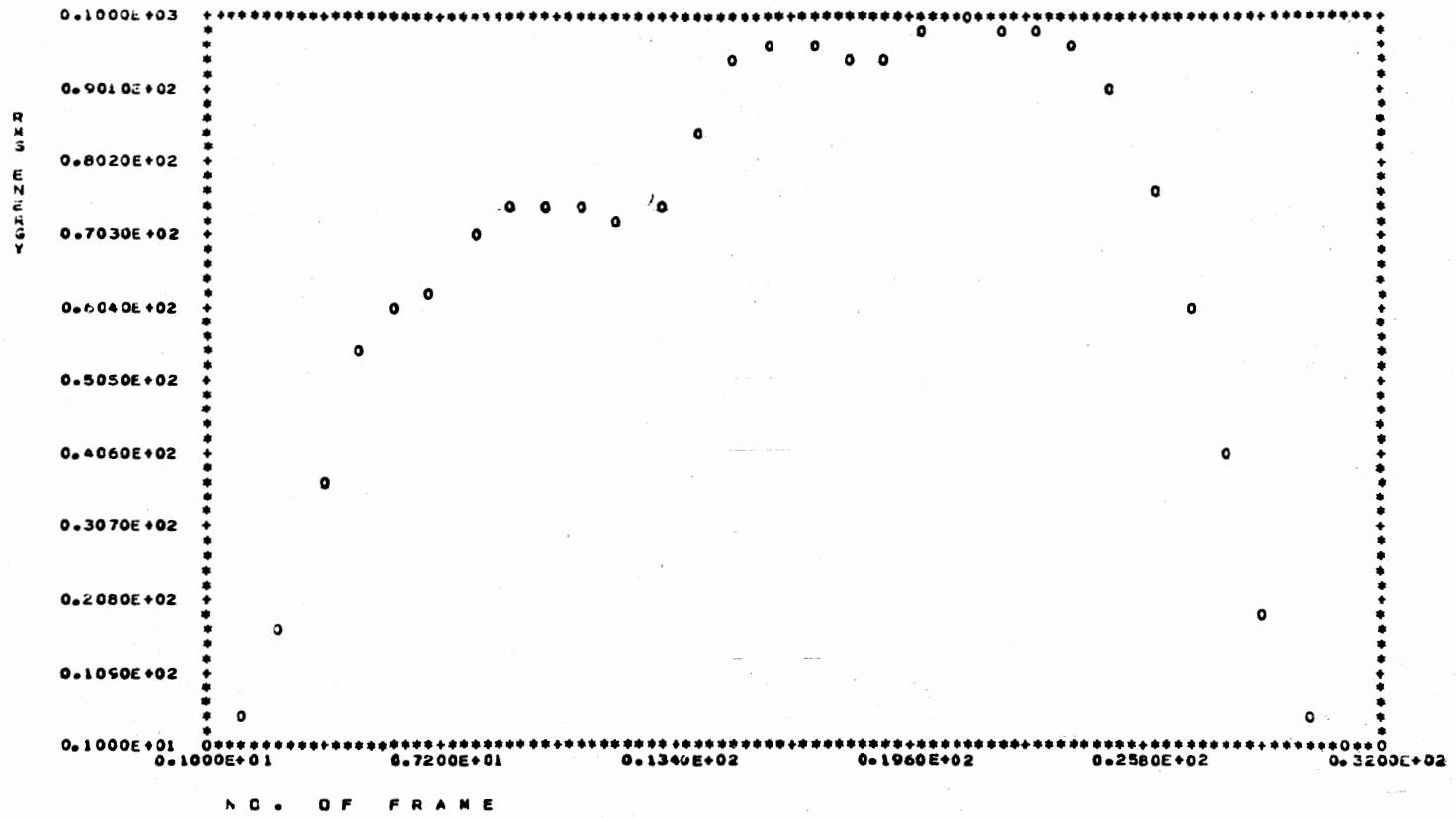


Figure 44. Smoothed and Quantized RMS Energy Contour for Digit Five, i.e. /faIv/ Spoken by International Speaker

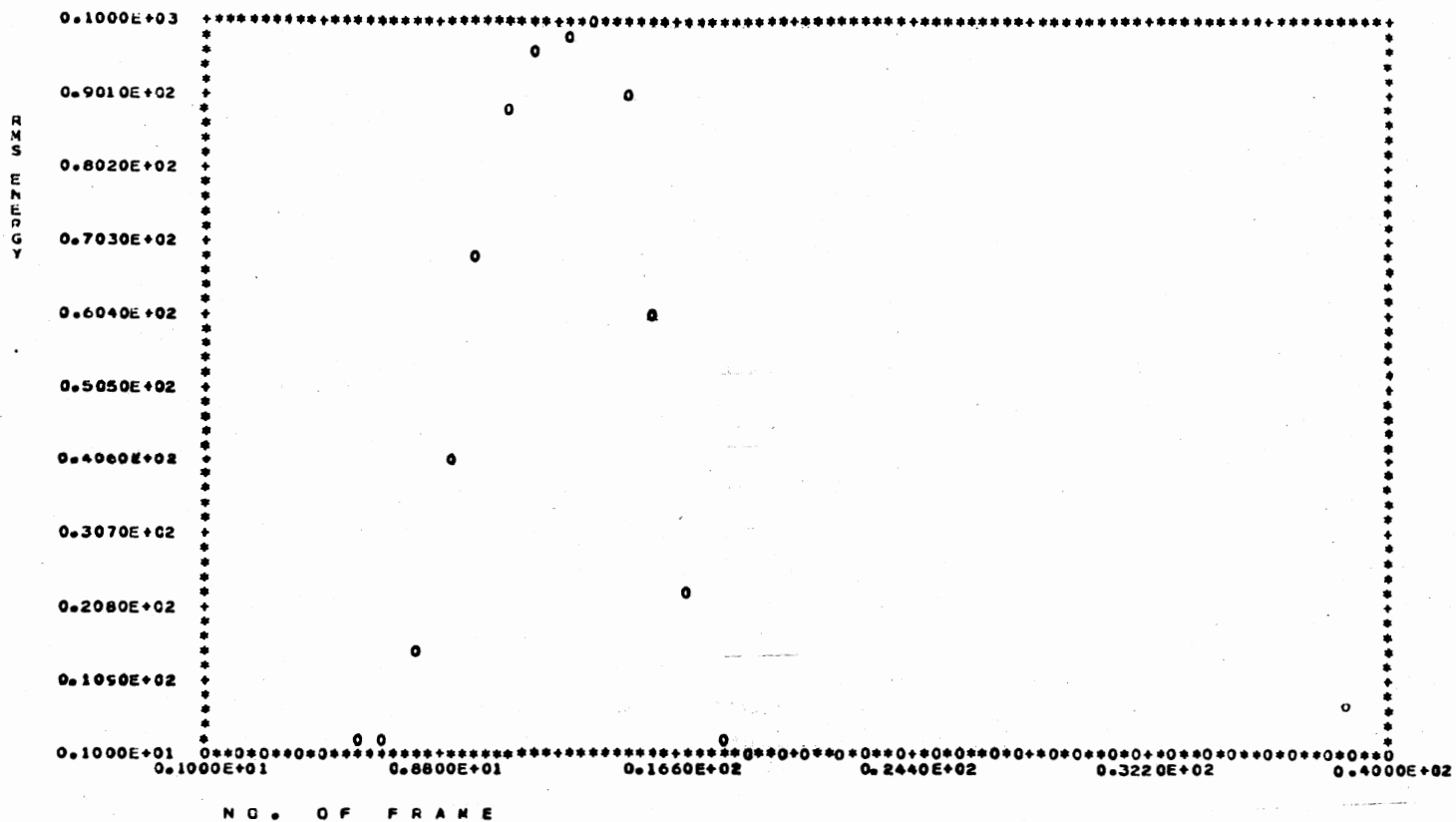


Figure 45. Smoothed and Quantized RMS Energy Contour for Digit Six, i.e. /sIks/ Spoken by International Speaker

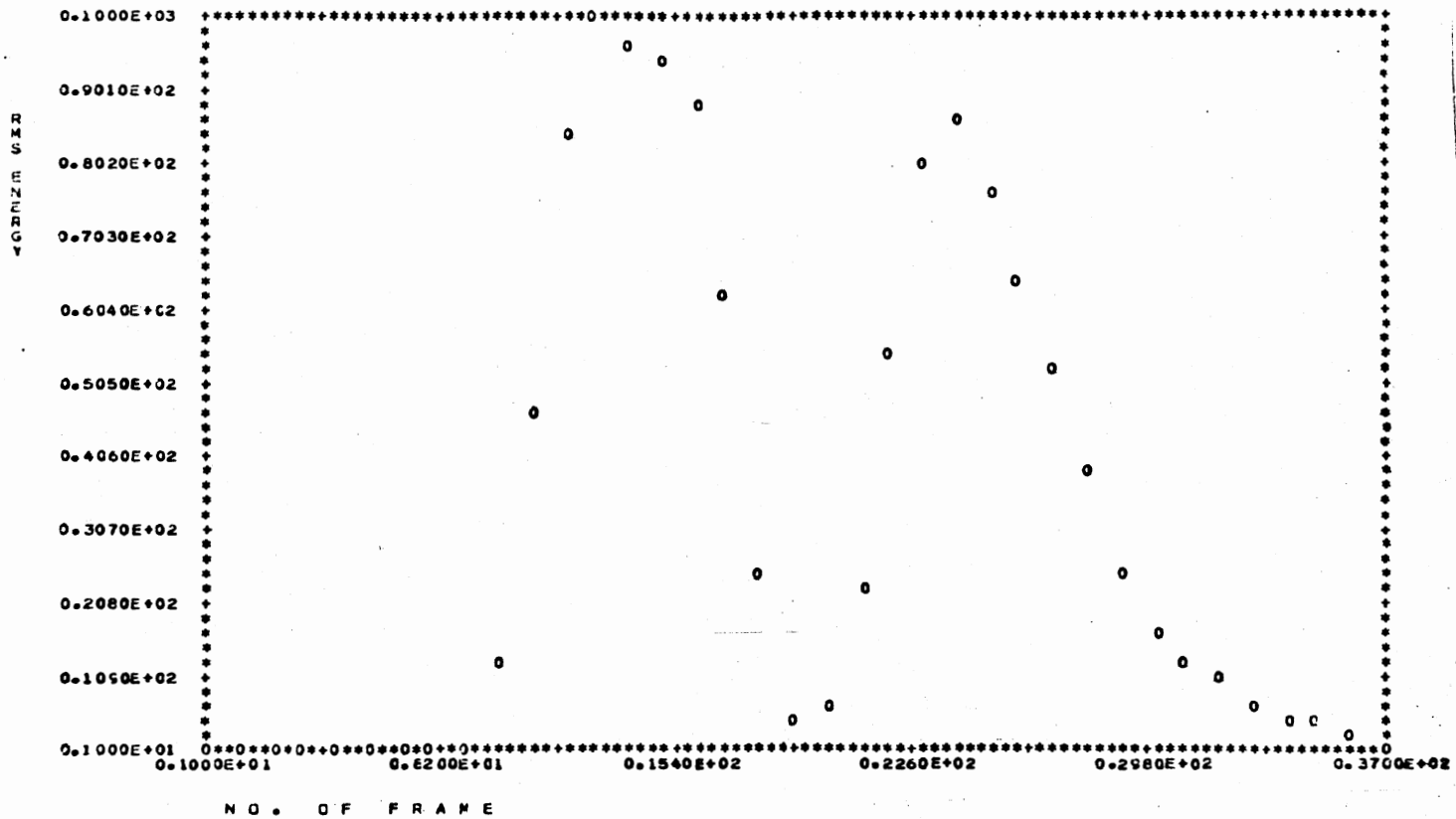


Figure 46. Smoothed and Quantized RMS Energy Contour for Digit Seven, i.e. /seven/ Spoken by International Speaker

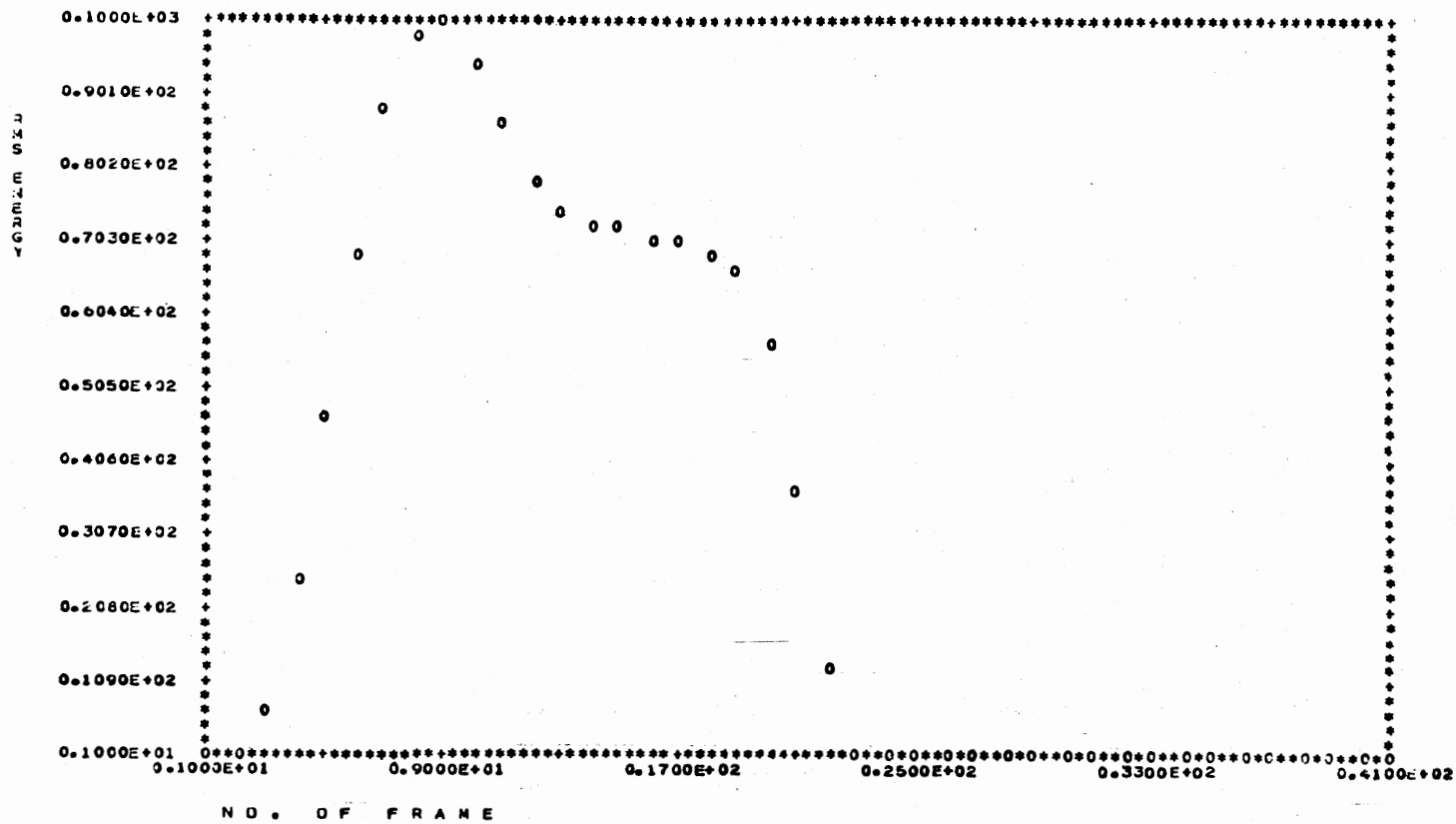


Figure 47. Smoothed and Quantized RMS Energy Contour for Digit Eight, i.e. /eIt/ Spoken by International Speaker

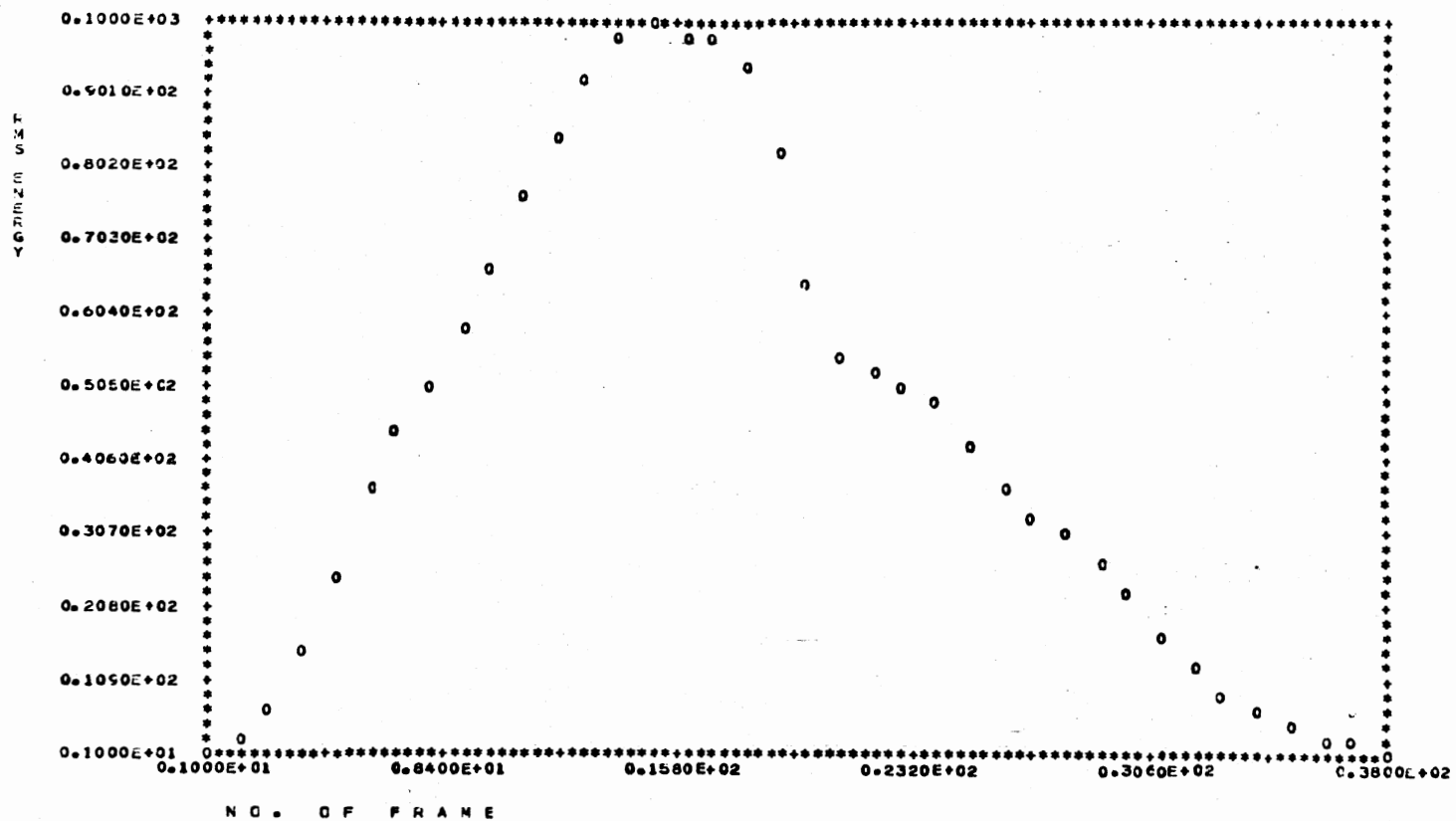


Figure 48. Smoothed and Quantized RMS Energy Contour for Digit Nine, i.e. /naIn/ Spoken by International Speaker

TABLE V
SMOOTHED RMS ENERGY PEAKS AND RATIOS FOR DIGITS SPOKEN IN
AMERICAN ENGLISH FOR FIXED MEAN AND VARYING PEAKS

Digit	P ₁	P ₂	Ratio
/zIro/	15.05	12.33	0.82
/wAn/	4.19	11.71	0.36
/tu/	8.64	5.14	0.595
/θri/	10.86	14.46	0.75
/fɔr/	12.80		
/faIv/	14.28	14.24	1.0
/sIks/	56.79		
/seven/	36.04	22.84	0.63
/eIt/	22.72		
/naIn/	23.76	10.28	0.43

TABLE VI
SMOOTHED AND QUANTIZED RMS ENERGY PEAKS AND THEIR RATIOS FOR
DIGITS SPOKEN IN AMERICAN ENGLISH

Digit	P_1	P_2	Ratio
/zIro/	100	91	0.91
/wAn/	100	72	0.72
/tu/	100	97	0.97
/θri/	48	100	0.48
/fɔr/	100	63	0.63
/faIV/	96	100	0.96
/sIks/	100		
/seven/	100	86	0.86
/eIt/	100		
/naIn/	100		

Following the block diagram in Figure 28, the second path after end-point detection is through Hamming window. As discussed in Chapter III, the window length is taken as 150 points. From the windowed data, a 14th order LPA model is computed. From this model, the BTR, CTR, FTR, SFBR, smoothed BTR, smoothed CTR, and smoothed FTR are computed for each frame. These plots are given in Figures 49-58. Plots of smoothed BTR, smoothed CTR, and smoothed FTR was obtained, in order to determine which parameter would give the best phonemic feature, to separate vowels from consonants. A threshold level is set to distinguish between vowel, vowel-like, and non-vowel segments. Using the plots in Figures 49-58, it can be seen that all these parameters have some useful phonetic and acoustical features. In addition, the BTR and CTR indicate better acoustical configuration than the FTR. Furthermore, the BTR is found to give moderately better results than the CTR. However, some modification of the CTR range is required prior to future applications.

The parameters derived are used in the phonetic feature detection stage. For phoneme segmentation and other classifications, the RMS dip-classification algorithm is used, which was discussed in Chapter III. The algorithm detects three types of dips or valleys in the smoothed RMS energy contour according to the signs of the functions Z_1 and Z_2 as a vowel, vowel-like, and non-vowel. Based upon the discussion in Chapter III, the vowel, vowel-like and non-vowel decision is made for the entire RMS energy contour. Also the vowel, vowel-like and non-vowel decision is made using the smoothed BTR contour. Then an "OR" decision is made using the above two classification

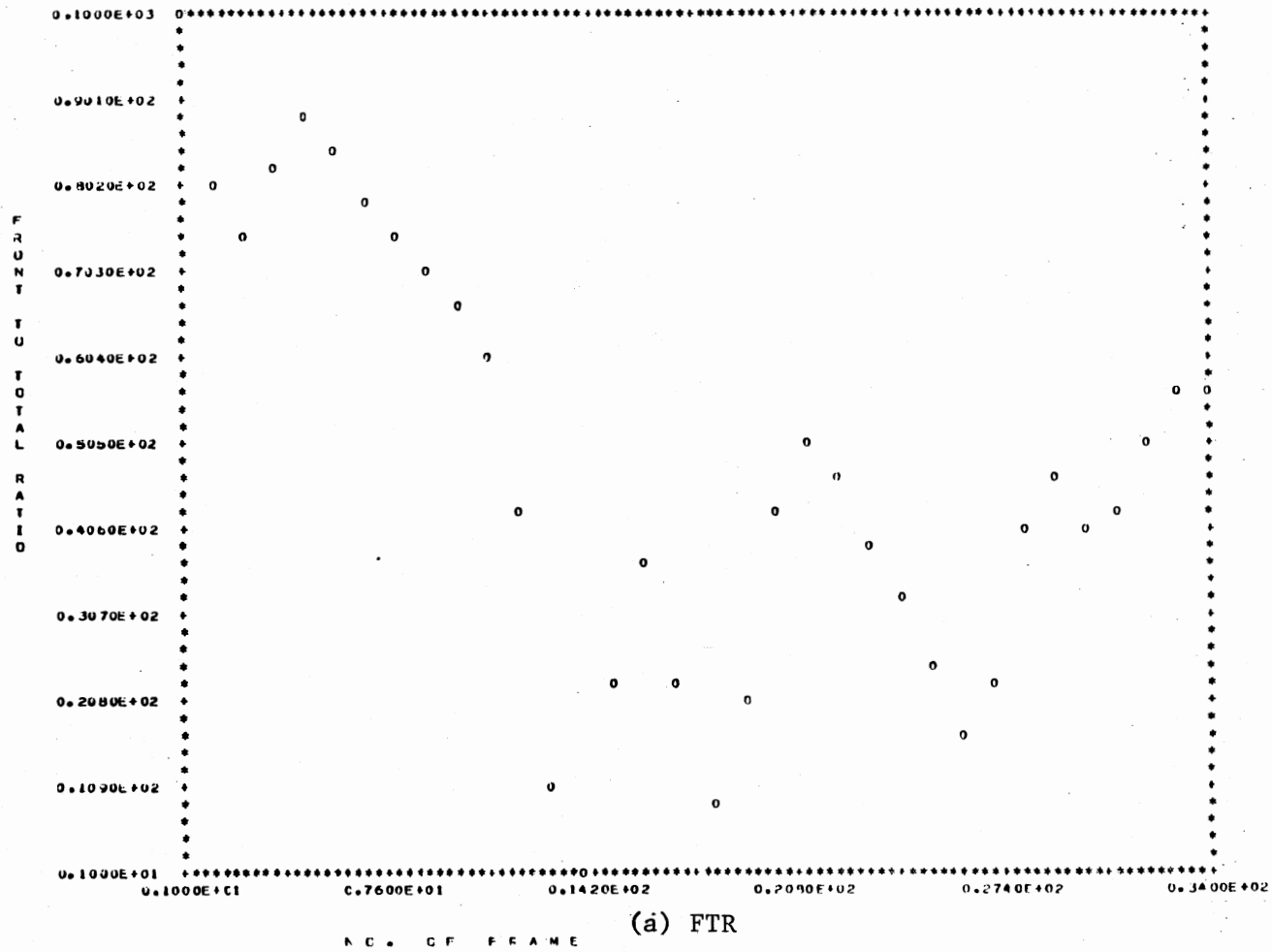


Figure 49. Smoothed and Quantized Feature Parameters for Digit Zero Spoken in American English

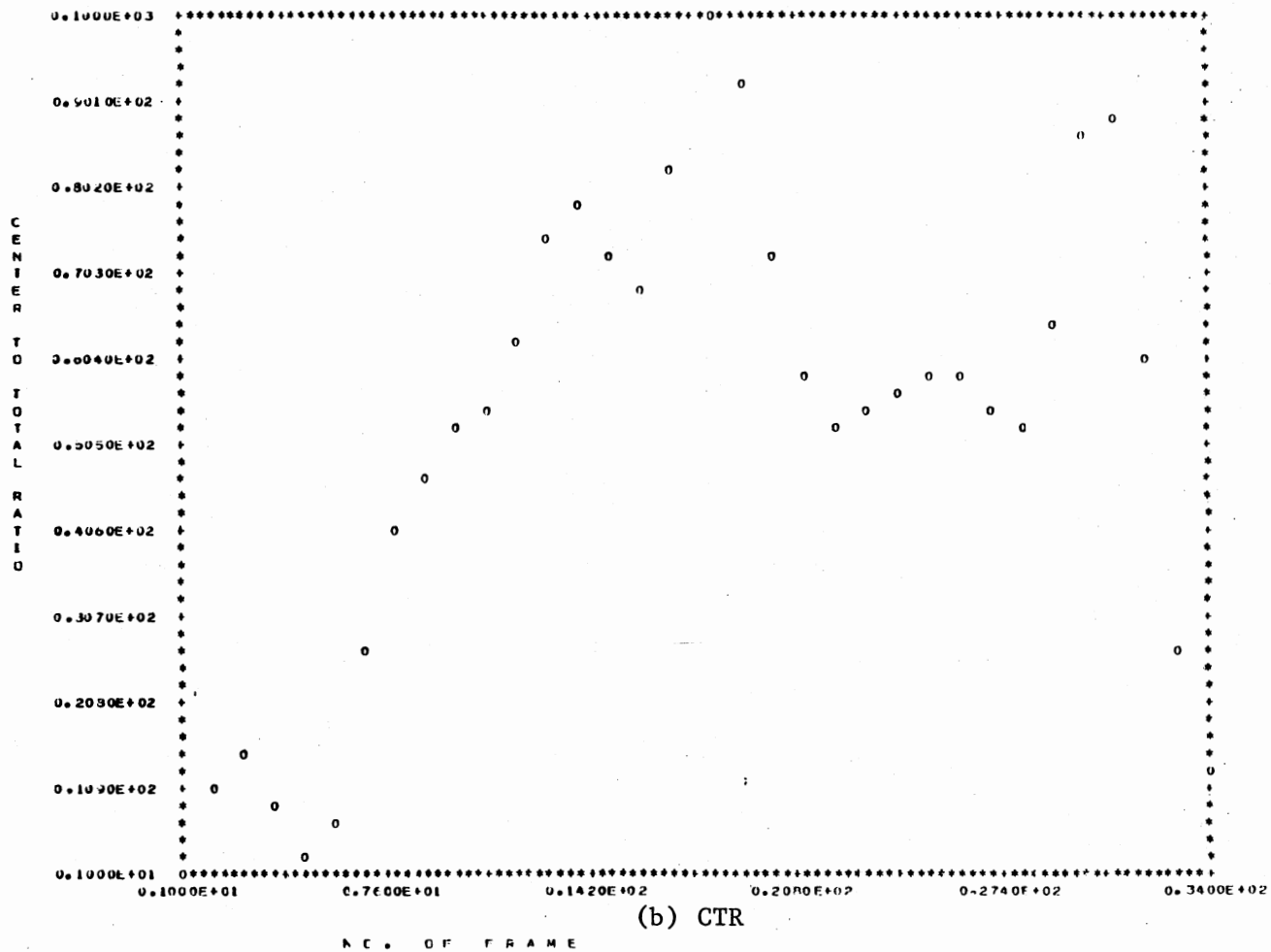


Figure 49. (Continued)

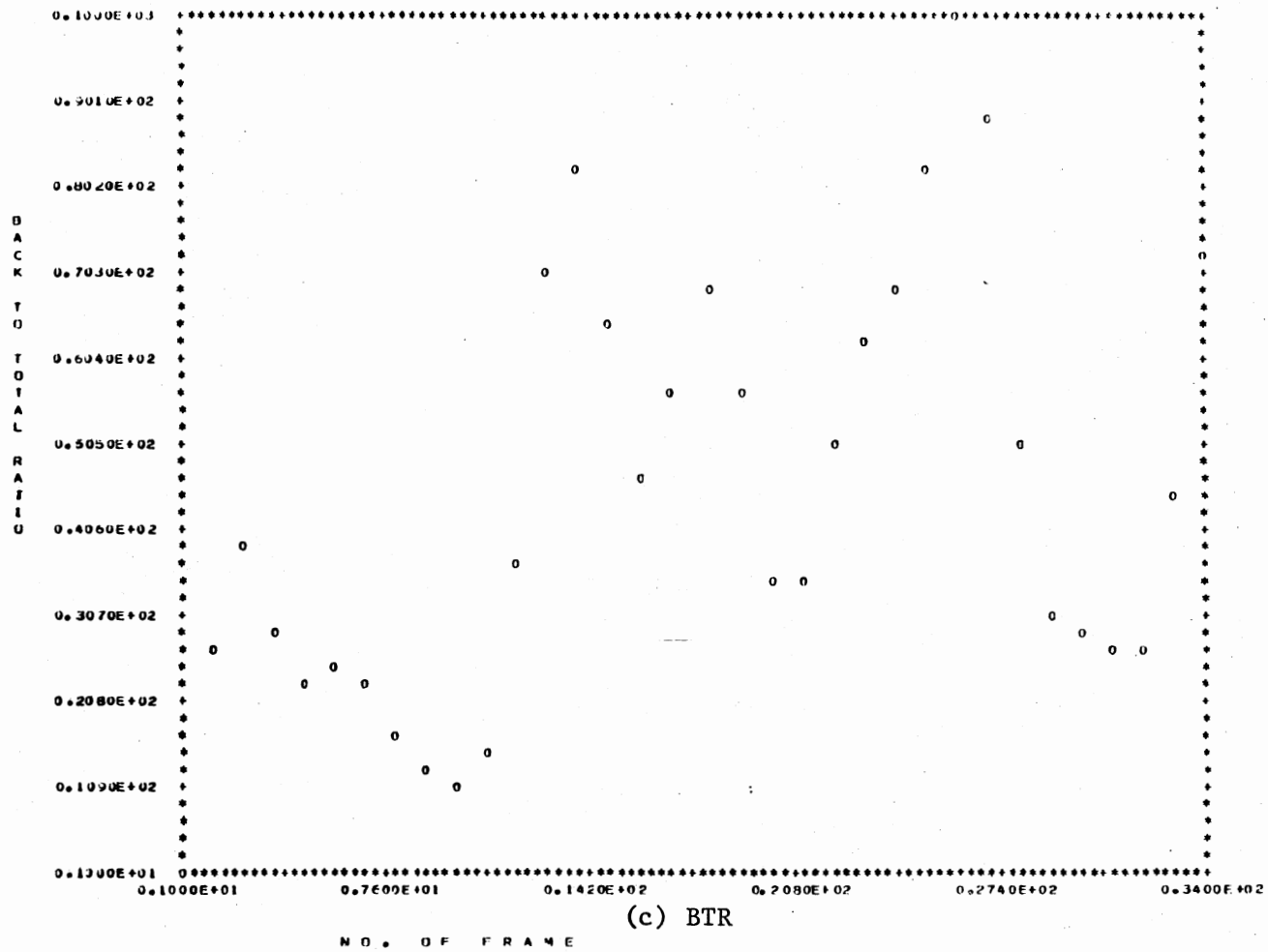


Figure 49. (Continued)

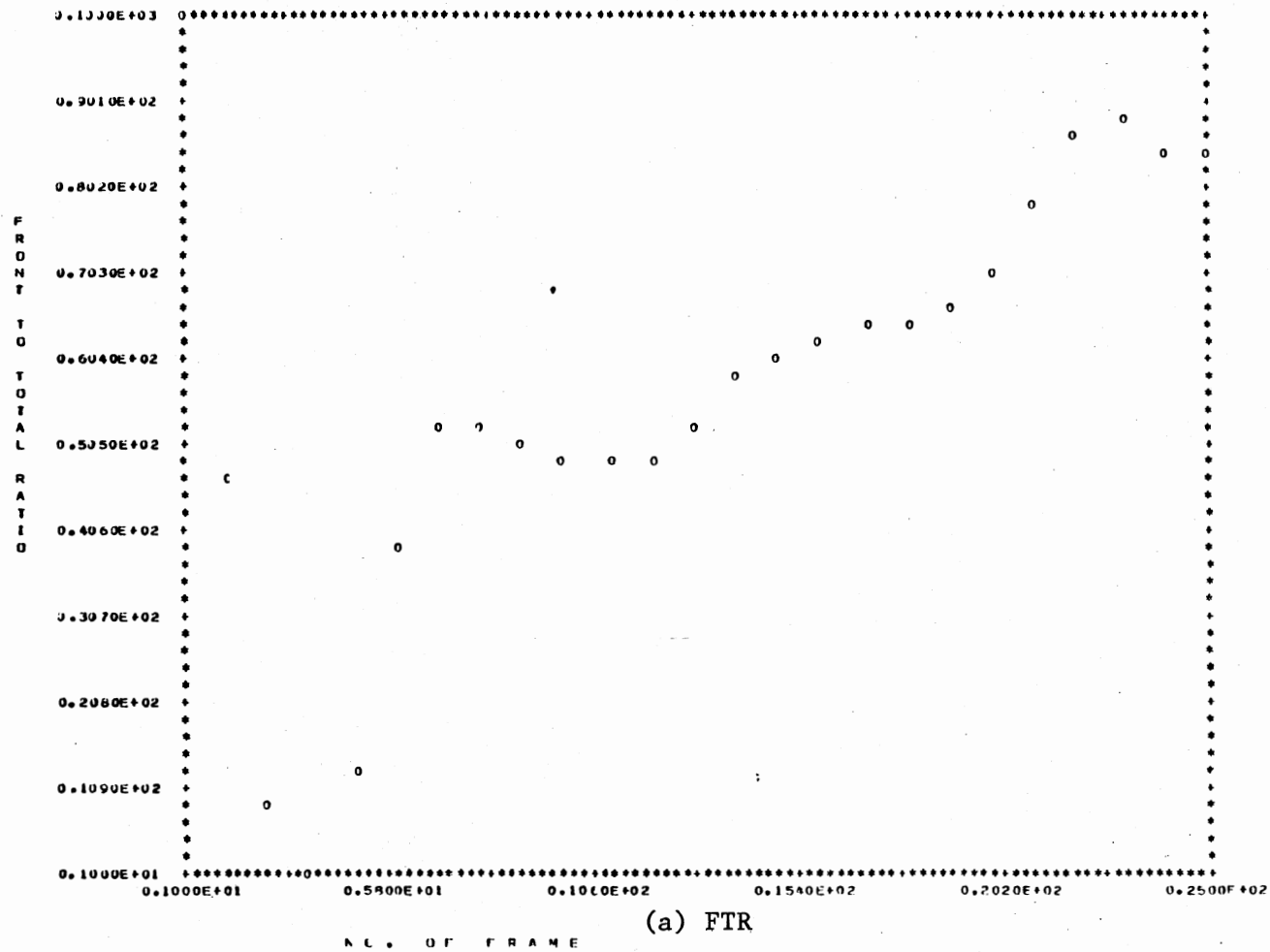


Figure 50. Smoothed and Quantized Feature Parameters for Digit One Spoken in American English

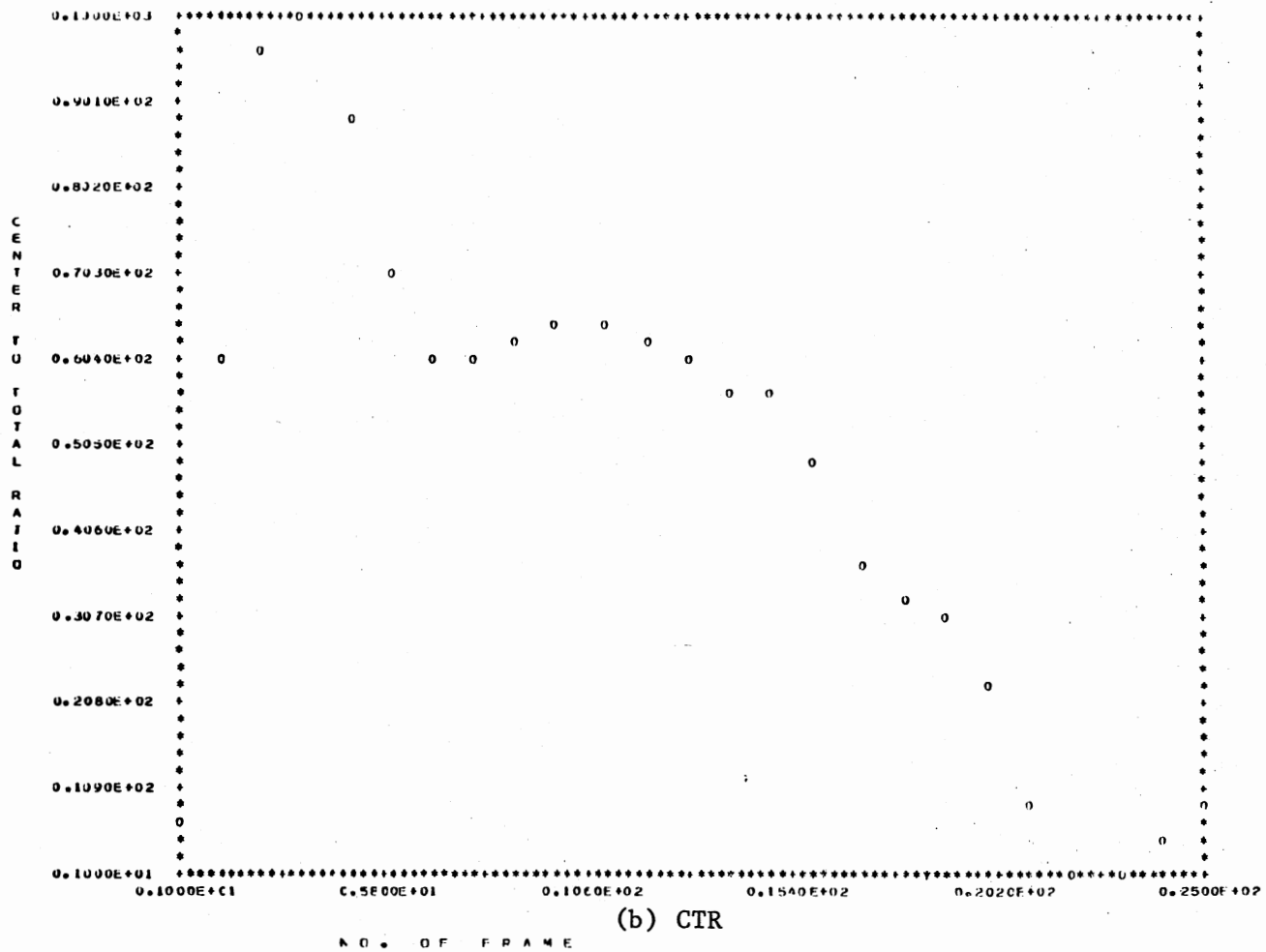


Figure 50. (Continued)

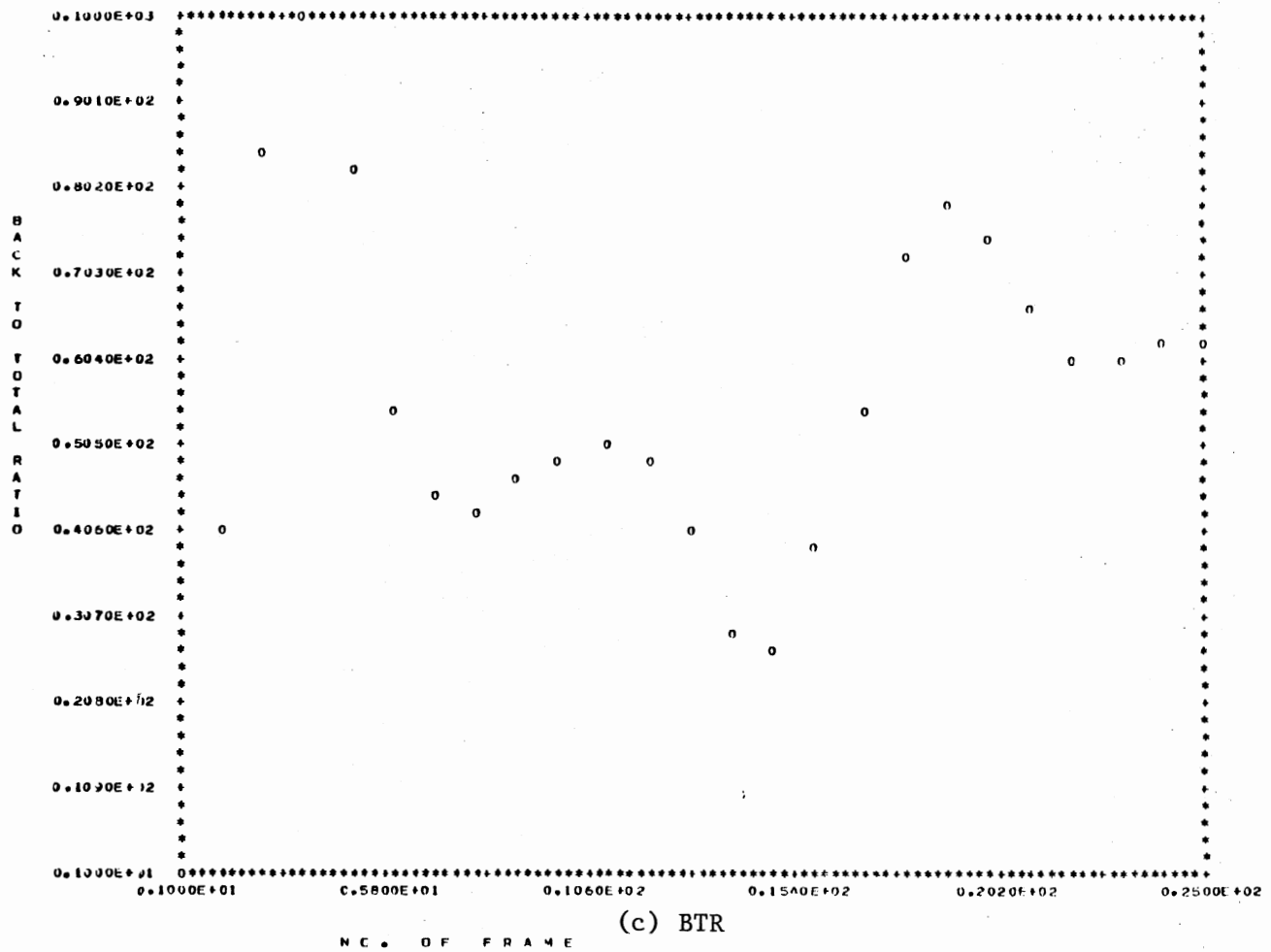


Figure 50. (Continued)

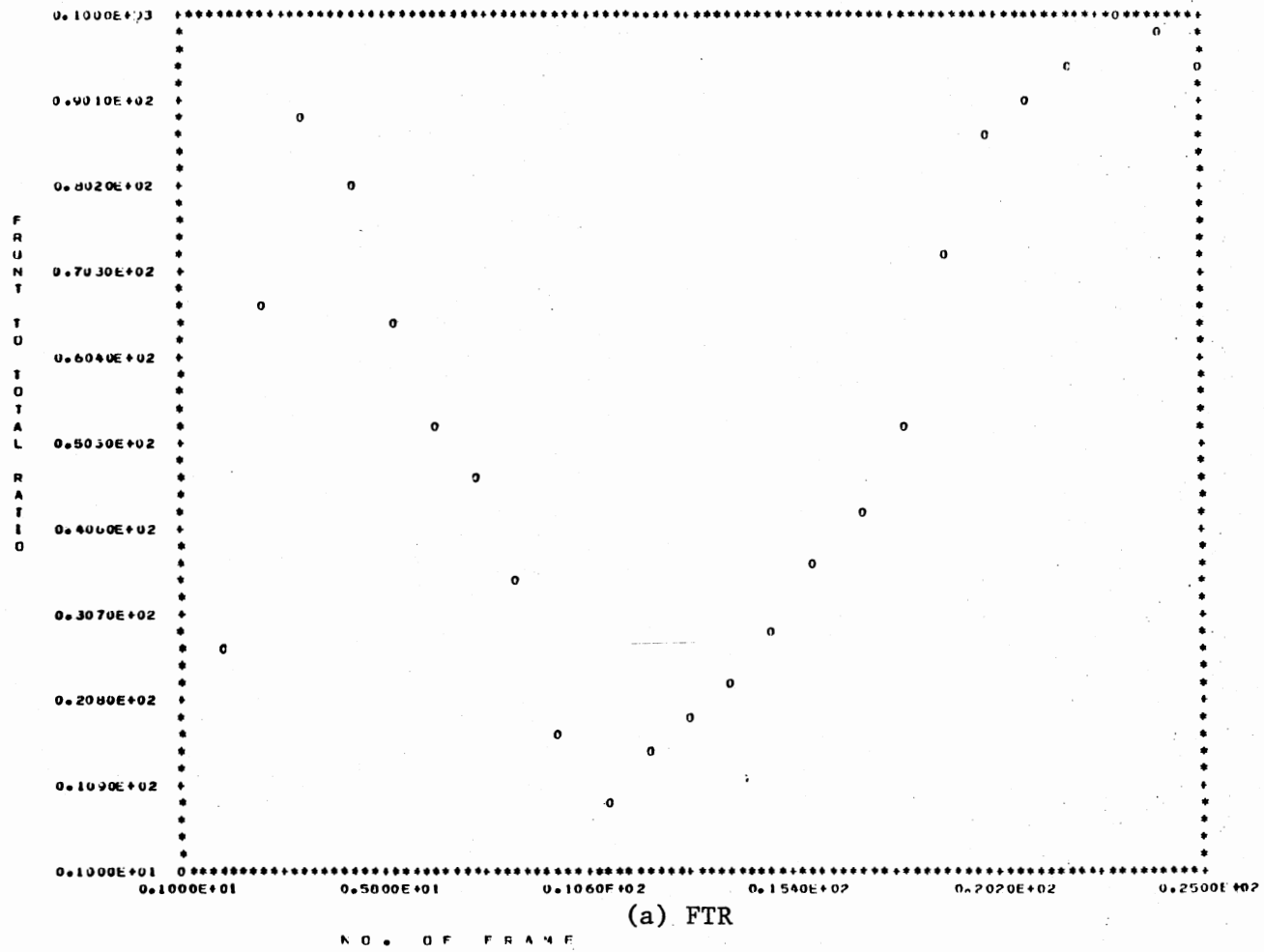


Figure 51. Smoothed and Quantized Feature Parameters for Digit Two Spoken in American English

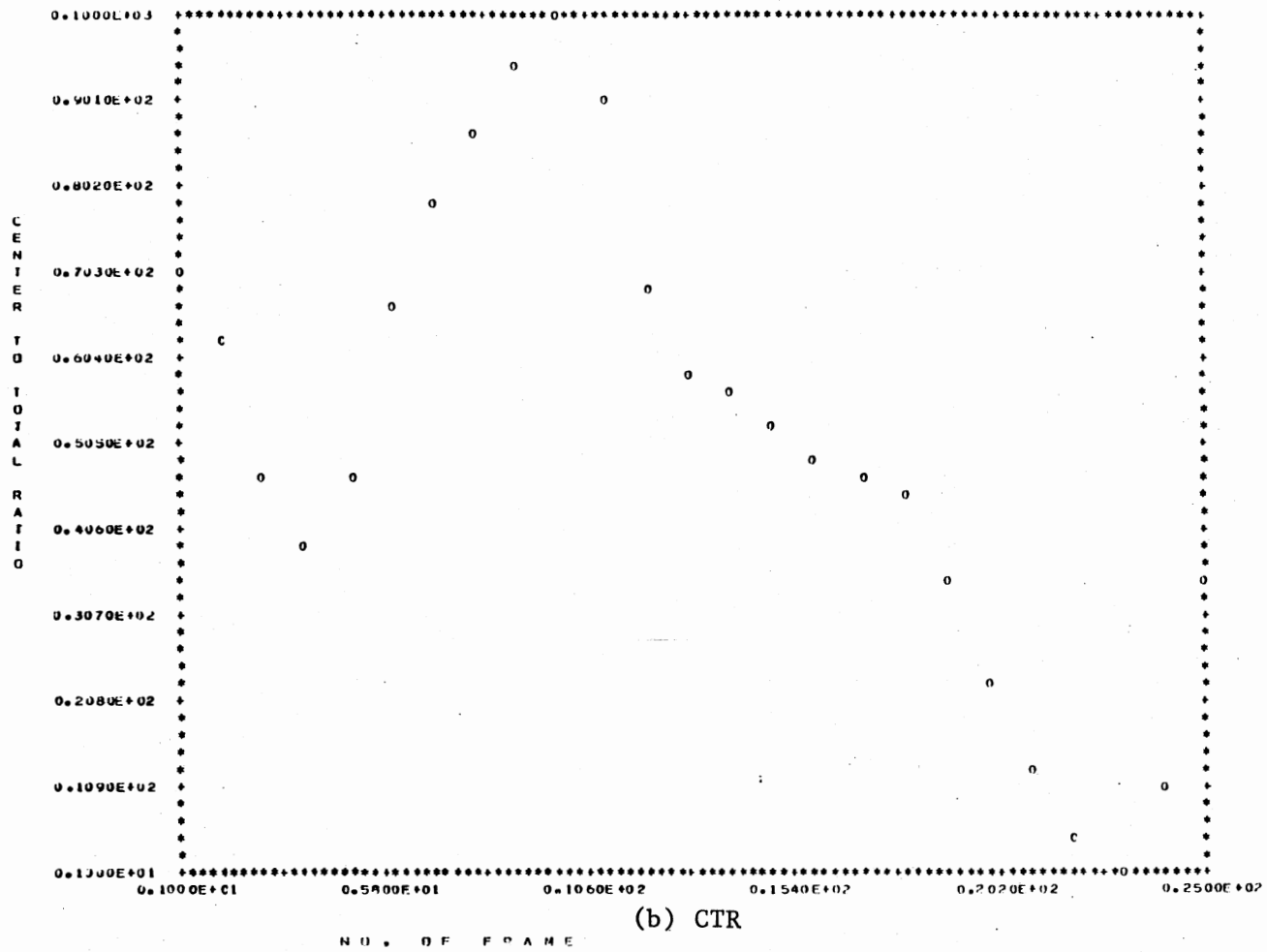


Figure 51. (Continued)

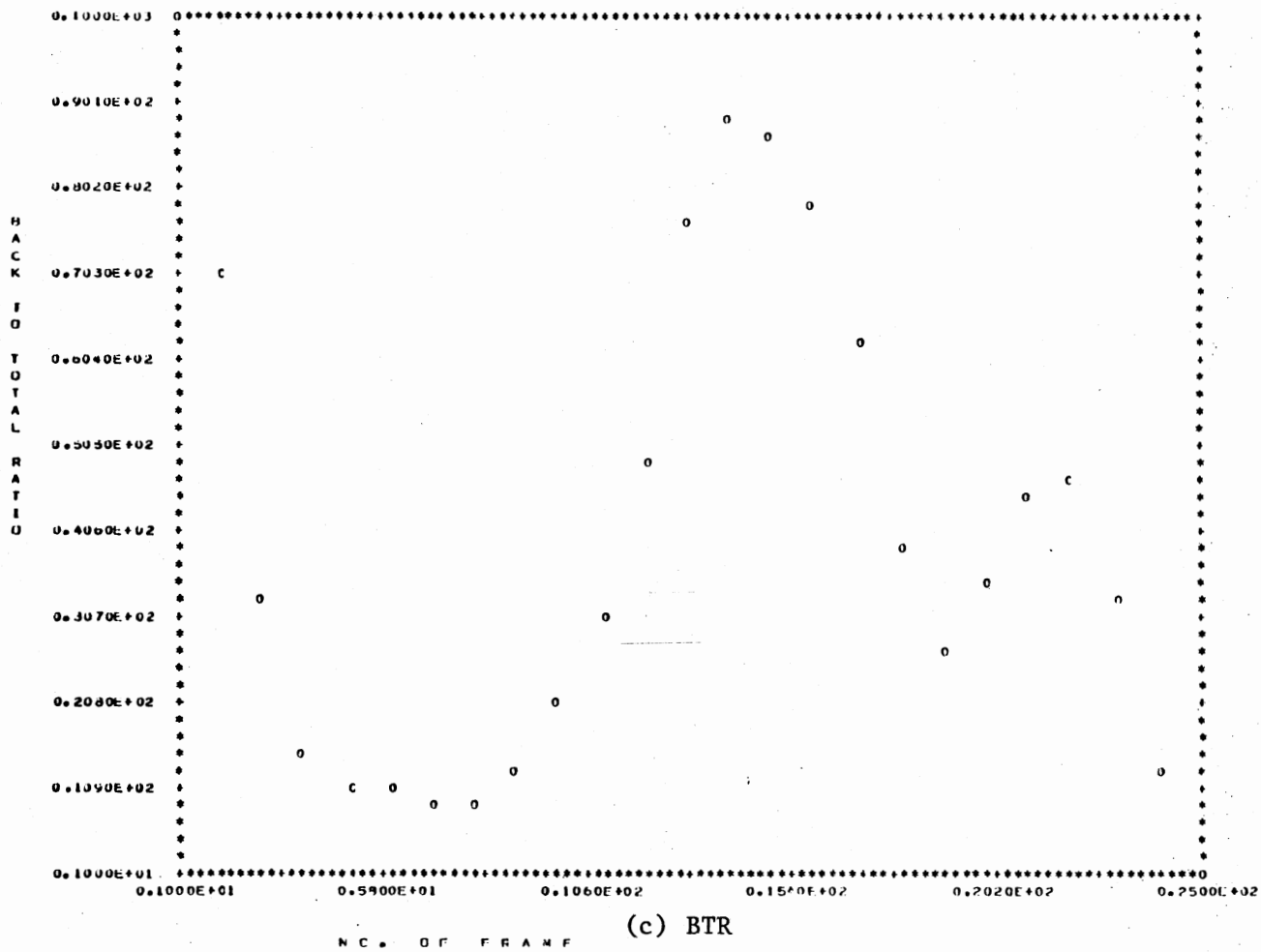


Figure 51. (Continued)

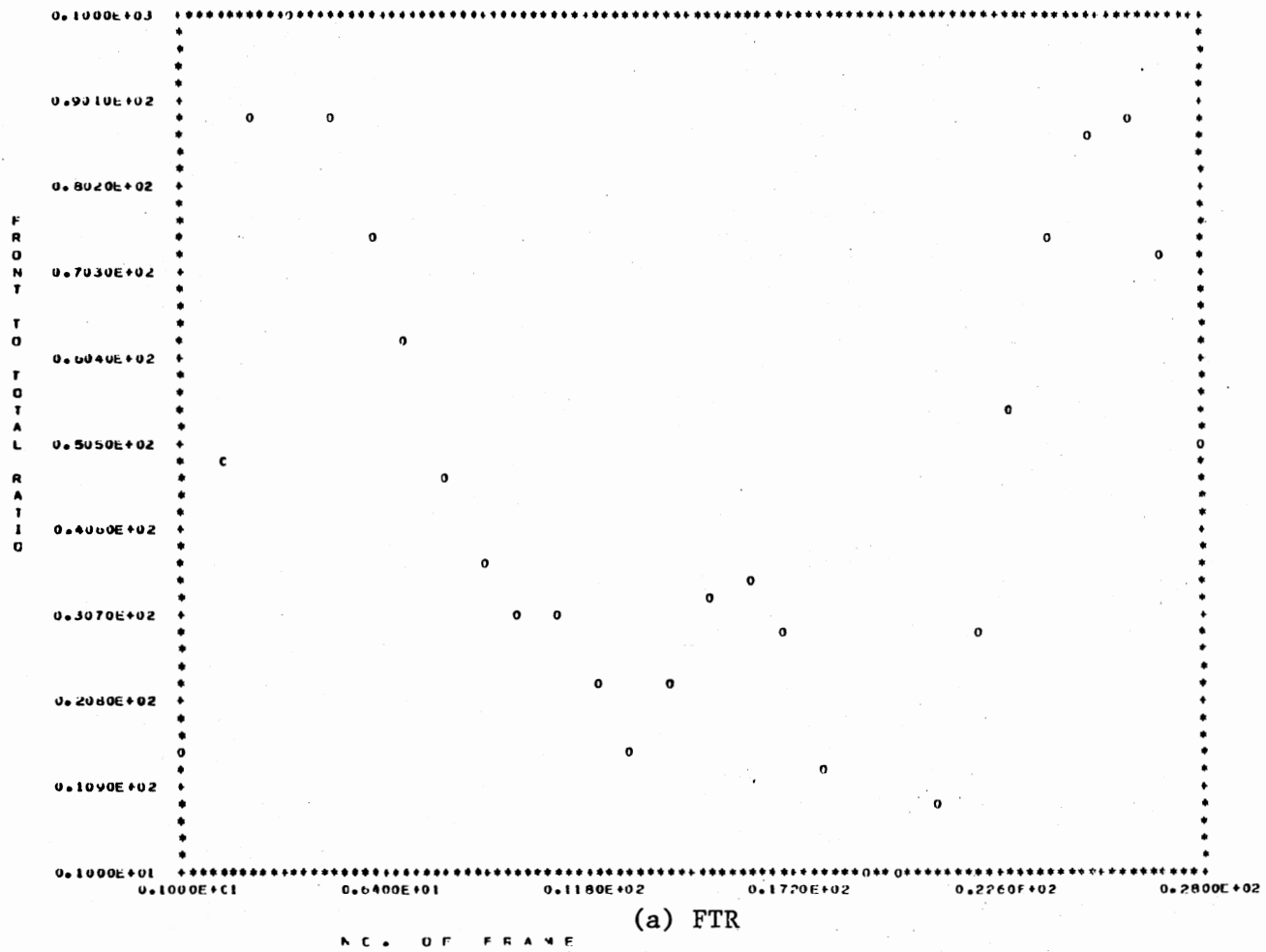


Figure 52. Smoothed and Quantized Feature Parameters for Digit Three Spoken in American English

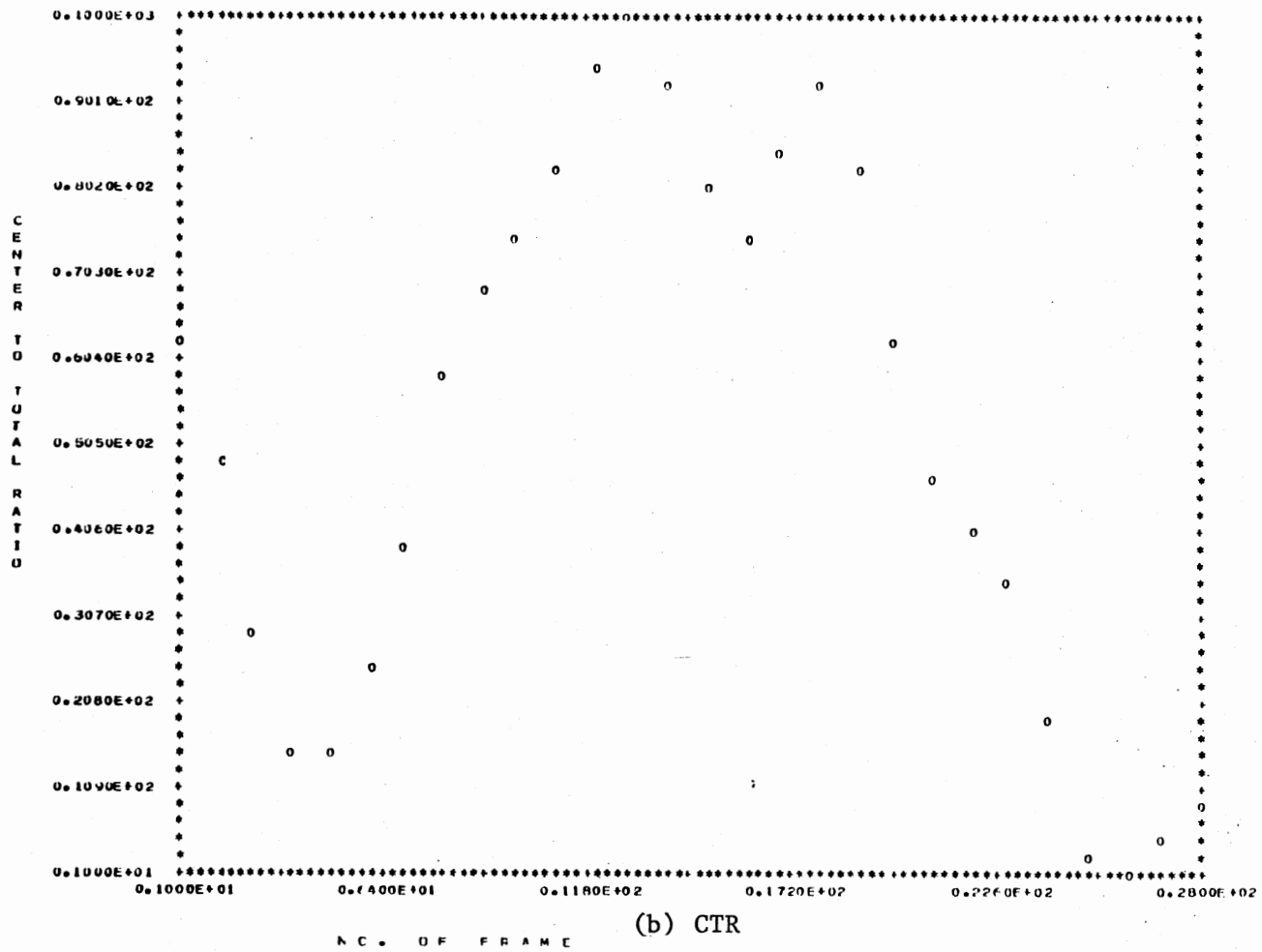


Figure 52. (Continued)

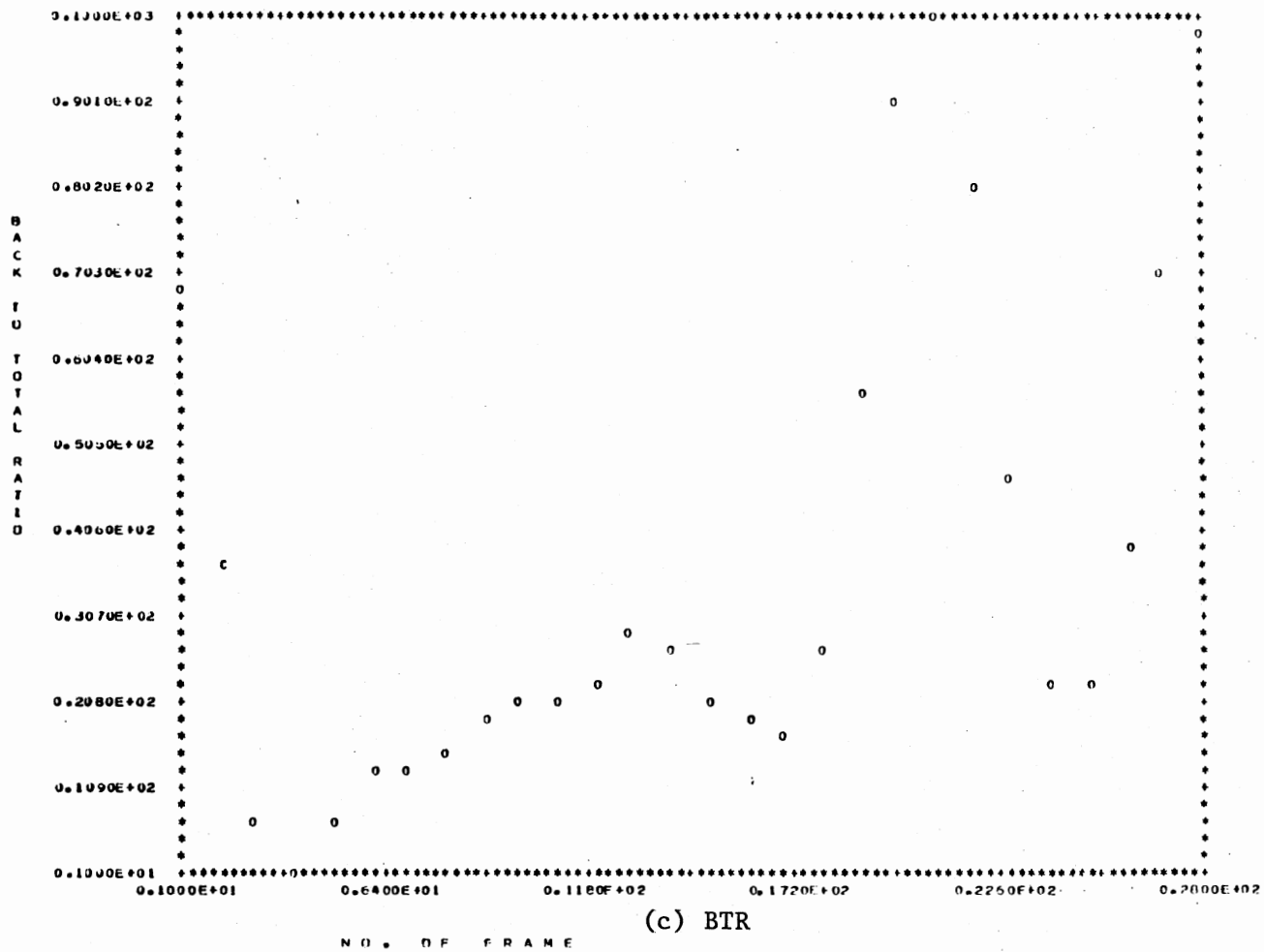


Figure 52. (Continued)

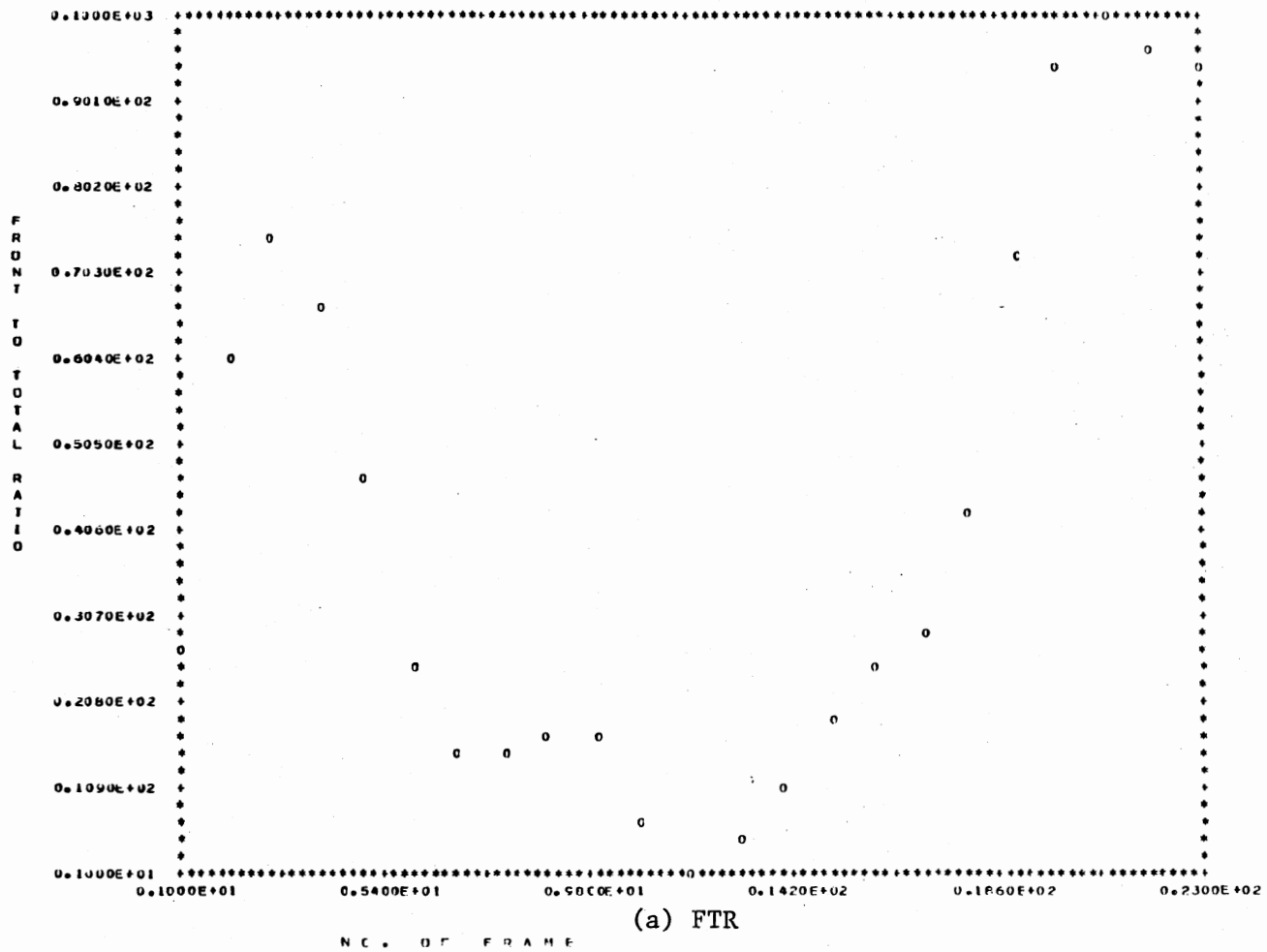


Figure 53. Smoothed and Quantized Feature Parameters for Digit Four Spoken in American English

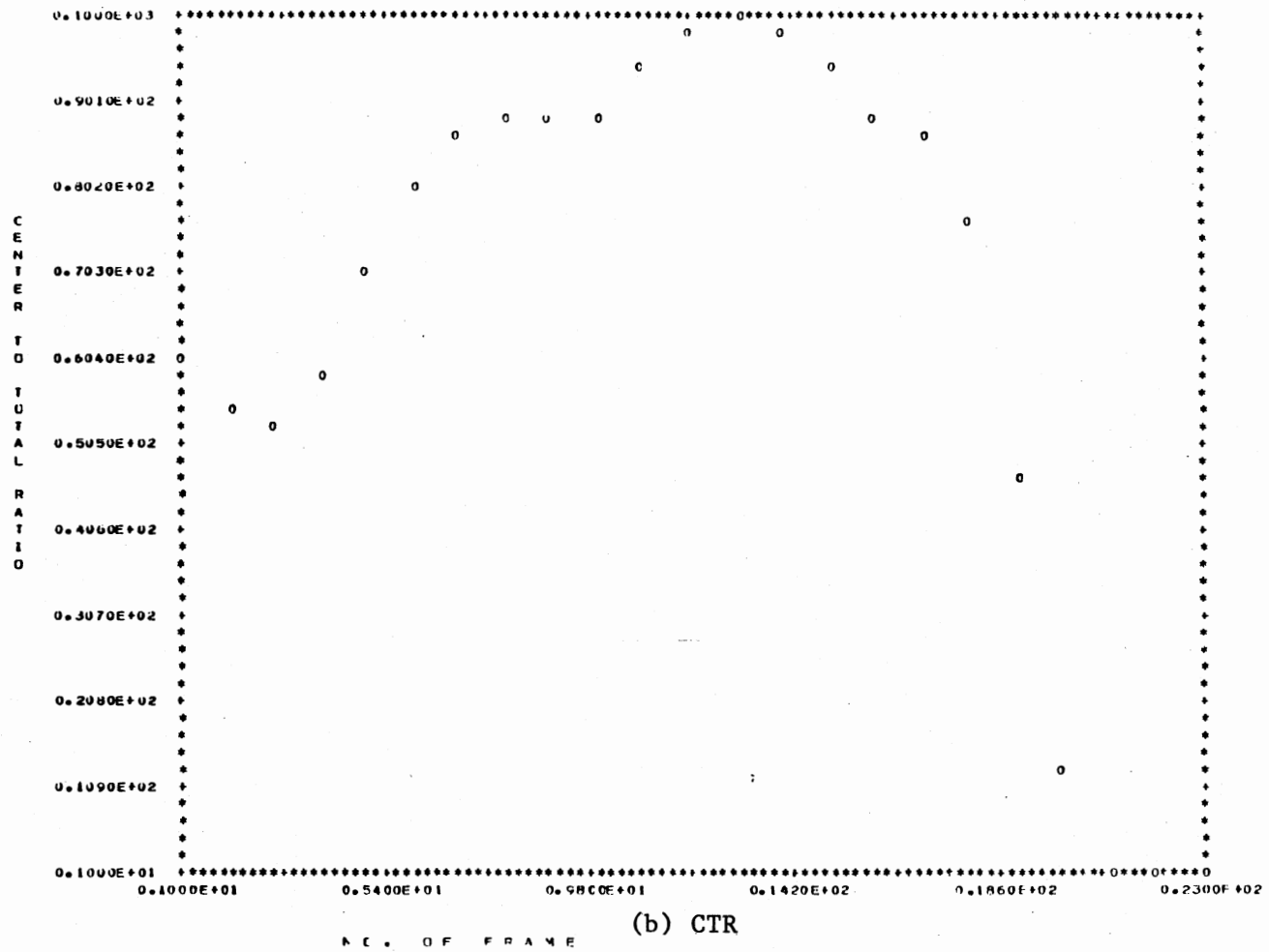


Figure 53. (Continued)

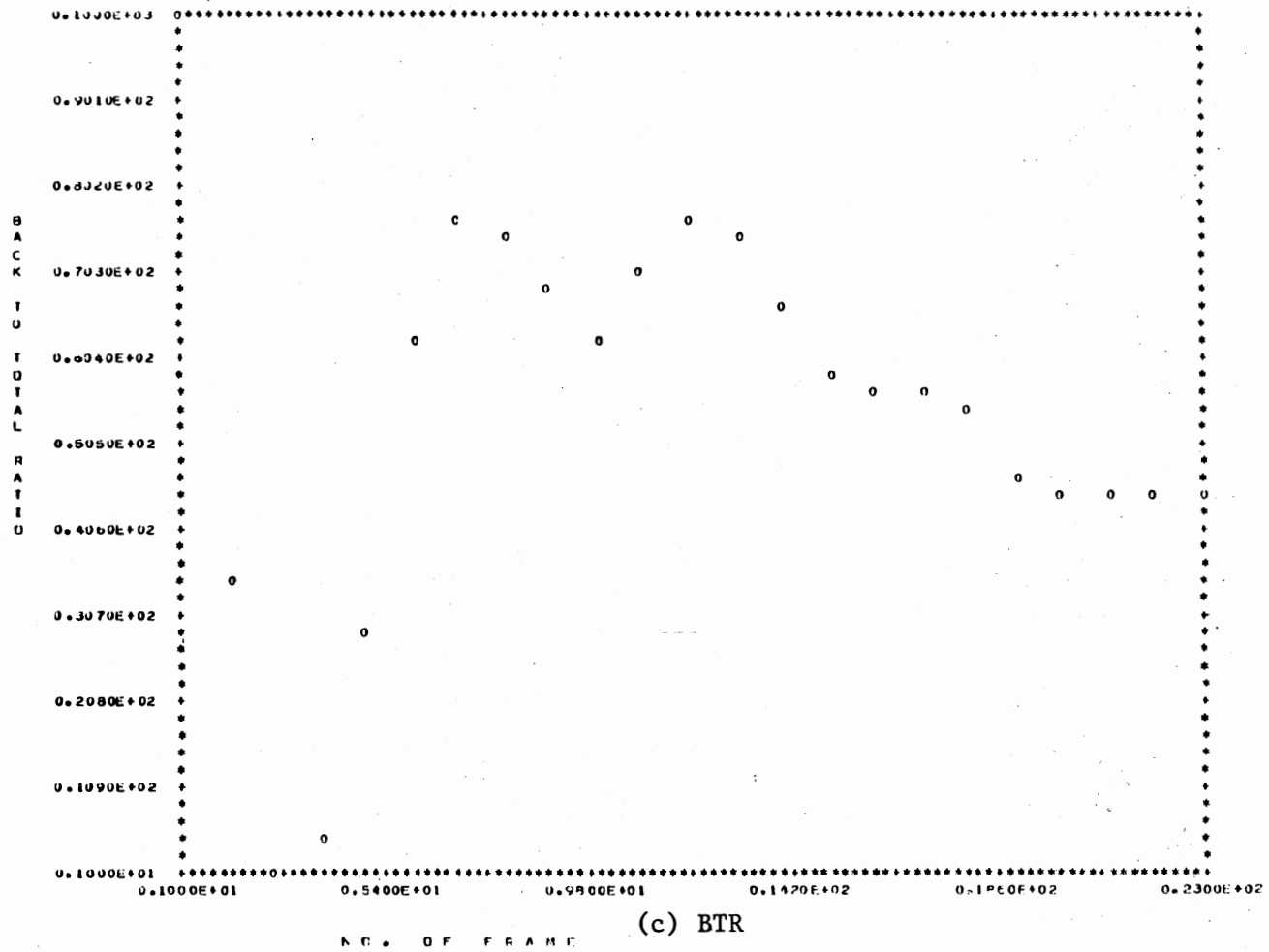


Figure 53. (Continued)

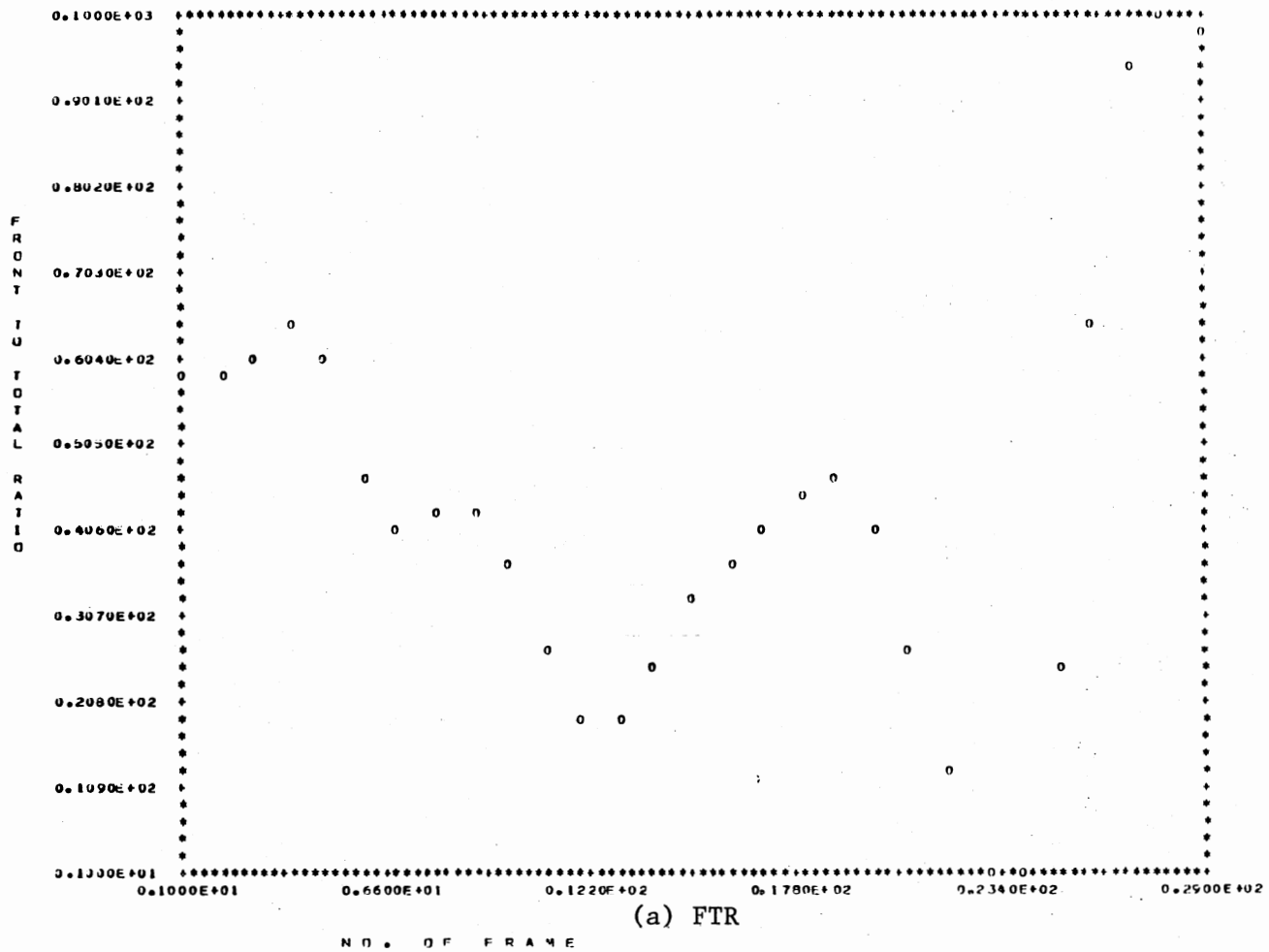


Figure 54. Smoothed and Quantized Feature Parameters for Digit Five Spoken in American English

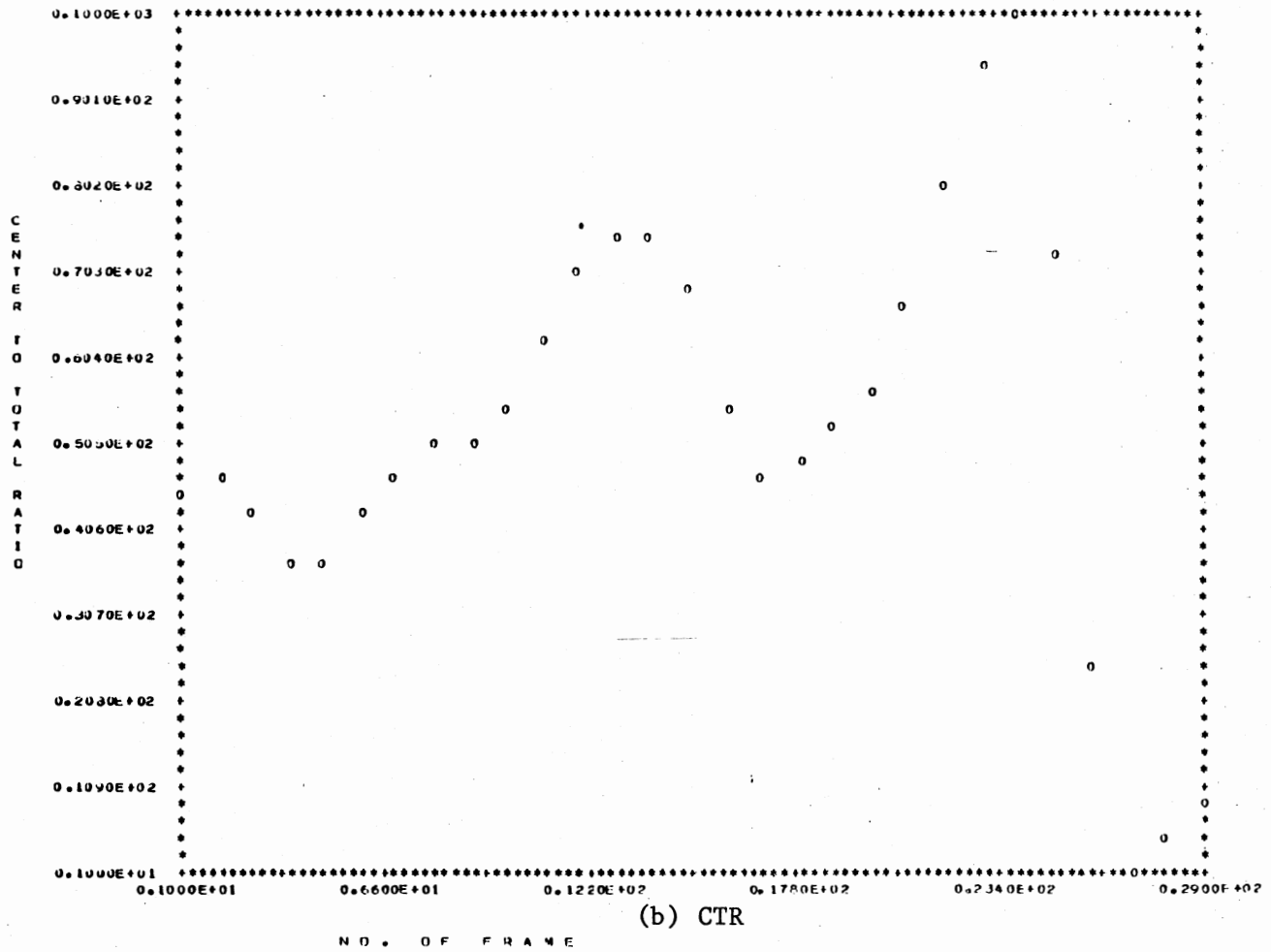


Figure 54. (Continued)

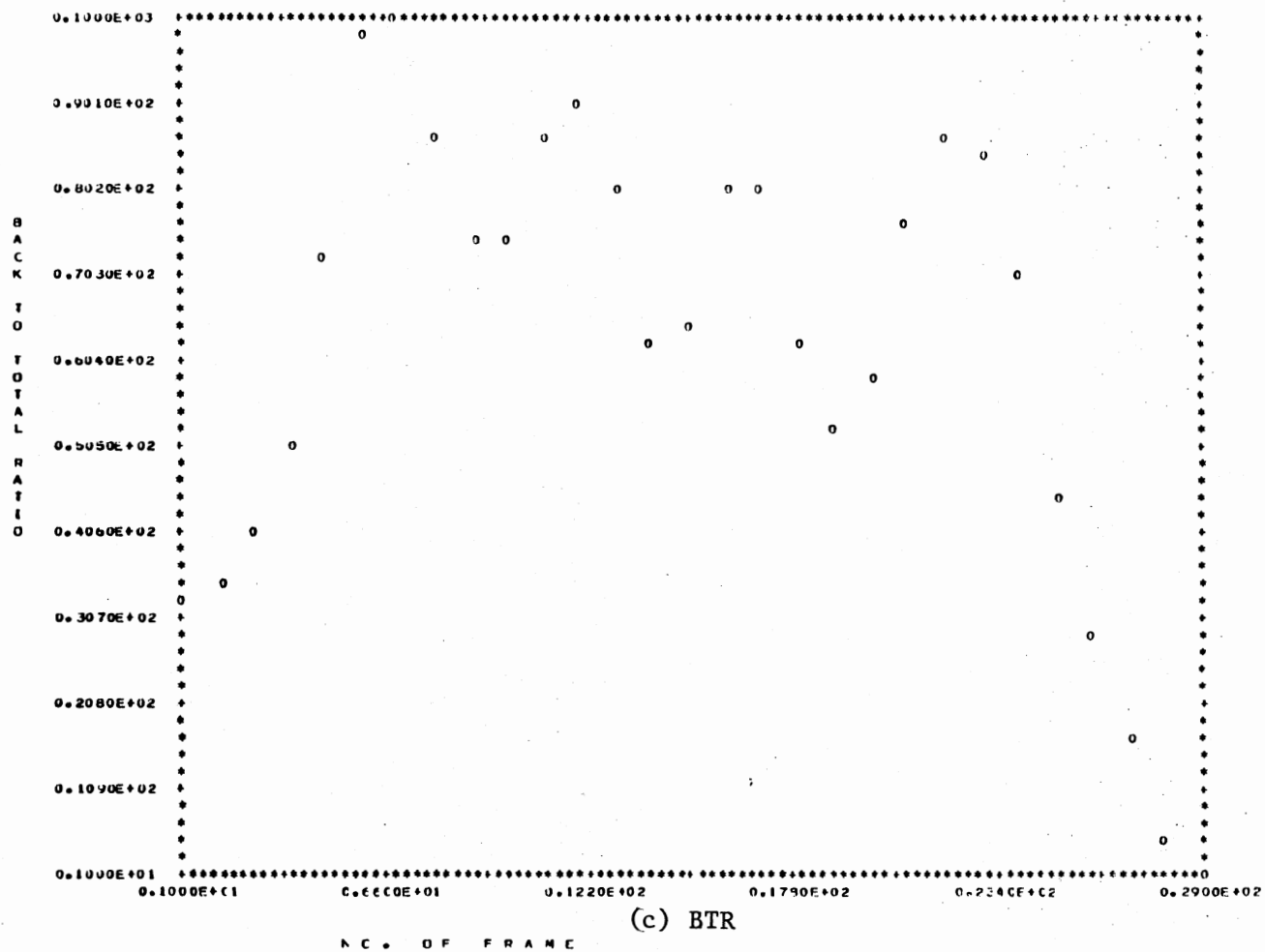


Figure 54. (Continued)

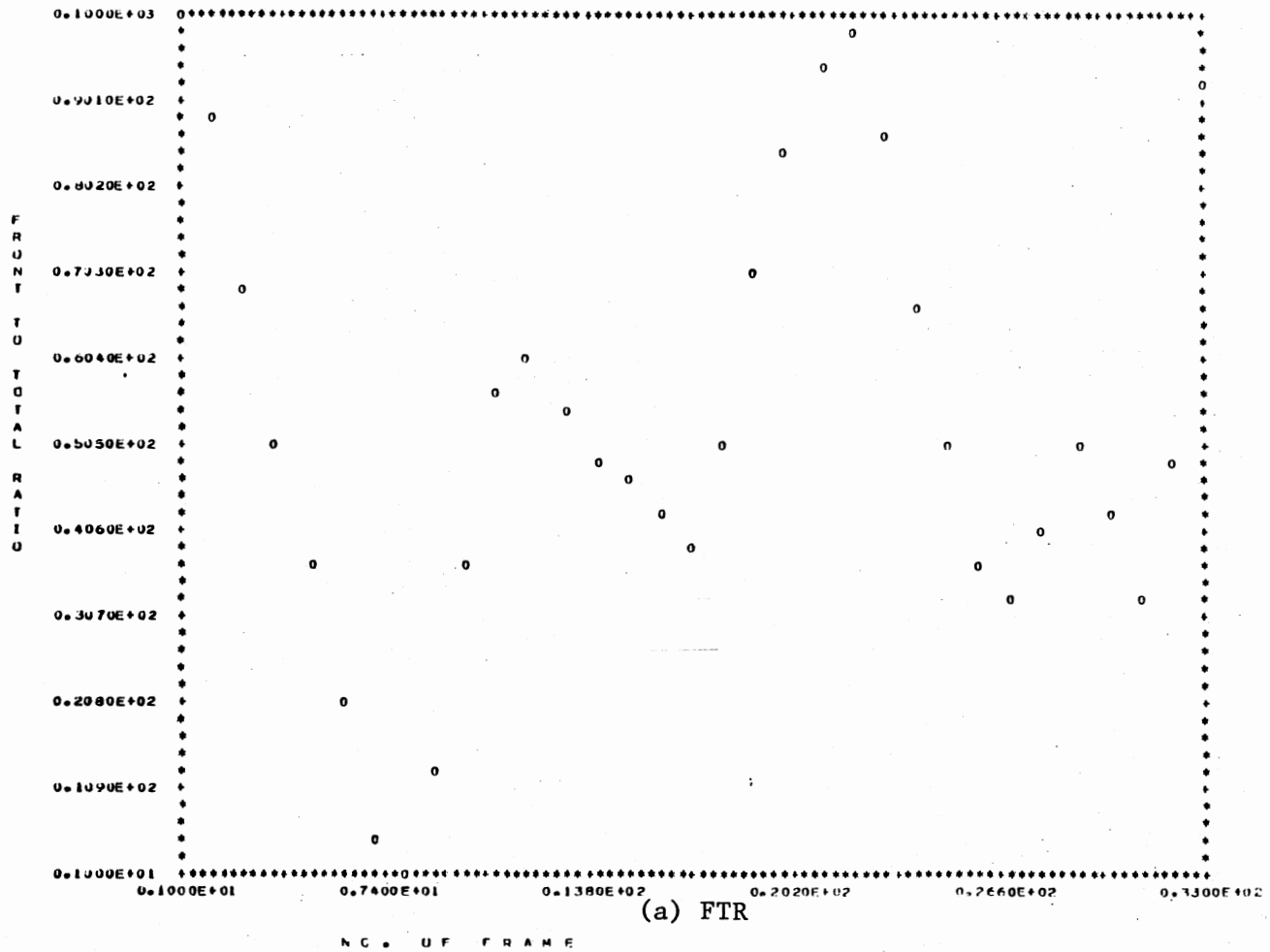


Figure 55. Smoothed and Quantized Feature Parameters for Digit Six Spoken in American English

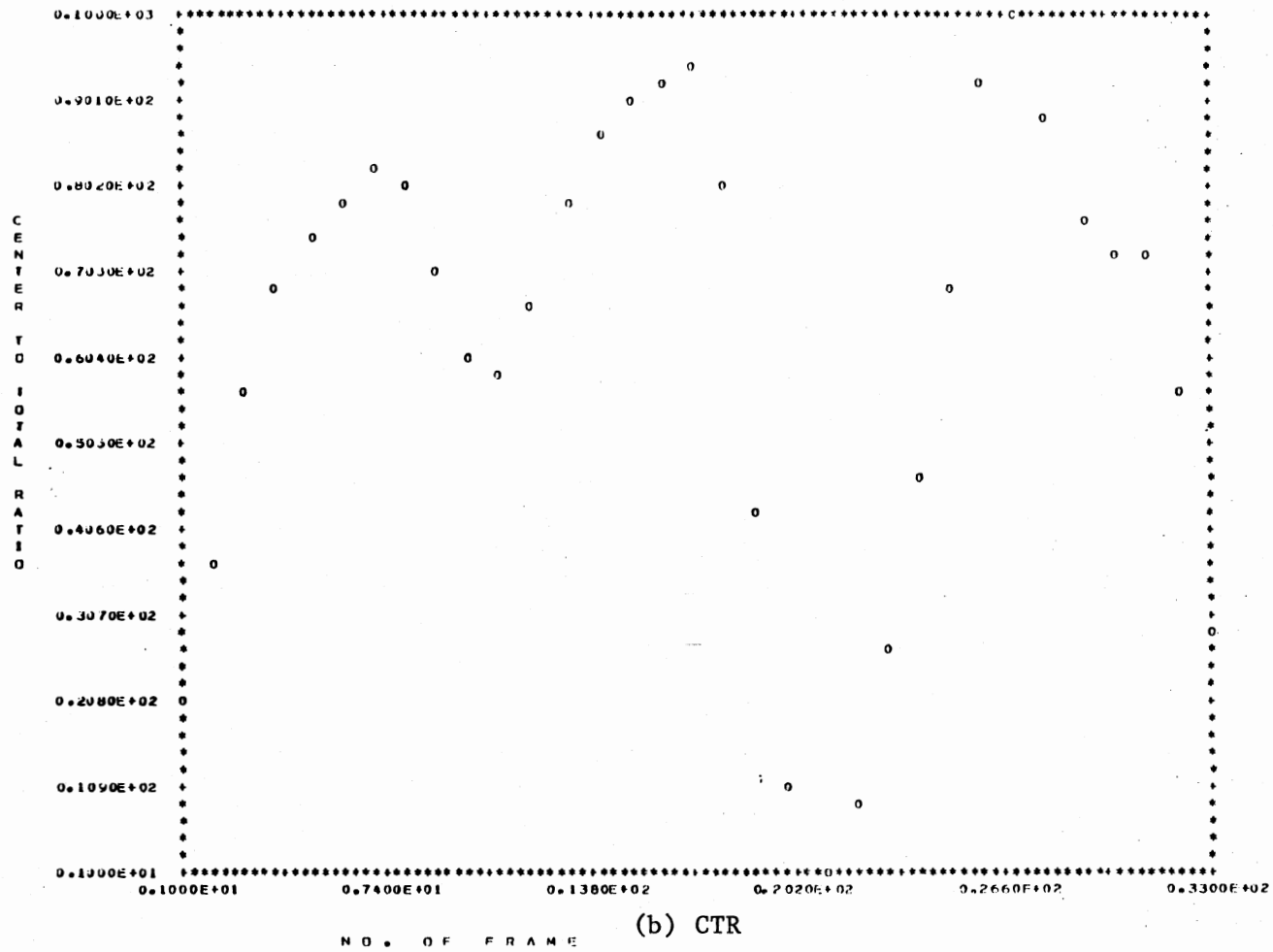


Figure 55. (Continued)

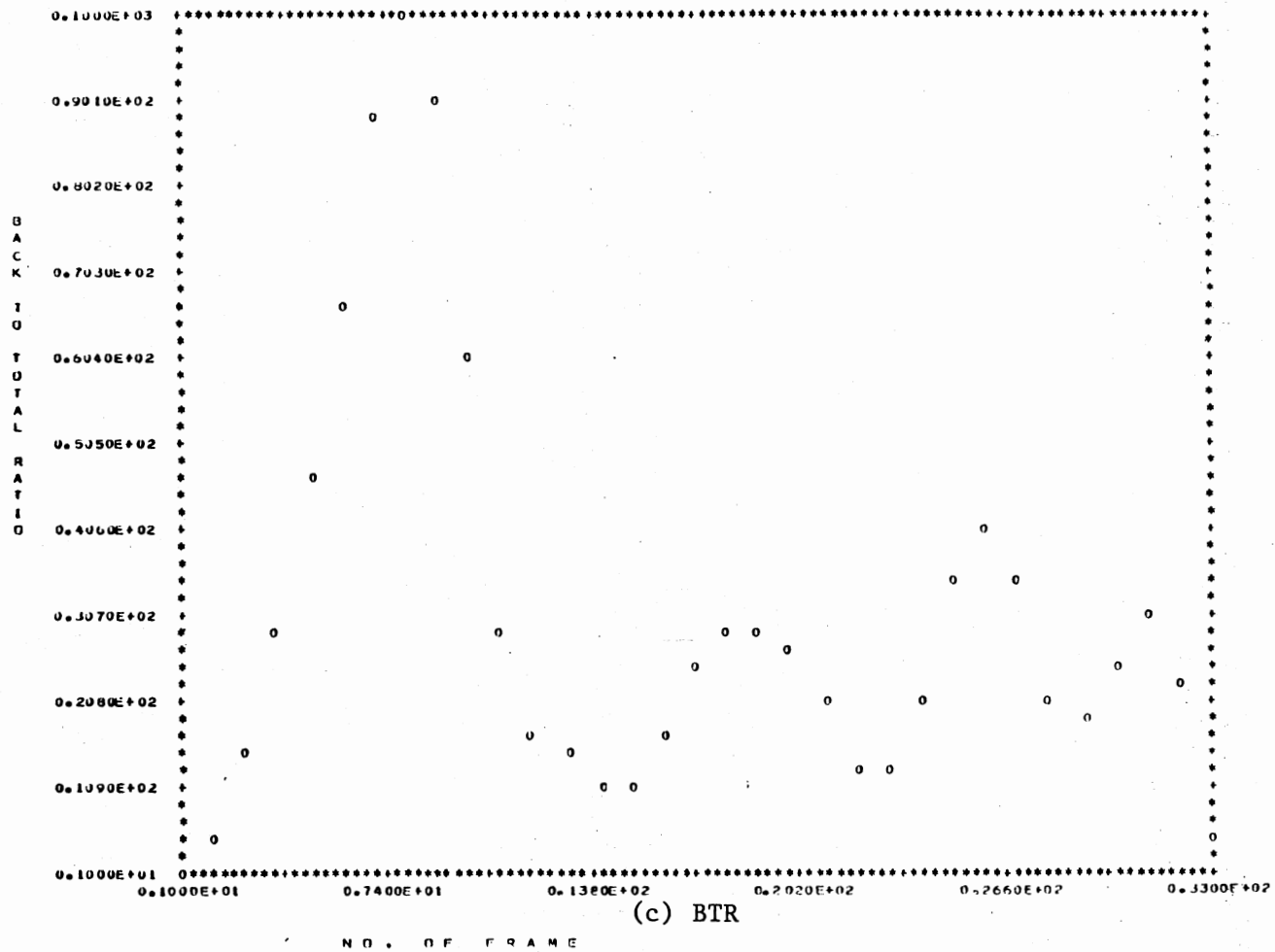


Figure 55. (Continued)

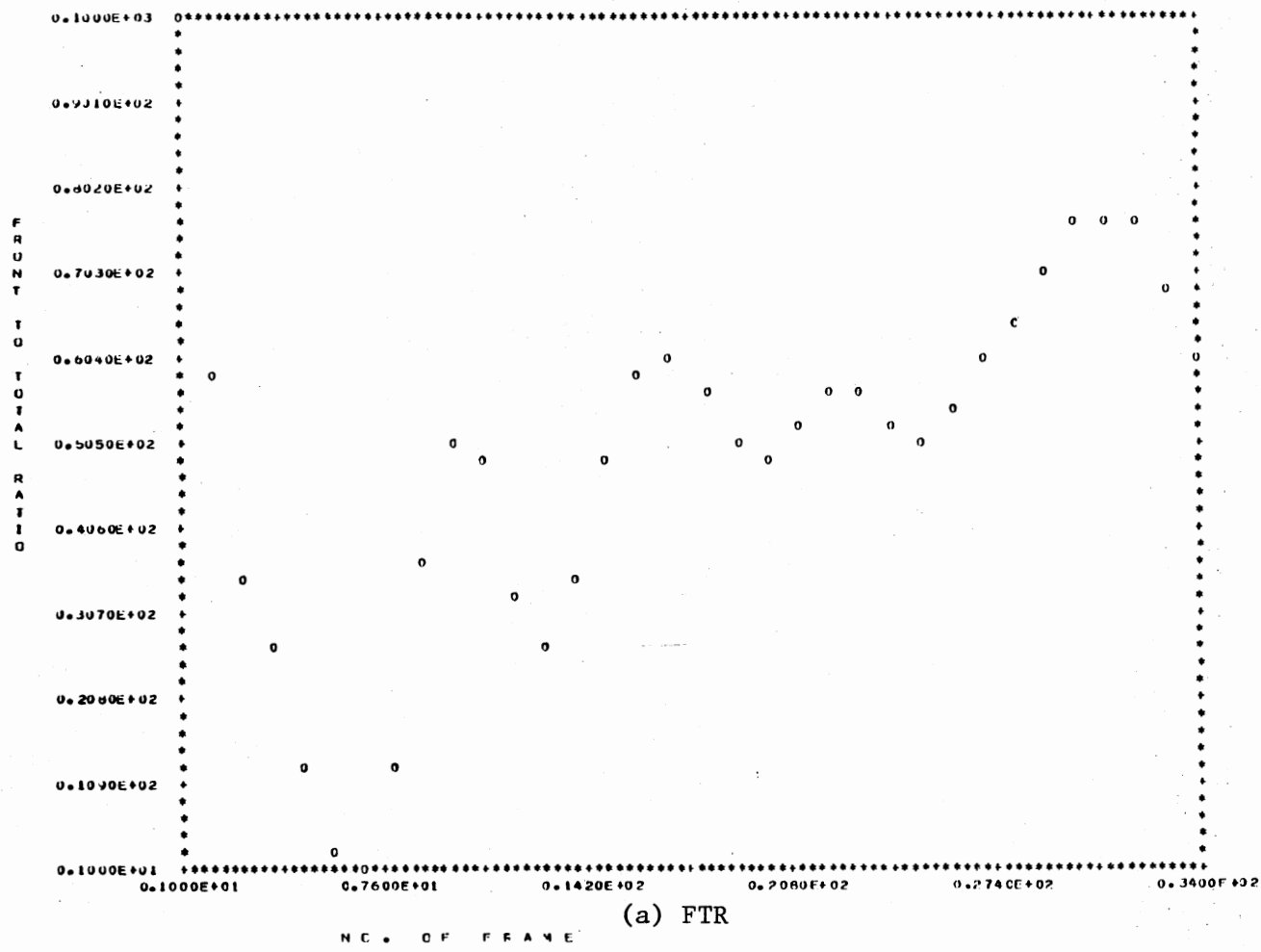


Figure 56. Smoothed and Quantized Feature Parameters for Digit Seven Spoken in American English

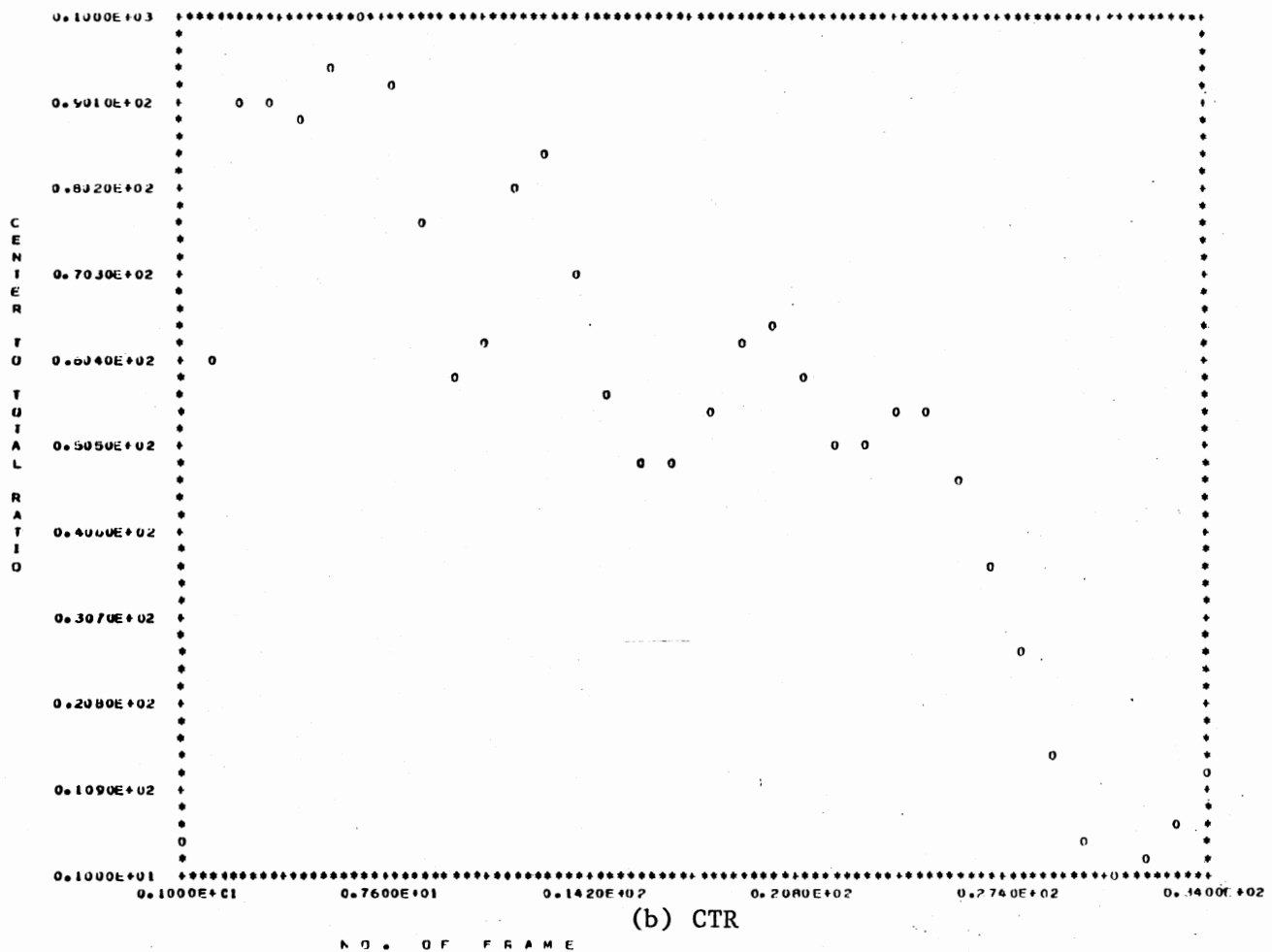


Figure 56. (Continued)

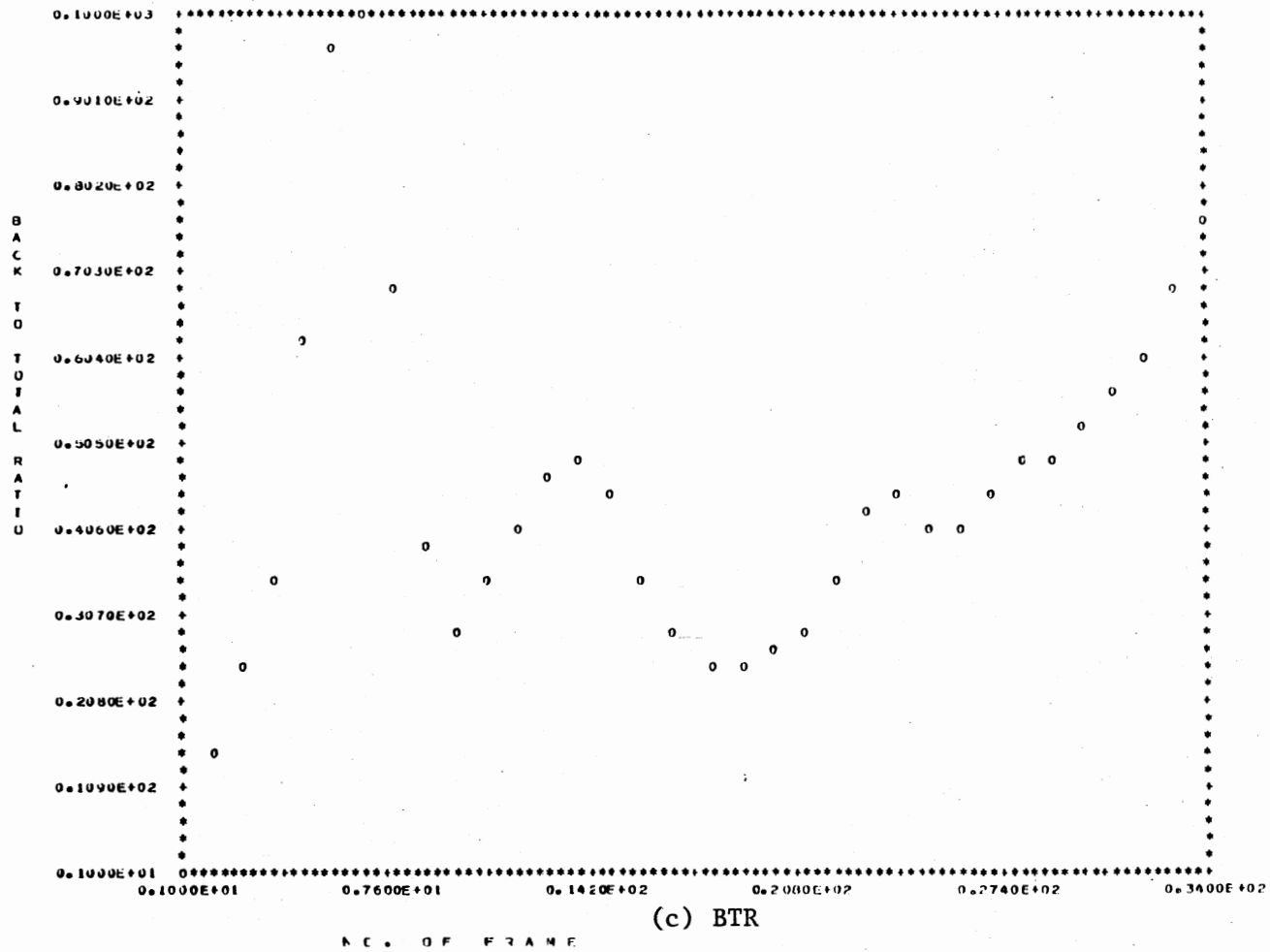


Figure 56. (Continued)

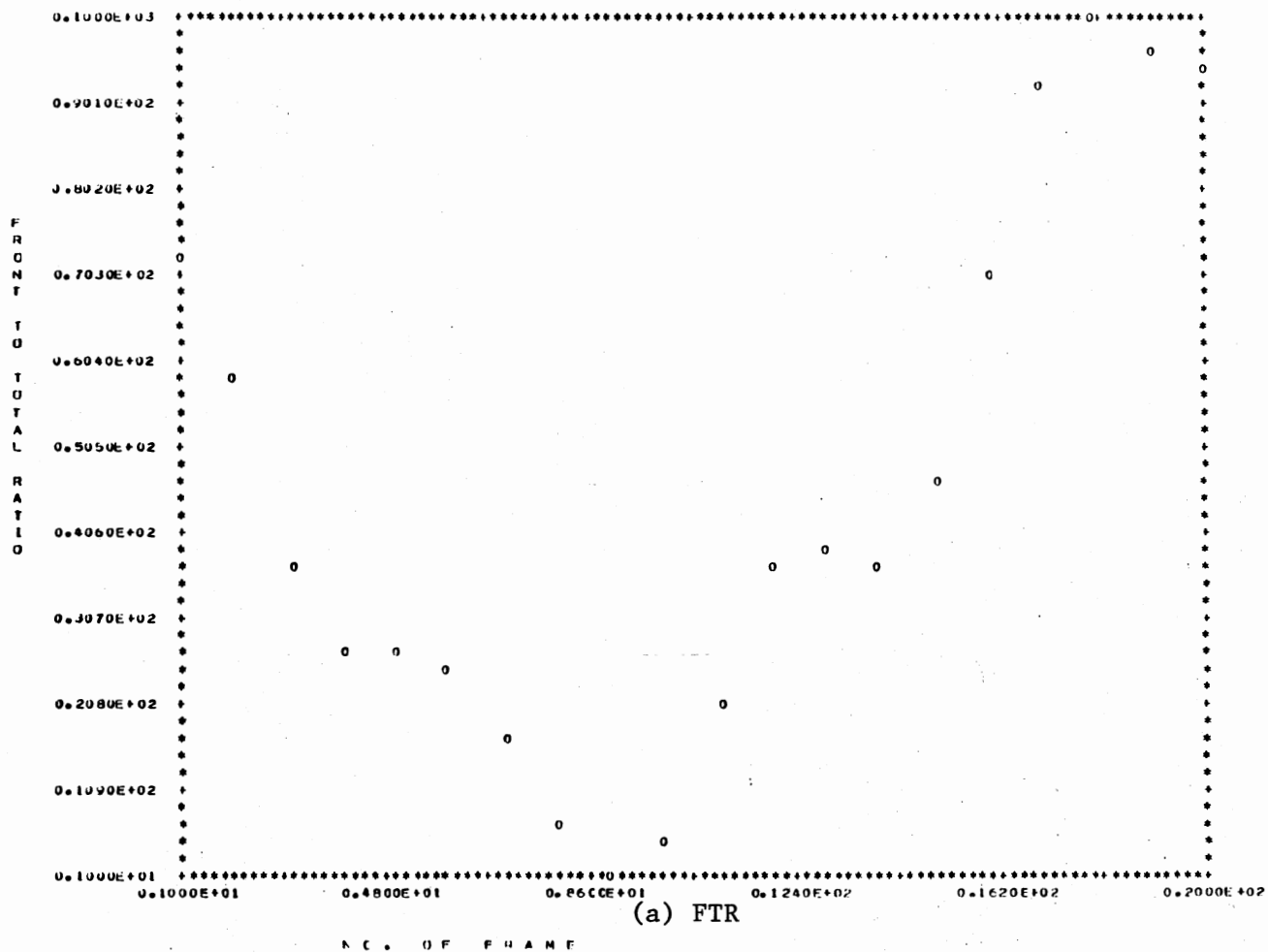


Figure 57. Smoothed and Quantized Feature Parameters for Digit Eight Spoken in American English

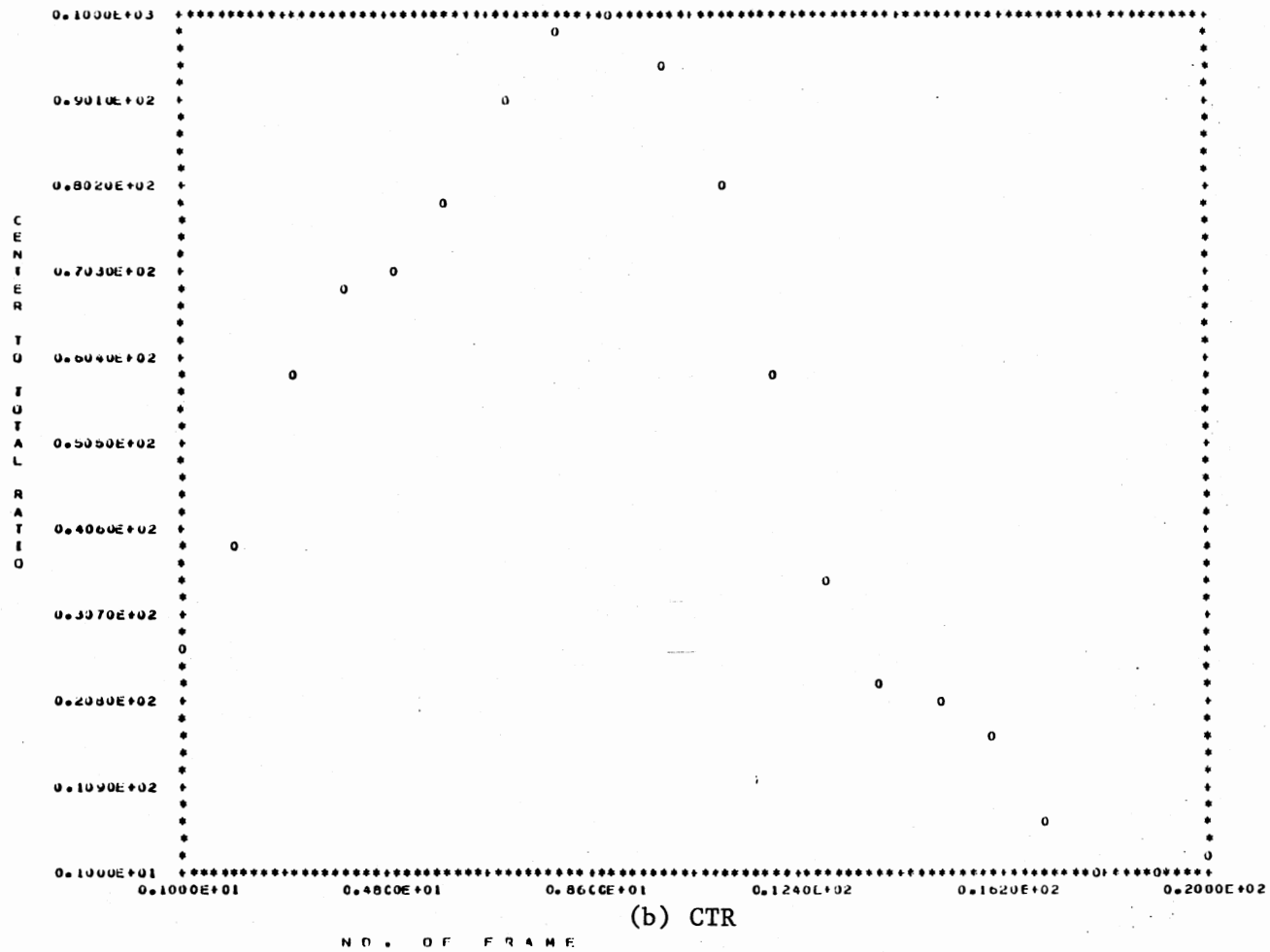


Figure 57. (Continued)

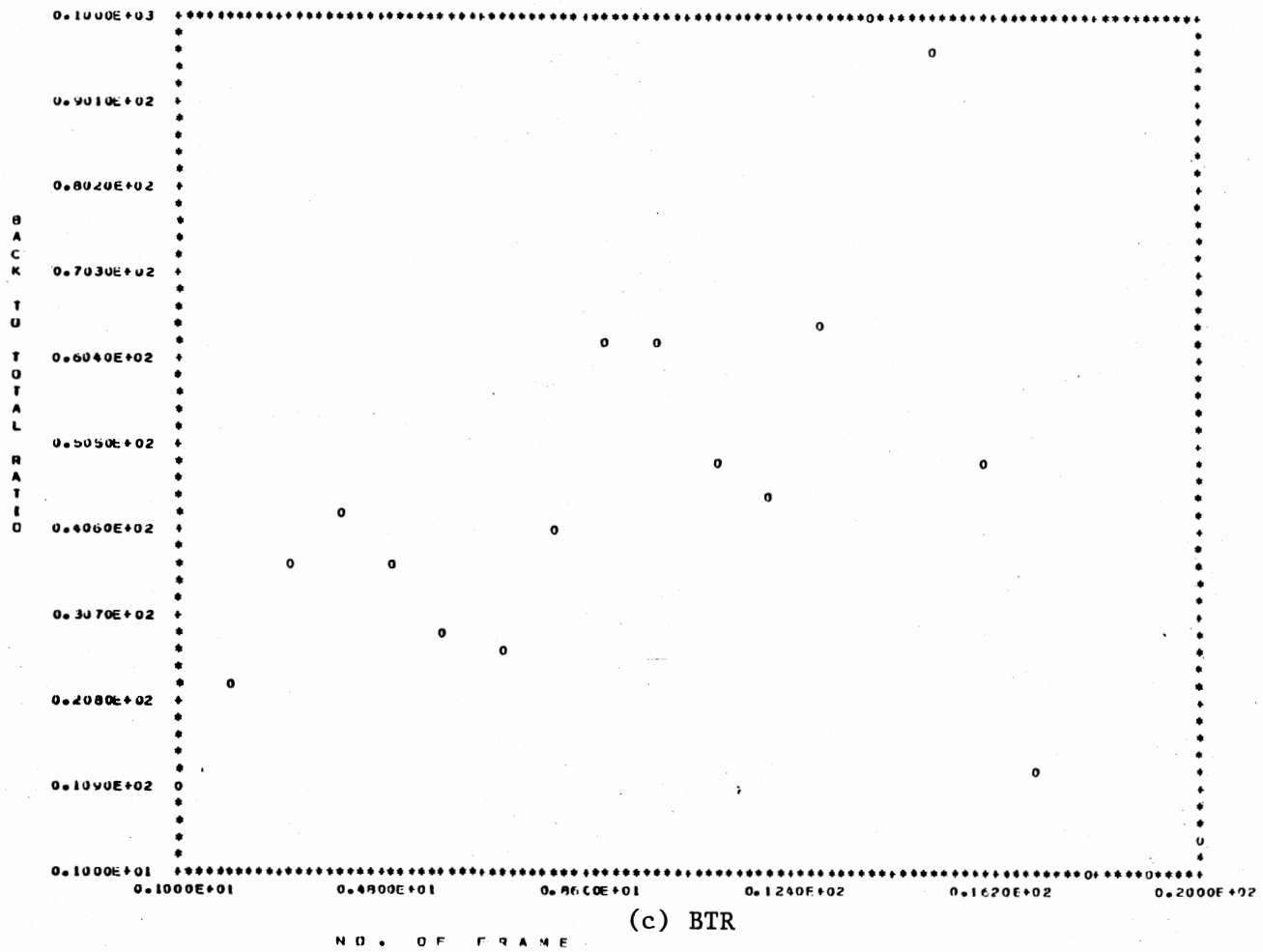


Figure 57. (Continued)

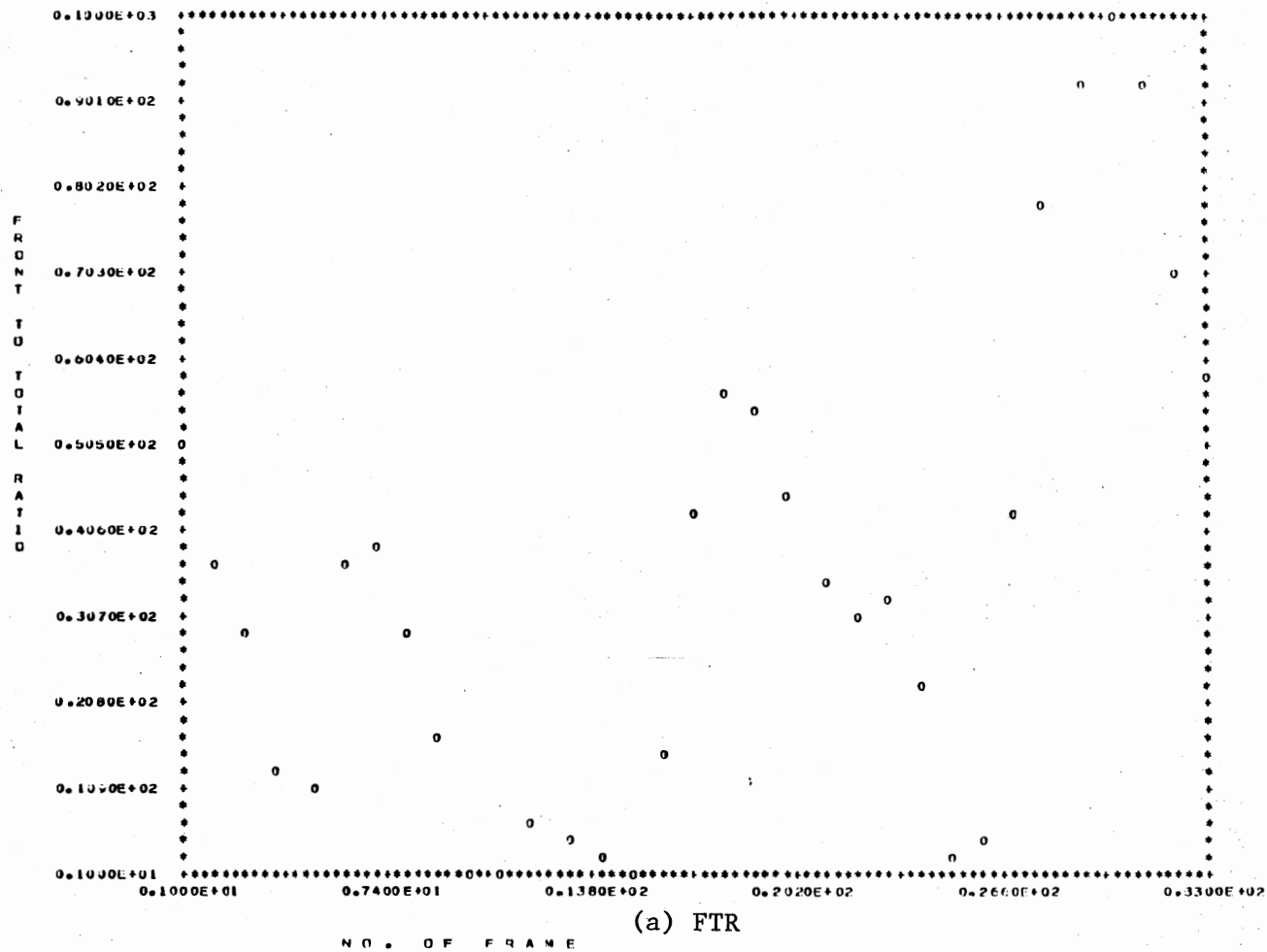


Figure 58. Smoothed and Quantized Feature Parameters for Digit Nine Spoken in American English

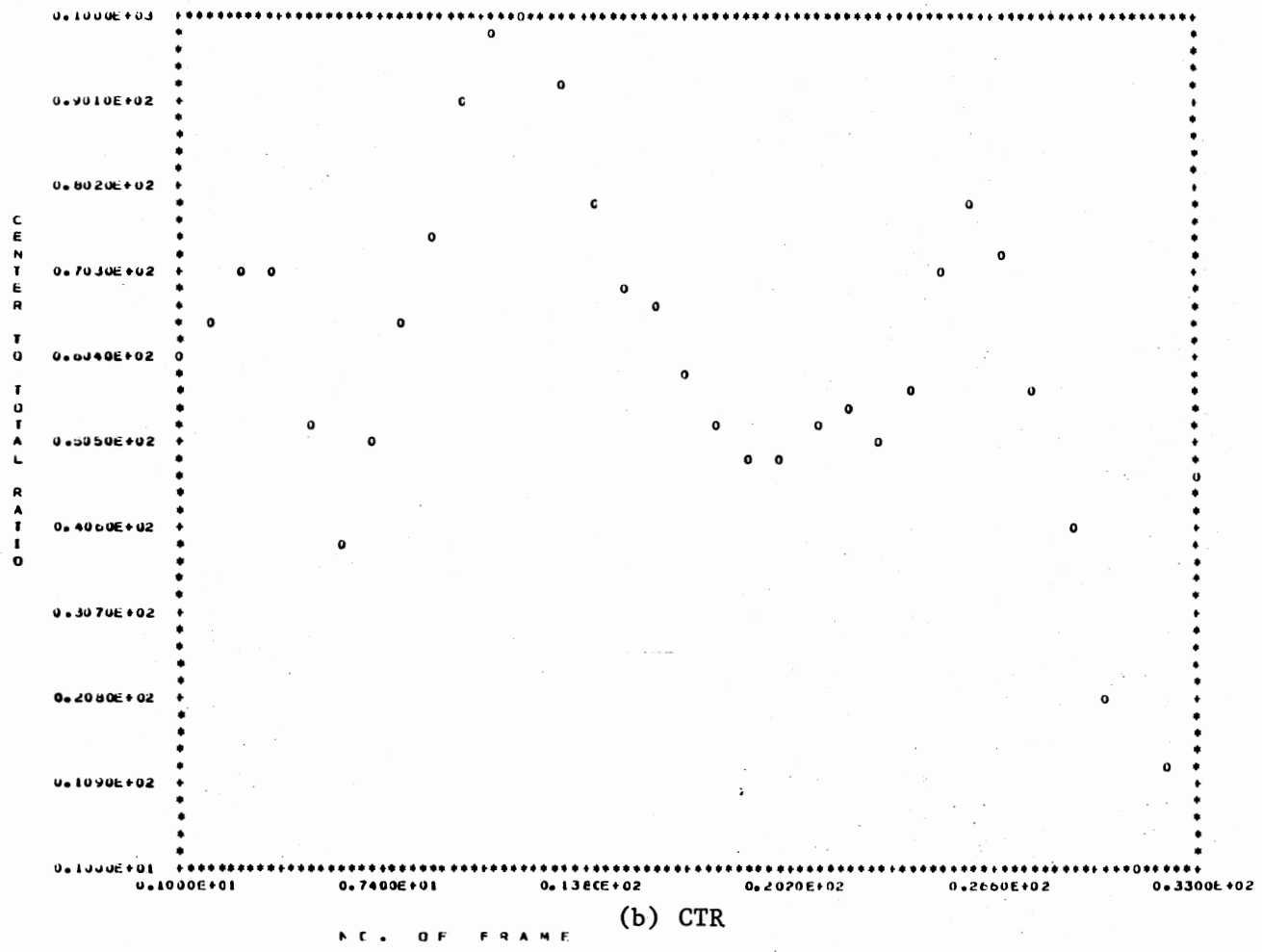


Figure 58. (Continued)

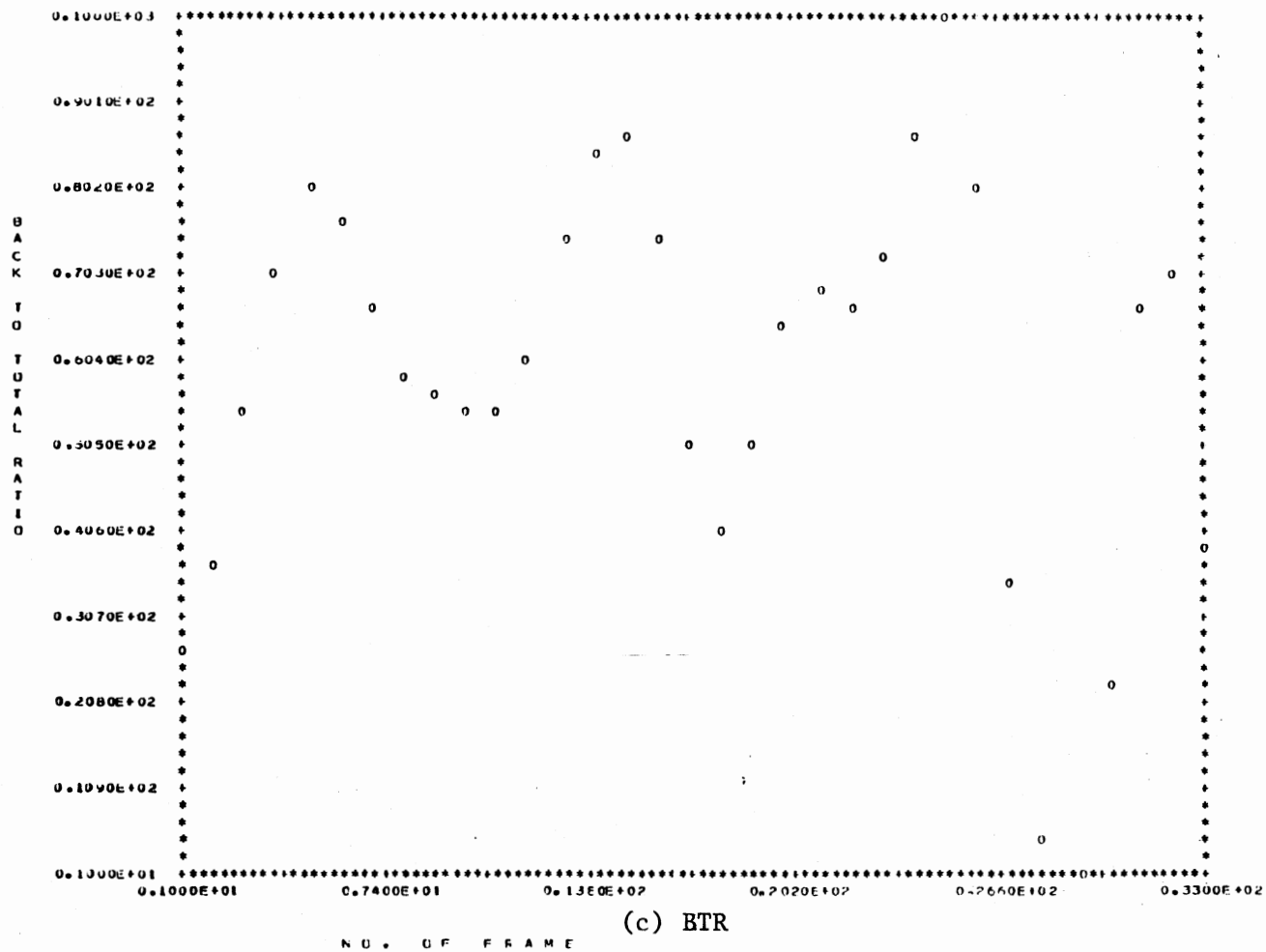


Figure 58. (Continued)

schemes. In addition, the BTR is used to separate nasals from vowel segments, and the SFBR is used to separate turbulence noise segments from back vowel segments. The classification of each segment is stored for further usage.

The third stage in Figure 28 corresponds to the actual digit recognition. The discussion on peak ratios was presented earlier. A summary of the beginning, middle and ending sound classification of digits spoken in American English is tabulated (Table VII). This Table is derived from Table II. The uniqueness aspects of the classifications are used in the final decision. The decision algorithm is based upon using Table VII and sound classification at the beginning and the end region of a particular digit. If a decision cannot be made, then the middle region is referred for identification. For example, the digits zero, six, and eight have unique classification for beginning and ending sounds. These three digits are decided on the beginning and ending sound. It is interesting to point out that the digit seven classified as having two front vowels in the middle region. The digits one and nine have the same sound patterns for end regions. However, the middle region is different. In a similar manner, the sound classification of the remaining digits can be discussed.

Figure 59 shows the flow chart for the digit recognition as based upon the uniqueness aspects of Table VII. The beginning region is tested first for "VL" indication to separate zero and one and nine from the remaining digits. If the beginning sound is vowel-like, then the middle region is tested. A middle vowel to front vowel indication is made to distinguish between one and nine. Middle

TABLE VII

SEQUENCE OF SOUND CLASS REGIONS OF DIGITS SPOKEN IN AMERICAN ENGLISH

Digit	RMS and BTR Decision		
	Beginning	Middle	End
/zIro/	VL	FV/BV	V
/wAn/	VL	MV	VL
/tu/	NV	FV	V
/θri/	NV	VL	V
/fɔr/	NV	BV/MV	VL
/faIv/	NV	MV/FV	VL
/sIks/	NV	FV	NV
/seven/	NV	FV/FV	VL
/eIt/	V		NV
/naIn/	VL	MV/FV	VL

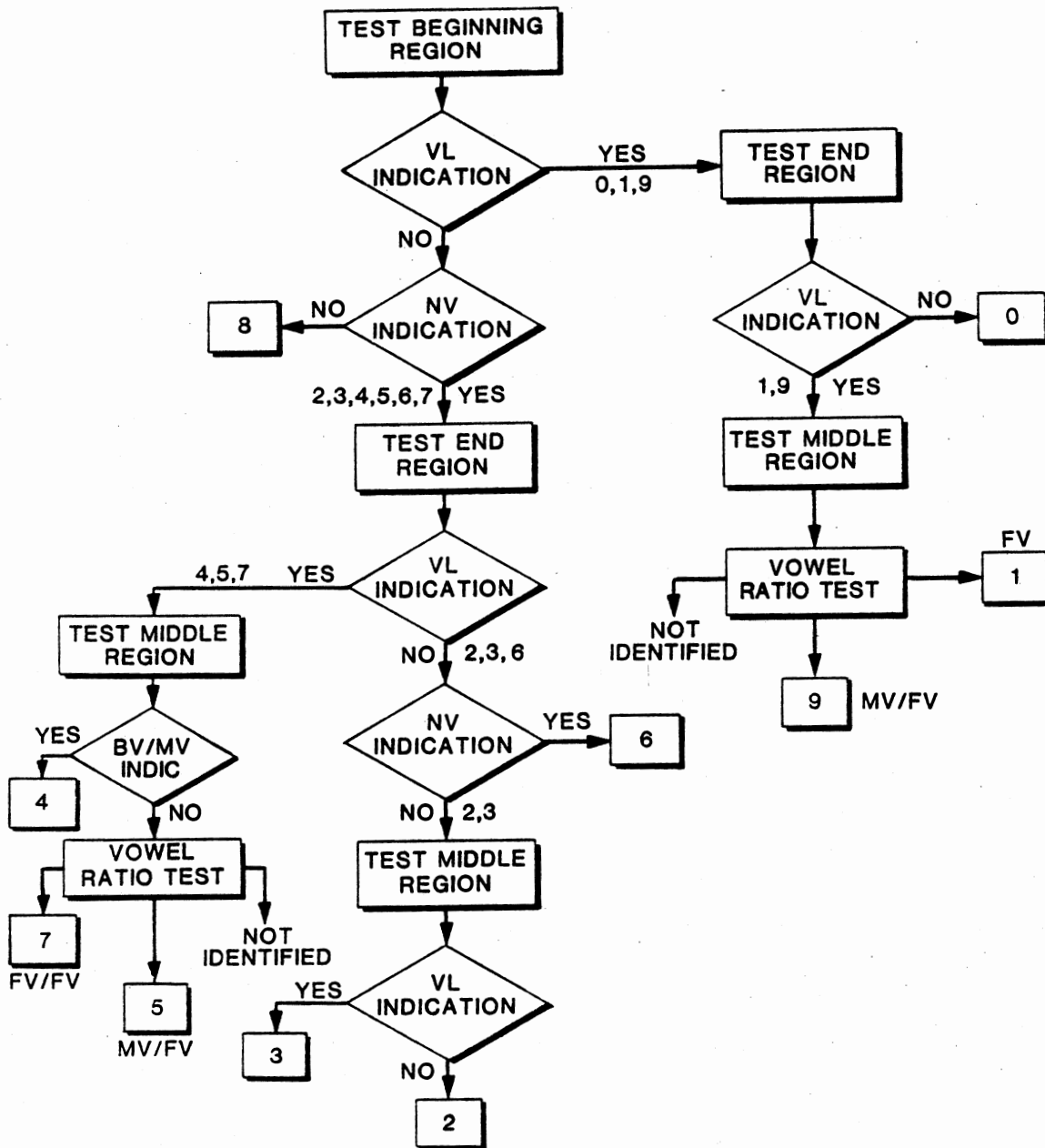


Figure 59. Decision Tree for Digit Identification Scheme for American English Language

vowel to front vowel indicated that digit nine is spoken. A front vowel indicates that digit one was spoken. For the other digits, the beginning region is tested for non-vowel "NV" segments. The remaining flow chart can be identified in a similar manner.

The above algorithm was tested for seven speakers both male and female. The accuracy rate was about 75 percent. Some of the speakers were international students. The results were much better for Americans and are accurate to about 95 percent, especially when the spoken digits were intact. By this method the digits two and four have the major problem in recognition. The reason for this is that both of these digits have only one essential vowel in the center region giving almost similar energy contour for most speakers. The digit four has a back vowel in the center region followed by mid-vowel in the end region since the semivowel /r/ is not pronounced distinctly by the Americans and is stressed by internationals. If both paths in Figure 28 give the same digit result, then the recognition is completed. If the two paths give different results, or that no decision has been made, then the following procedure is used.

Pattern Recognition Scheme

The final step in the recognition scheme is based on the cross-correlation coefficients in Equation (43) in Chapter III. The correlation of digits i and j are considered for the digits zero through nine. For the seven speakers, correlation tables are constructed. Some are given in Tables VIII - XI. Since the seven sets of digits are spoken by seven people having different dialects and accents, the

TABLE VIII
 RMS ENERGY CORRELATION TABLE FOR DIGITS SPOKEN IN AMERICAN ENGLISH
 (MAXIMUM NUMBER OF FRAMES USED)

SET NUMBER :		RMS CORRELATION								
ZERO	ONE	TWO	THREE	FOUR	FIVE	SIX	SEVEN	EIGHT	NINE	
1.000	0.206	0.389	0.003	0.224	0.033	0.724	-0.004	-0.085	-0.016	
0.206	1.000	0.280	0.186	0.464	0.003	0.391	0.017	-0.116	0.410	
0.389	0.280	1.000	0.151	0.693	0.292	0.614	0.395	0.622	0.605	
0.003	0.186	0.151	1.000	0.253	0.039	-0.057	0.115	0.147	0.076	
0.224	0.464	0.693	0.253	1.000	0.147	0.677	0.679	0.212	0.800	
0.033	0.003	0.292	0.039	0.147	1.000	0.156	-0.057	0.249	0.447	
0.724	0.391	0.614	-0.057	0.677	0.156	1.000	0.384	0.062	0.444	
-0.004	0.017	0.395	0.115	0.679	0.057	0.384	1.000	0.239	0.617	
-0.085	-0.116	0.622	0.147	0.212	0.249	0.062	0.239	1.000	0.287	
-0.016	0.410	0.605	0.076	0.800	0.447	0.444	0.617	0.287	1.000	

TABLE IX

RMS ENERGY CORRELATION TABLE FOR DIGITS IN AMERICAN ENGLISH
(MAXIMUM NUMBER OF FRAMES USED)

SET NUMBER :		5							
		RMS CORRELATION							
ZERO	ONE	TWO	THREE	FOUR	FIVE	SIX	SEVEN	EIGHT	NINE
1.000	-0.095	0.131	0.786	0.400	0.569	0.559	0.745	0.573	0.599
-0.095	1.000	0.650	0.047	0.428	-0.155	0.553	0.173	0.493	0.651
0.131	0.650	1.000	0.044	0.584	0.147	0.690	0.280	0.776	0.573
0.786	0.047	0.044	1.000	0.379	0.762	0.394	0.581	0.410	0.715
0.400	0.428	0.584	0.379	1.000	0.414	0.225	0.174	0.481	0.613
0.569	-0.155	-0.147	0.762	0.414	1.000	0.091	0.419	0.092	0.457
0.559	0.553	0.690	0.394	0.225	0.091	1.000	0.598	0.747	0.665
0.745	0.173	0.280	0.581	0.174	0.419	0.598	1.000	0.300	0.408
0.573	0.493	0.776	0.410	0.491	0.092	0.747	0.300	1.000	0.635
0.599	0.651	0.573	0.715	0.613	0.457	0.665	0.408	0.635	1.000

TABLE X

BTR CORRELATION TABLE FOR DIGITS SPOKEN IN AMERICAN ENGLISH
(MAXIMUM NUMBER OF FRAMES USED)

SET NUMBER :		4							
ZERO	ONE	TWO	THREE	FOUR	FIVE	SIX	SEVEN	EIGHT	NINE
1.000	-0.103	0.027	0.157	-0.333	-0.125	-0.347	-0.002	-0.119	0.192
-0.103	1.000	-0.639	-0.298	-0.724	0.489	0.143	-0.239	-0.405	0.428
0.027	-0.639	1.000	0.045	0.423	0.196	-0.465	-0.524	0.320	0.139
0.157	-0.298	0.045	1.000	-0.123	-0.249	-0.224	-0.437	-0.532	-0.159
-0.333	-0.724	0.423	-0.123	1.000	0.648	0.244	-0.138	0.199	0.134
-0.125	0.489	0.196	-0.249	0.648	1.000	0.360	0.006	0.351	0.519
-0.347	0.143	-0.465	-0.224	0.244	0.360	1.000	0.429	0.165	0.145
-0.002	-0.239	-0.524	-0.437	-0.138	0.006	0.429	1.000	0.003	0.071
-0.119	-0.405	0.320	-0.532	0.199	0.351	0.165	0.003	1.000	0.353
0.192	0.428	0.139	-0.159	0.134	0.519	0.145	0.071	0.353	1.000

TABLE XI
 BTR CORRELATION TABLE FOR DIGITS SPOKEN IN AMERICAN ENGLISH
 (MAXIMUM NUMBER OF FRAMES USED)

SET NUMBER :		BTR CORRELATION							
ZERO	ONE	TWO	THREE	FOUR	FIVE	SIX	SEVEN	EIGHT	NINE
1.000	-0.400	-0.547	-0.003	0.318	-0.067	0.032	-0.062	0.158	-0.241
-0.400	1.000	0.339	0.161	0.568	0.354	0.426	-0.100	0.338	-0.182
-0.547	0.339	1.000	-0.099	0.787	0.532	0.264	-0.167	0.835	0.225
-0.003	0.161	-0.099	1.000	0.007	-0.268	0.133	0.085	0.153	-0.496
0.318	0.568	0.787	0.007	1.000	0.922	0.508	-0.034	0.767	0.224
-0.067	0.354	0.532	-0.268	0.922	1.000	0.444	-0.362	0.698	-0.021
0.032	0.426	0.264	0.133	0.508	0.444	1.000	0.284	0.212	-0.108
-0.062	-0.100	-0.167	0.085	-0.034	-0.362	0.284	1.000	0.144	0.111
0.158	0.338	0.835	0.153	0.767	0.698	0.212	0.144	1.000	0.363
-0.241	-0.182	0.225	-0.496	0.224	-0.021	-0.108	0.111	0.363	1.000

correlation data gives in some way, the general pattern for the correlation of the spoken digits. Tables VIII - XI indicate that, some digits are heavily correlated.

In the following, pattern matching is used to classify the digits that have not been recognized by the procedure discussed earlier. Recognition of these spoken digits is achieved by cross-correlating the pattern of the unknown spoken digit with the stored test patterns of the digits zero through nine, as shown by Figures 60 and 61. The standard patterns used are that of the RMS and BTR parameters. The digit which has the highest correlation is selected. This system utilizes no linguistic information, but uses procedures of acoustic characteristics imbedded in the RMS energy and BTR parameter derived from the equivalent area function. This completes the digit recognition scheme.

The above method was applied for sample digits and the results were very good as shown by Tables XII and XIII. However, more sample digits have to be tested before the recognition rate can be given. Figures 60 and 61 can be used for further research for testing a digit recognition scheme.

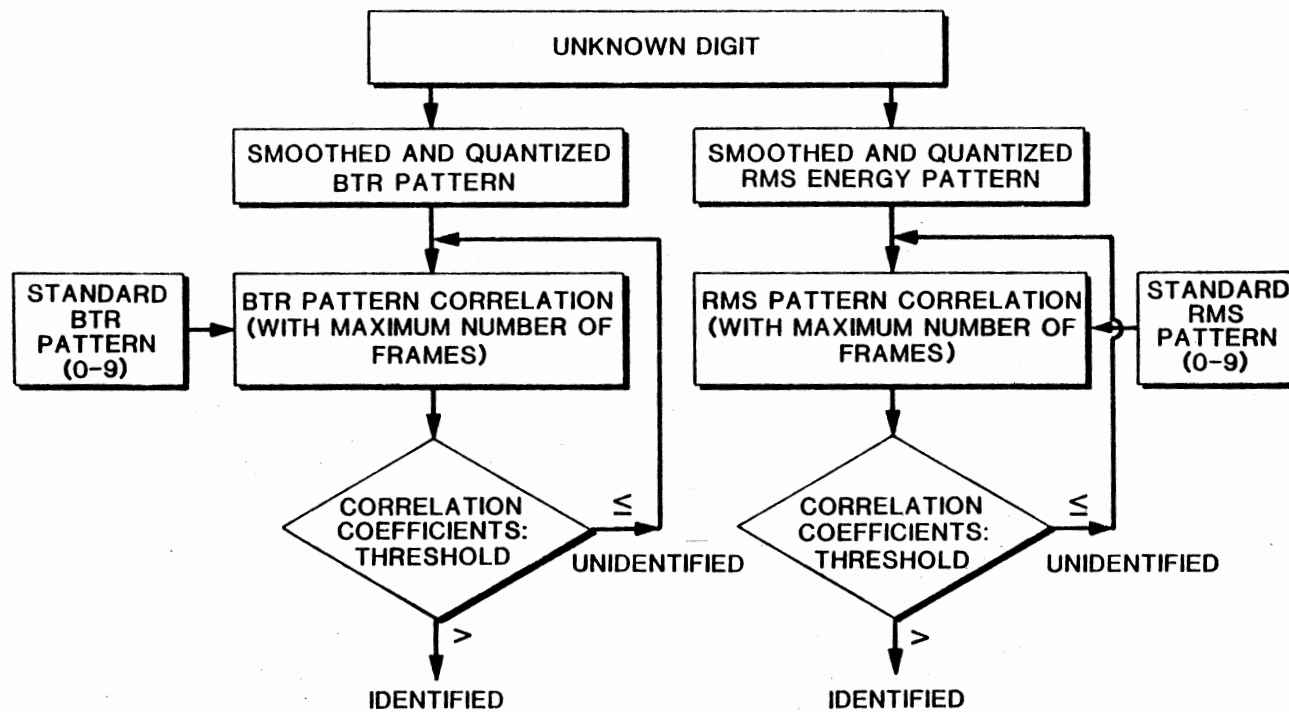


Figure 60. RMS and BTR Correlation Flow Diagram for Digit Recognition

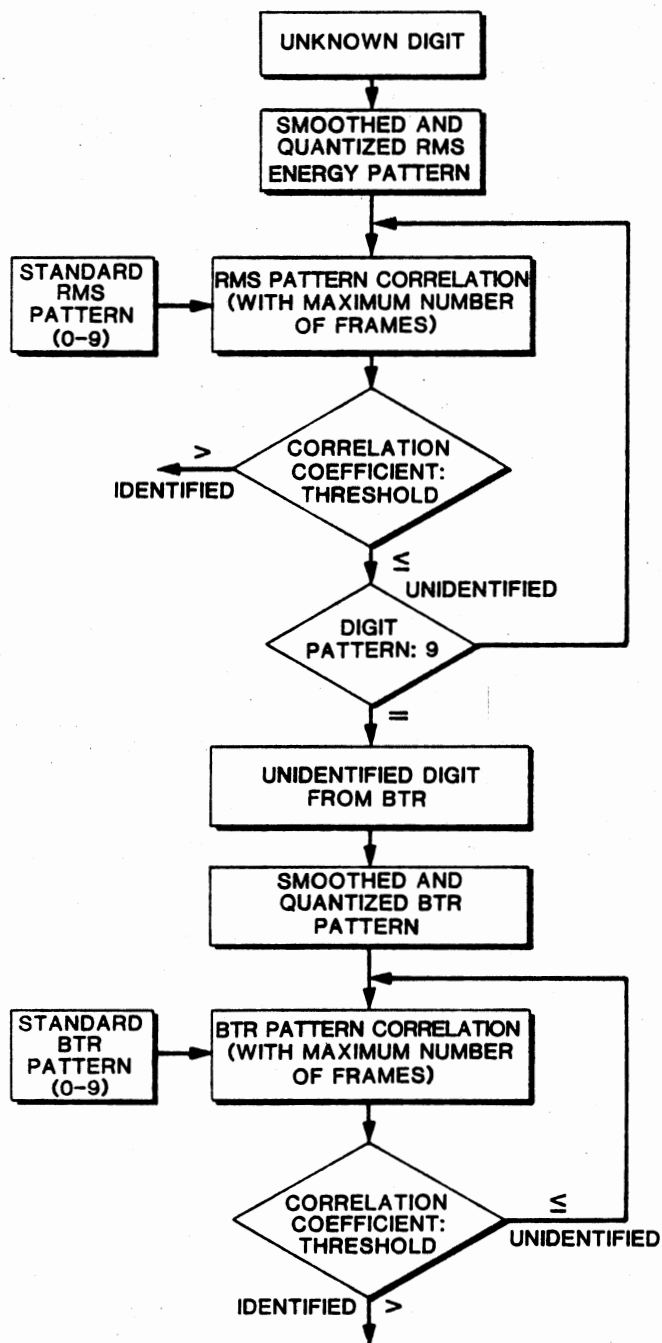


Figure 61. Future RMS and BTR Correlation Flow Diagram for an Efficient Digit Recognition System

TABLE XII

RMS ENERGY CORRELATION TABLE FOR TWO SETS OF DIGITS SPOKEN IN AMERICAN
ENGLISH BY TWO DIFFERENT SPEAKERS

ZERO	ONE	TWO	RMS CORRELATION						
			THREE	FOUR	FIVE	SIX	SEVEN	EIGHT	NINE
0.791	-0.107	0.076	0.149	0.379	-0.040	0.212	0.471	-0.130	0.251
-0.107	0.981	0.172	0.584	0.155	0.616	0.355	0.229	0.901	-0.005
0.076	0.172	0.981	0.841	0.476	0.781	0.988	0.901	0.045	0.604
0.149	0.584	0.841	0.949	0.504	0.856	0.887	0.818	0.489	0.571
0.379	0.155	0.476	0.504	0.855	0.185	0.614	0.805	-0.166	0.828
-0.040	0.616	0.781	0.856	0.185	0.960	0.818	0.698	0.663	0.242
0.212	0.355	0.938	0.887	0.614	0.818	0.987	0.858	0.287	0.470
0.471	0.229	0.901	0.818	0.805	0.698	0.858	0.793	0.342	0.451
-0.130	0.901	0.045	0.489	-0.166	0.863	0.287	0.342	0.972	-0.077
0.251	-0.005	0.604	0.571	0.828	0.242	0.470	0.451	-0.077	0.954

TABLE XIII

BTR CORRELATION TABLE FOR THE SAME TWO SETS OF DIGITS USED IN TABLE XII

ONE	TWO	THREE	BTR CORRELATION						
			FOUR	FIVE	SIX	SEVEN	EIGHT	NINE	TEN
0.662	0.078	-0.491	-0.079	-0.678	-0.553	-0.375	-0.400	0.088	-0.403
0.078	0.745	-0.224	-0.267	0.165	0.144	-0.136	-0.021	0.677	0.283
-0.491	-0.224	0.758	0.395	0.630	0.143	0.530	0.338	0.310	-0.399
-0.070	-0.267	0.385	0.713	0.153	-0.082	0.620	-0.022	0.309	-0.569
-0.678	0.165	0.630	0.153	0.801	-0.047	0.305	0.580	0.336	-0.110
-0.553	0.144	0.143	-0.082	-0.047	0.832	0.526	0.160	0.650	0.550
-0.375	-0.136	0.530	0.620	0.305	0.526	0.709	-0.216	0.582	0.465
-0.400	-0.021	0.338	-0.022	0.580	0.160	-0.216	0.656	0.307	-0.143
0.088	0.677	0.310	0.309	0.336	0.650	0.582	0.307	0.762	0.377
-0.403	0.283	-0.399	-0.559	-0.110	0.550	0.465	-0.143	0.377	0.756

CHAPTER V

ACOUSTIC PHONEMIC DIGIT RECOGNITION SCHEME FOR DIGITS SPOKEN IN ARABIC

Introduction

In the last chapter, the phonemic digit recognition was discussed based upon acoustic patterns of phonemes. Obviously, a similar approach can be used in recognizing digits in other languages. The only difference is that the phonemes for a given digit in two languages will generally be different. Also, some of the phonemes used in some languages may not be in other languages. For example, the glottal phoneme in Arabic is not in English. In this chapter, the ideas in Chapter IV are extended for digits spoken in Arabic and it is shown that the phonemic Arabic digit recognition can be accomplished in the general framework of speech processing.

Arabic Phonemes

There are about 57 phonemes in Arabic Language, but only 24 phonemes are used in the digits 0 to 9 in Arabic. Table XIV shows the Arabic alphabet and its equivalent translation into American English [70]. Also Figure 9 and Table IV classify the phonemes used in digits spoken in Arabic. There are ten vowels, three of which are short vowels and three can be long vowels, where it is really pronounced long as in /θalâθðh/ (three), i.e. the phoneme /â/ is a long vowel. The short

TABLE XIV

TRANSLITERATION OF ARABIC WORDS AND NAMES

ا	} Consonantal sound	a	ط	t
ء			ظ	z
ا	Long vowel *	a	ع	' (Inverted apostrophe)
ب		b	ف	f
ت		t	ق	q
ث		th	ك	k
ج		j	ل	l
ح		h	م	m
خ		kh	ن	n
د		d	ه	h
ذ		z	و	consonant	w
ر		r	و	long vowel *	u
ز		z	و	diphthong	au
س		s	ي	consonant	y
ش		sh	ي	long vowel *	i
ص		s	ي	diphthong	ai
ض		dh			
Short vowels:		/ (fatḥa)	a		
		ـَ (kasra)	i		
		ـُ (dhamma)	u		

vowels are /a/, /i/, and /u/. For example /a/ as in /wahid/, (one). The long vowels are /â/, /û/, and /î/, as shown by Figure 9 and 11. Diphthongs do not appear in Arabic digits.

There are nine basic sounds which fundamentally differ from any sound in American English. These are /h/, /kh/, /s/, /dh/, /t/, /Z/, /p/ or /A/, /g/, and /q/. Since these phonemes are very difficult to be produced by a non-Arabic speaking person, a brief description of the articulatory problems will be given below for interest.

The Emphatics

This group of Arabic sounds share some common features: 1) they are all produced with the back of the tongue raised towards the back roof of the mouth; 2) all of them have non-emphatic counterparts, from which they should be clearly distinguished; and 3) in producing echo, these sounds produce more echo and 'thickness of voice' than their counterparts.

Both vowel and consonant productions may involve the simultaneous activity of many articulators. For vowels, the tongue, velum, lips, and larynx are all in operation. It is difficult to say that one articulator's movements is more important than another's. In languages other than American English, the lips and tongue may form simultaneous vocal-tract stoppages [19][21]. Multiple articulation is important in Arabic.

Velarized Phonemes [67]

The Arabic /s/ is a velarized alveolar /s/. It is alveolar, which means that it is produced not at the teeth, like the /s/ as in 'seen',

but further back in the mouth. It is velarized, which means that the back part of the tongue is made tense, with some raising up toward the soft palate or velum, as illustrated by Figure 62, giving the /s/ sound a velar effect. Thus the phoneme /s/ is quite different from the phoneme /s/.

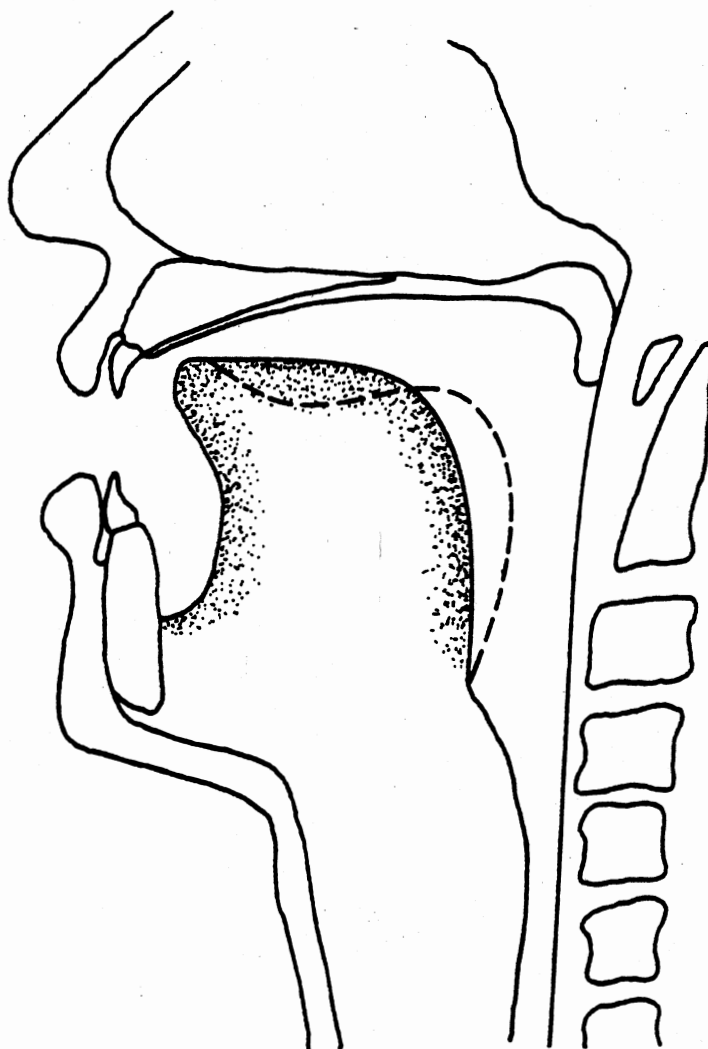
The Arabic phoneme /dh/ or /d/ pronounced as 'daad', has no American English equivalent. It is a velarized voiced alveolar stop. From another view point, it is like the phoneme /d/, with the addition of velarization [67]. It should be pointed out that one side of the tongue gets closer to the side of the mouth than touching the molars.

The phoneme /t/, pronounced as 'taa' is the emphatic counterpart of the phoneme /t/. The consonant /t/ is a voiceless alveolar stop. It differs from the plain American English /t/, which is often aspirated (i.e., it is produced with a slight puff of breath); while /t/ is not aspirated [69].

The phoneme /ṯ/ pronounced as 'thaa' is a velarized voiced interdental fricative; it is like the American phoneme /ṯ/ as in that, with the addition of velarization (the tensing up of the base of the tongue).

Throat and Back-of-Mouth Sounds

The phoneme /ʁ/, i.e. /A/ pronounced as 'Aayn' is not generally used in the English Language. It is a voiced pharyngeal fricative. The 'Aayn' is similar to the sound one produces sometimes to express feeling of being strangled. This sound is produced deep in the throat contracting the muscles in there and forcing the air through, like



[s] ————— } AS SPOKEN BY A
 [ʕ] - - - - - } SPEAKER OF ARABIC

Figure 62. Tongue Position for the Emphatic-Nonemphatic [s] vs [ʕ] Fricative Consonants in Arabic

'arbaḅḏh' in Arabic for digit four. It is important to distinguish this sound from 'Haa' as it causes some echo in the mouth, which could be clearly noticed if one closes one's ears when producing the sound 'Aayn'.

The phoneme /g/, pronounced as 'gayn' is very close in pronunciation to standard French 'r'. It is like the noise one makes when gargling without water. It is a voiced velar fricative. This sound should be distinguished from American English 'g' for which the air stream is completely interrupted in the area of the back roof of the mouth.

The phoneme /q/, pronounced as 'qaaf' is produced in the area deeper and more to the back than k. It has a very clear echo in the mouth, as distinct from k or g in American English. The closest sound in English is heard in pronunciation of words like 'caught', but it is still deeper.

The phoneme /h/ like /wâhid/ for Arabic digit one; it is a voiceless pharyngeal fricative. It is produced with the base of the tongue near the back of the pharynx (throat) and the pharynx walls strongly constricted. There is no contact whatsoever between the base of the tongue and the velum.

The phoneme /hk/ is a voiceless velar fricative. It is produced by narrowing the passageway between the back of the tongue and the velum, so creating friction as the air passes through; the vocal cords are at rest. The digit (five) in Arabic is /khḷmaḅh/ has the phoneme /kh/ which is sometimes written in the form /x/ which has nothing to do with the American English sound x. The phoneme transcribed as /kh/

sound is produced with a continuous stream of air, as opposed to k which has an interrupted stream of air in pronunciation. It is closer to German ch in 'nacht', Scotch 'ch' in 'loch' and Latin American pronunciation of 'j' in 'mejor'.

Minor Differences

The phoneme /r/, pronounced as 'raa' should be distinguished from its counterpart in American English. In Arabic, it is produced by the tip of the tongue rapidly touching the alveolar ridge (or the area just over the gums of the upper front teeth). Occasionally, we hear a very similar sound in the pronunciation of the t and d in American English 'latter' and 'ladder'. It is also similar to British English r in 'very' and Spanish r too. The tip of the tongue is not curled when pronouncing the phoneme /r/ as in /ʃɛfr/.

The phoneme /l/, pronounced 'laam' as in digit /θalāʔθh/ in Arabic, meaning three, has a unique case in Arabic. Usually it is like in western languages in general and American English /l/ in words like 'leaf', 'lip', but not 'lot' or 'low' .

There is an important difference between the Arabic sound /l/ and the American English /l/. Most speakers of American English pronounce the phoneme /l/ with back of the tongue raised somewhat toward the velum, resulting in a velarized /l/. This velarized quality of the English /l/ is especially noticeable at the end of the word. For example as in feel or bell. The Arabic /l/ has a non-velarized or clear sound. The clear /l/ results when the back of the tongue is relaxed and not raised.

The Arabic phoneme /s/ corresponds to the American English /s/ as in 'see'. The Arabic /s/ is dental pronounced with the tongue tip at the upper teeth as in the Arabic spoken digit /khΛmsðh/, i.e. five, other examples are shown in Table III and IV. It is interesting to point out that the American English /s/ is alveolar, pronounced slightly behind the teeth, giving slightly lower-pitched /s/.

The Arabic phoneme /h/ as in /θaλa^ˆθðh/, is like the American English /h/ as in 'hat' and 'heat'; it is a voiceless glottal fricative. It is generally considered as whisper. The Arabic /h/ differs from the American English /h/ in the following ways: 1) it is pronounced with more force than in the English /h/; 2) it can be pronounced at the end of a syllable word; 3) also it can be held twice as long.

The Arabic phoneme /θ/, as in /θaλa^ˆθðh/, digit three in Arabic, is like the American phoneme /θ/, as in digit /θri/ in American English. The only difference is that the phoneme /θ/ is followed by a front vowel in digit /θri/ spoken in English, while it is followed by a back vowel in the same digit spoken in Arabic. Similarly the phonemes /f/, /ω/, /d/, /n/, /m/, /b/ and /y/ have no noticeable differences when used only in digits spoken in Arabic.

Finally an important feature in Arabic phonology is that of germination, which means that the consonant is doubled in pronunciation as in the digit /sitt^ˆðh/, i.e. six. Furthermore, the feature impeded in the spoken digit depends considerably on the degree of constriction and tongue hump position, as shown by Table IV and discussed in detail in Chapter II. In addition a simple comparison spoken by American English and Arabic is shown in Table XV.

TABLE XV

COMPARISON BETWEEN BEGINNING AND END OF ENGLISH AND ARABIC DIGITS

Digit	Start	End	Digit	Start	End
/zIro/	Voiced f	B.V	/serf/	UV.f	Semi Vowel
/wAn/	Semi Vowel	Nasal	/wâhid/	Semi Vowel	Voiced Stop
/tu/	UV.f	FV-BV	/iθnân/	F.V	Nasal
/θri/	UV.f	F.V	/θalâθðh/	UV.f	M.V*
/fɔr/	UV.f	Semi Vowel	/arbaðð/	M.V	M.V*
/faIv/	UV.f	Voiced f	/xλmsðh/	Fricative	M.V*
/sIks/	UV.f	UV.f	/sIttðh/	UV.f	M.V*
/seven/	UV.f	Nasal	/sλbρðh/	UV.f	M.V*
/eIt/	Diphthong	UV. Stop	/θamânyðh/	UV.f	M.V*
/naIn/	Nasal	Nasal	/tIspðh/	UV. Stop	M.V*

* These mid-vowels (M.V) are followed by a whisper phoneme /h/ as shown in Table IV.

Segmentation

The segmentation scheme, used previously in Chapter IV, is utilized for locating the boundaries of digits spoken in Arabic, because it is based on silence, voiced and unvoiced decisions. In addition, the end-point detection scheme discussed earlier is based on the assumption that all digits spoken in American English have no silence intervals among them, with the exception of digits zero and seven. However, the segmentation algorithm will account for the number of boundaries in digits spoken in Arabic provided that the number of spoken digits are known. In addition the number of phonemes for each digit must be known. This is because seven of the digits spoken in Arabic have silence interval in the middle region and the rest have voiced interval in the middle region. Table IV illustrates the sequence of sound classes for the digits spoken in Arabic. The middle region dips, for the digits that have silence region, can be seen from the smothered RMS energy plot. Successful segmentation is achieved provided that the ZCR and energy threshold for silence, unvoiced and voiced segments was properly set. This method of segmentation is discussed in detail in Chapters III and IV.

Digit Recognition Flow Chart

Once the digit boundaries are located, the RMS energy per frame is computed for the spoken digit as discussed earlier and shown by the flow chart of Figure 28. The difference between English and Arabic digit recognition is that the digit phoneme tree is different. The

RMS energy is smoothed and then quantized to a maximum level of 100, which essentially normalizes the data. Following the flow chart, the smoothed quantized RMS peaks are computed and stored for the uttered digit. RMS plots are shown in Figures 63-83. For each uttered digit the ratio of the largest two peaks is calculated, so that it is always less than unity. As a first step towards phonemic digit recognition, the computed values for the peak ratios for each uttered digit, zero through nine is then stored. A threshold level is adjusted according to an empirical value. Tables XVI and XVII give the range for the largest two peaks and the value of their ratio. The computed ratio is first compared to the empirically established threshold. If the peak ratio of the unknown spoken digit does not match the established range, then the first idea based on RMS energy does not work. It has been found from measurements, based on five sets of data, that the actual digit is recognized more than half the time. The other part of the time it indicates that no digit is identified. It can be concluded that the peak ratio value is speaker dependent. This completes the first digit recognition idea.

Proceeding in the second path of the flow chart in Figure 28, the signal is Hamming windowed, where the window length is assumed to be 150 points. Fourteenth order LPA model is computed using 128 points per frame. As before, the parameters signed FBR, the smoothed BTR, CTR, and FTR are computed. Plots are then obtained for these parameters, as shown by Figures 84-93. The RMS dip-classification is applied in the phonetic feature detection section followed by the two vowel, vowel-like and non-vowel decisions, which are based on the RMS and BTR contours. An 'OR' decision is finally made using the above two classification

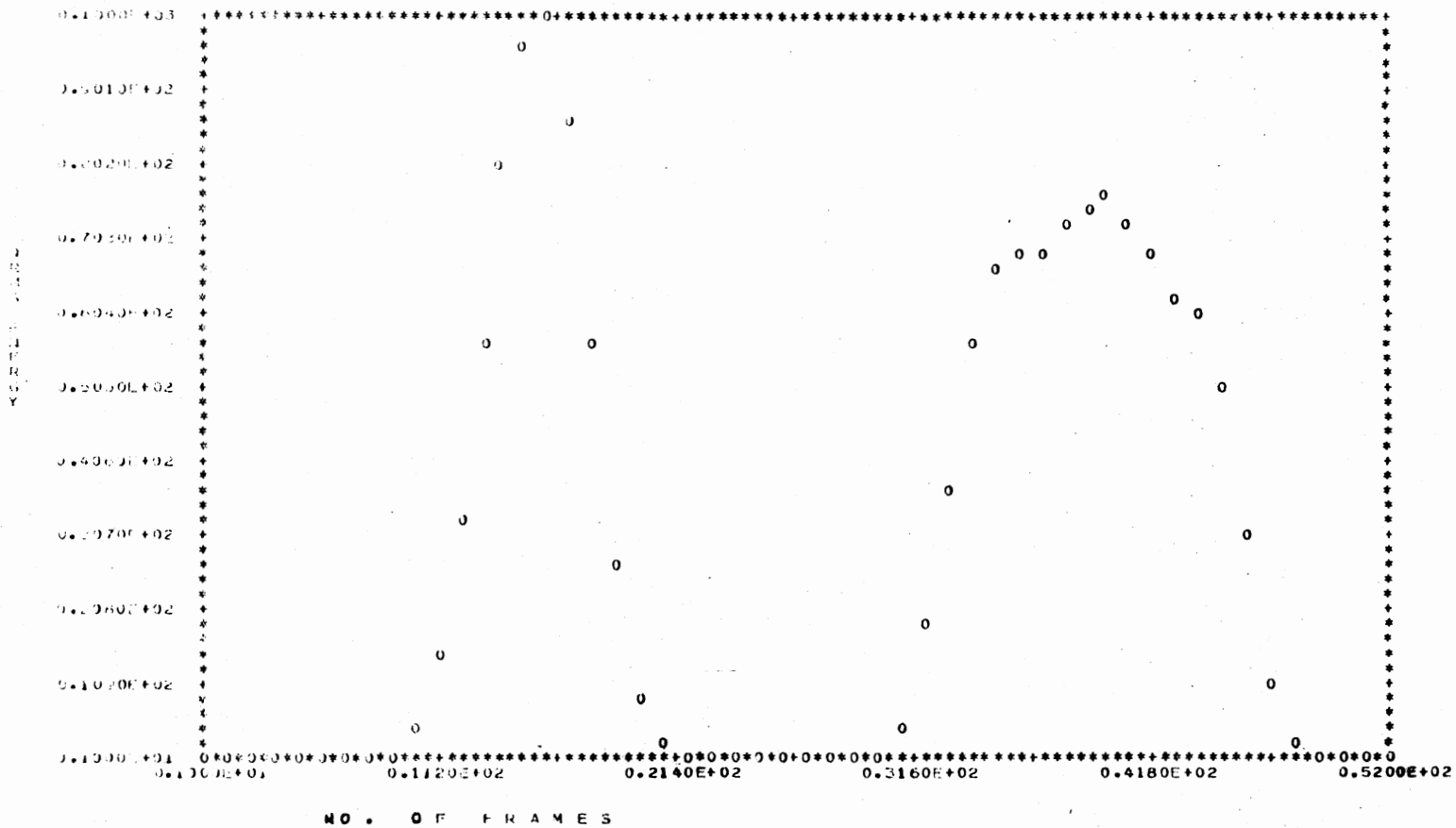


Figure 63. Smoothed and Quantized RMS Energy Contour for Digit Zero, i.e. /sefr/ Spoken in Arabic, Sample 1

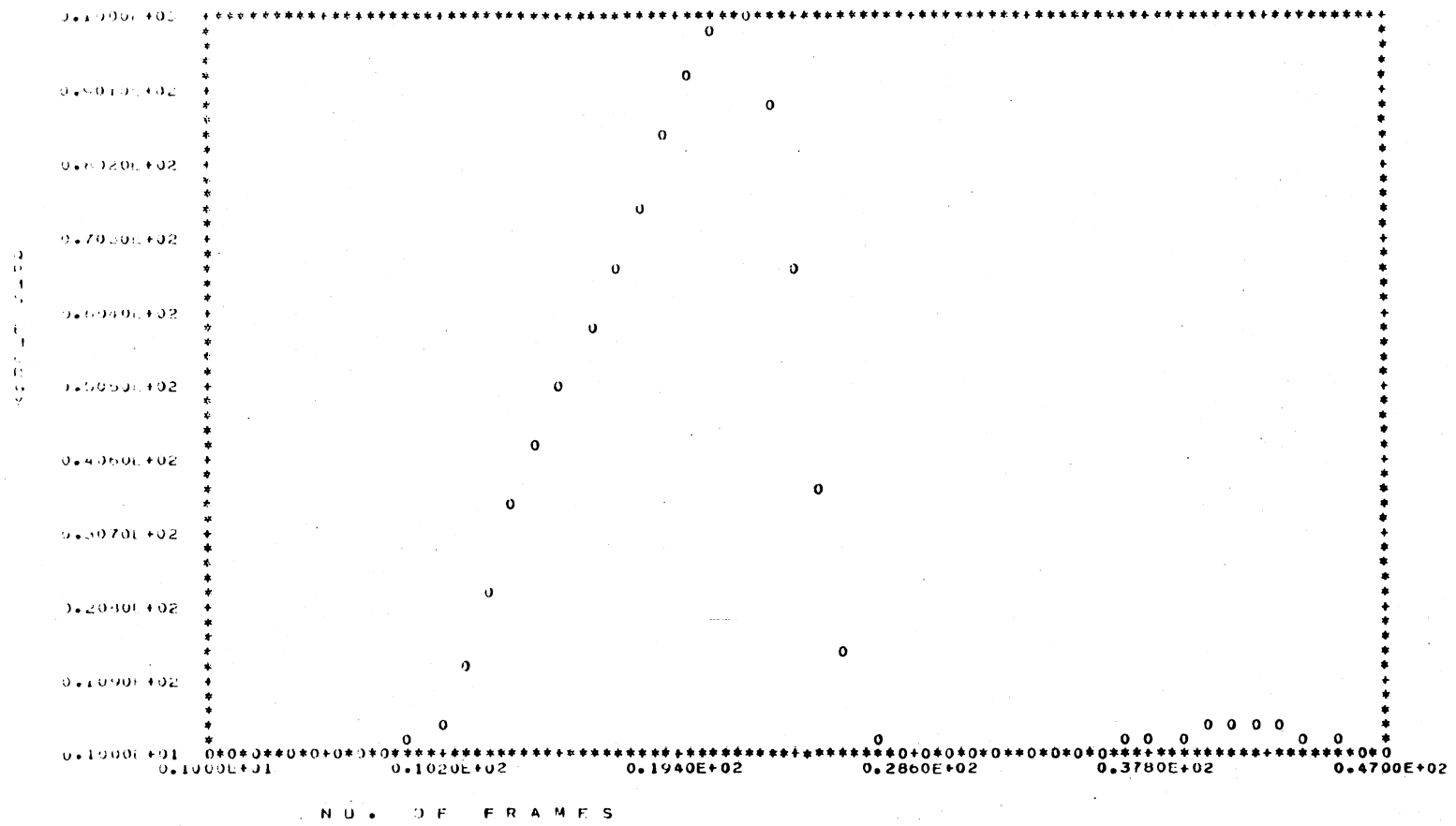


Figure 64. Smoothed and Quantized RMS Energy Contour for Digit One, i.e. /wāhid/ Spoken in Arabic, Sample 1

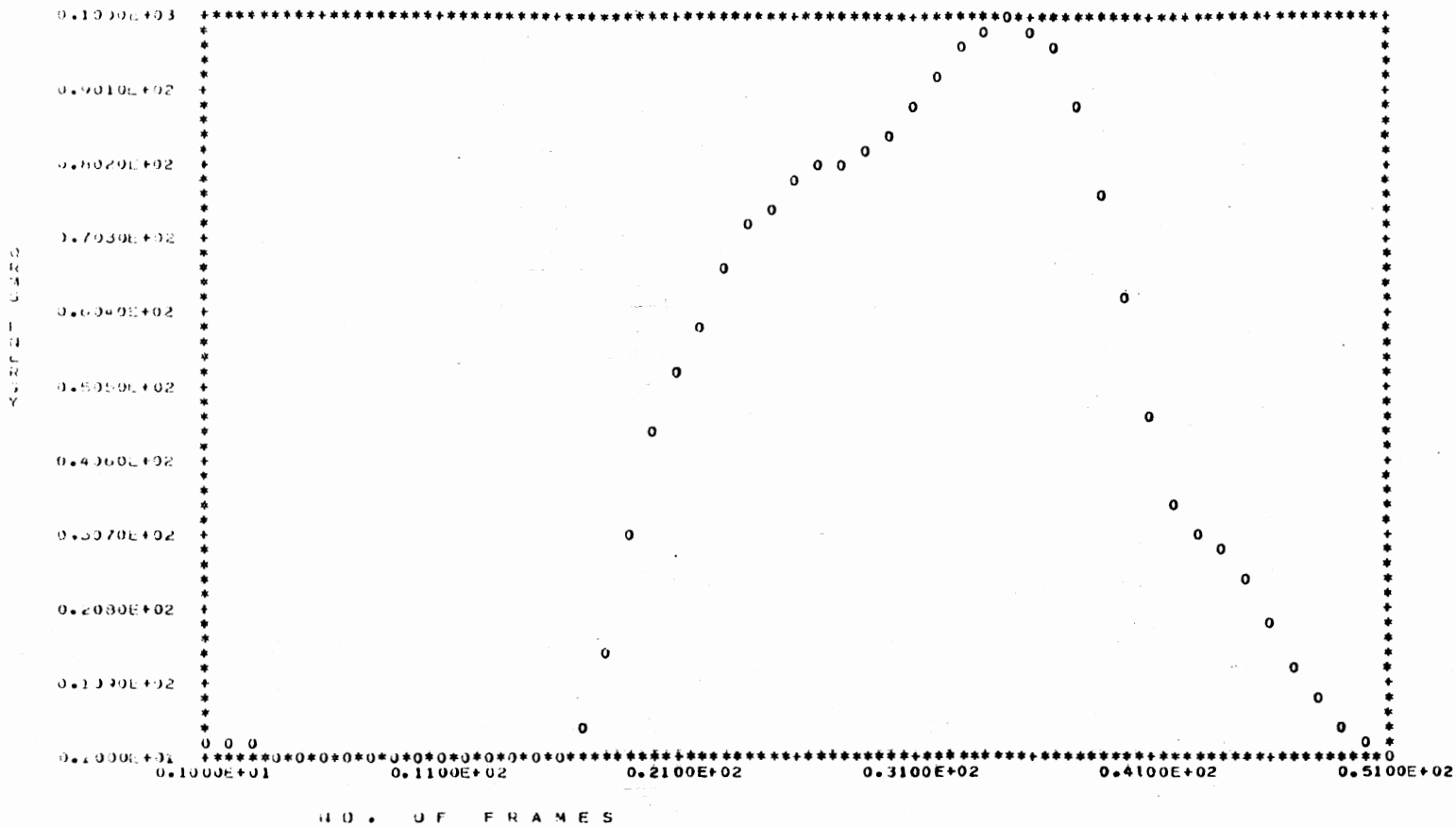


Figure 65. Smoothed and Quantized RMS Energy Contour for Digit Two, i.e. /iθnân/ Spoken in Arabic, Sample 1

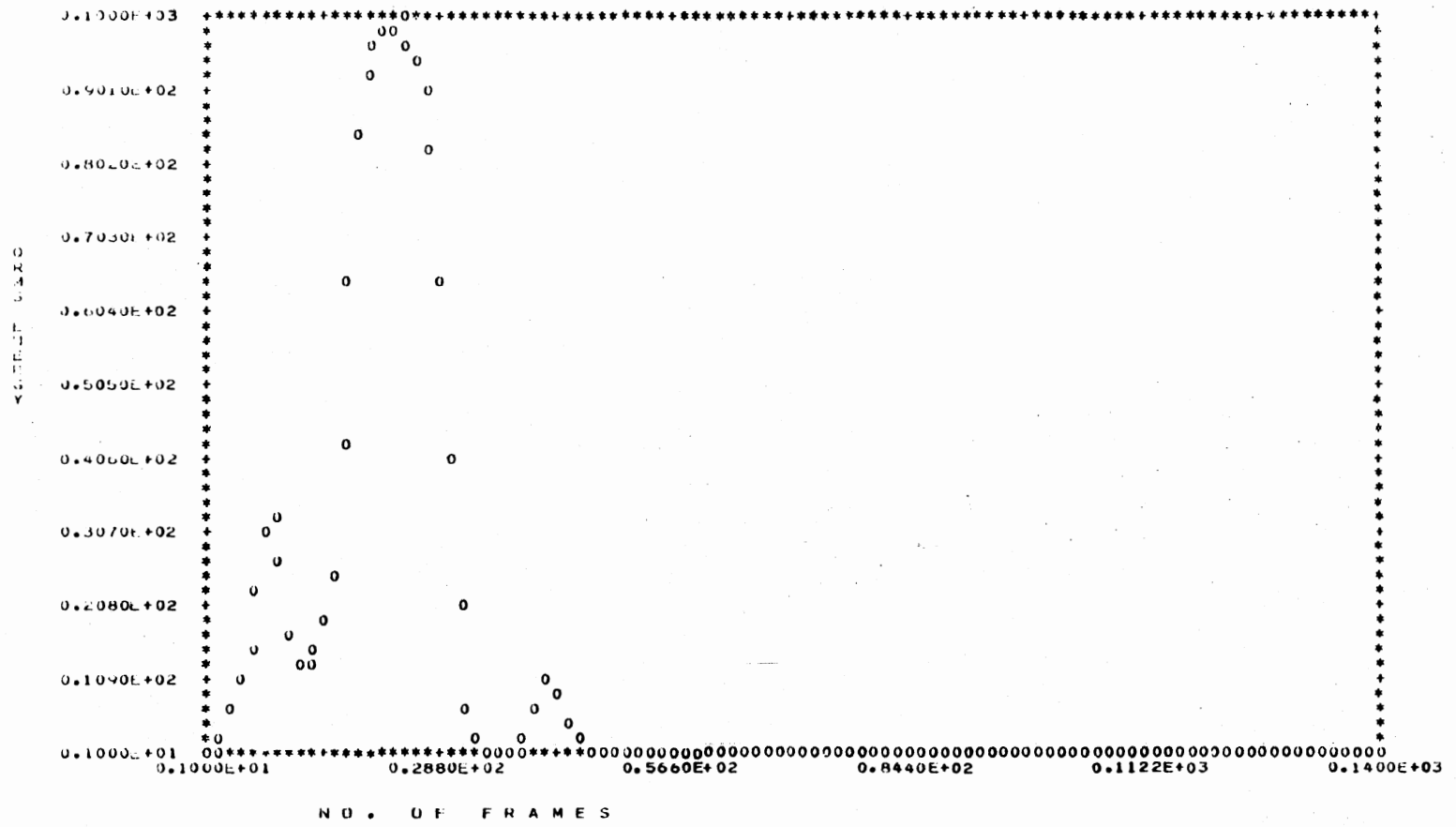


Figure 66. Smoothed and Quantized RMS Energy Contour for Ditit Three, i.e. /θalâṯh/ Spoken in Arabic Showing Wrong End Point Detection, Sample 1

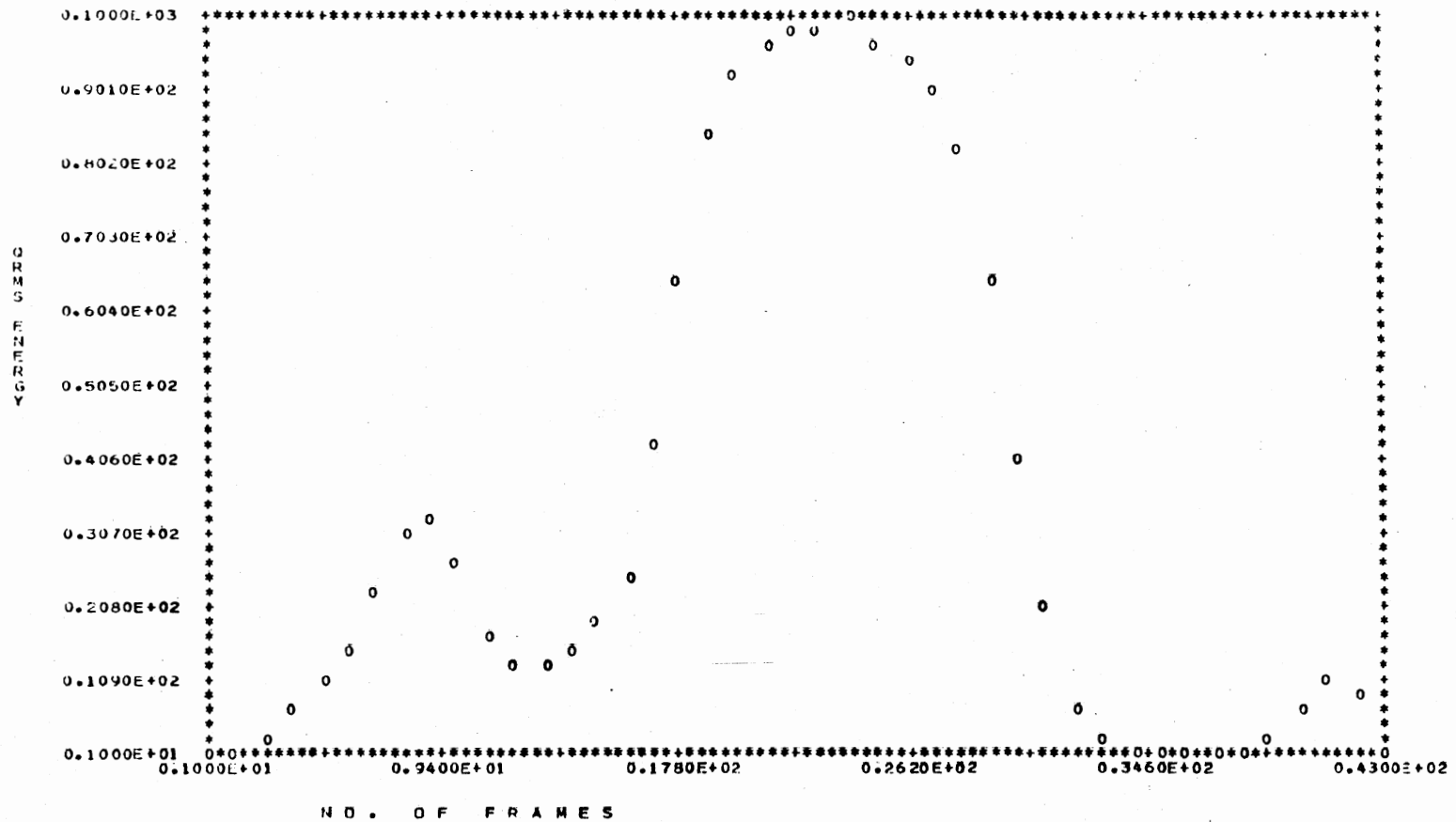


Figure 67. Smoothed and Quantized RMS Energy Contour for Digit Three, i.e. /θalâθh/ Spoken in Arabic, Sample 1

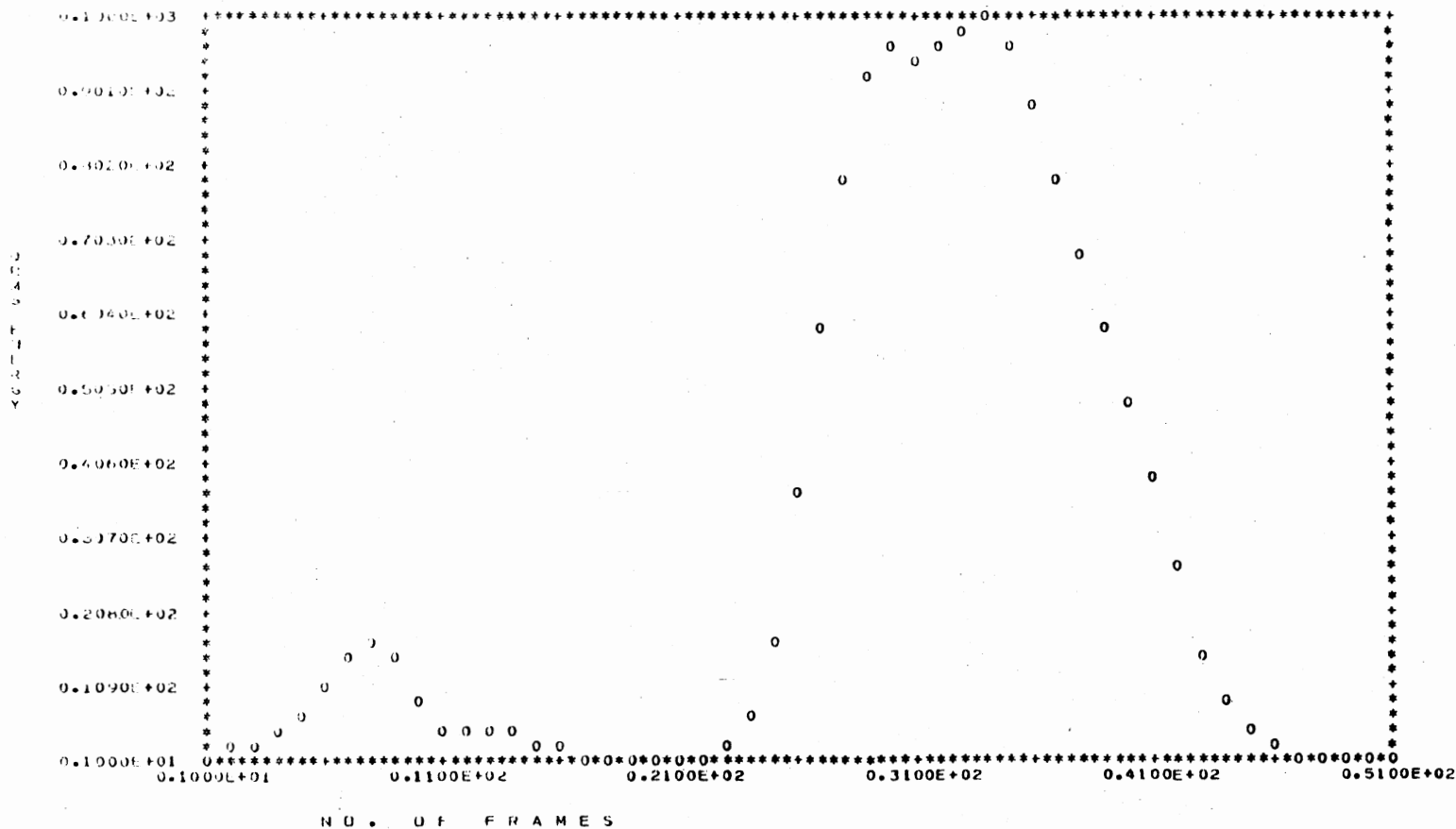


Figure 68. Smoothed and Quantized RMS Energy Contour for Digit Four, i.e. /arbaʔh/ Spoken in Arabic, Sample 1

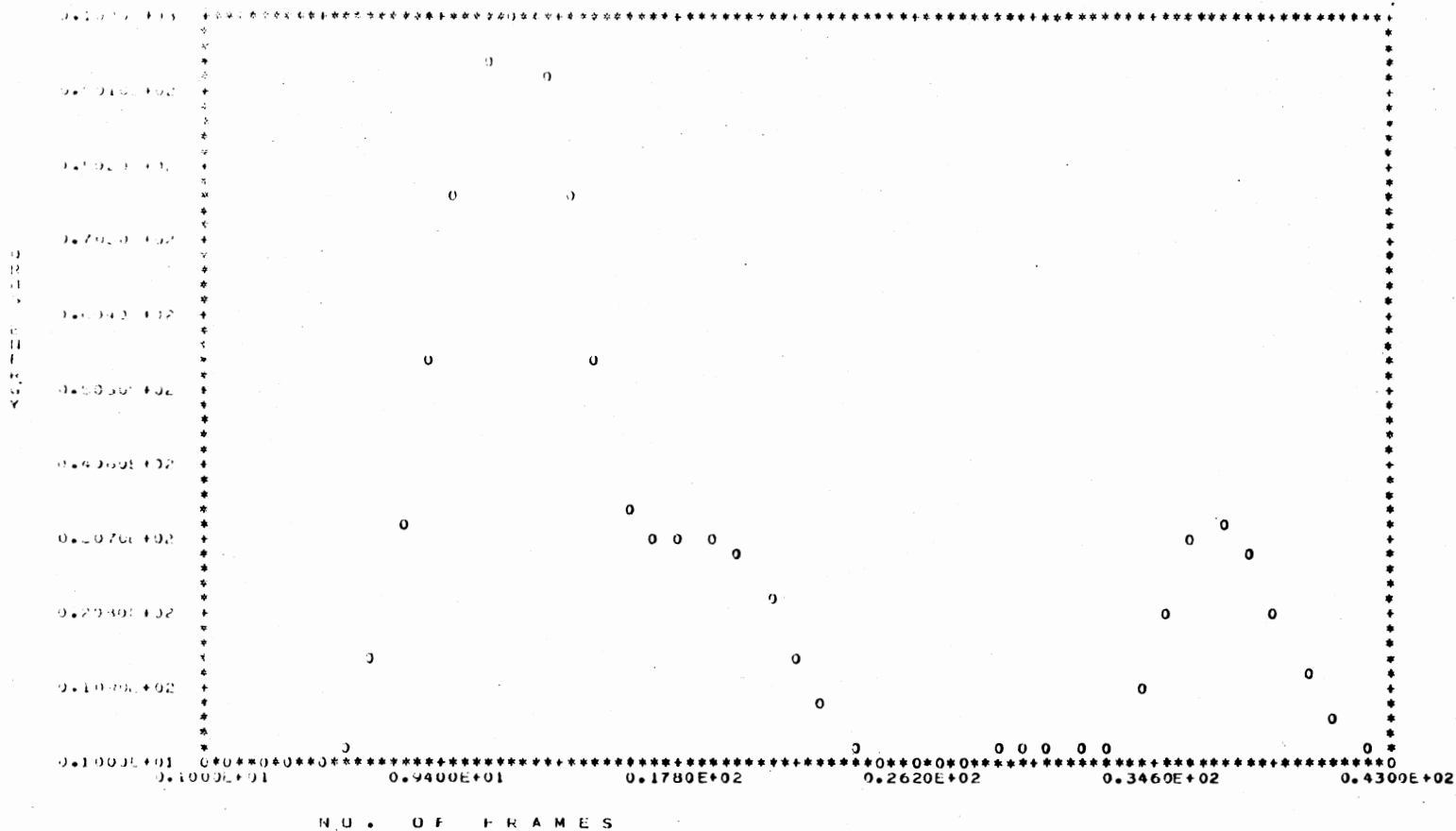


Figure 69. Smoothed and Quantized RMS Energy Contour for Digit Five, i.e. /khΛmsðh/ Spoken in Arabic, Sample 1

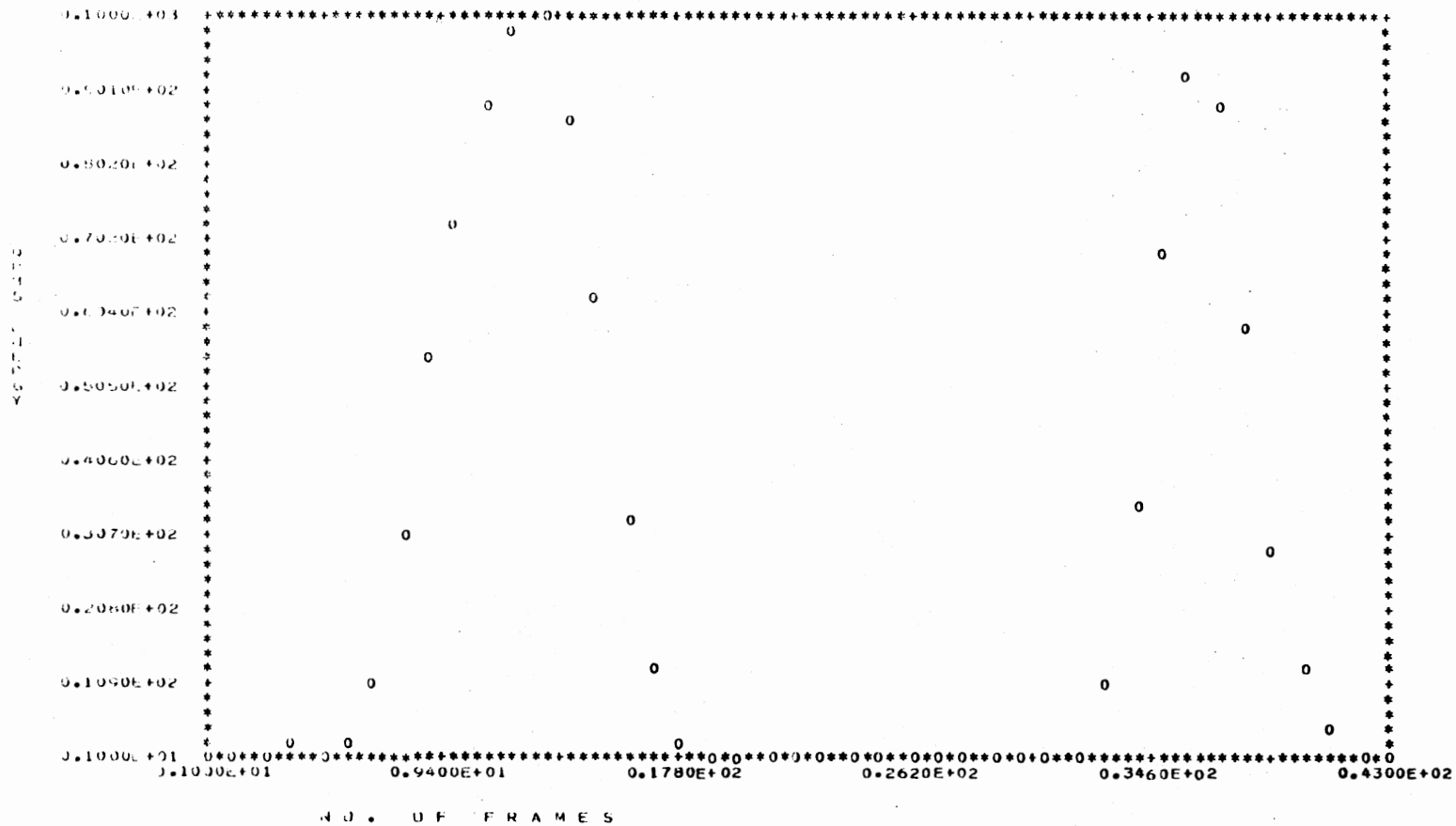


Figure 70. Smoothed and Quantized RMS Energy Contour for Digit Six, i.e. /sIttθh/ Spoken in Arabic, Sample 1

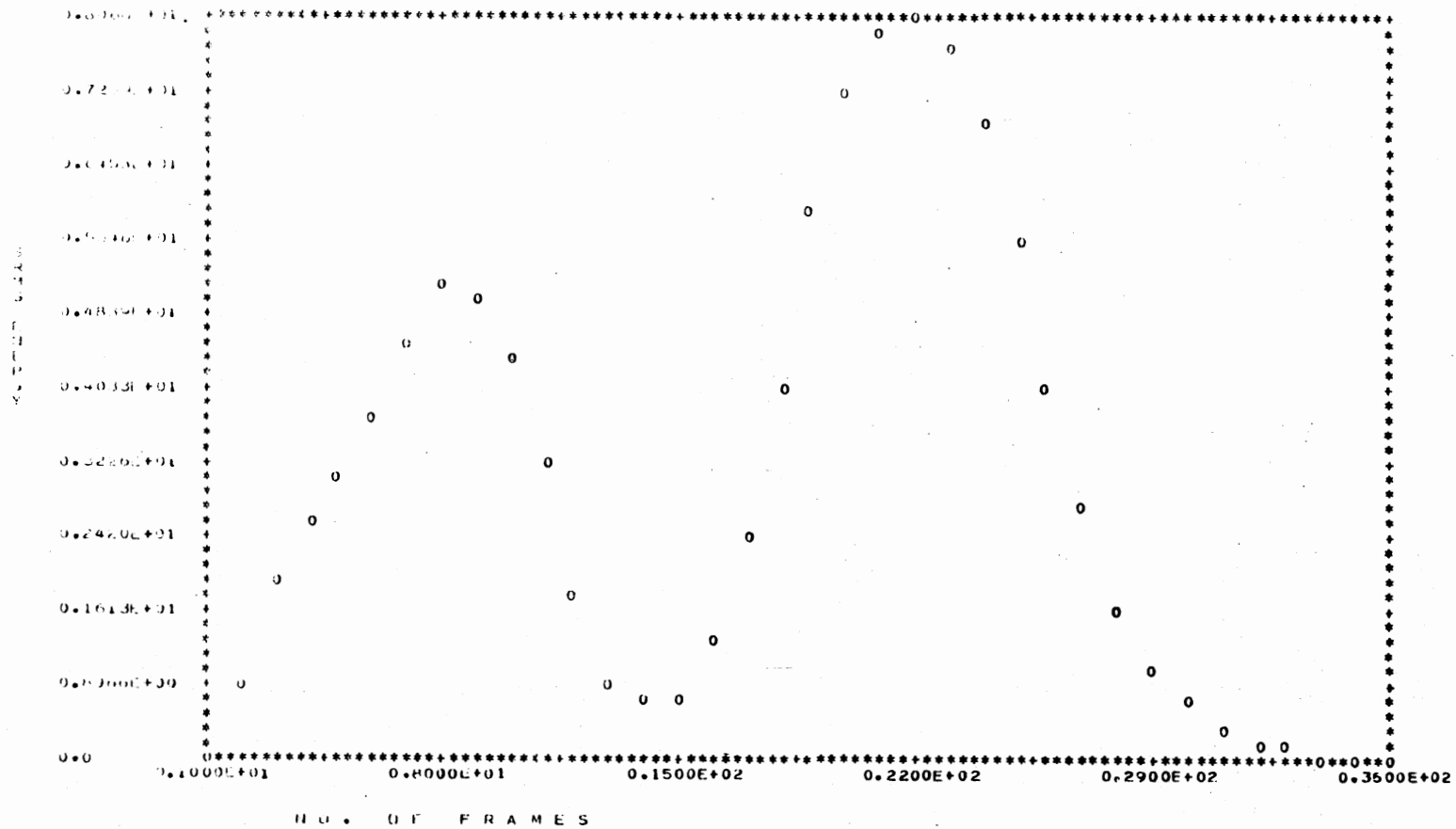


Figure 71. Smoothed and Quantized RMS Energy Contour for Digit Seven, i.e. /sʌbɔðh/ Spoken in Arabic, Sample 1

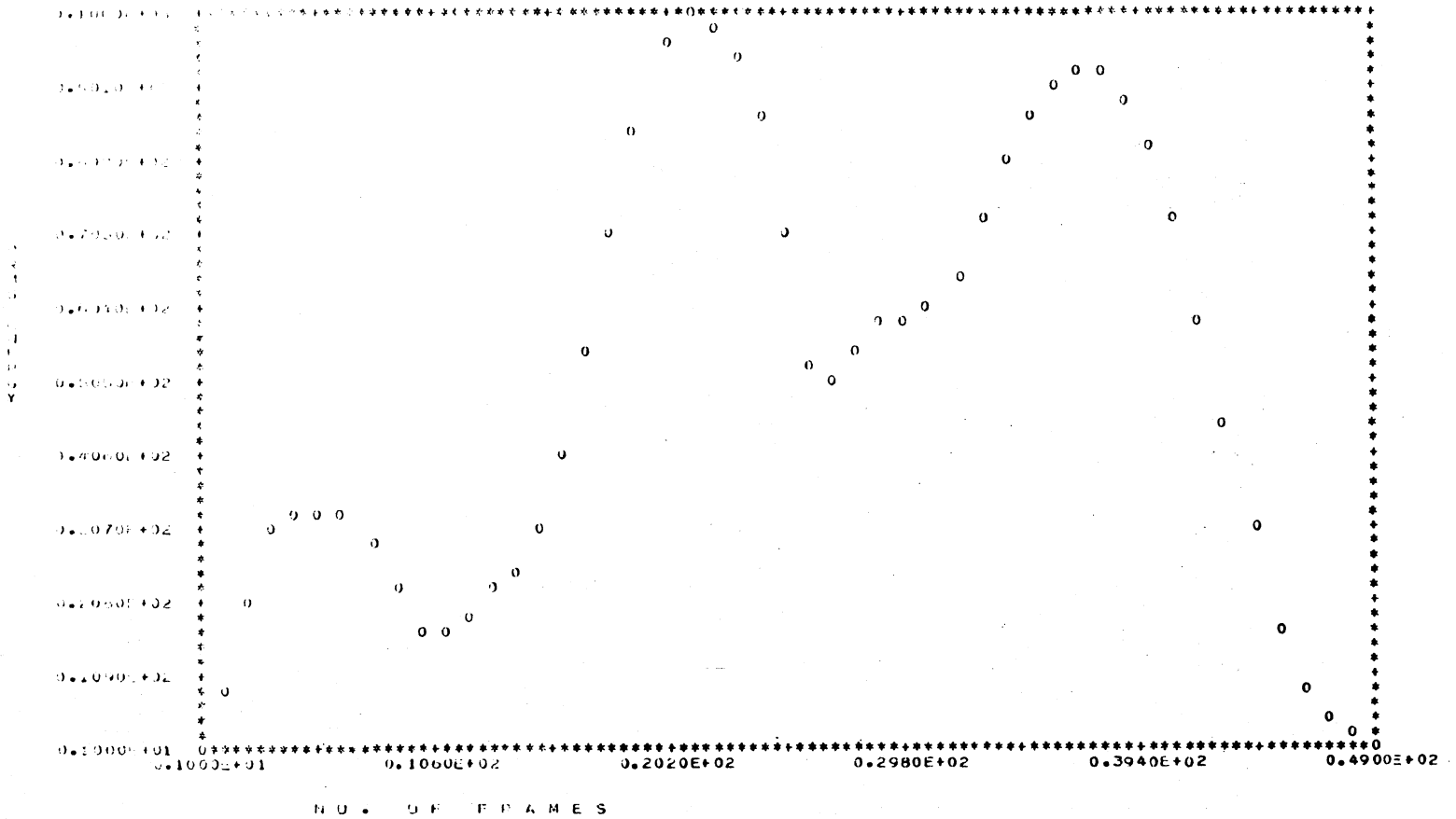


Figure 72. Smoothed and Quantized RMS Energy Contour for Digit Eight, i.e. /θamānyðh/ Spoken in Arabic, Sample 1

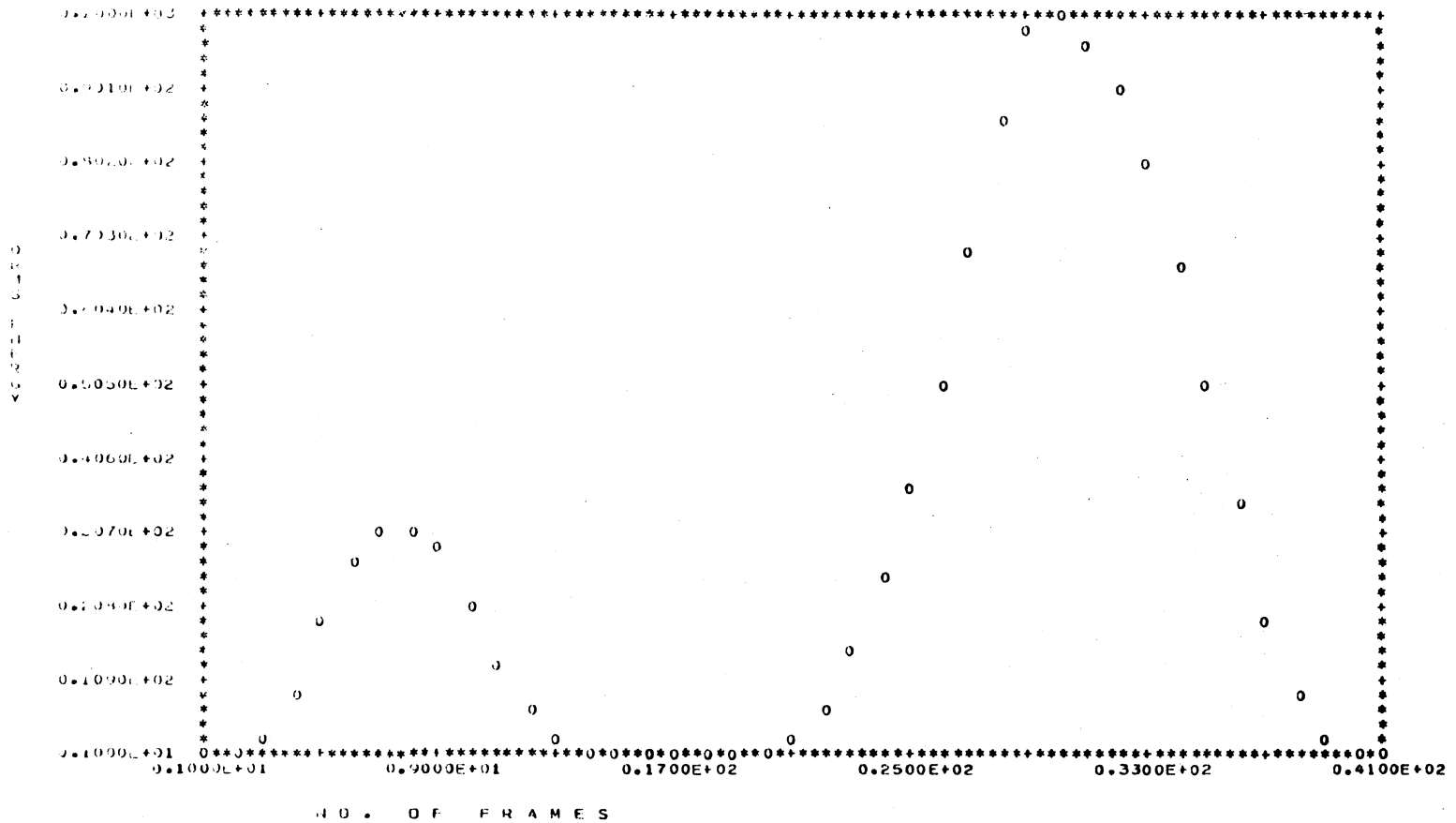


Figure 73. Smoothed and Quantized RMS Energy Contour for Digit Nine, i.e. /tIsəðh/ Spoken in Arabic, Sample 1

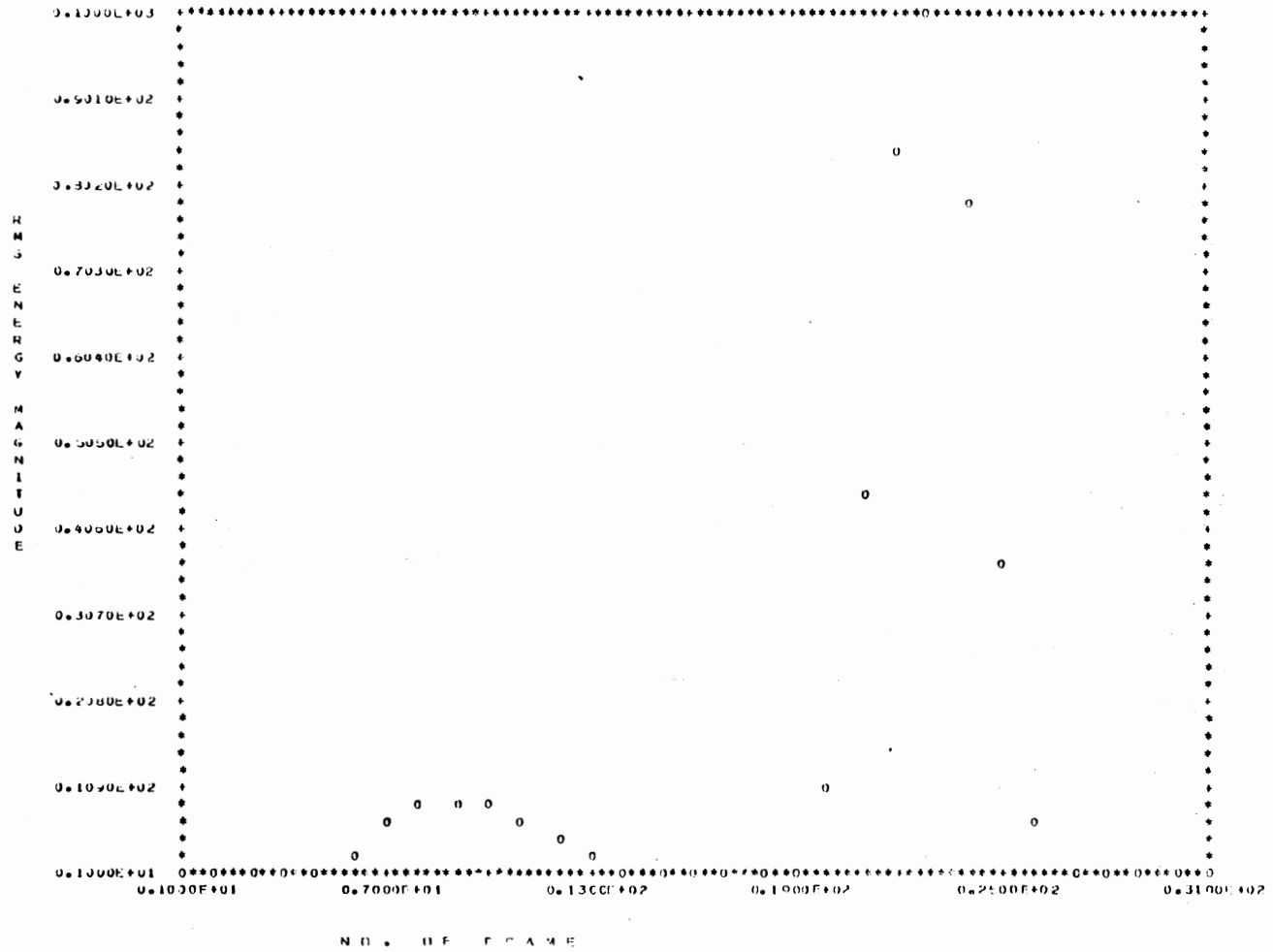


Figure 74. Smoothed and Quantized RMS Energy Contour for Digit Zero, i.e. /sefr/, Spoken in Arabic, Sample 2

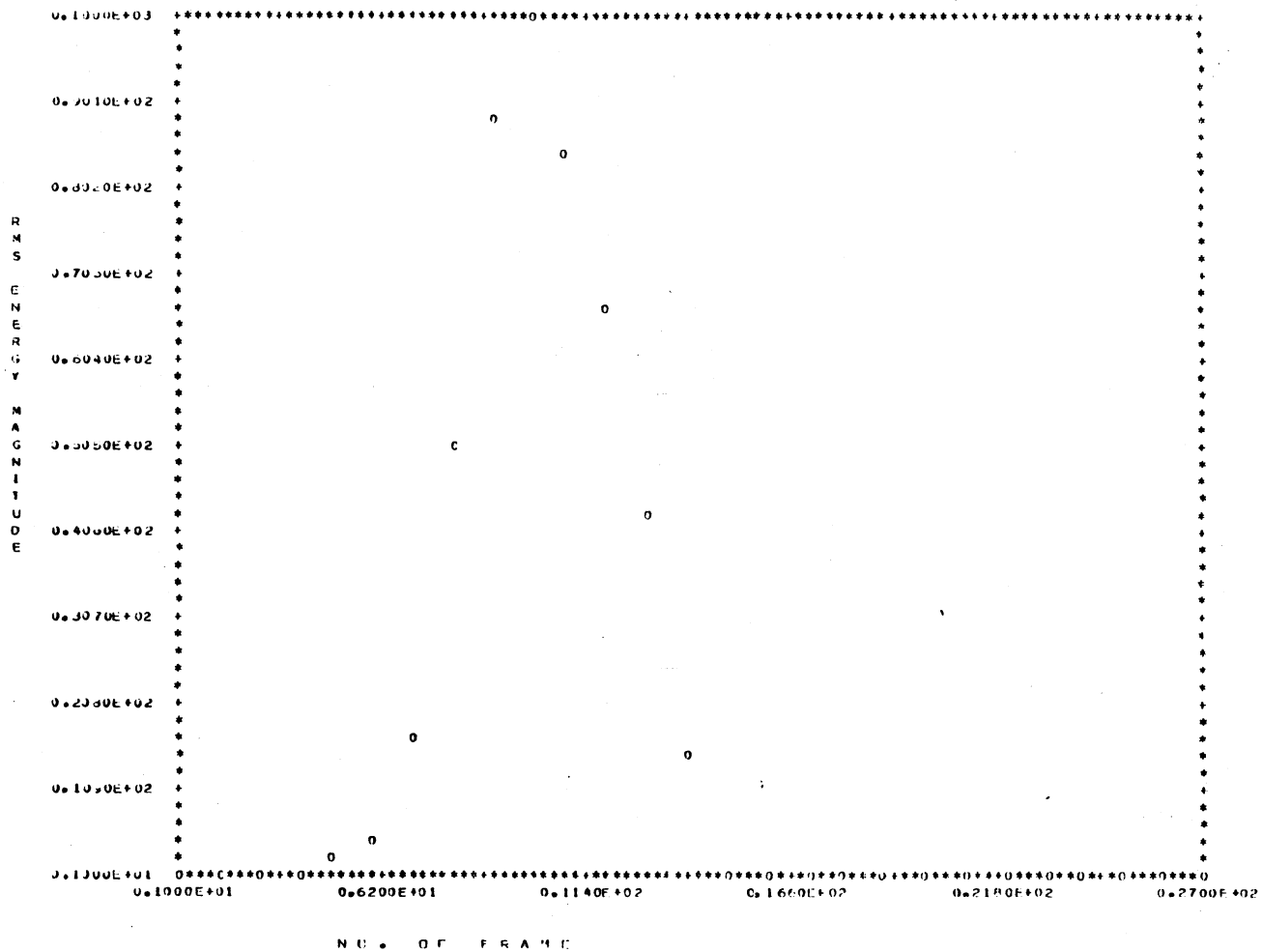


Figure 75. Smoothed and Quantized RMS Energy Contour for Digit One, i.e. /wâhid/ Spoken in Arabic, Sample 2

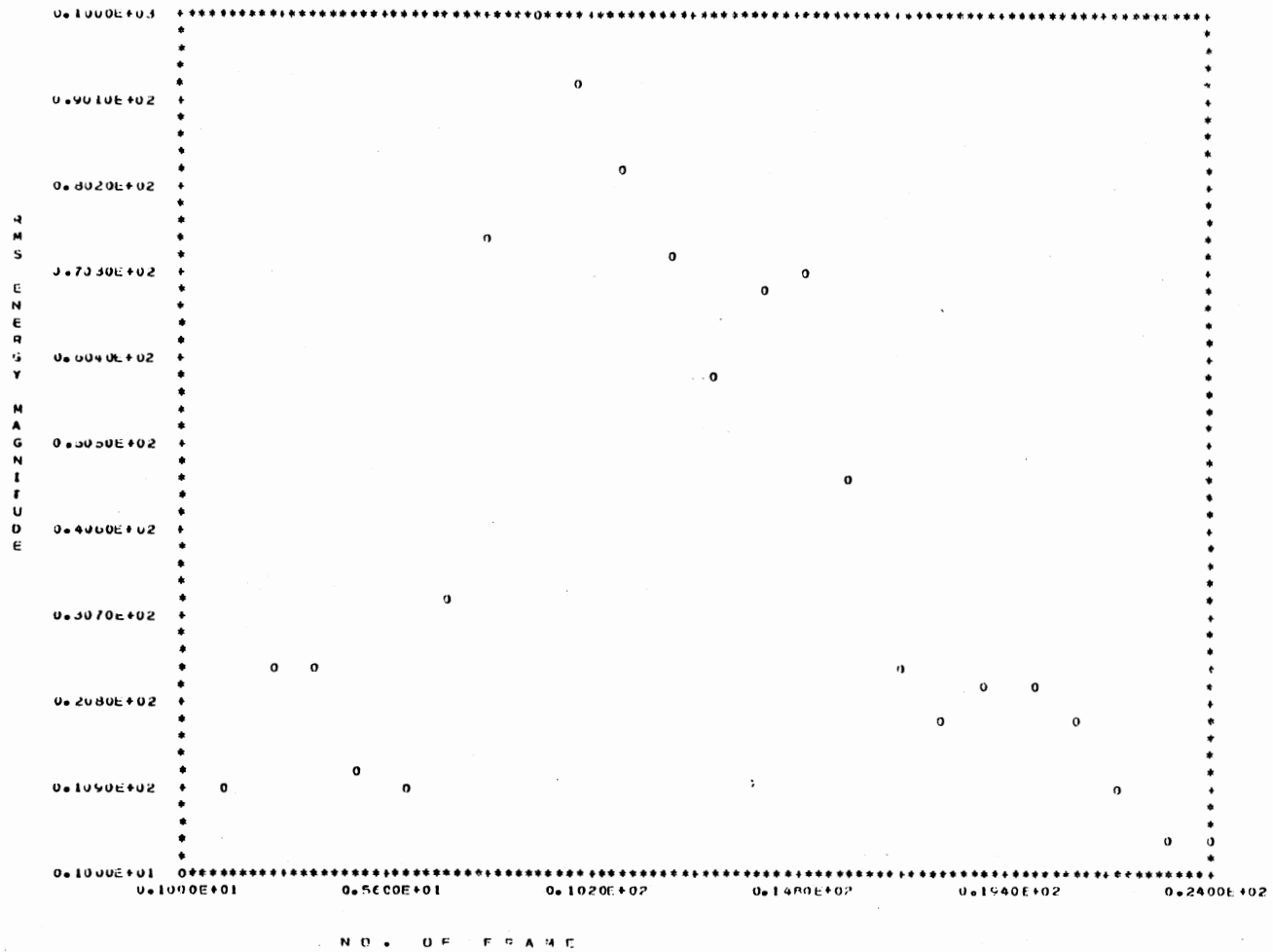


Figure 76. Smoothed and Quantized RMS Energy Contour for Digit Two, i.e. /iθnân/, Spoken in Arabic, Sample 2

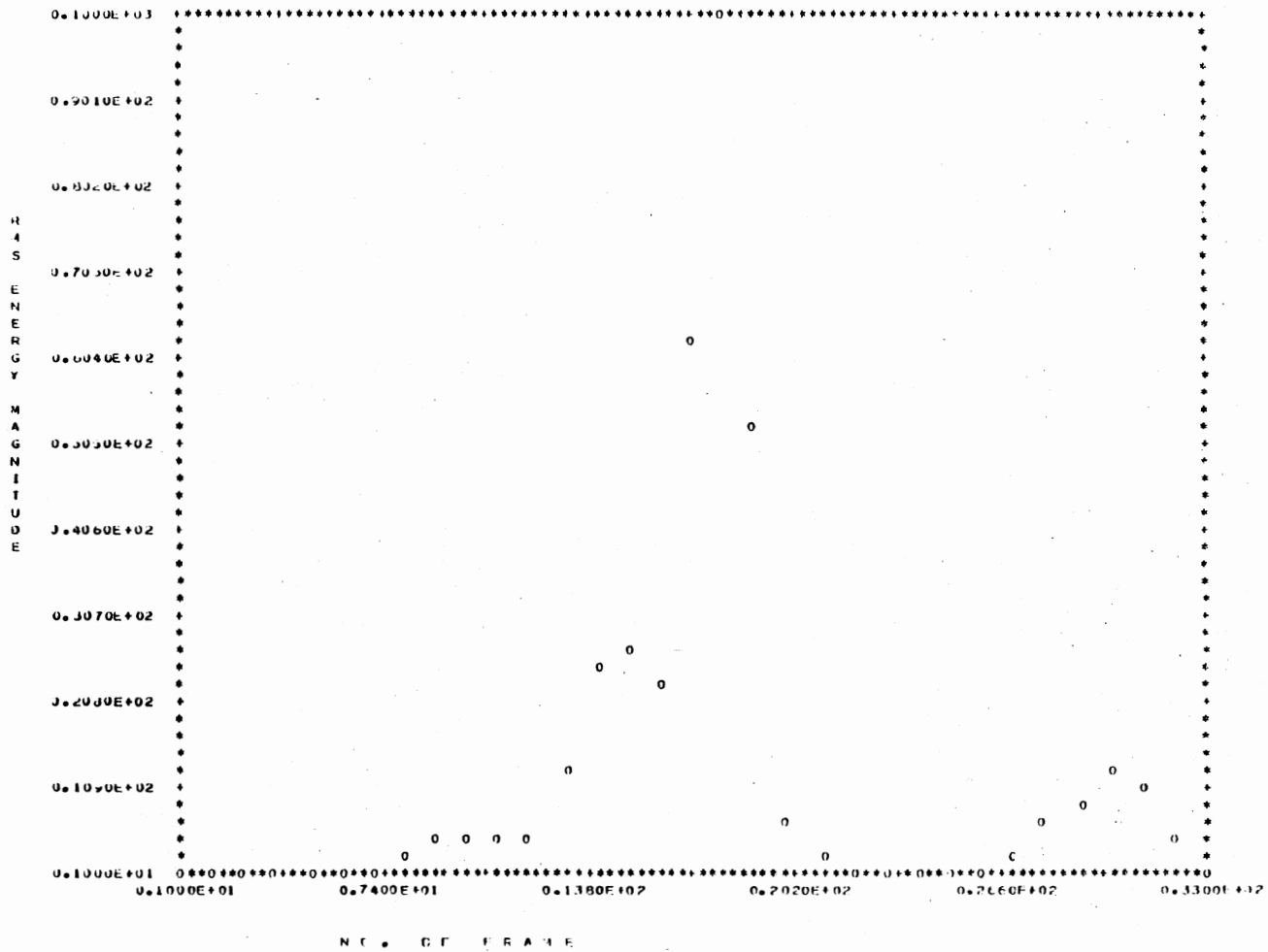


Figure 77. Smoothed and Quantized RMS Energy Contour for Digit Three, i.e. /θalāθθh/ Spoken in Arabic, Sample 2

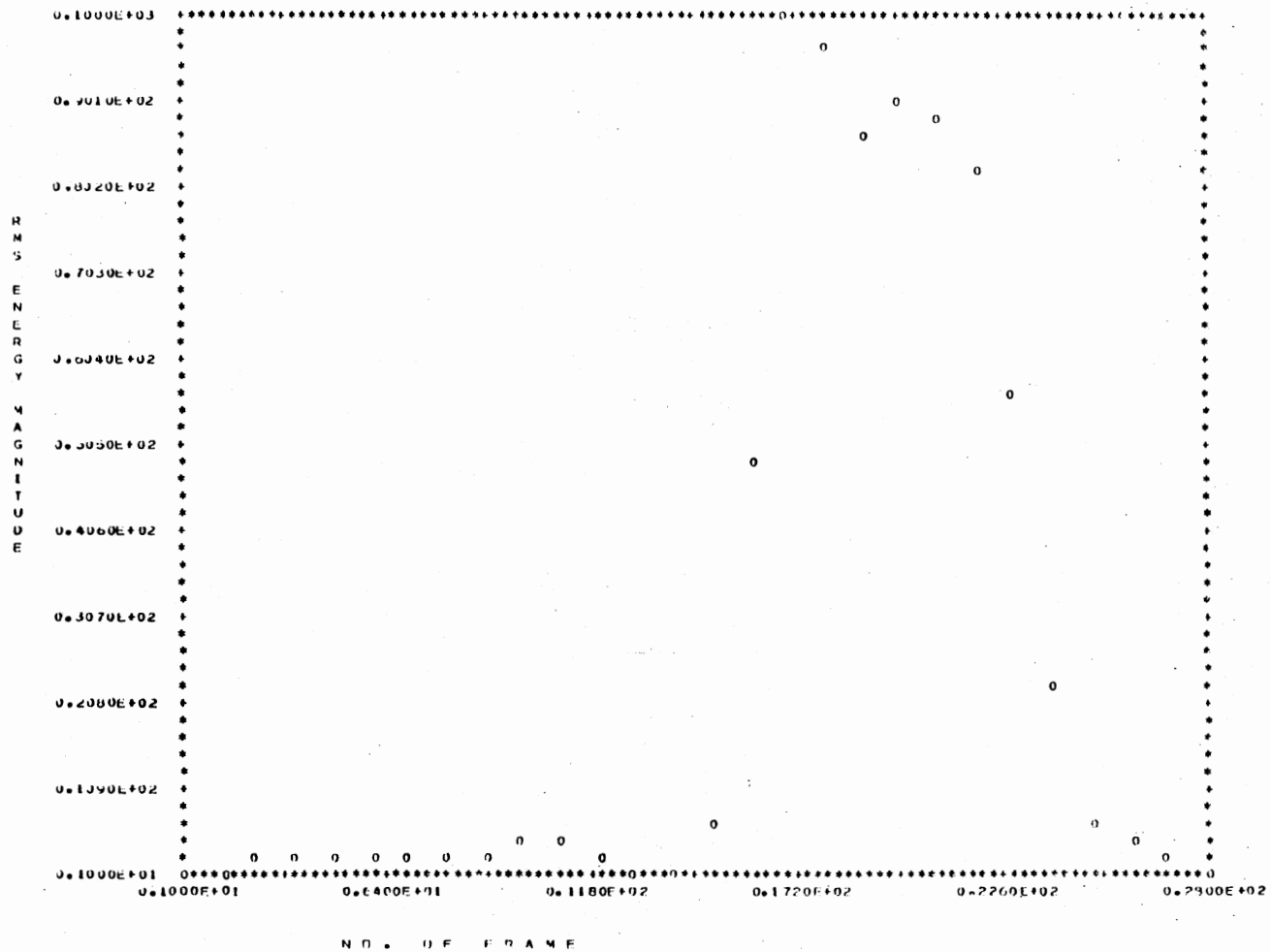


Figure 78. Smoothed and Quantized RMS Energy Contour for Digit Four, i.e. /arbaʁəh/ Spoken in Arabic, Sample 2

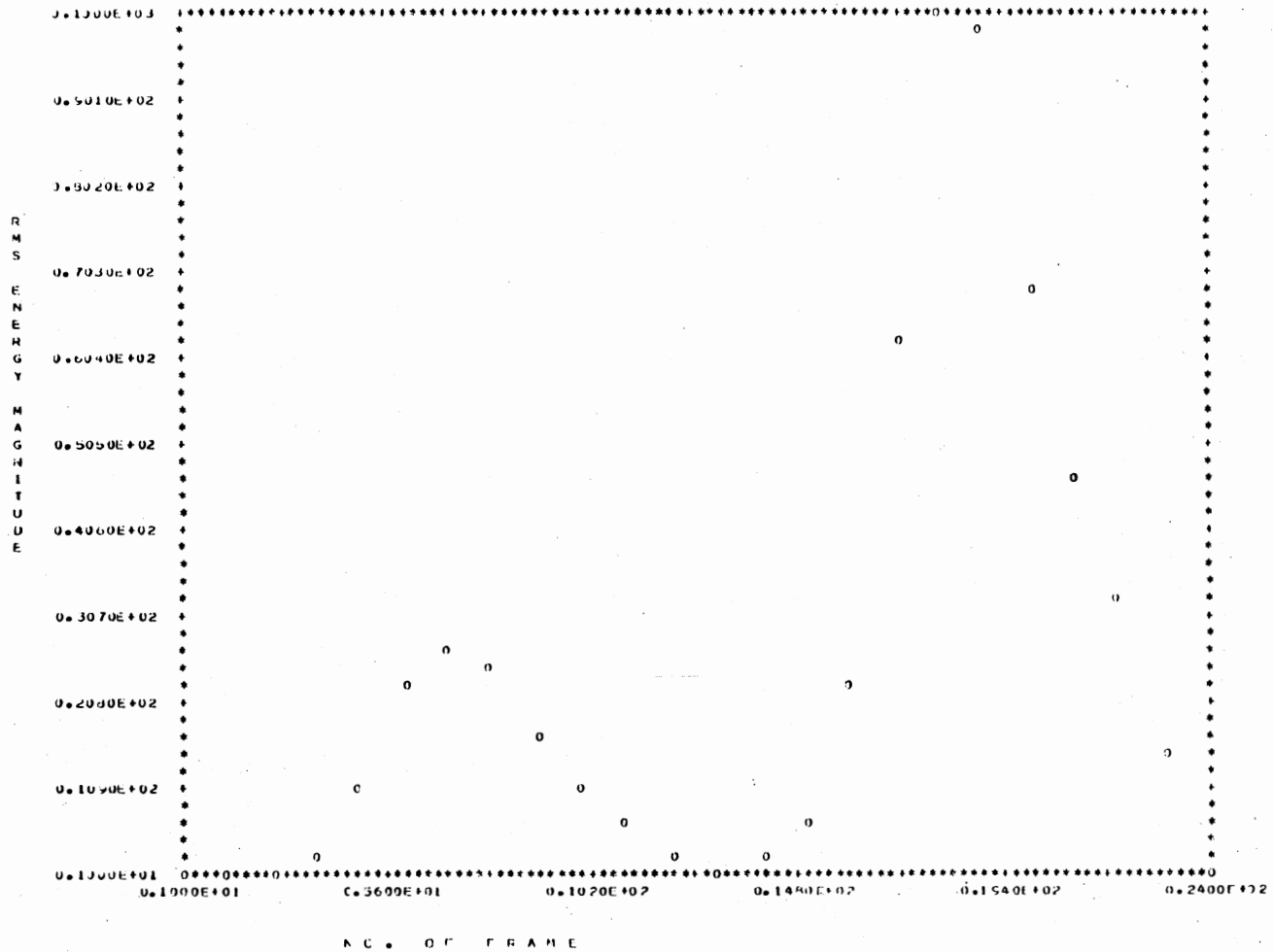


Figure 79. Smoothed and Quantized RMS Energy Contour for Digit Five, i.e. /khΛmsəh/ Spoken in Arabic, Sample 2

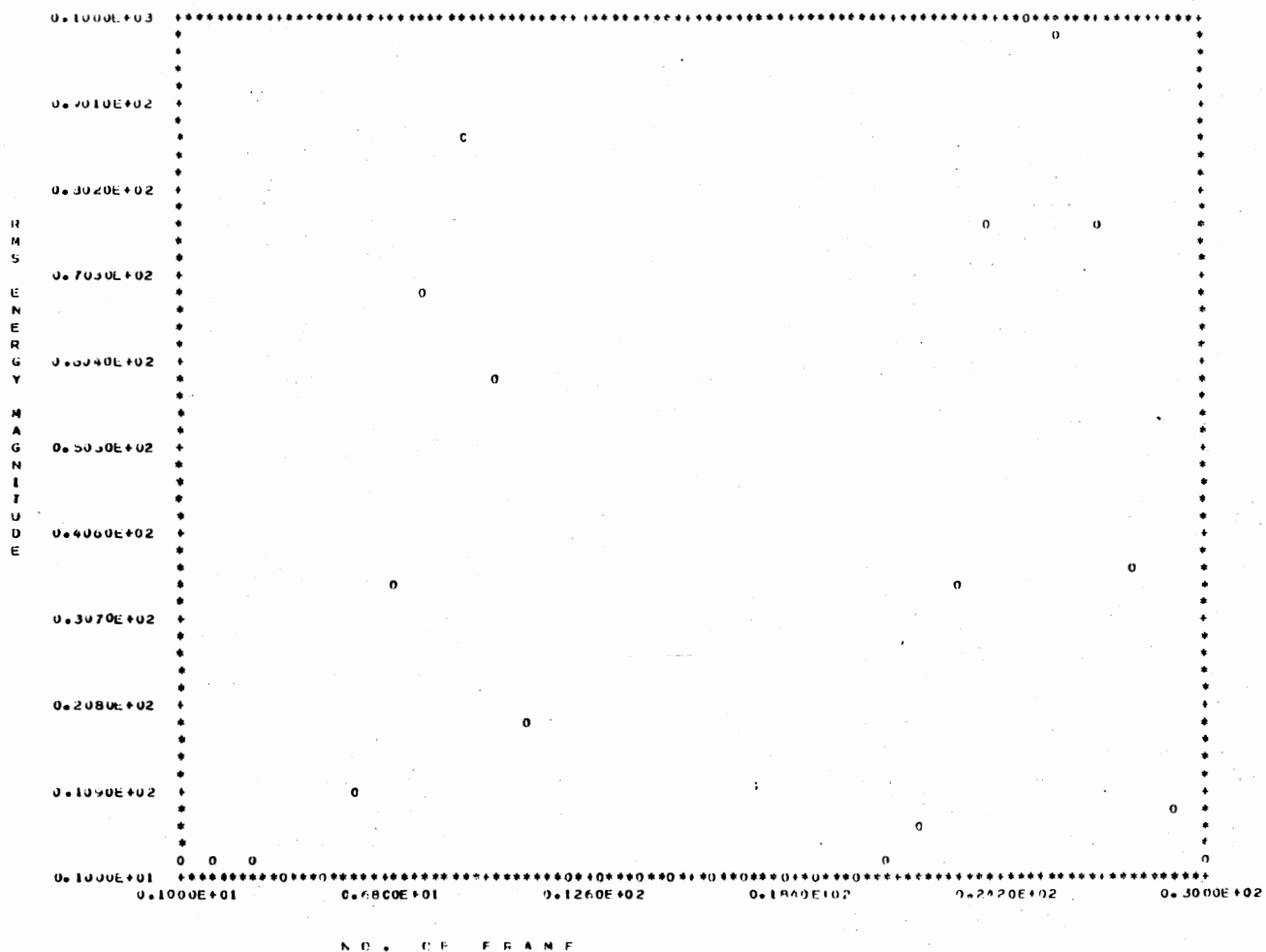


Figure 80. Smoothed and Quantized RMS Energy Contour for Digit Six, i.e. /sIttθh/ Spoken in Arabic, Sample 2

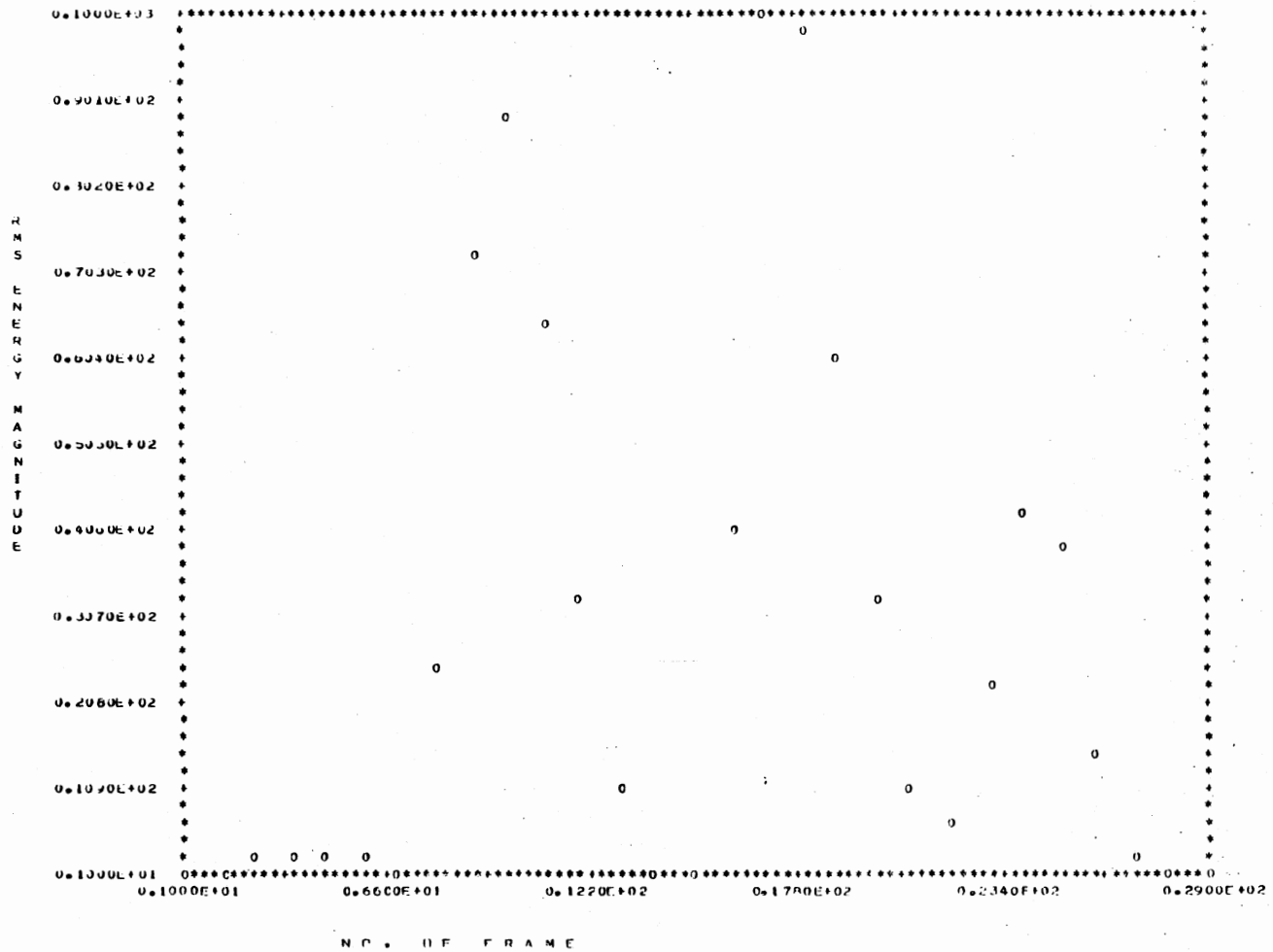


Figure 81. Smoothed and Quantized RMS Energy Contour for Digit Seven, i.e. /sʌbɒðh/ Spoken in Arabic, Sample 2

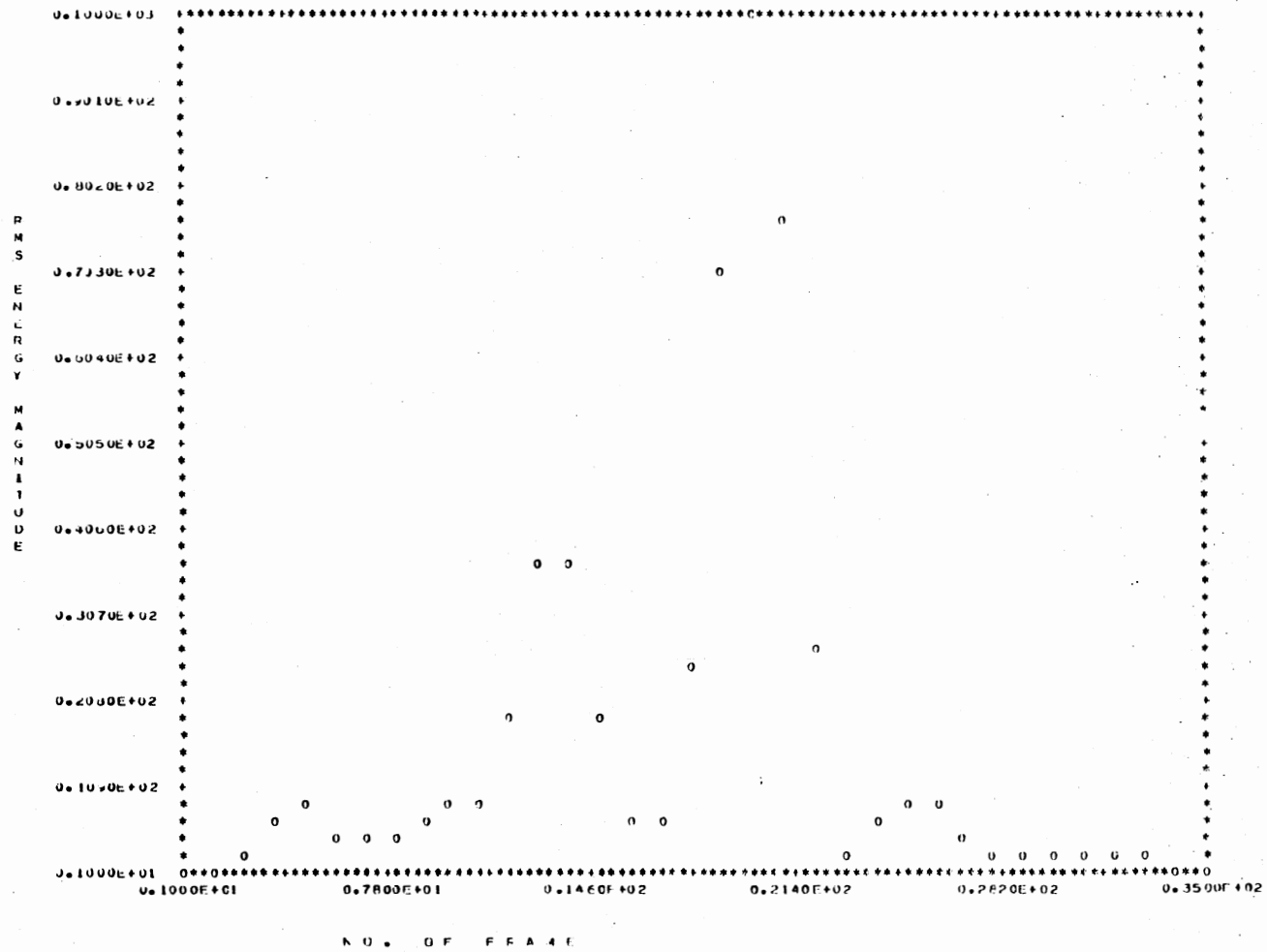


Figure 82. Smoothed and Quantized RMS Energy Contour for Digit Eight, i.e. /θamânyðh/ Spoken in Arabic, Sample 2

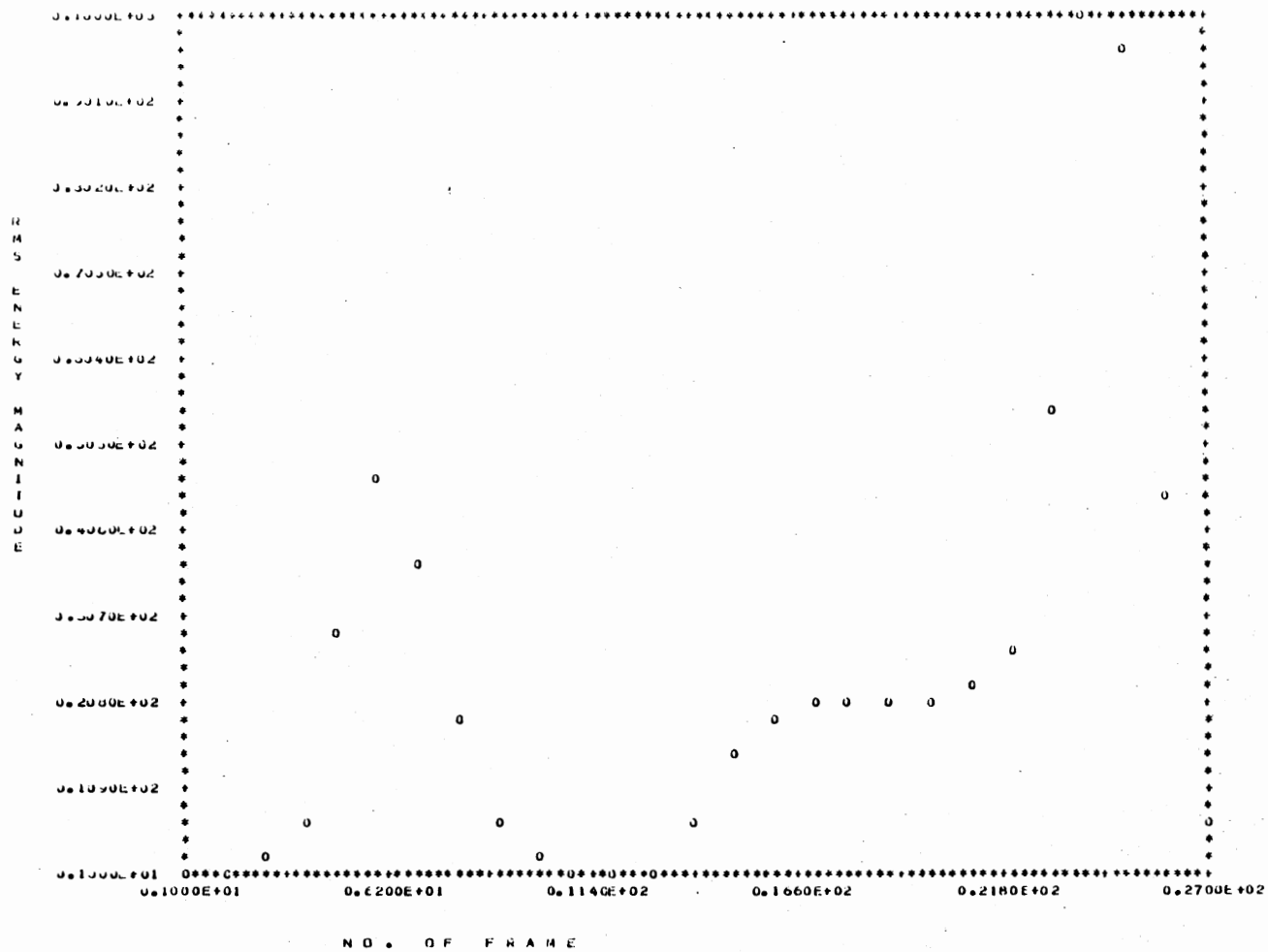


Figure 83. Smoothed and Quantized RMS Energy Contour for Digit Nine, i.e. /tIs ρ h/ Spoken in Arabic, Sample 2

TABLE XVI

SMOOTHED RMS ENERGY PEAKS AND THEIR RATIO'S FOR DIGITS SPOKEN
IN ARABIC WITH VARYING AMPLITUDE AND FIXED MEAN

Digit	P ₁	P ₂	Ratio
/sefr/	11.35	8.64	0.76
/wâhid/	12.39		
/iðnân/	23.11		
/ðalâððh/	27.38	100.67	0.27
/arbaððh/	1.52	9.48	0.16
/khΛmsðh/	46.98	15.97	0.34
/sIttðh/	34.85	32.41	0.93
/sΛbððh/	5.1	8.07	0.63
/ðamânyðh/	5.54	16.97	14.27
/tIspðh/	4.05	13.34	0.30

TABLE XVII
SMOOTHED AND QUANTIZED RMS ENERGY PEAK RATIO'S
FOR DIGITS SPOKEN IN ARABIC

Digit	P_1	P_2	Ratio
/sefr/	10.9	100	0.11
/w ^h ahid/	100		
/i ^h θnan/	70	100	0.7
/θa ^h la ^h θθh/	100	17	0.17
/arbaρθh/	100	90	0.9
/kh ^h Λmsθh/	27	100	0.27
/s ^h I ^h ttθh/	87	100	0.87
/s ^h Λbρθh/	89	100	0.89
/θa ^h manyθh/	37	100	0.37
/tI ^h spθh/	46	100	0.46

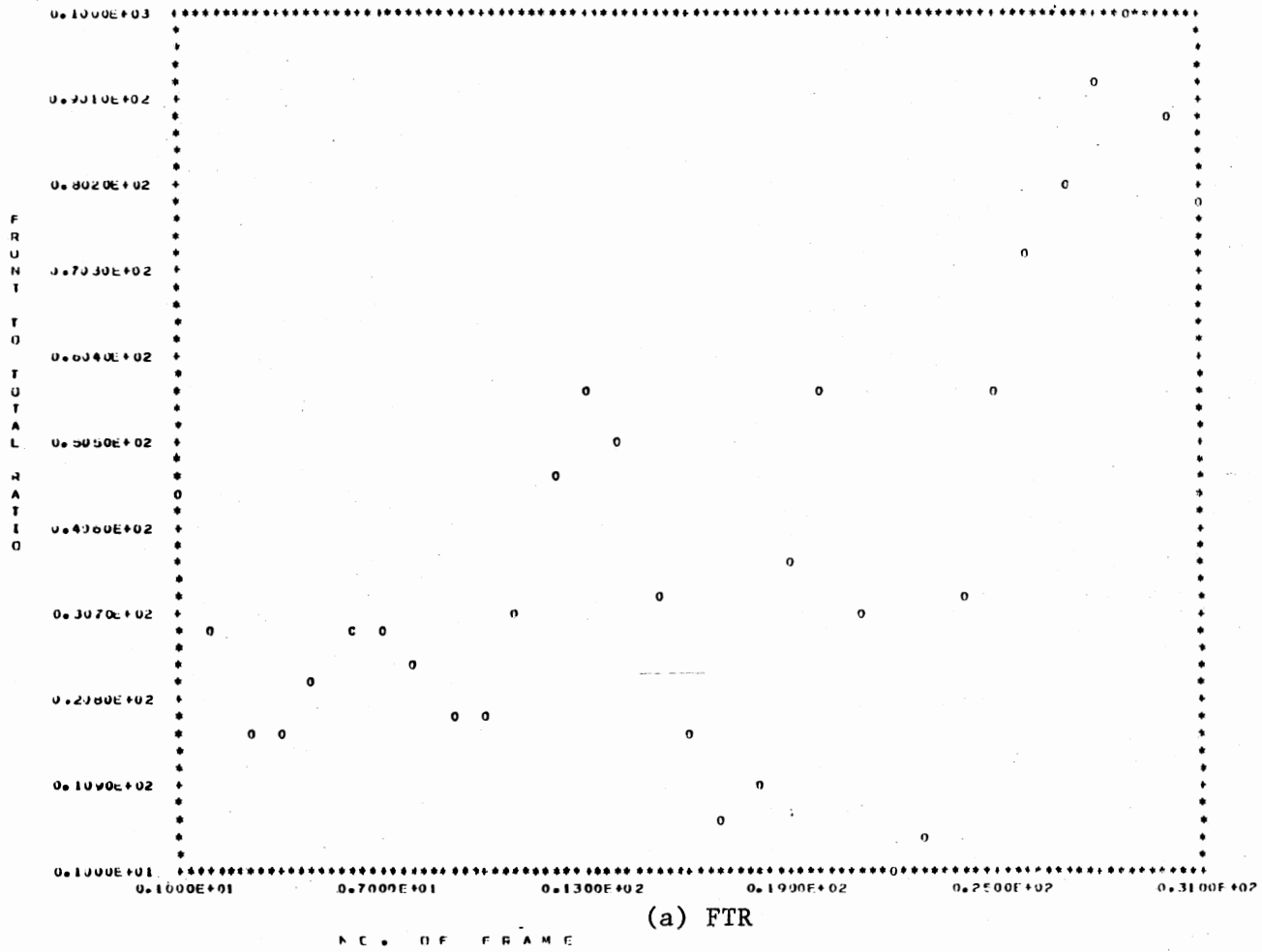


Figure 84. Smoothed and Quantized Feature Parameter for Digit Zero, i.e. /sefr/ Spoken in Arabic

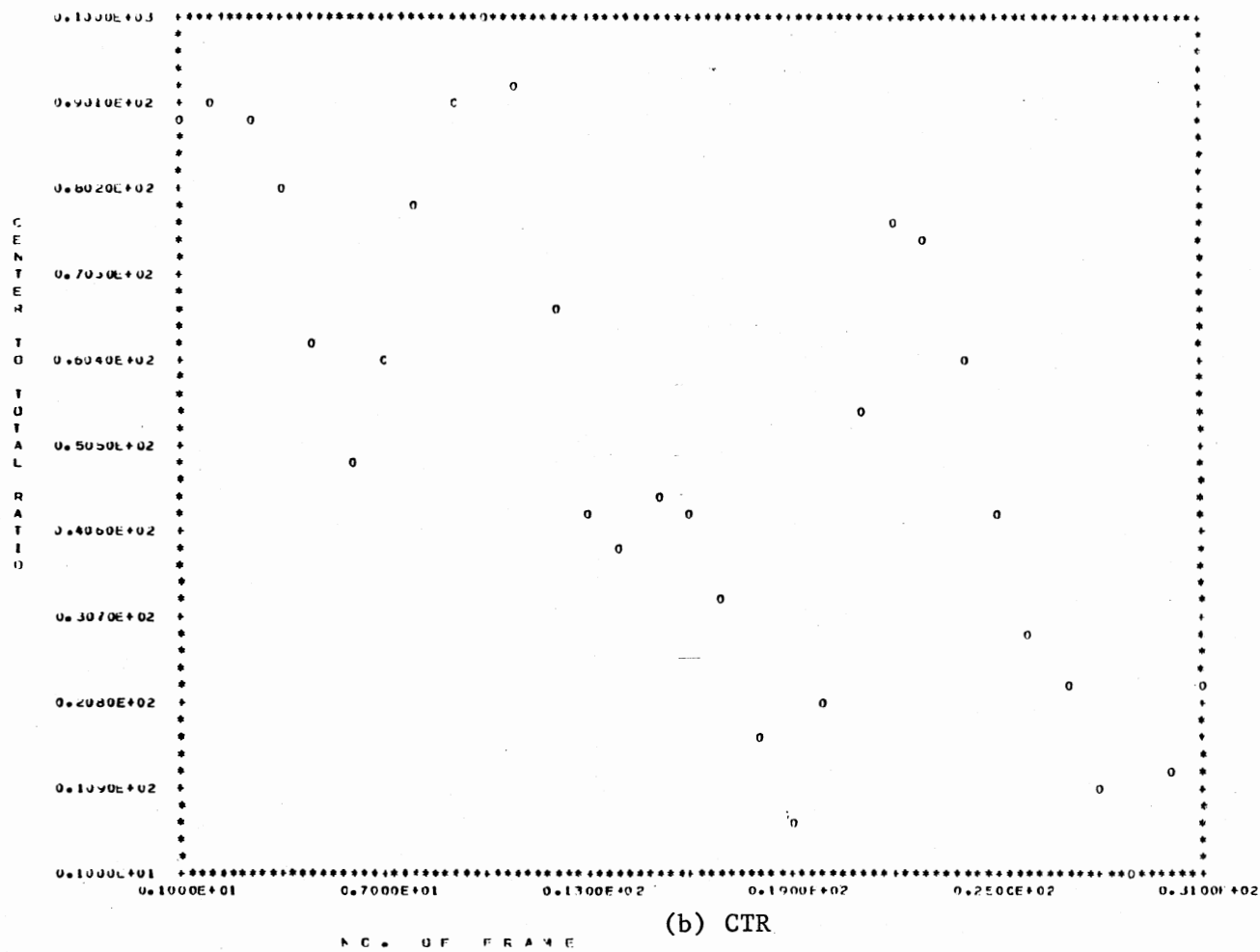


Figure 84. (Continued)

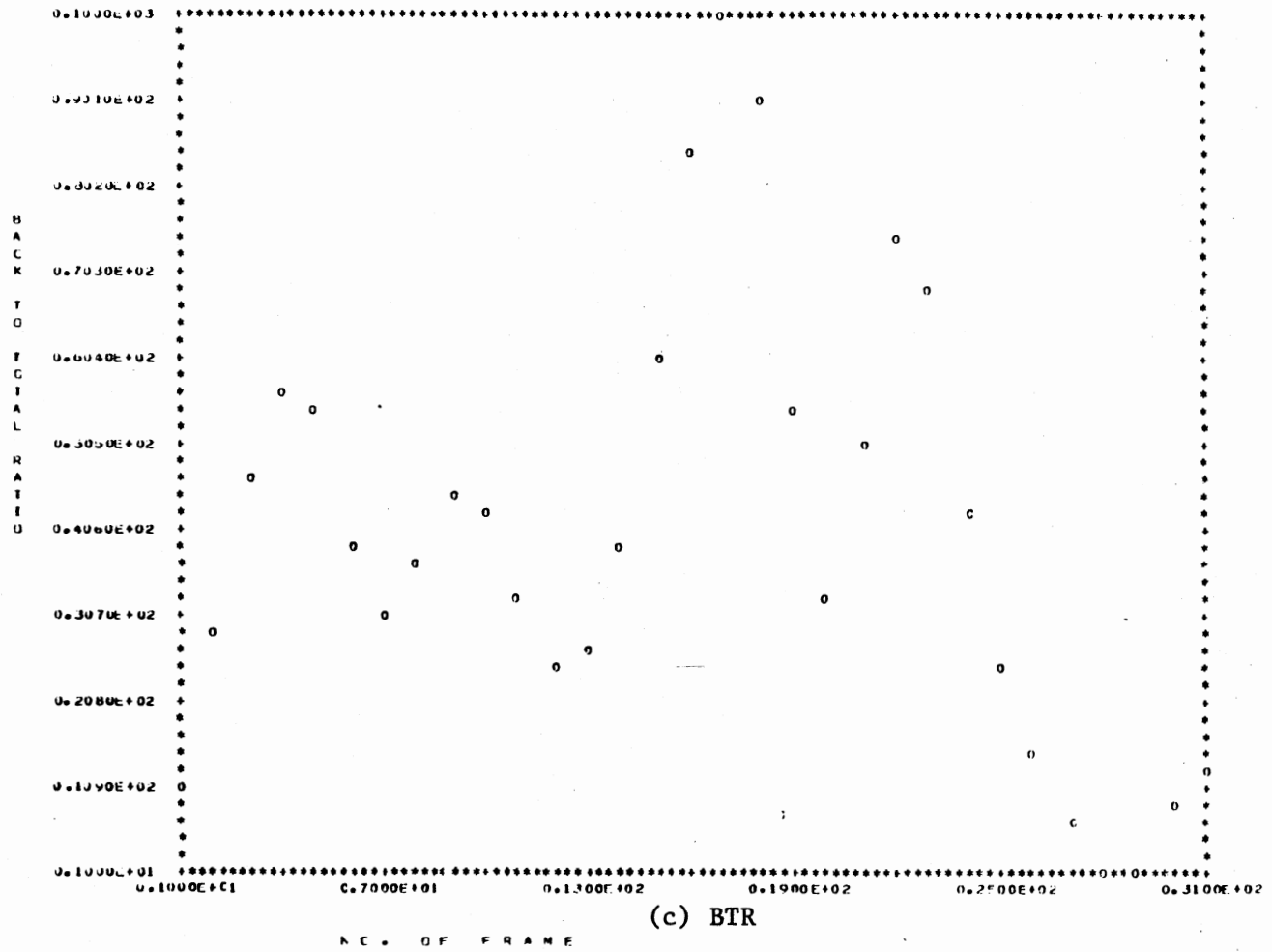


Figure 84. (Continued)

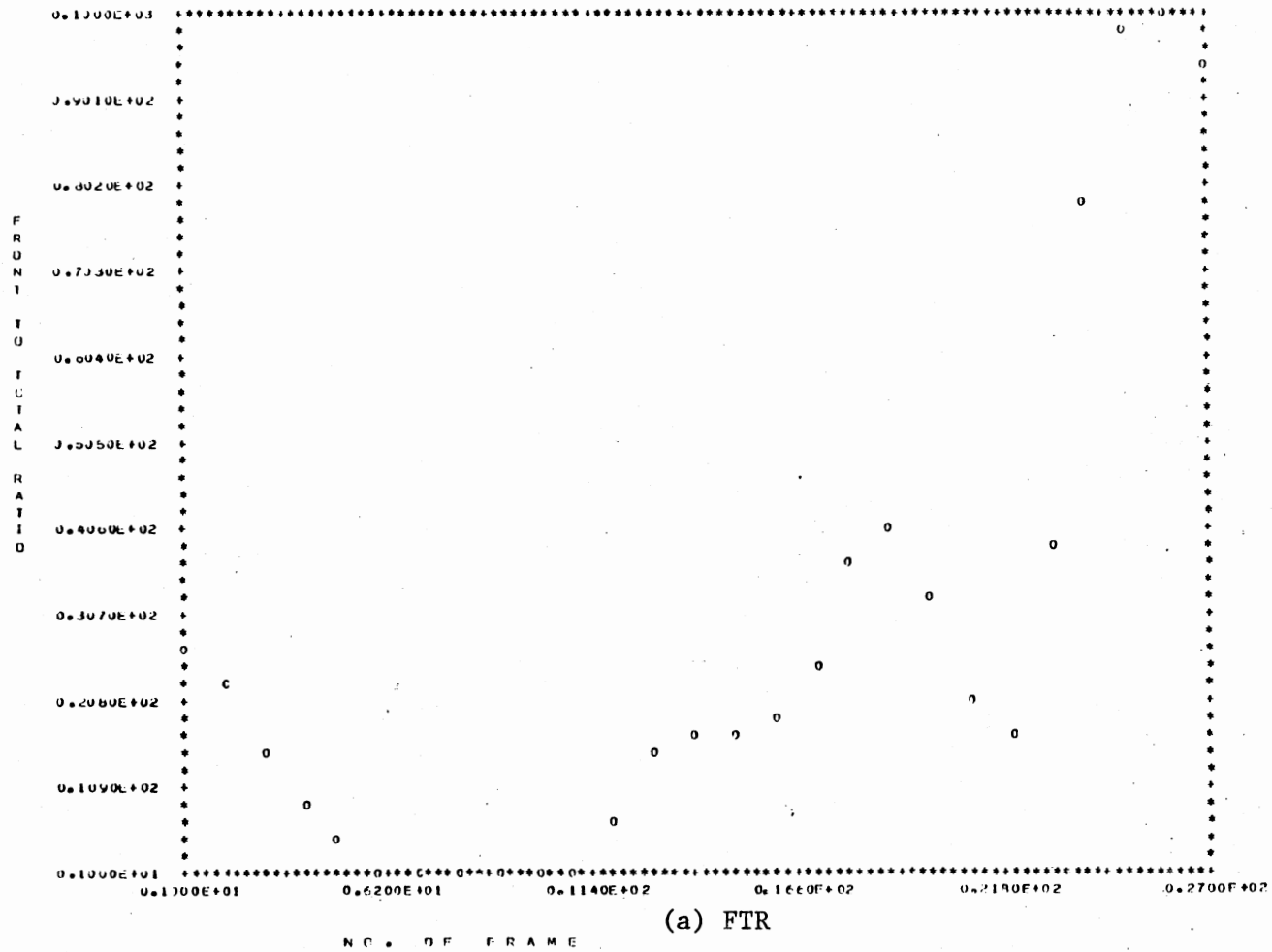


Figure 85. Smoothed and Quantized Feature Parameter for Digit One, i.e. /wâhid/ Spoken in Arabic

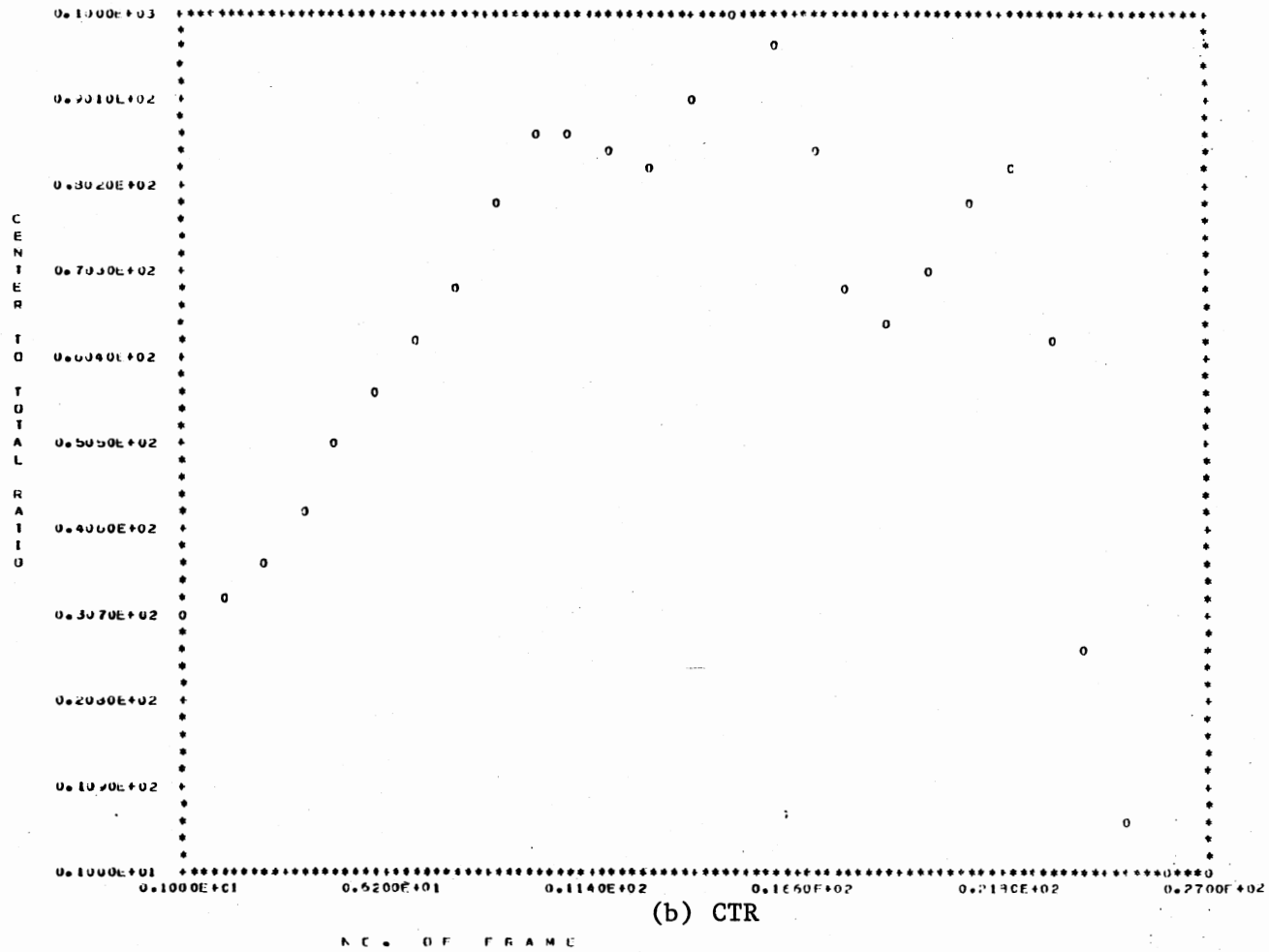


Figure 85. (Continued)

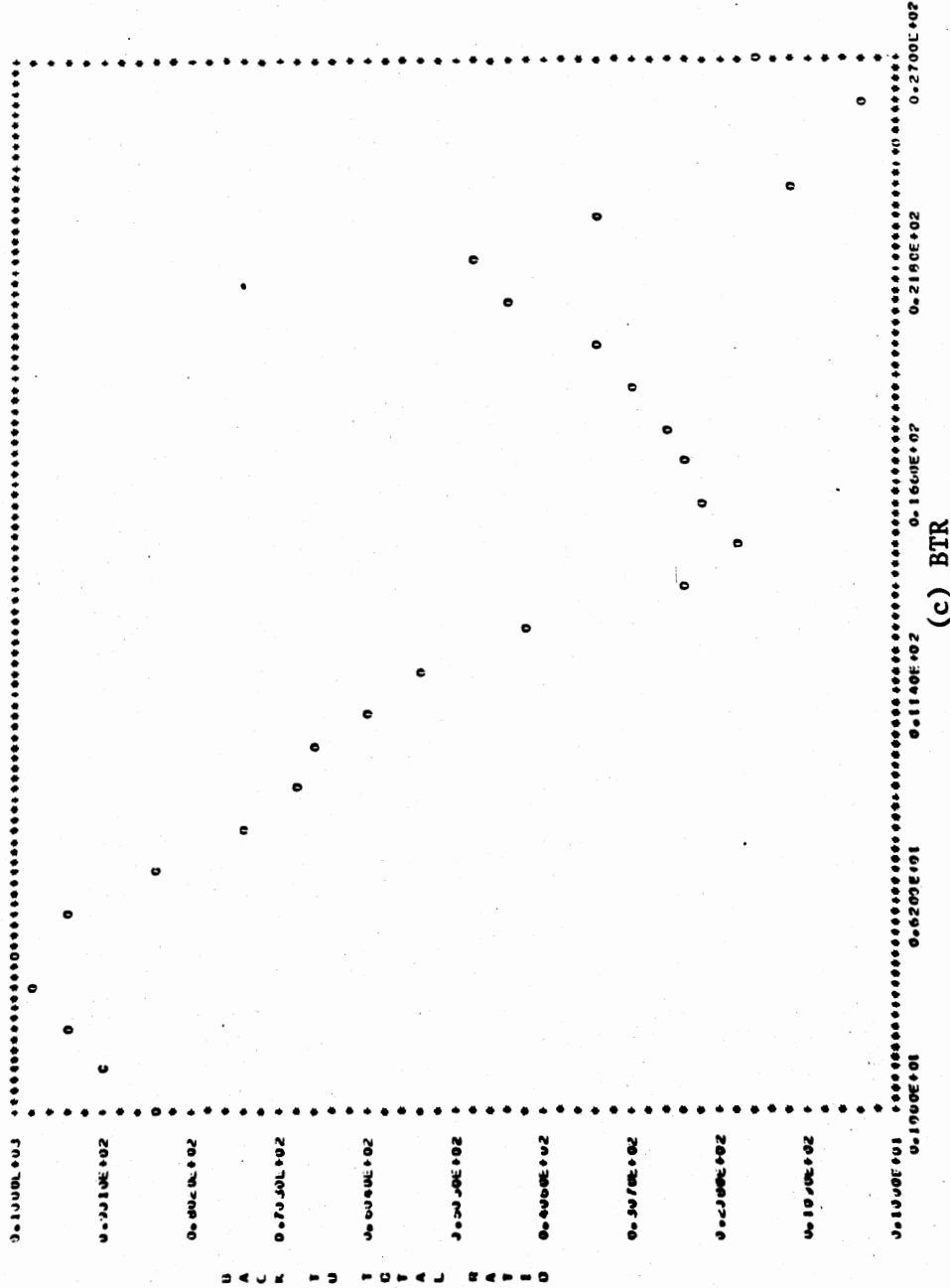


Figure 85. (Continued)

A C . O F F R A M E

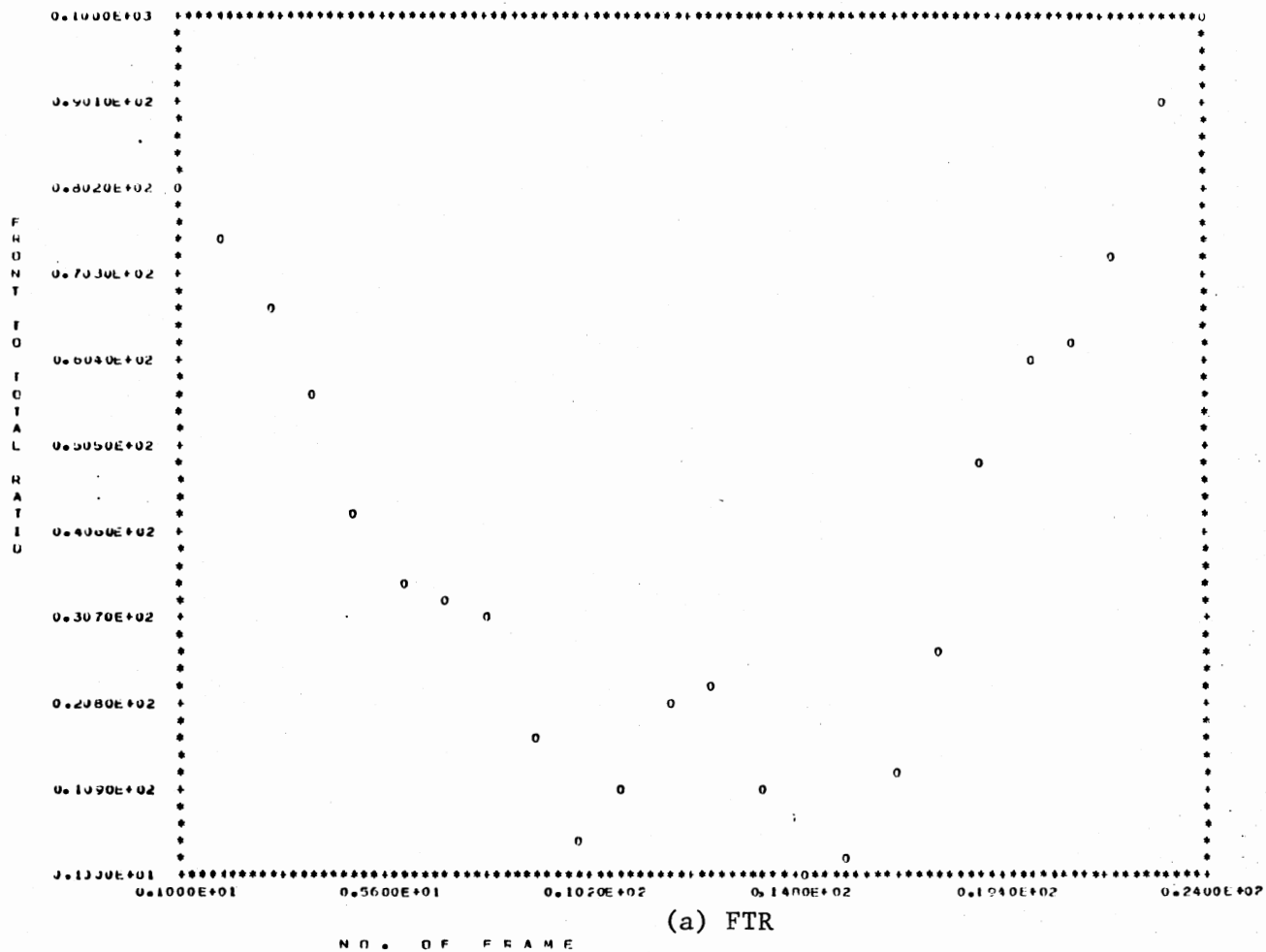


Figure 86. Smoothed and Quantized Feature Parameter for Digit Two, i.e. /iθnān/ Spoken in Arabic

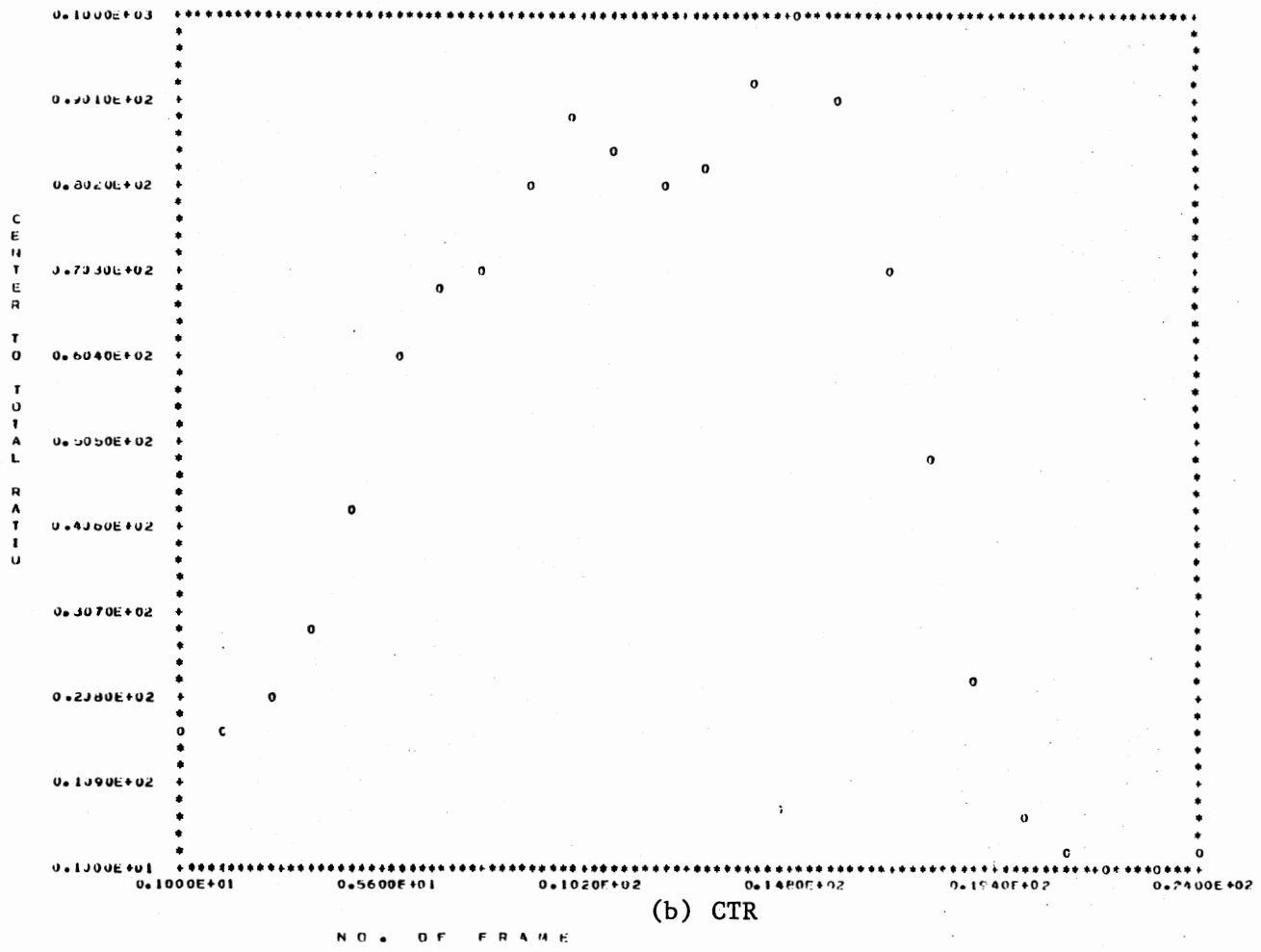


Figure 86. (Continued)

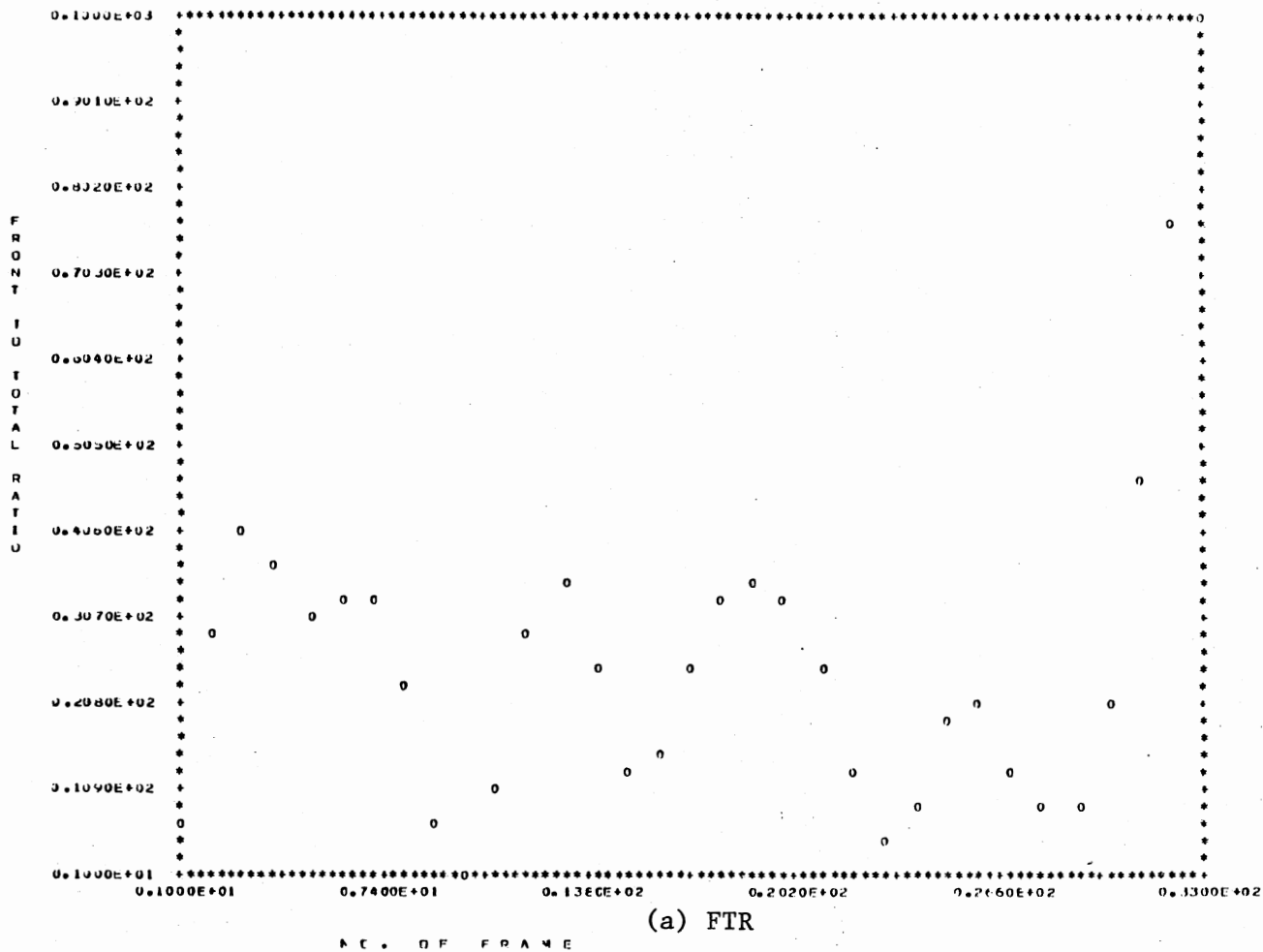
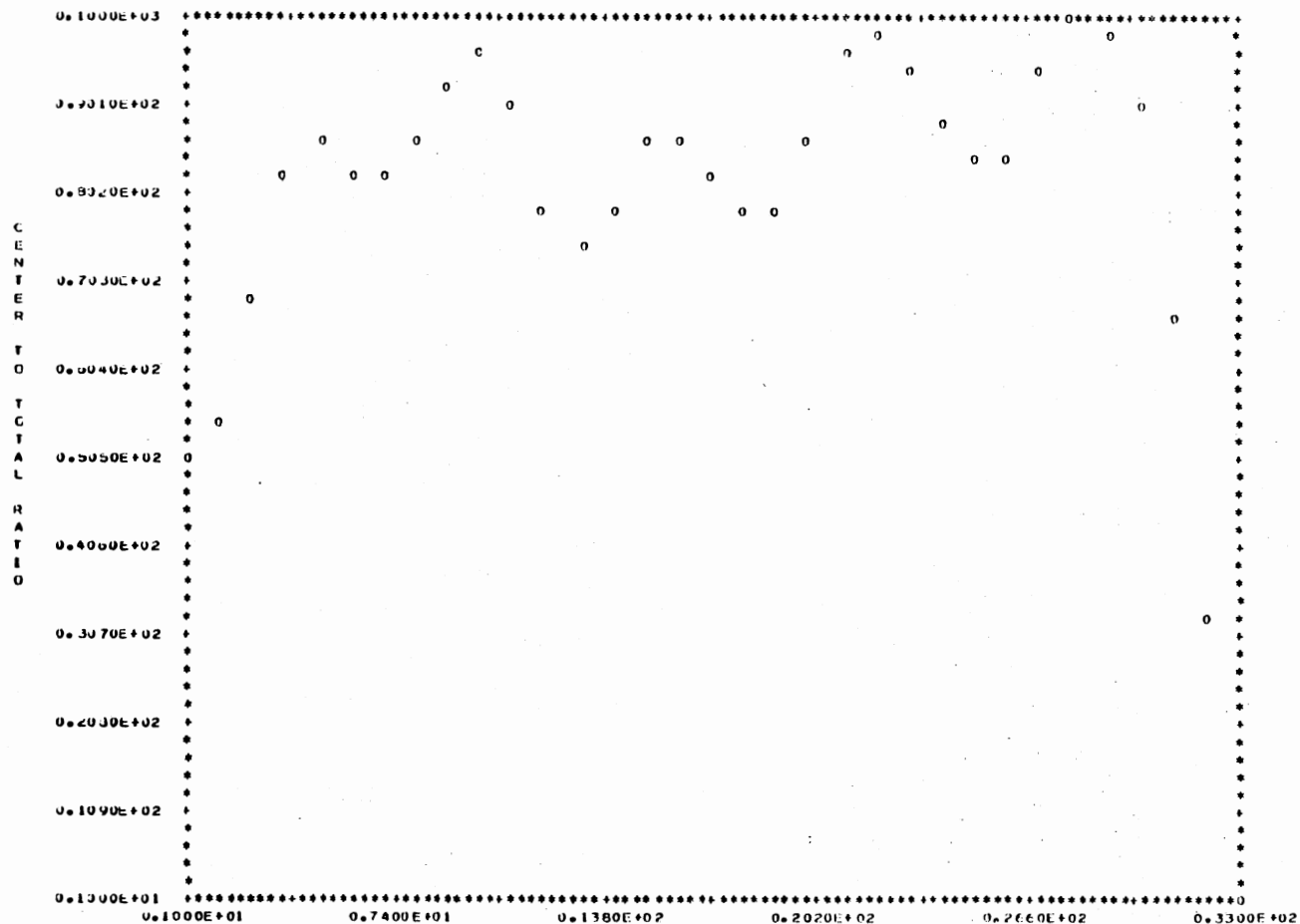


Figure 87. Smoothed and Quantized Feature Parameter for Digit Three, i.e./θaλaθθh/ Spoken in Arabic



(b) CTR

Figure 87. (Continued)

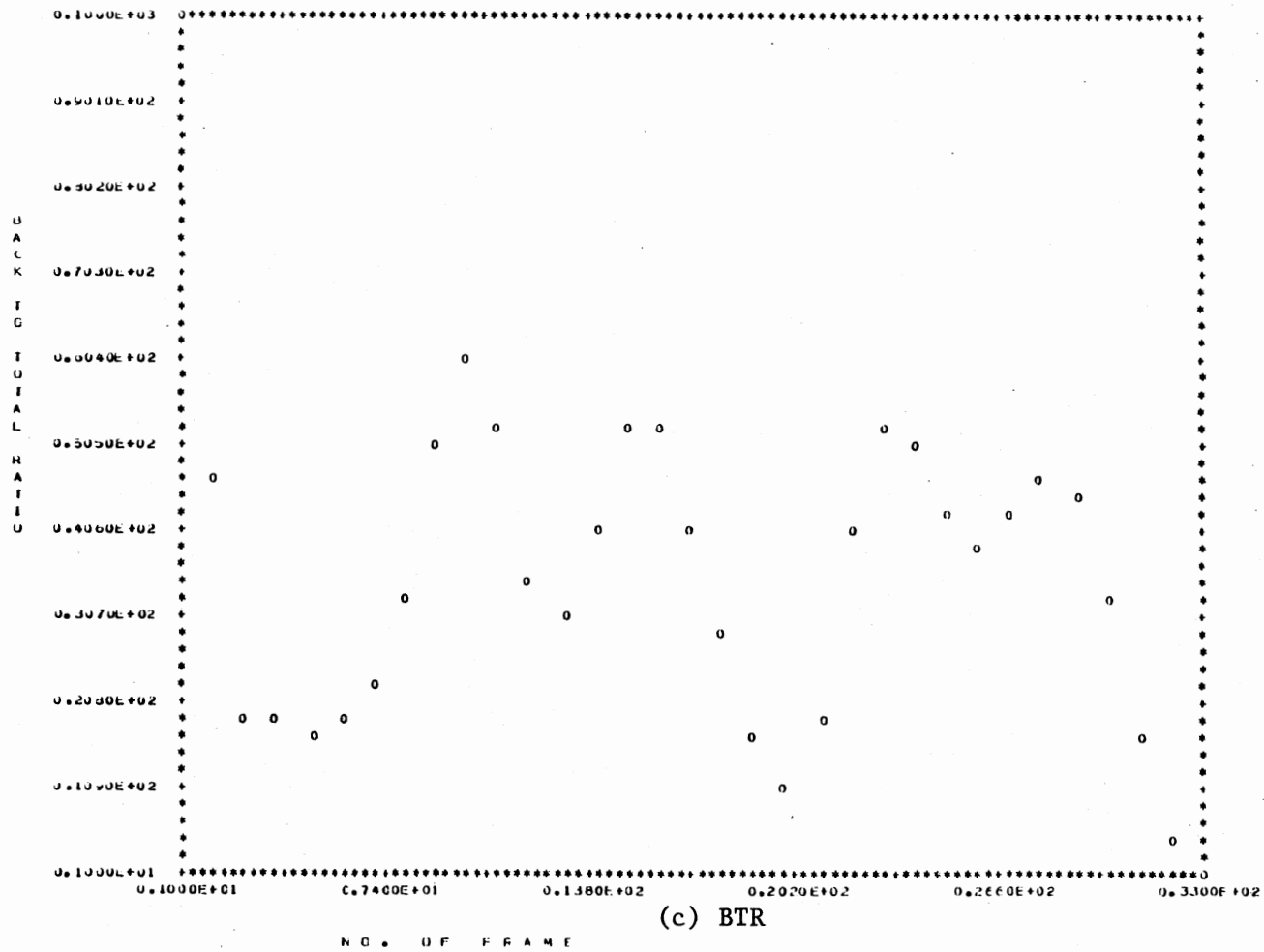


Figure 87. (Continued)

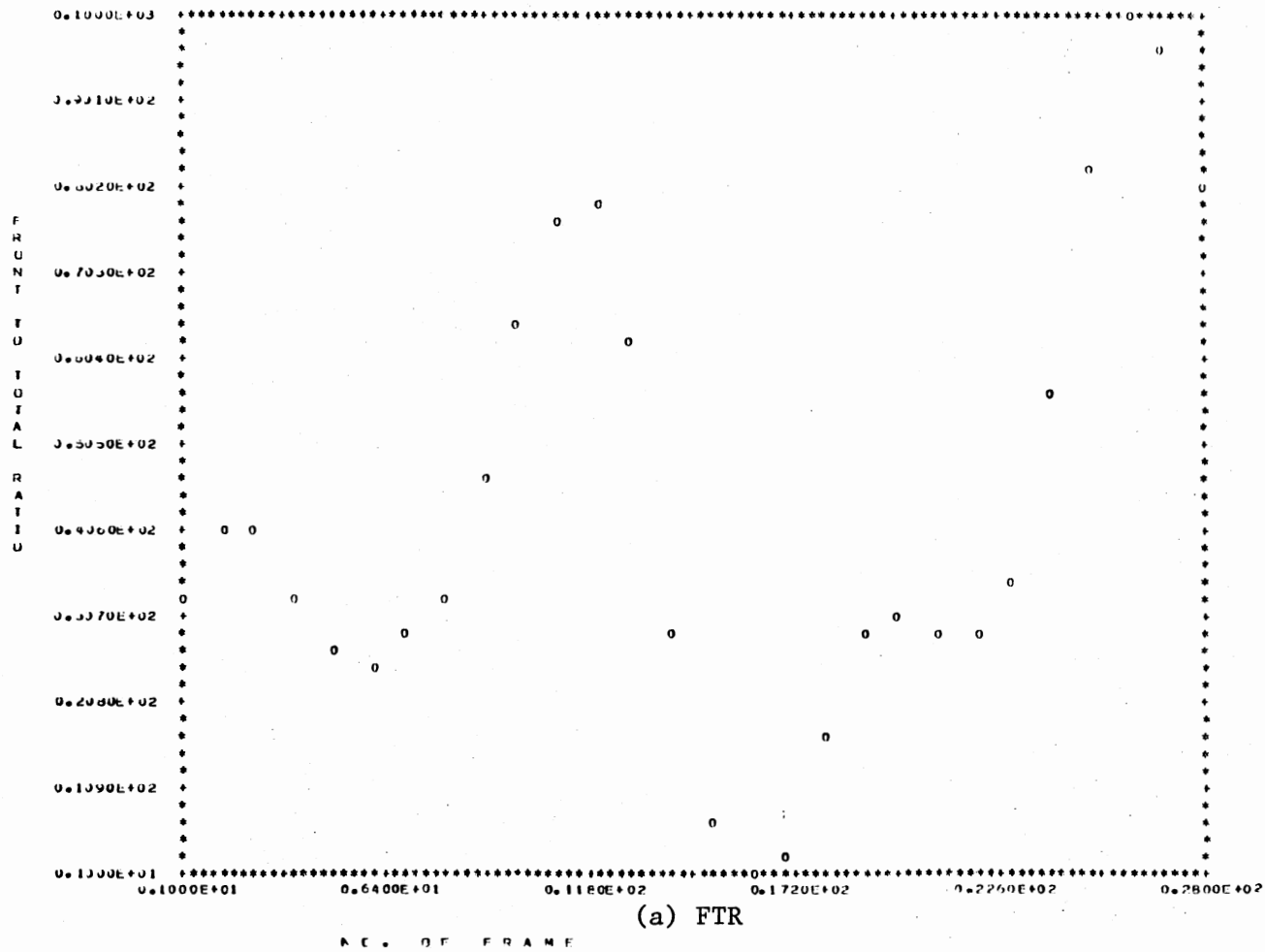


Figure 88. Smoothed and Quantized Feature Parameter for Digit Four, i.e. /arbaðh/ Spoken in Arabic

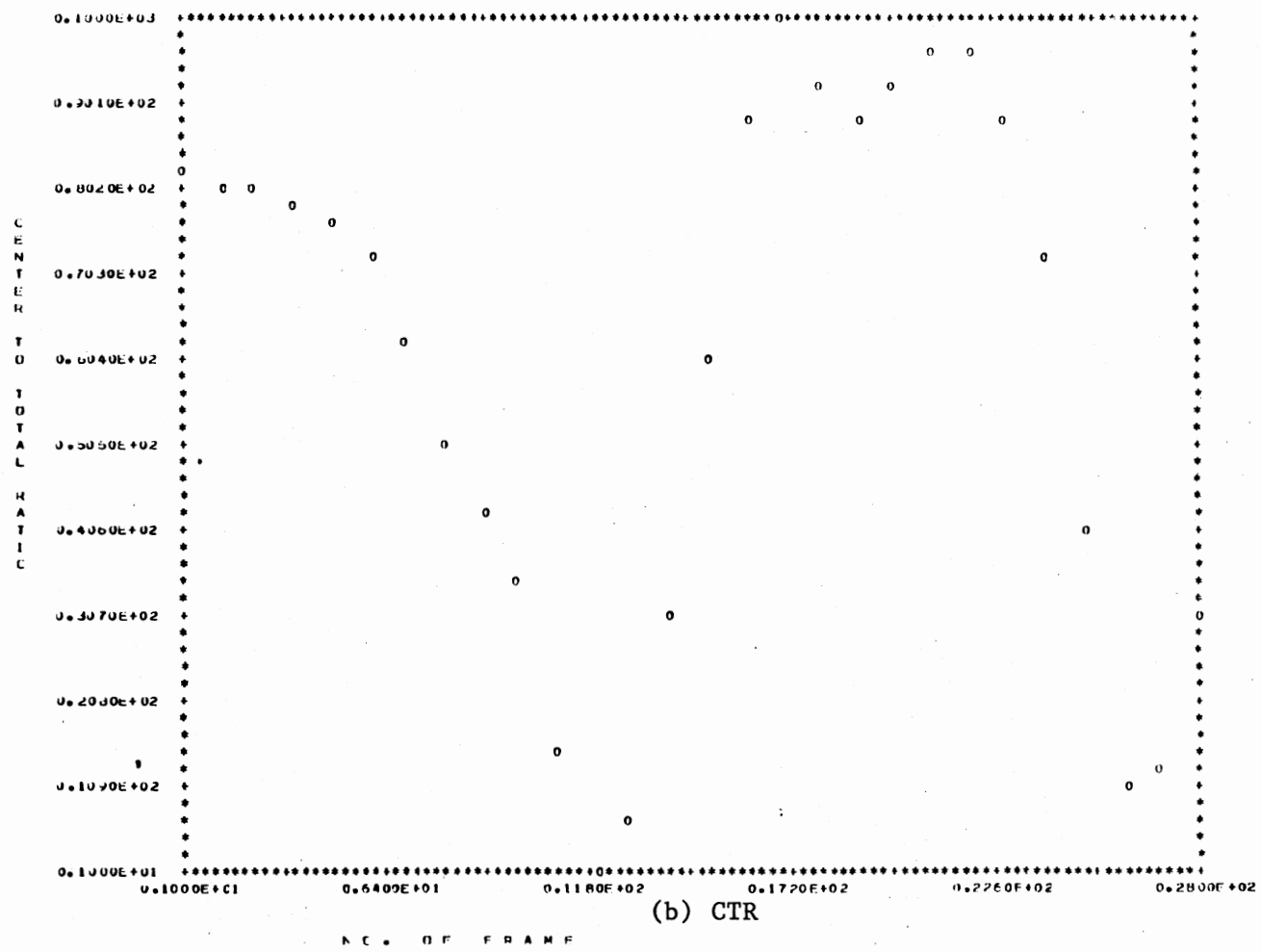


Figure 88. (Continued)

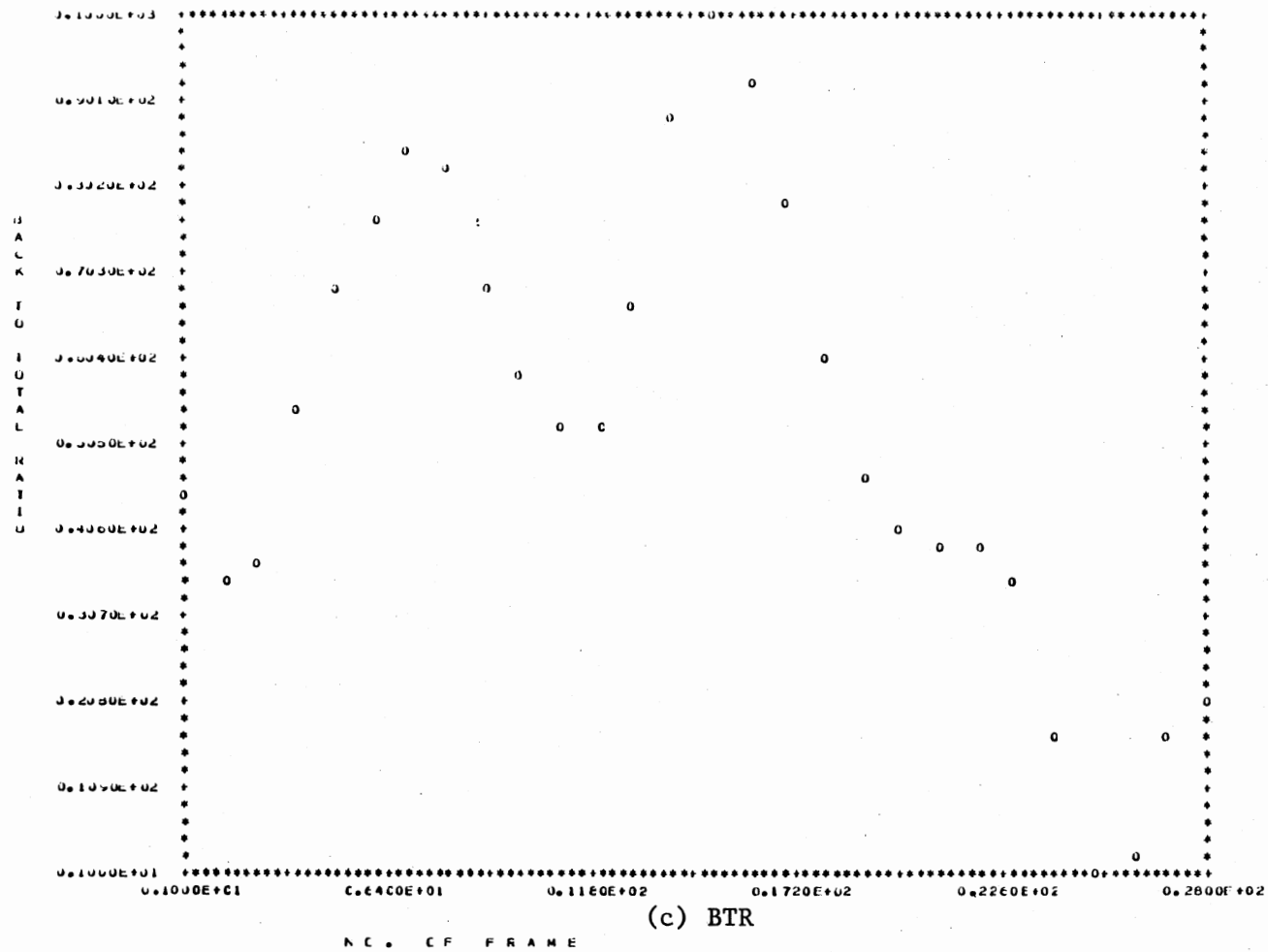


Figure 88. (Continued)

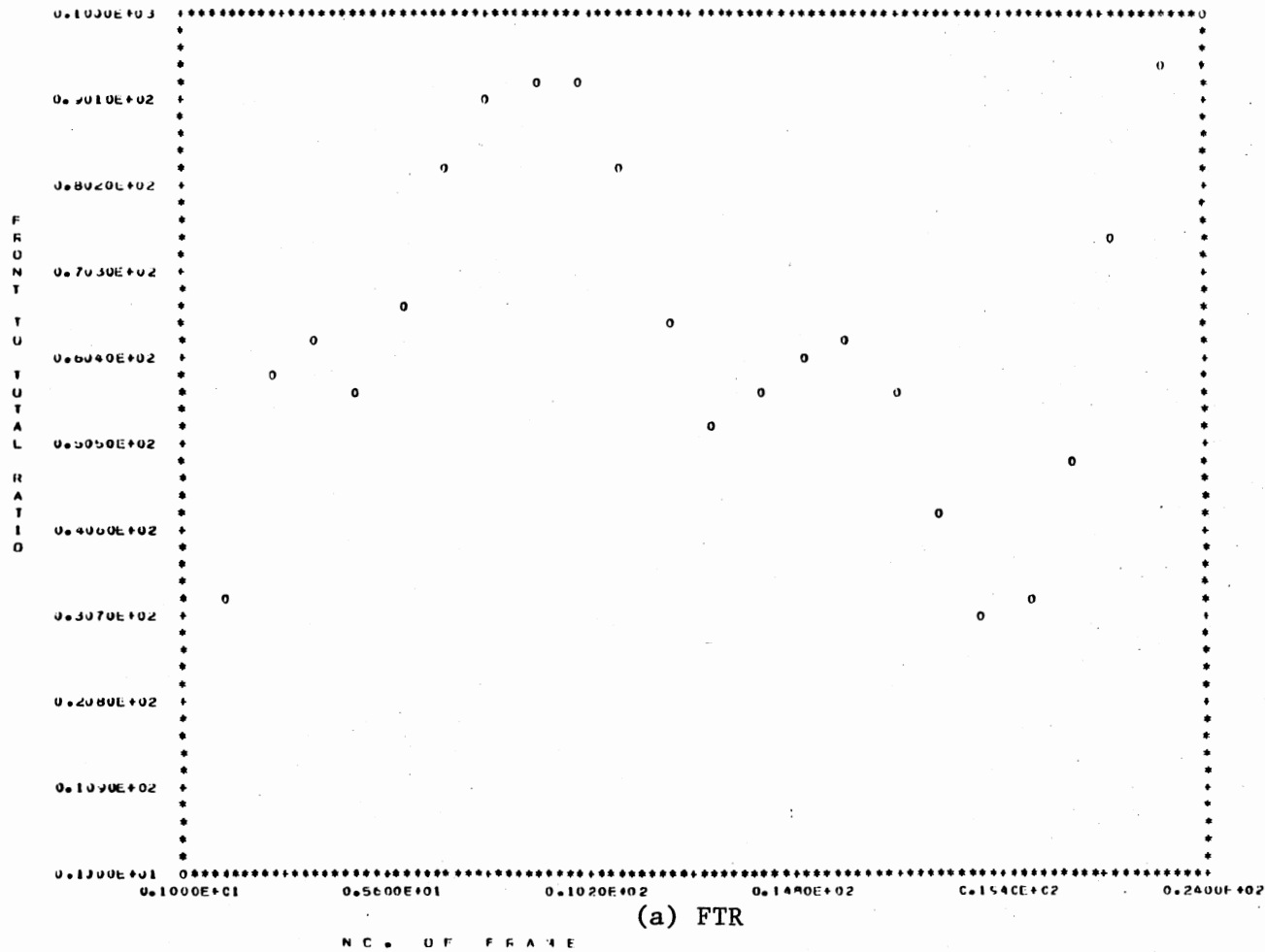


Figure 89. Smoothed and Quantized Feature Parameter for Digit Five, i.e. /khʌmsəh/ Spoken in Arabic

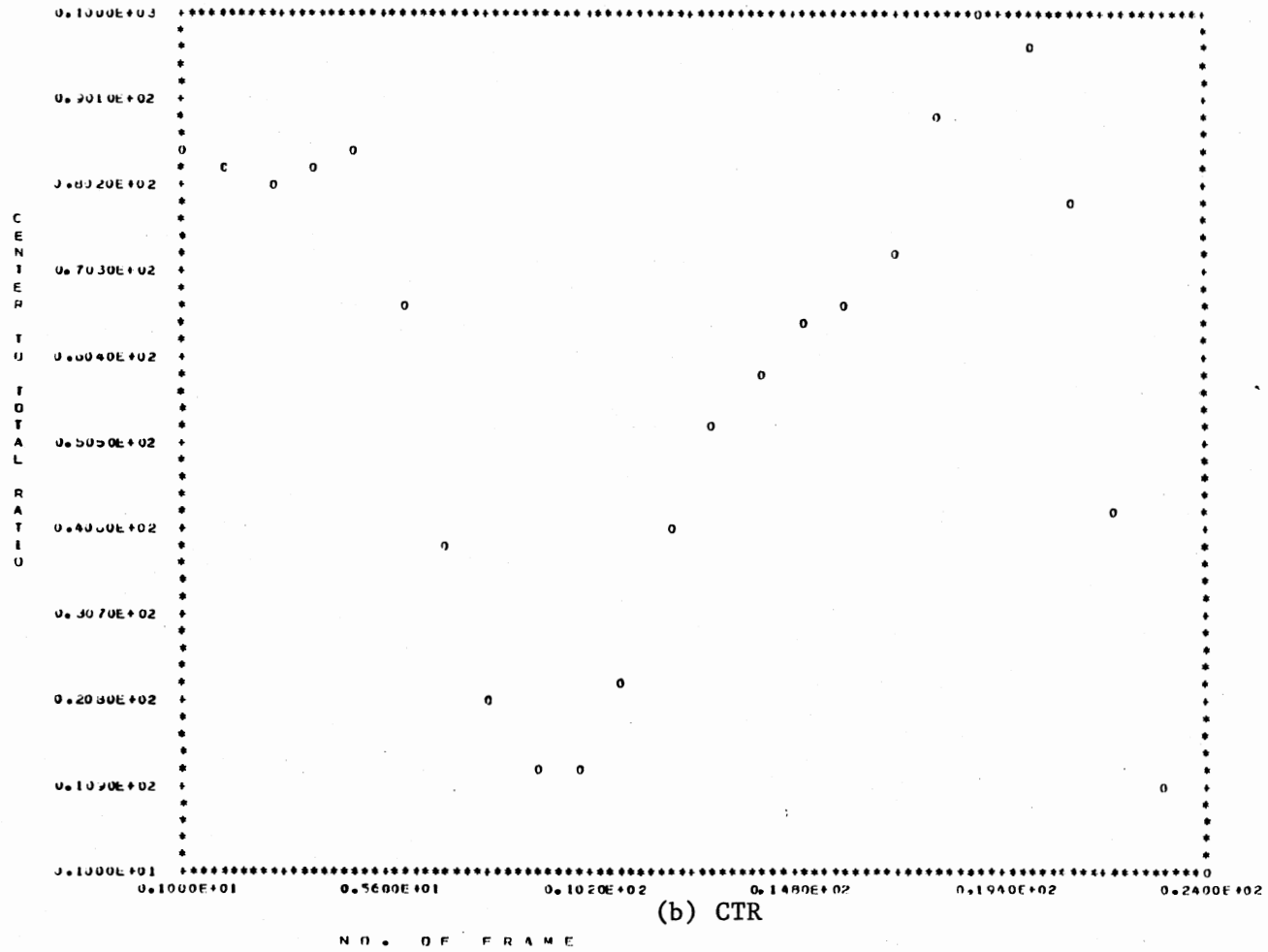


Figure 89. (Continued)

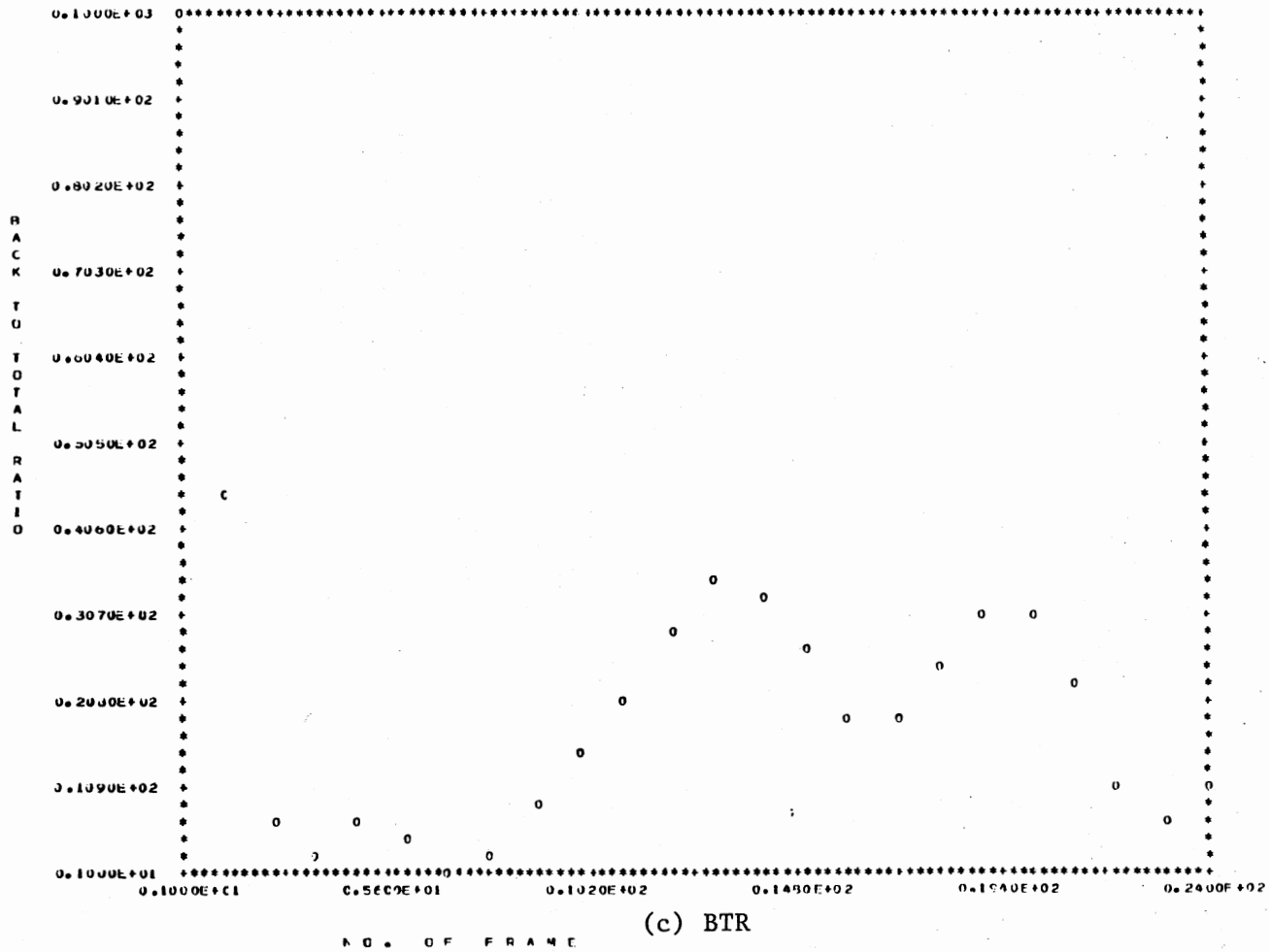


Figure 89. (Continued)

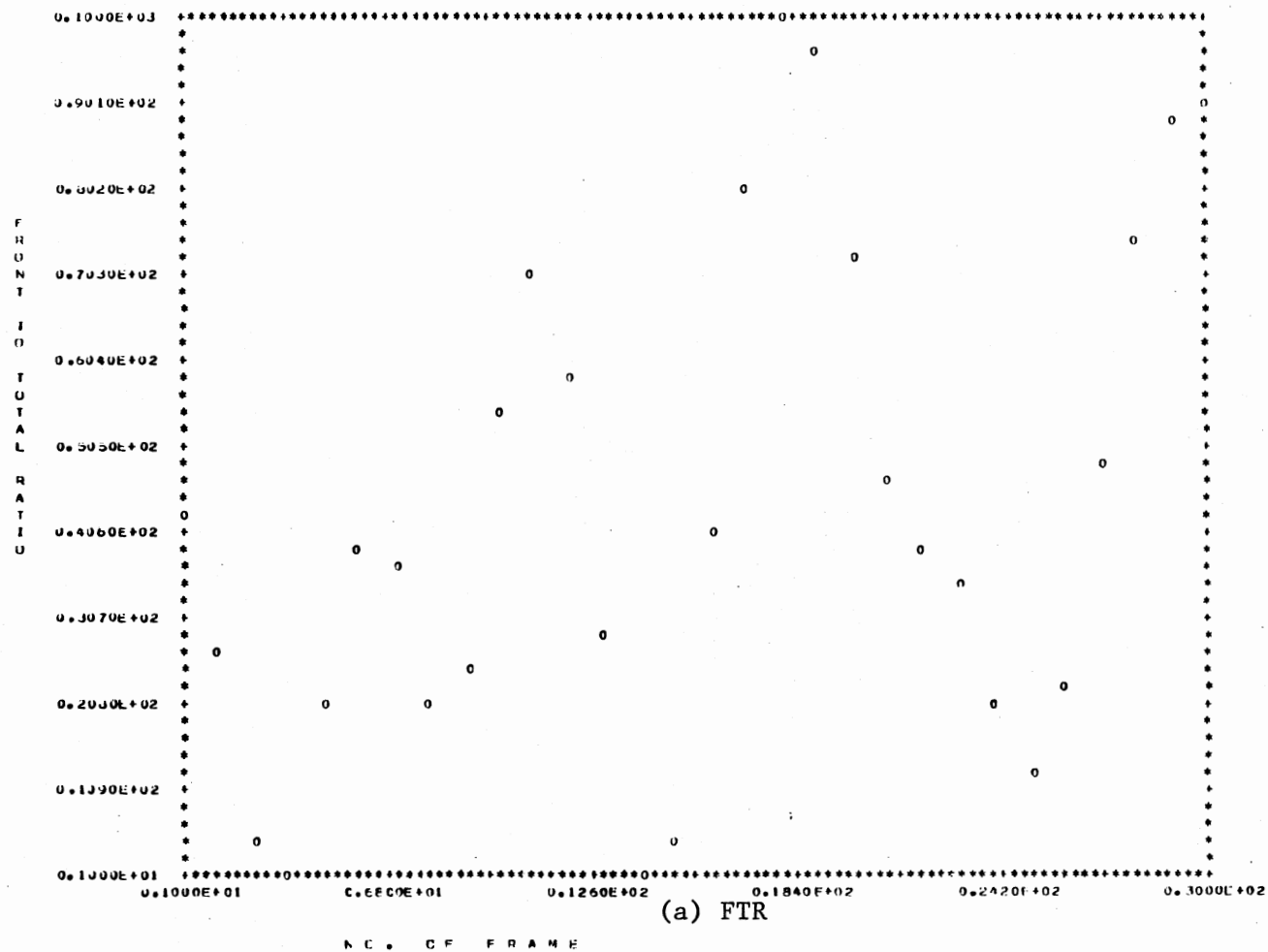


Figure 90. Smoothed and Quantized Feature Parameters for Digit Six, i.e. /sittəh/ Spoken in Arabic

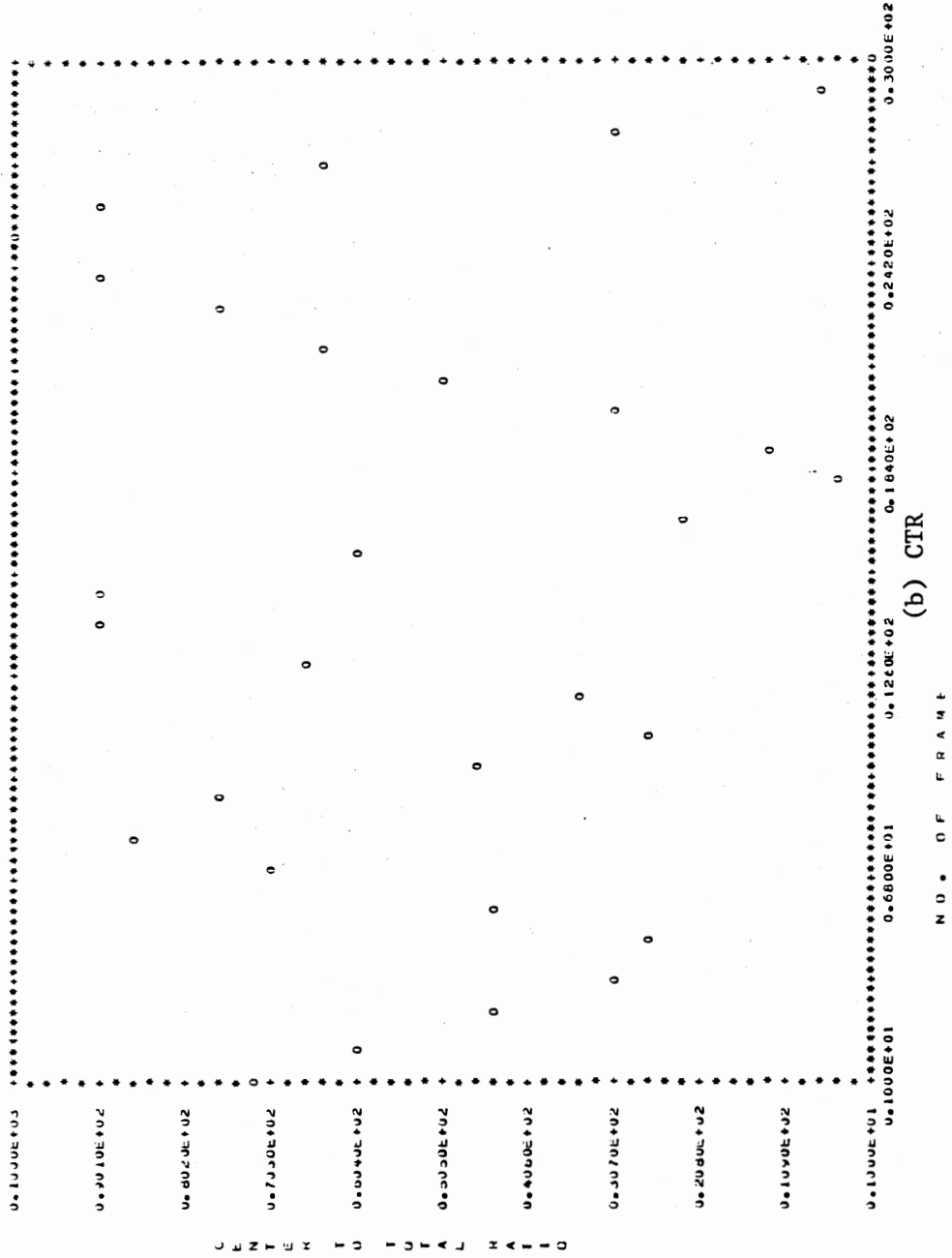
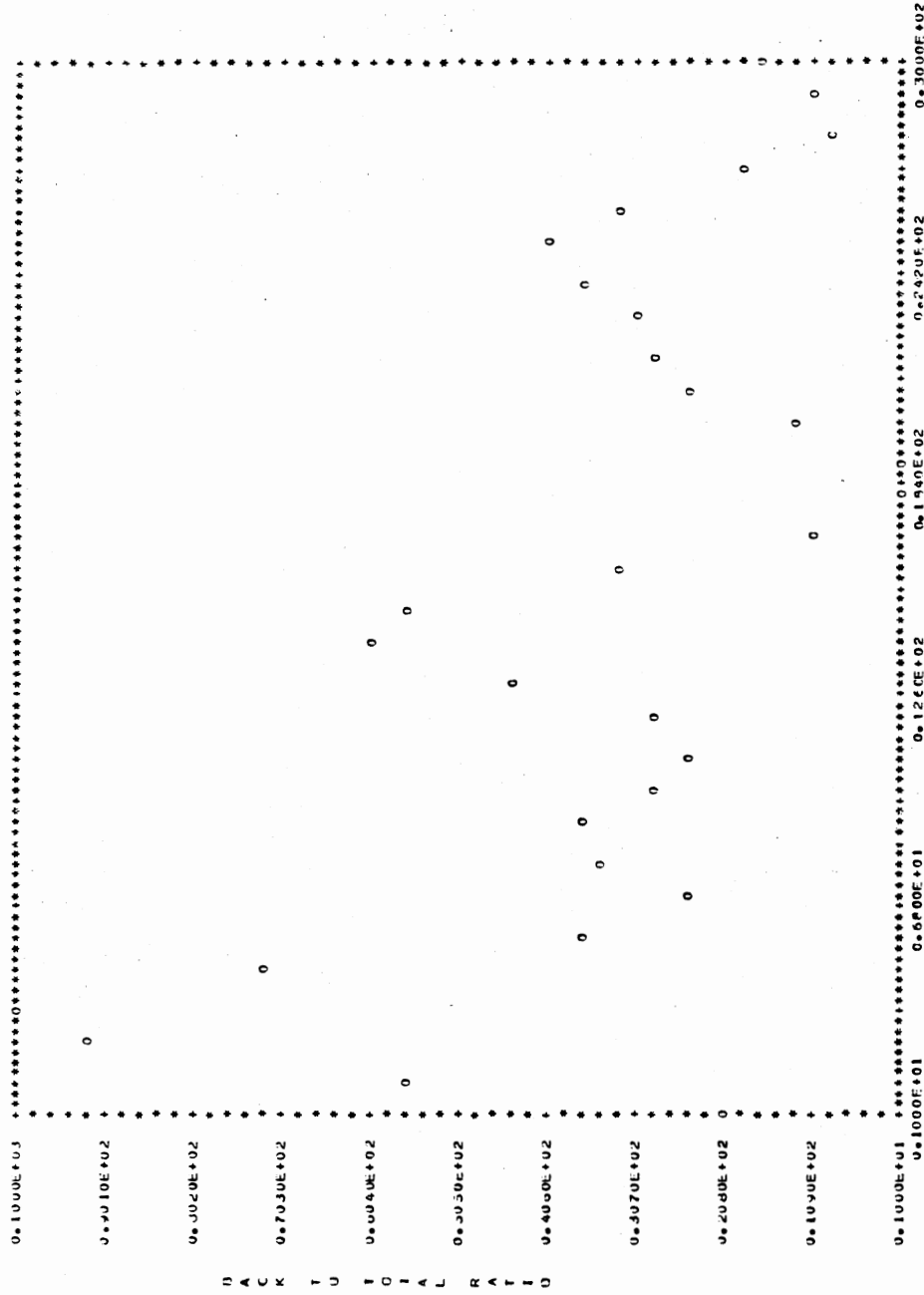


Figure 90. (Continued)



(c) BTR

Figure 90. (Continued)

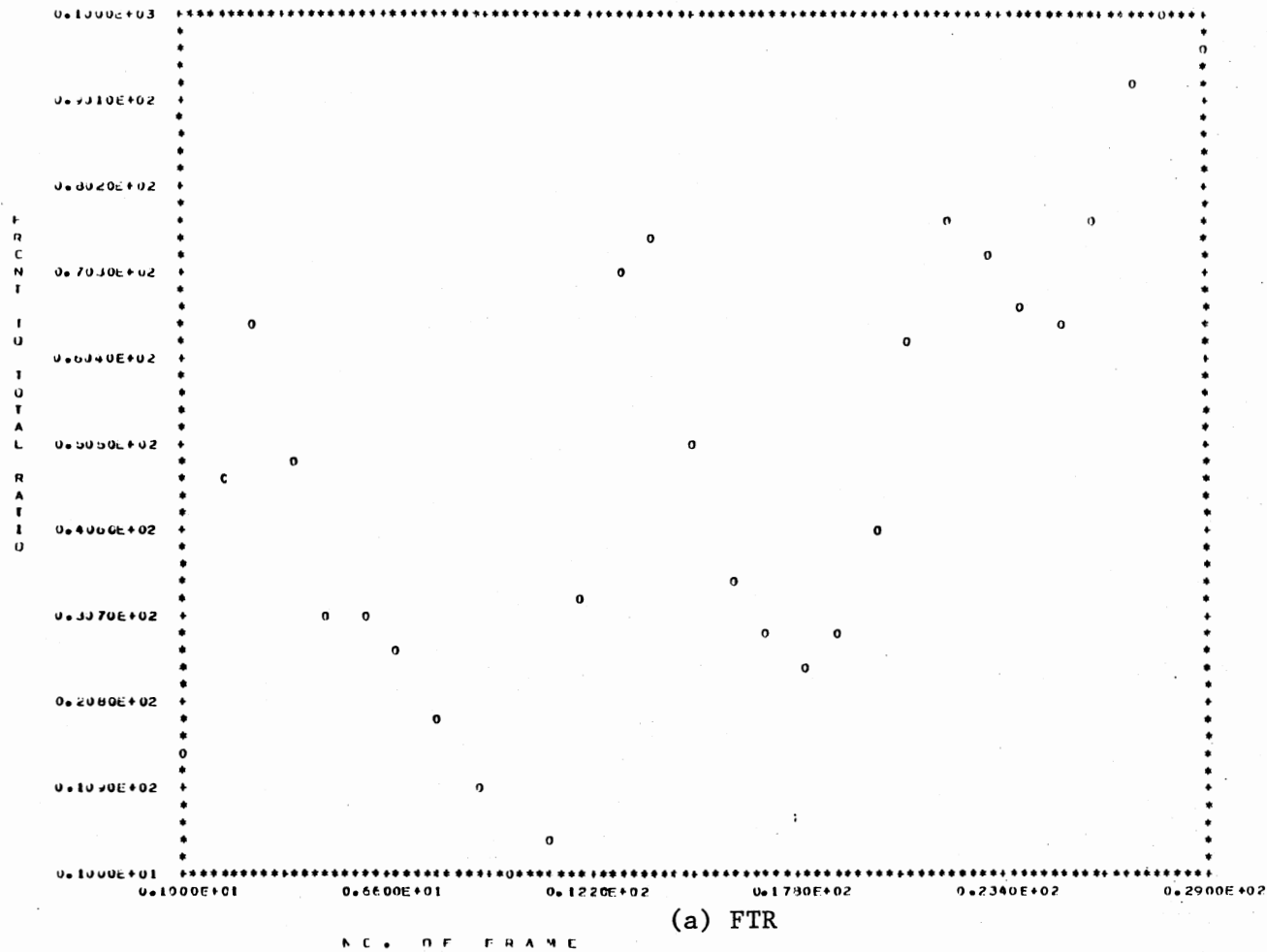
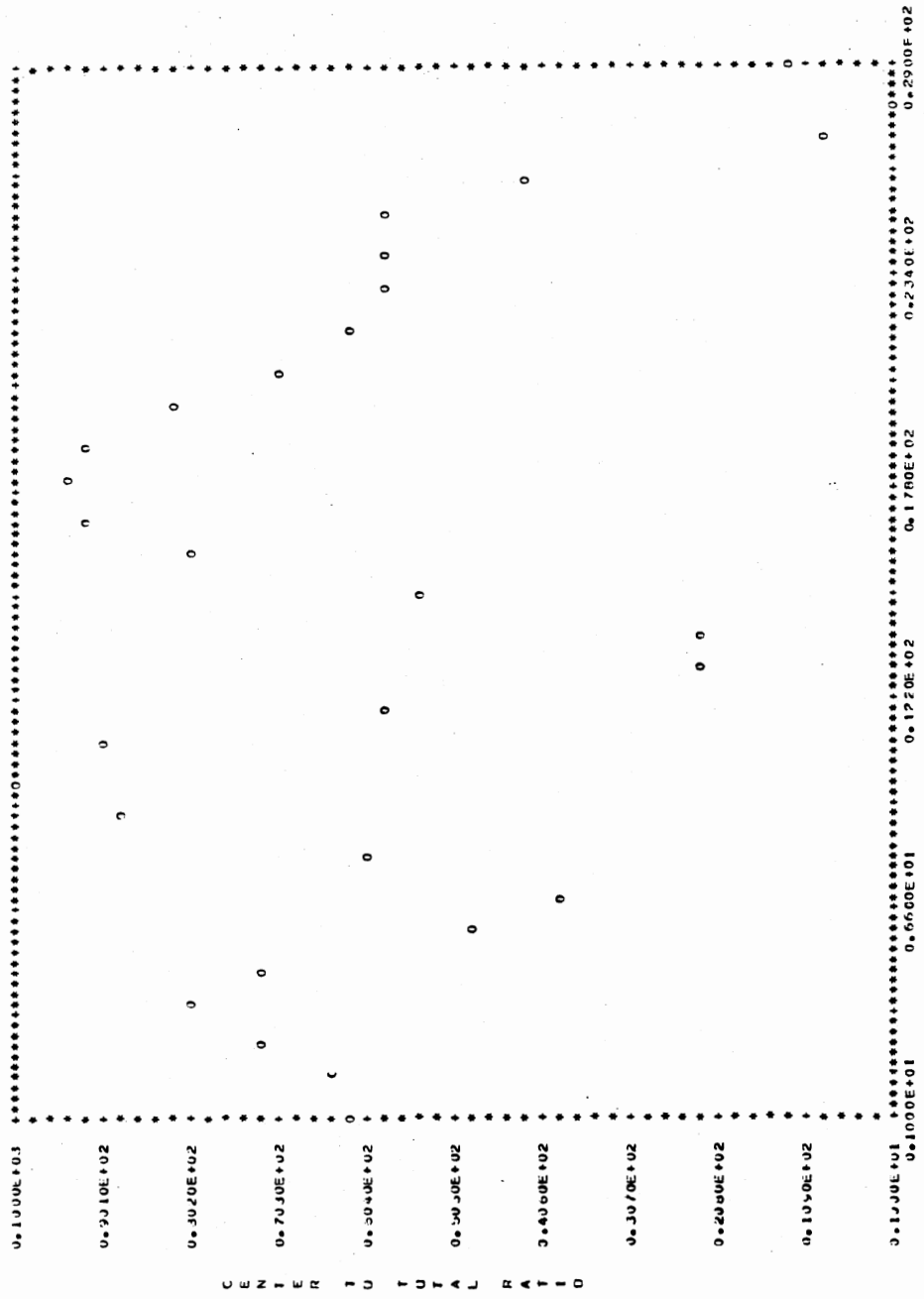


Figure 91. Smoothed and Quantized Feature Parameter for Digit Seven, /sʌbɔðh/ Spoken in Arabic



(b) CTR

Figure 91. (Continued)

NO. OF FRAME

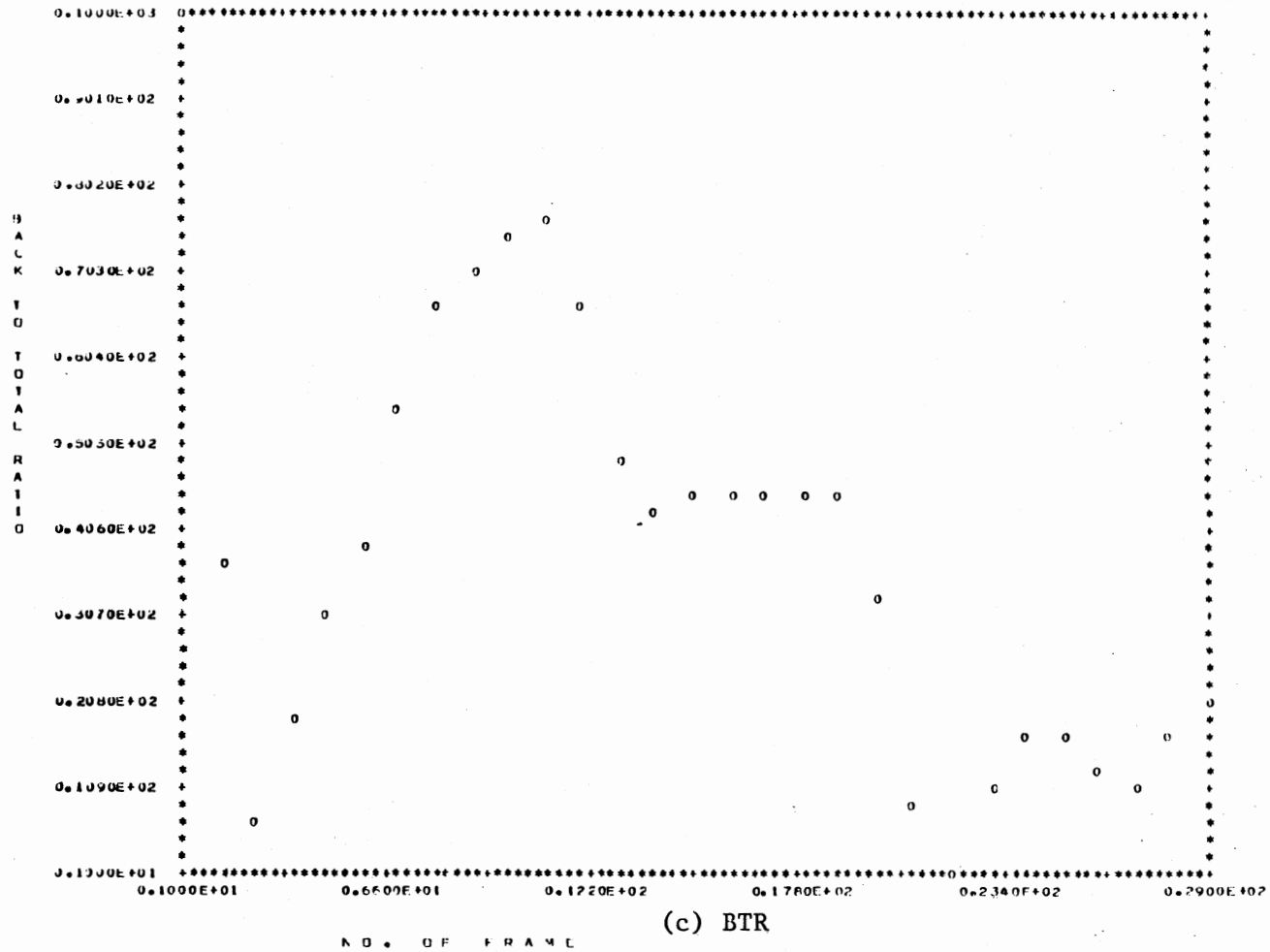


Figure 91. (Continued)

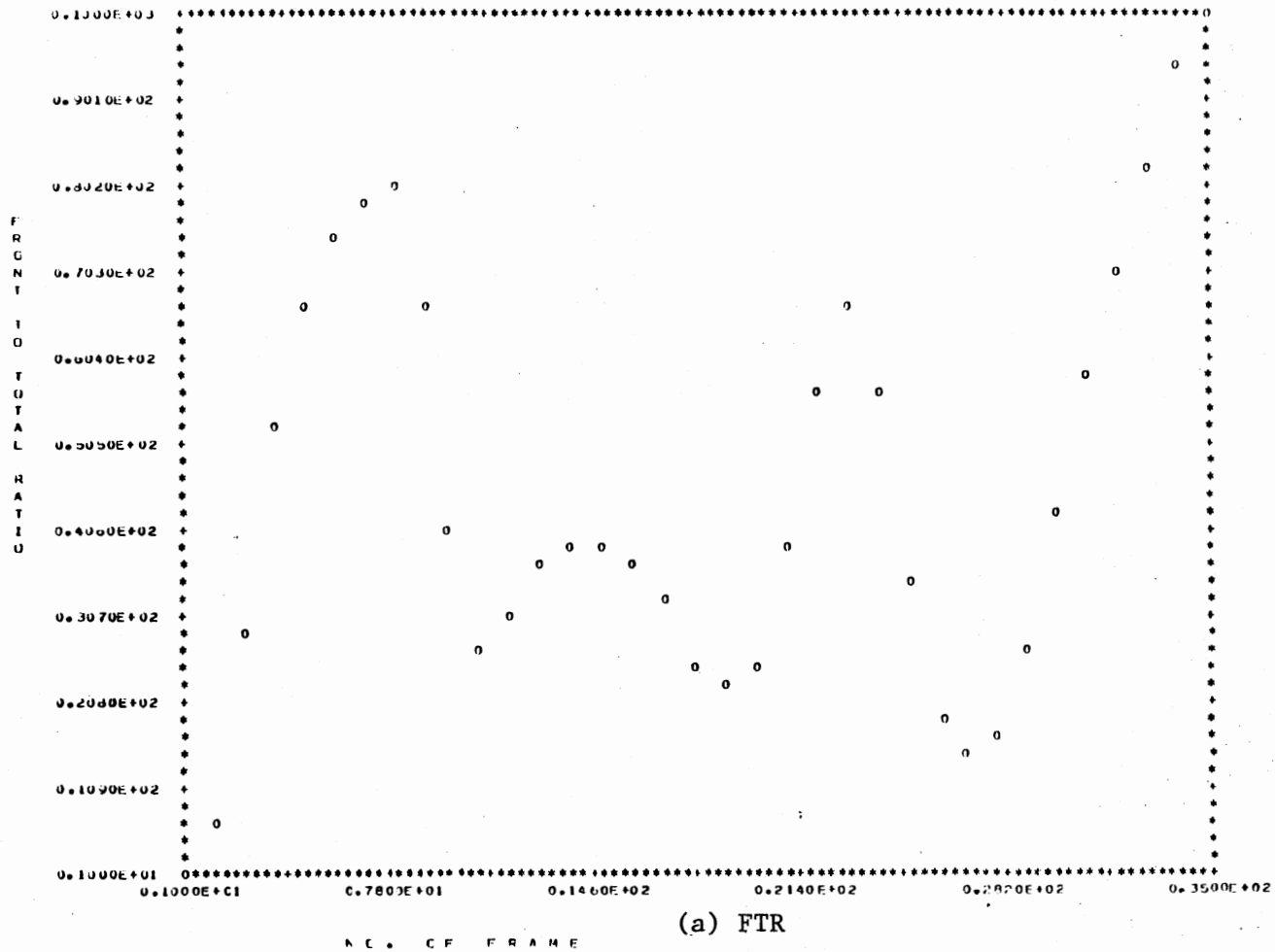


Figure 92. Smoothed and Quantized Feature Parameter for Digit Eight, i.e. /θamānyðh/ Spoken in Arabic

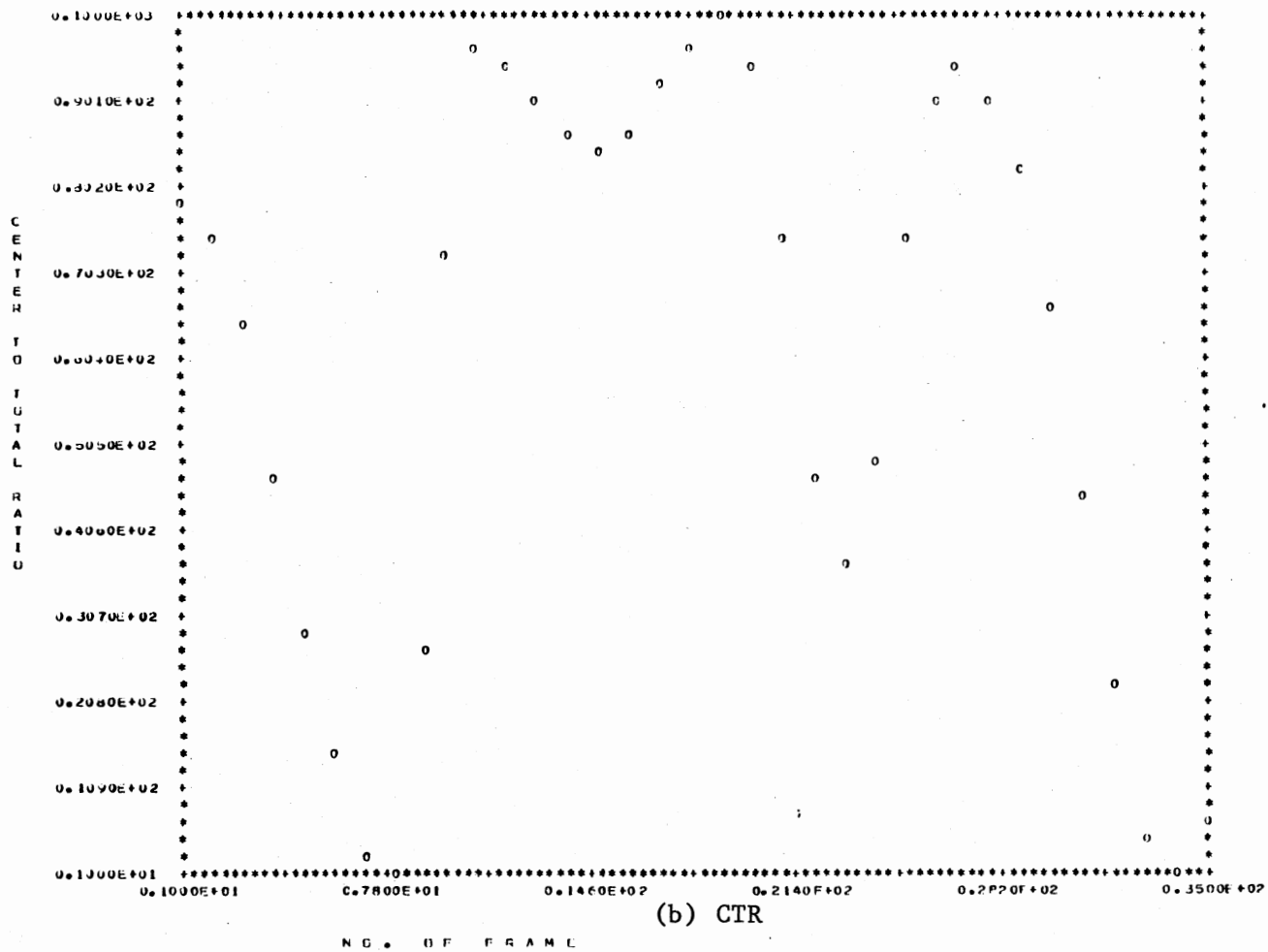


Figure 92. (Continued)

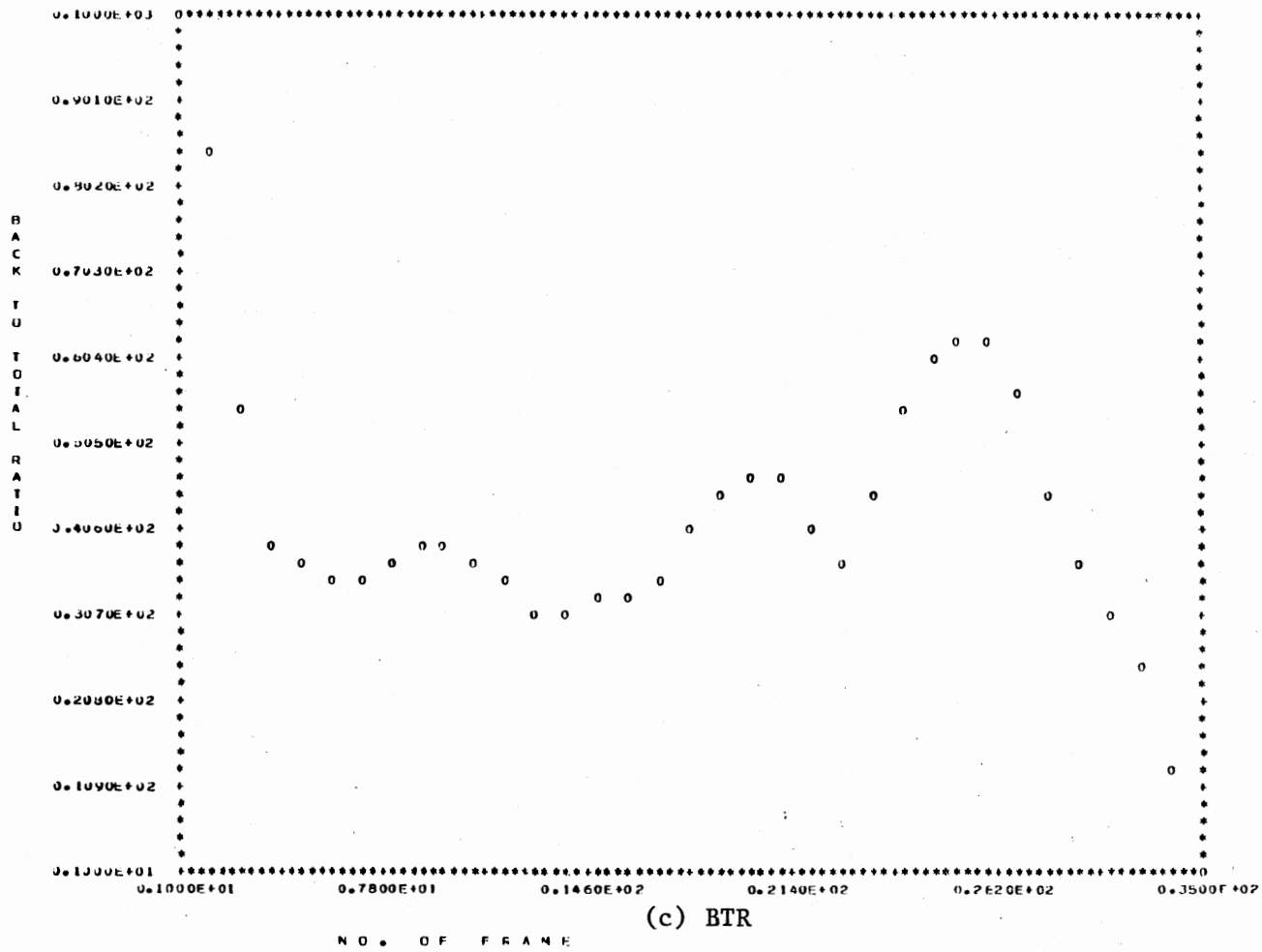


Figure 92. (Continued)

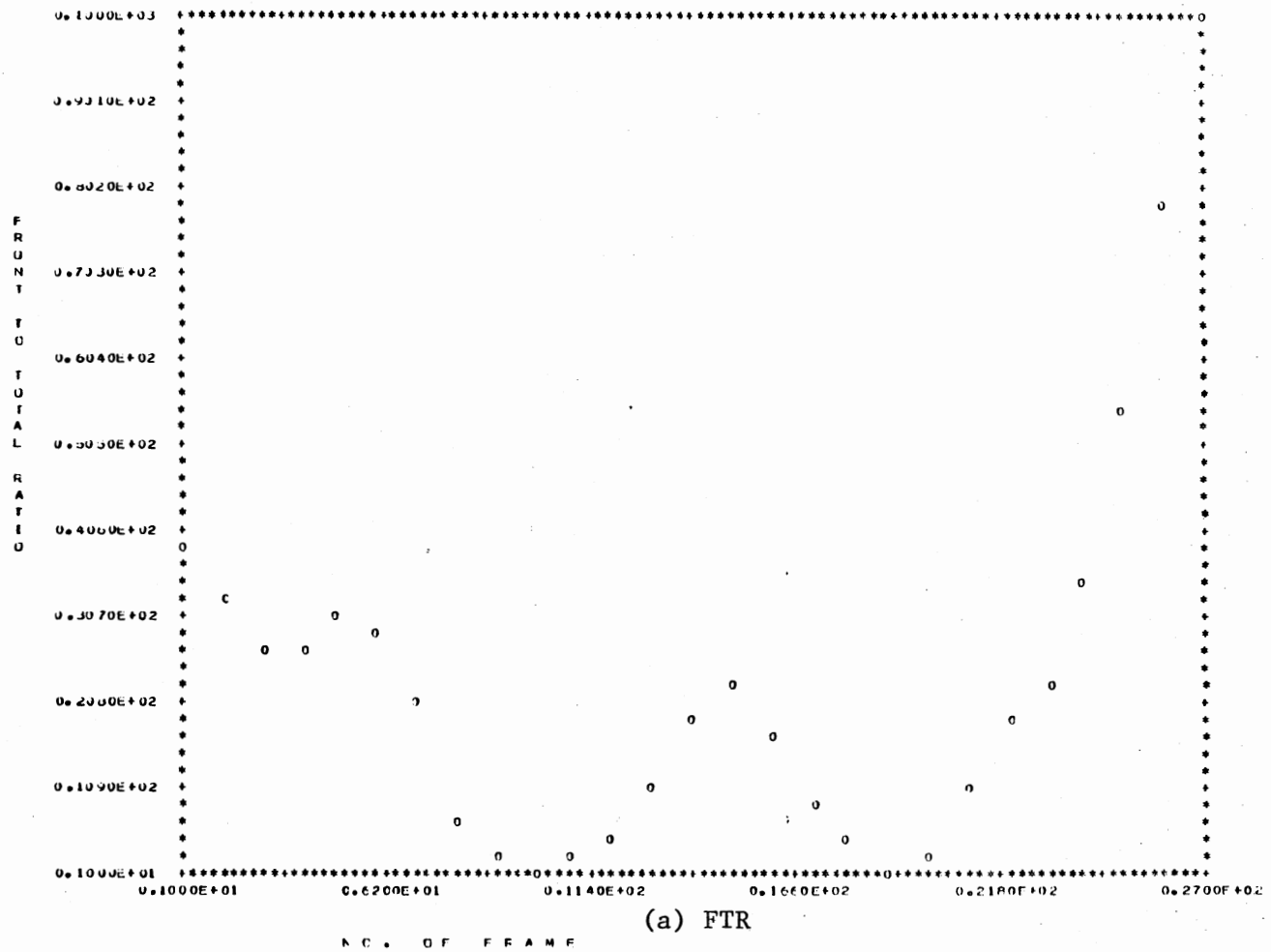


Figure 93. Smoothed and Quantized Feature Parameter for Digit Nine, i.e. /tIsɔðh/ Spoken in Arabic

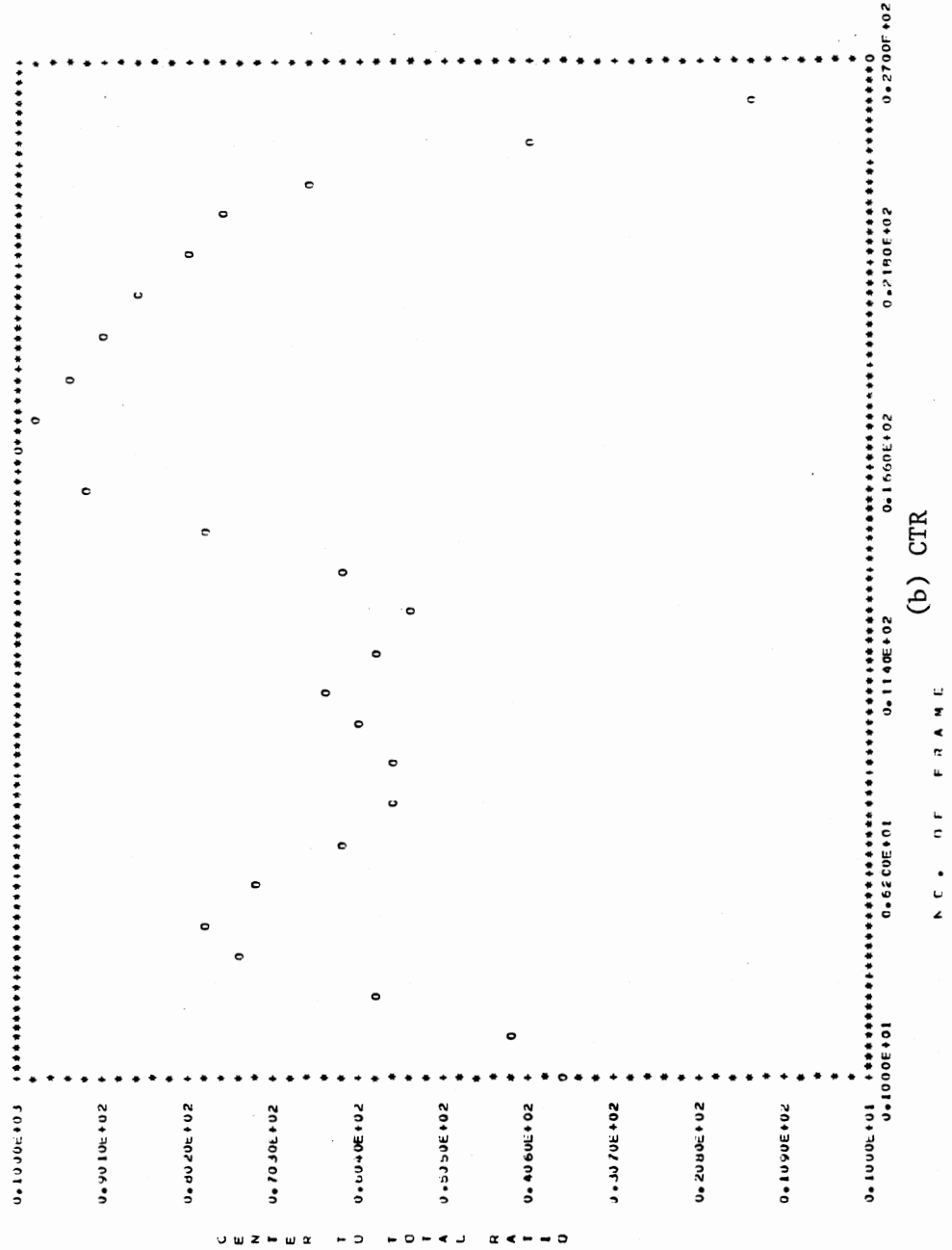


Figure 93. (Continued)

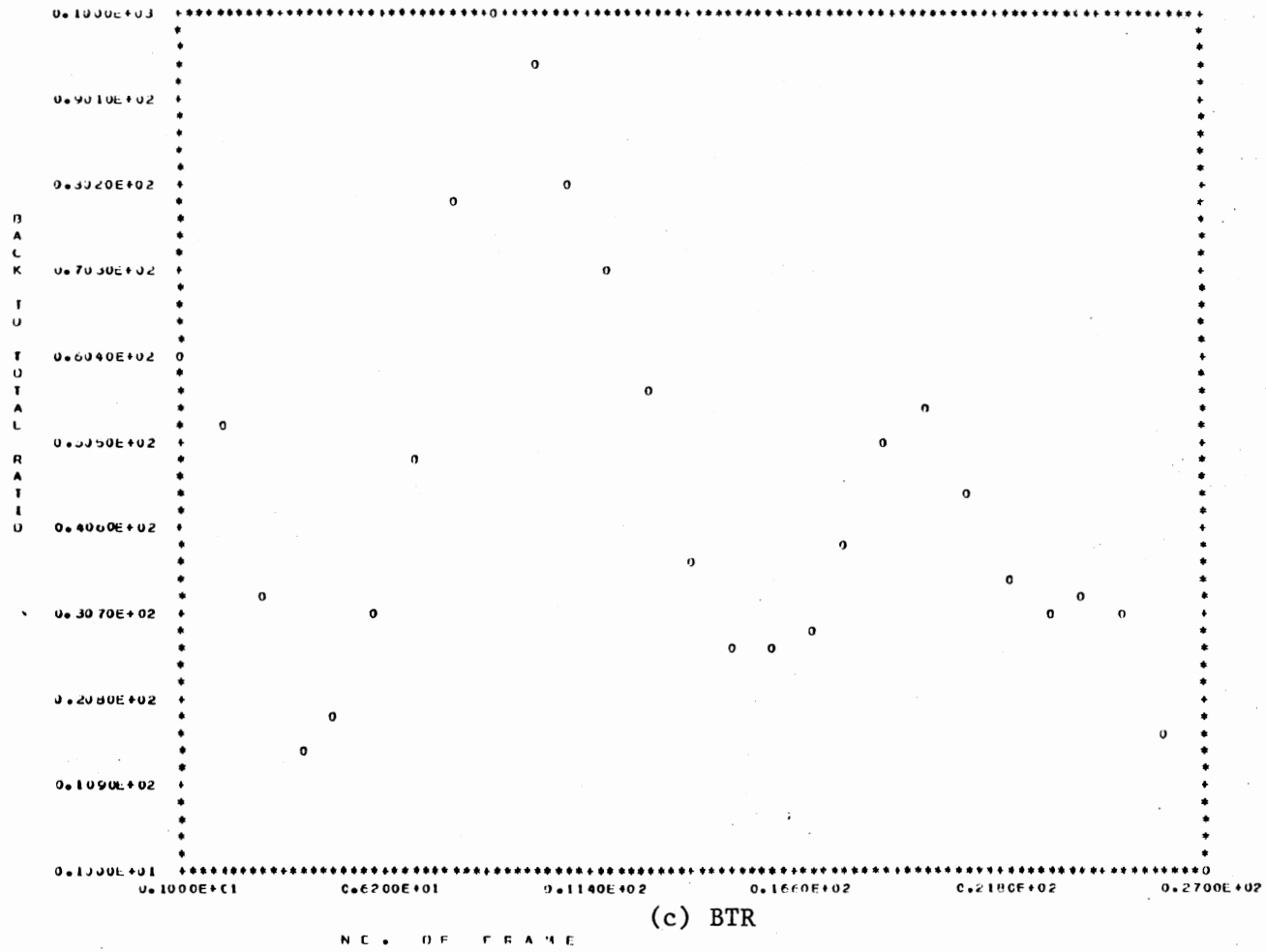


Figure 93. (Continued)

schemes. Furthermore, the SFBR parameter and the BTR parameter are used to separate turbulence noise segments from back vowel segments and nasals segments from vowel segments respectively. The classification of each segment is stored, so that it can be utilized in the third stage, for actual digit recognition.

The second path in Figure 28 is implemented as before. Table XVIII gives a summary of the beginning, middle, and ending sound classification of digits spoken in Arabic. This table is derived from Table IV. The decision algorithm is based upon the uniqueness aspects of the classifications of a particular digit. This is illustrated by Table XVIII. For example, the digits zero, one, two and four have unique classification for beginning and ending sounds. These four digits are therefore decided on the beginning and ending sounds. The digits seven and eight can be separated from the remaining digits using the non-vowel uniqueness aspects in the middle region. The remaining digits thru, five, six and seven have the same uniqueness aspects at front and end region, but the middle region is different. The flow chart in Figure 94 gives this procedure for recognizing the digits spoken in Arabic.

The above algorithm was tested for five male speakers whose native tongue is Arabic. The accuracy rate was about 70 percent. The speakers have different accents, which accounts for the inconsistency of the peak ratios among different speakers. The results were much better when the spoken digits were intact, and the accuracy ratio was about 95 percent. By this method, the digits six and nine have the major problem. The reason is due to the fact that different speakers tend to stress different phonemes in the spoken digit. This fact affects

TABLE XVIII
SEQUENCE OF SOUND CLASS REGIONS OF DIGITS SPOKEN IN ARABIC

Digit	RMS and BTR Decision		
	Beginning	Middle	End
/sefr/	NV	FV/NV	VL
/wâhid/	VL	BV/FV	NV
/iθnân/	FV	VL/BV	VL
/θalâθθh/	NV	VL/BV/NV	MV
/arbaρθh/	MV	VL/MV	MV
/khΛmsθh/	NV	MV/VL/NV	MV
/sIttθh/	NV	FV/NV	MV
/sΛbρθh/	NV	MV/VL	MV
/θamanyθh/	NV	MV /BV	MV
/tIspθh/	NV	VL/BV/VL	MV

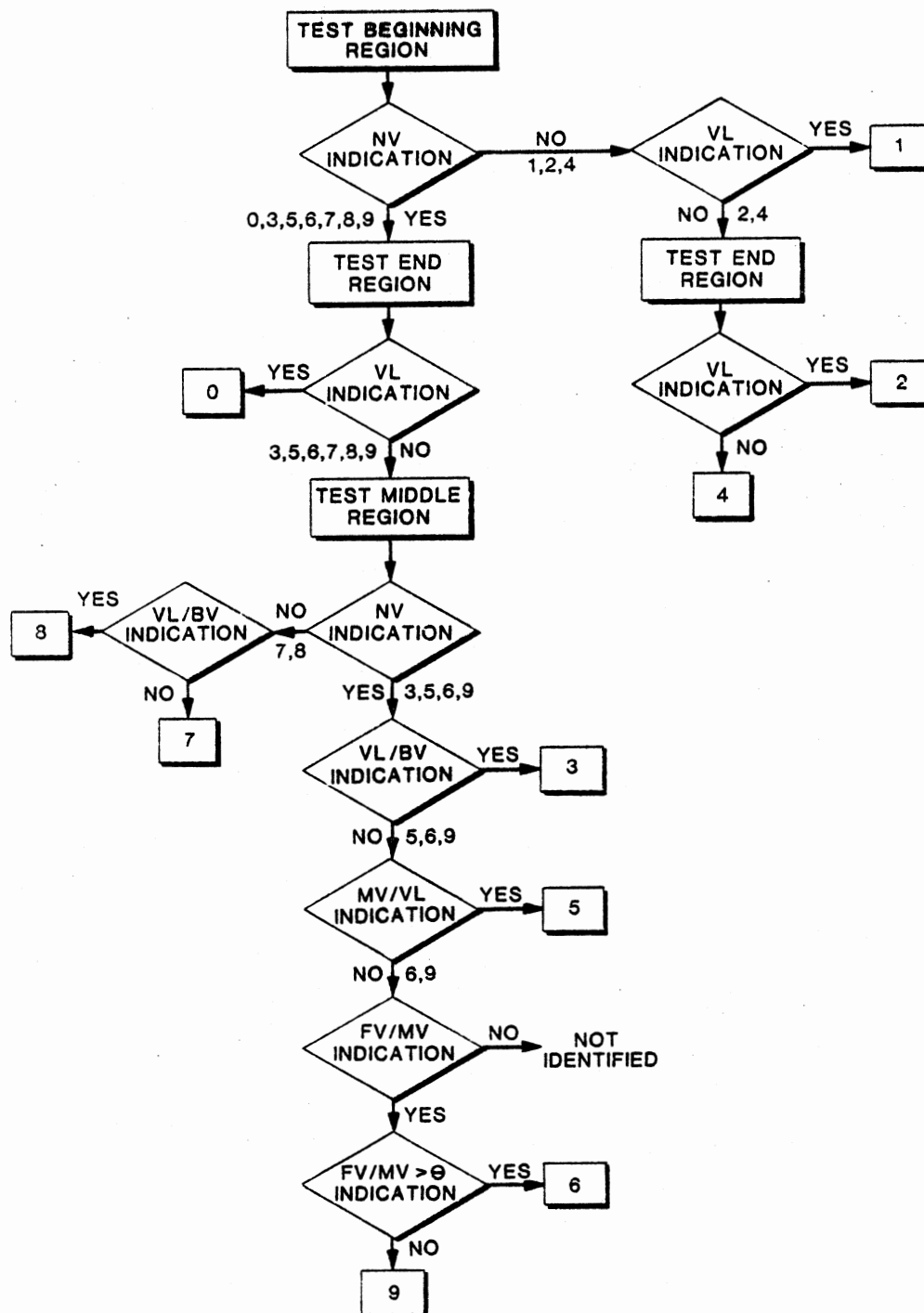


Figure 94. Decision Tree for Digit Identification for Arabic Language

the value range of FV-MV indication for different speakers. It is interesting to point out that if the spoken Arabic digit has a long vowel followed by a consonant (VVC) then that vowel is always stressed. It is pronounced louder than the other vowels; otherwise the first vowel of the digit is stressed. This can be seen from Table IV, which shows digit nine, /tIspəh/, have VCCV, where as digit six have either VccV or VCV, depending on how the /t/ is stressed. In addition, these two digits have almost similar characteristics in all the three regions. In fact the digit six, /sIttəh/, indicates that it has VCCV in the center because there is double /t/, i.e. the duration of the /t/ is doubled. Even though the second recognition idea is more reliable, the recognition result is based on the two paths for robustness. If the two paths give different results, or that no decision has been made, the pattern recognition procedure will be used. This is discussed next.

Pattern Recognition Scheme

In order to ascertain that the recognition scheme is robust, the cross-correlation coefficients are used in a pattern recognition scheme. The correlation coefficients is computed for the spoken digit using Equation (43). The flow chart of Figure 60 shows that a standard RMS and BTR pattern of the digits zero to nine are first selected using empirical rule, then stored. The RMS pattern correlation coefficients of the unknown digit is computed with maximum number of frames, then compared with the standard pattern. The correlation coefficient threshold is chosen empirically. If no recognition results, then the unknown digit is correlated with the

patterns of the remaining digit. The digit which has the highest correlation is selected. The same procedure is followed using the BTR pattern. The BTR pattern of the spoken digit is matched by the cross-correlation process. The RMS energy and BTR correlations for digits spoken in Arabic are shown in Tables XIX-XXIV. The cross-correlation process involves selecting the best match by locating the reference pattern which gives the greatest correlation coefficients with the pattern to be recognized. This system utilizes no linguistic information, but procedures of acoustic characteristics impeded in the RMS energy and the BTR parameter as discussed in the previous chapter.

The above scheme was applied for sample digits. The results, as shown in Tables XXV and XXVI, were good. However, before the recognition rate can be given, more digit patterns have to be listed. The computer programs are given in detailed form in Appendix B. This completes the discussion on the digit recognition scheme for digits spoken in Arabic.

TABLE XIX

RMS ENERGY CORRELATION TABLE FOR DIGITS SPOKEN IN ARABIC
(MAXIMUM NUMBER OF FRAMES USED)

SET NUMBER :		1								
RMS CORRELATION										
ZERO	ONE	TWO	THREE	FOUR	FIVE	SIX	SEVEN	EIGHT	NINE	
1.000	-0.304	-0.399	-0.385	-0.276	-0.213	-0.286	-0.015	-0.329	-0.291	
-0.304	1.000	-0.089	-0.258	0.421	-0.226	0.132	0.721	0.183	0.927	
-0.399	-0.089	1.000	0.305	-0.253	0.306	-0.231	-0.116	-0.305	-0.115	
-0.385	-0.258	0.305	1.000	0.251	0.425	-0.145	-0.394	0.194	-0.294	
-0.276	0.421	-0.253	0.251	1.000	0.153	0.109	0.267	0.329	0.399	
-0.213	-0.226	-0.306	0.425	0.153	1.000	0.588	-0.135	0.710	-0.284	
-0.286	0.132	-0.231	-0.145	0.109	0.588	1.000	0.297	0.795	0.043	
-0.015	0.721	-0.116	-0.394	0.267	-0.135	0.297	1.000	0.418	0.725	
-0.329	0.183	-0.305	0.194	0.329	0.710	0.795	0.418	1.000	0.121	
-0.291	0.927	-0.115	-0.294	0.399	-0.284	0.043	0.725	0.121	1.000	

TABLE XX

RMS ENERGY CORRELATION TABLE FOR DIGITS SPOKEN IN ARABIC
(MAXIMUM NUMBER OF FRAMES USED)

SET NUMBER :		2		RMS CORRELATION						
ZERO	ONE	TWO	THREE	FOUR	FIVE	SIX	SEVEN	EIGHT	NINE	
1.000	-0.137	-0.236	-0.223	0.355	0.007	0.244	0.008	0.110	0.615	
-0.137	1.000	0.845	-0.144	-0.339	-0.214	0.260	0.442	-0.050	-0.367	
-0.236	0.845	1.000	0.053	-0.255	-0.261	0.023	0.364	0.127	-0.498	
-0.223	-0.144	0.053	1.000	0.598	0.740	-0.291	0.610	0.245	-0.066	
0.355	-0.339	-0.255	0.598	1.000	0.855	-0.305	0.464	0.634	0.123	
0.007	-0.214	-0.261	0.740	0.855	1.000	-0.291	0.546	0.632	-0.017	
0.244	0.260	0.023	-0.291	-0.305	0.291	1.000	0.166	-0.229	0.503	
0.008	0.442	0.364	0.610	0.464	0.546	0.166	1.000	0.163	0.076	
0.110	-0.050	0.127	0.245	0.634	0.632	-0.229	0.163	1.000	0.037	
0.615	-0.367	-0.498	-0.066	0.123	-0.017	0.503	0.076	0.037	1.000	

TABLE XXI

RMS ENERGY CORRELATION TABLE FOR DIGITS SPOKEN IN ARABIC
(MAXIMUM NUMBER OF FRAMES USED)

SET NUMBER :		RMS CORRELATION							
ZERO	ONE	TWO	THREE	FOUR	FIVE	SIX	SEVEN	EIGHT	NINE
1.000	-0.253	-0.166	-0.234	0.137	-0.150	-0.132	-0.419	0.033	-0.440
-0.253	1.000	0.104	0.795	-0.187	0.843	-0.156	0.103	-0.068	0.501
-0.166	0.104	1.000	0.149	0.817	-0.018	-0.189	0.139	0.449	0.197
-0.234	0.795	0.149	1.000	-0.104	0.928	-0.290	-0.040	0.103	0.202
0.137	-0.187	0.817	-0.104	1.000	-0.253	-0.124	0.109	0.213	-0.009
-0.150	0.843	-0.018	0.928	-0.253	1.000	-0.124	-0.064	0.050	0.202
-0.132	-0.156	-0.189	-0.290	-0.124	-0.124	1.000	0.000	-0.038	0.182
-0.419	0.103	0.139	-0.040	0.109	-0.064	0.000	1.000	-0.367	0.773
0.033	-0.068	0.449	0.103	0.213	0.050	-0.038	-0.367	1.000	-0.348
-0.440	0.501	0.197	0.202	-0.009	0.202	0.182	0.773	-0.348	1.000

TABLE XXII

BTR CORRELATION TABLE FOR DIGITS SPOKEN IN ARABIC
(MAXIMUM NUMBER OF FRAMES USED)

SLT NUMBER :		BTR CORRELATION							
ZERO	ONE	TWO	THREE	FOUR	FIVE	SIX	SEVEN	EIGHT	NINE
1.000	-0.262	-0.123	-0.227	-0.294	-0.158	-0.235	-0.212	-0.138	0.084
-0.262	1.000	0.470	-0.409	0.317	-0.092	0.263	0.136	0.559	0.485
-0.123	0.470	1.000	-0.306	0.268	-0.213	0.287	0.155	0.130	0.291
-0.227	-0.409	-0.306	1.000	0.196	-0.185	-0.166	0.022	-0.009	-0.021
-0.294	0.317	0.268	0.196	1.000	0.582	0.567	0.472	0.475	0.163
-0.158	-0.092	-0.213	-0.185	0.582	1.000	0.429	0.324	0.085	-0.354
-0.235	0.263	0.287	-0.166	0.567	0.429	1.000	-0.050	0.617	-0.205
-0.212	0.136	0.155	0.022	0.472	0.324	-0.050	1.000	-0.002	0.335
-0.138	0.559	0.130	-0.009	0.475	0.085	0.617	-0.002	1.000	0.172
0.084	0.485	0.291	-0.021	0.163	-0.354	-0.205	0.335	0.172	1.000

TABLE XXIII

BTR CORRELATION TABLE FOR DIGITS SPOKEN IN ARABIC
(MAXIMUM NUMBER OF FRAMES USED)

SET NUMBER :		2		BTR CORRELATION						
ZERO	ONE	TWO	THREE	FOUR	FIVE	SIX	SEVEN	EIGHT	NINE	
1.000	0.172	0.653	0.064	0.564	0.025	0.149	0.055	-0.106	0.148	
0.172	1.000	0.190	0.073	0.323	0.109	0.526	0.347	0.237	0.294	
0.653	0.190	1.000	0.034	0.827	-0.185	0.092	0.503	-0.201	0.407	
0.064	0.073	0.034	1.000	0.171	0.530	0.044	0.549	0.621	0.352	
0.564	0.323	0.827	0.171	1.000	0.133	0.223	0.529	-0.120	0.276	
0.025	0.109	-0.185	0.530	0.133	1.000	-0.105	0.570	0.512	0.289	
0.149	0.526	0.092	0.044	0.223	-0.105	1.000	-0.147	0.152	-0.052	
0.055	0.347	0.503	0.549	0.528	0.570	-0.147	1.000	0.234	0.741	
-0.106	0.237	-0.201	0.621	-0.120	0.512	0.152	0.234	1.000	0.143	
0.148	0.294	0.407	0.352	0.276	0.289	-0.052	0.741	0.143	1.000	

TABLE XXIV

BTR CORRELATION TABLE FOR DIGITS SPOKEN IN ARABIC
(MAXIMUM NUMBER OF FRAMES USED)

SET NUMBER :		BTR CORRELATION							
ZERO	ONE	TWO	THREE	FOUR	FIVE	SIX	SEVEN	EIGHT	NINE
1.000	0.527	0.205	0.041	0.127	0.117	0.027	0.292	-0.128	0.283
0.527	1.000	-0.380	0.480	0.105	0.413	0.269	0.566	0.459	0.371
0.205	-0.380	1.000	-0.340	0.023	-0.132	-0.072	0.086	-0.520	0.057
0.041	0.480	-0.340	1.000	0.694	-0.104	0.097	-0.397	0.551	0.206
0.127	0.105	0.023	0.694	1.000	-0.002	0.512	0.128	0.256	0.227
0.117	0.413	-0.132	-0.104	-0.002	1.000	0.520	0.551	0.756	0.585
0.027	0.269	-0.072	0.097	0.512	0.520	1.000	0.750	0.597	0.157
0.292	0.566	0.086	-0.397	0.128	0.551	0.750	1.000	0.372	0.206
-0.128	0.459	-0.520	0.551	0.256	0.756	0.597	0.372	1.000	0.702
0.283	0.371	0.057	0.206	0.227	0.585	0.157	0.206	0.702	1.000

TABLE XXV

RMS ENERGY CROSS-CORRELATION TABLE FOR TWO SETS OF DIGITS
SPOKEN IN ARABIC BY TWO DIFFERENT SPEAKERS

ZERO	ONE	TWO	RMS CORRELATION						
			THREE	FOUR	FIVE	SIX	SEVEN	EIGHT	NINE
0.929	0.394	-0.039	-0.360	0.243	-0.065	0.140	0.681	0.328	-0.444
0.394	0.965	0.720	-0.131	-0.107	0.340	0.005	0.399	0.470	-0.297
-0.039	0.720	0.861	0.007	-0.164	0.502	-0.026	0.196	0.250	-0.316
-0.360	-0.131	0.007	0.916	0.358	-0.144	-0.205	-0.012	0.626	0.810
0.243	-0.107	-0.164	0.358	0.886	-0.020	0.143	0.533	0.225	0.369
-0.065	0.340	0.502	-0.144	-0.020	0.820	-0.094	0.159	-0.119	-0.130
0.140	0.005	-0.026	-0.205	0.143	-0.094	0.967	0.413	-0.525	-0.385
0.681	0.399	0.196	-0.012	0.533	0.159	0.413	0.575	0.580	-0.108
0.328	0.470	0.250	0.626	0.225	-0.119	-0.525	0.580	0.795	0.735
-0.444	-0.297	-0.316	0.810	0.369	-0.130	-0.385	-0.108	0.735	0.911

TABLE XXVI

BTR CROSS-CORRELATION TABLE FOR THE SAME DIGITS
USED IN TABLE XXV

ZERO	ONE	TWO	BTR CORRELATION						
			THREE	FOUR	FIVE	SIX	SEVEN	EIGHT	NINE
0.567	0.009	0.708	0.141	0.694	-0.163	0.359	0.112	-0.350	0.297
0.009	0.721	0.389	-0.344	0.114	-0.143	0.269	0.508	0.080	0.751
0.708	0.389	0.798	0.026	0.340	-0.150	0.097	0.264	-0.165	0.292
0.141	-0.344	0.026	0.686	-0.134	-0.077	-0.235	-0.399	0.330	0.001
0.694	0.114	0.340	-0.134	0.821	0.104	-0.051	-0.038	-0.060	0.212
-0.163	-0.143	-0.150	-0.077	0.104	0.556	0.105	0.141	-0.193	-0.109
0.359	0.269	0.097	-0.235	-0.051	0.105	0.454	-0.132	-0.079	0.031
0.112	0.508	0.264	-0.399	-0.038	0.141	-0.132	0.493	0.064	0.407
-0.350	0.080	-0.165	0.330	-0.060	-0.193	-0.079	0.064	0.443	-0.047
0.297	0.751	0.292	0.001	0.212	-0.109	0.031	0.407	-0.047	0.554

CHAPTER VI

SUMMARY AND SUGGESTION FOR FURTHER STUDY

Summary

In this thesis a Robust Phonetic Digit Recognition System for digits spoken in American English and Arabic has been developed. A three-stage recognition scheme has been developed based on the RMS energy, the equivalent area functions and pattern analysis. The first stage uses rectangular window for segmenting the spoken digit into 128 data points per frame. The RMS energy is computed per frame, followed by a smoothing routine, to remove any undesirable ripples without losing any useful features. The smoothed RMS energy is then quantized to a maximum level of 100. This form of normalization made the RMS peaks to be more consistent with different speakers uttering the same digit than using a fixed mean and varying amplitudes. The later type of data scaling is avoided, because the amplitude value depends considerably on the number of data points which depends on the end point detection algorithm. The largest two peaks are computed from the smoothed RMS energy and their ratio is computed, which is assumed to be less than unity. An empirical threshold value is established from seven different speakers for American English digits and five different speakers for Arabic digits.

In the second stage, the data is first Hamming windowed, then a first order low-frequency pre-emphasis is used with a gain of (-1.0). A 14th order linear prediction analysis model is applied to the windowed signal. Useful measures are extracted from the smoothed and quantized RMS energy functions for effectively distinguishing vowel intervals from non-vowel intervals. The vowel, vowel-like and non-vowel intervals detection is based on the RMS dip-classification. In addition, the parameters FTR, CTR, and BTR are computed from the equivalent vocal-tract area via LPA. These parameters have some useful measures, acoustical and phonemic features. However, the BTR is not only found to be reliable, but sufficient for distinguishing vowel, and non-vowel intervals than the CTR. An 'OR' decision control is used based on the RMS dip-classification and BTR parameters for correctly identifying vowel, vowel-like, and non-vowel segments. Two different digit recognition trees are then used, based on vowel, vowel-like and non-vowel decision scheme for each language.

In the first stage, the digit recognition rate is about 60 percent (55 percent) for digits spoken in American English (Arabic), whereas in the second stage the digit recognition rate is about 75 percent (70 percent) for digits spoken in American English (Arabic). However, a digit recognition rate of 95 percent is obtained with American English, speaking in a sound proof chamber. The low recognition accuracy for digits spoken in both languages is due to that all the speakers have different accents. In addition most international speakers tend to stress a vowel or a consonant with a spoken digit. Another important reason is that some of the spoken digits are not intact. A lower recognition rate is obtained for digits spoken in

Arabic, because some Arabic digits have silence interval in the middle region of spoken digits. Wrong boundary detection resulted whenever the digit is not intact.

In the final stage, the digit recognition algorithm, based on the correlation coefficients of the RMS and BTR parameters, is used if no digit is identified by both procedures. A 97 percent recognition rate is obtained for digits spoken in both languages. The number of speakers used in this research is rather limited, indicating that the recognition scheme cannot be completely classified as speaker independent. The cross-correlation scheme improved the recognition rate accuracy considerably. In addition, it furnished a better understanding of the feature uniqueness among the digits. For example, it has been found that the digit seven is not heavily correlated with the same digit in another set. This is because international speakers tend to either emphasize the semi-vowel /v/ or mis-pronounce this vowel.

Similar recognition rate is achieved for digits spoken in Arabic. The RMS and BTR patterns are used so that, the spoken digit will always be identified, either by the BTR or by both parameters.

Suggestions for Further Research

Some extensions to the present research are proposed below.

The parameters derived from the equivalent area functions discussed in Chapter III need to be further studied. That is parameter normalization, where no prior knowledge about the input speakers is required, needs to be further investigated. Especially the vocal-tract

length and area functions may be first estimated from the acoustic speech waveform, and then the area function is normalized to an acoustic tube of the same shape having a certain reference length.

The two dimensional plots of BTR versus CTR may give a unique pattern for a given spoken digit after parameter normalization. Furthermore, three dimensional plots of BTR, CTR and FTR may give additional insight into the phonemic features of spoken words. Definite pattern uniqueness for the digits zero to nine may result after parameter normalization and boundary limit modification. An overlap limit for the front, central and back cavity may modify these parameters to be very effective for discriminating between front, middle, and back parameters.

In this study, spectrum analysis is avoided, in order to have an efficient recognition system. However, studying consonant features using selective Linear Prediction Analysis may improve detection of nasals, voiced fricatives and plosives. But normalization of the frequency scale is of major importance. This may lead to speaker verification system. Additional research in terms of the radius r_0 at the constriction and its distance from the glottis may allow new results in word recognition.

This research has opened up a new area in speech and speaker recognition in Arabic. Obviously, this area is wide open.

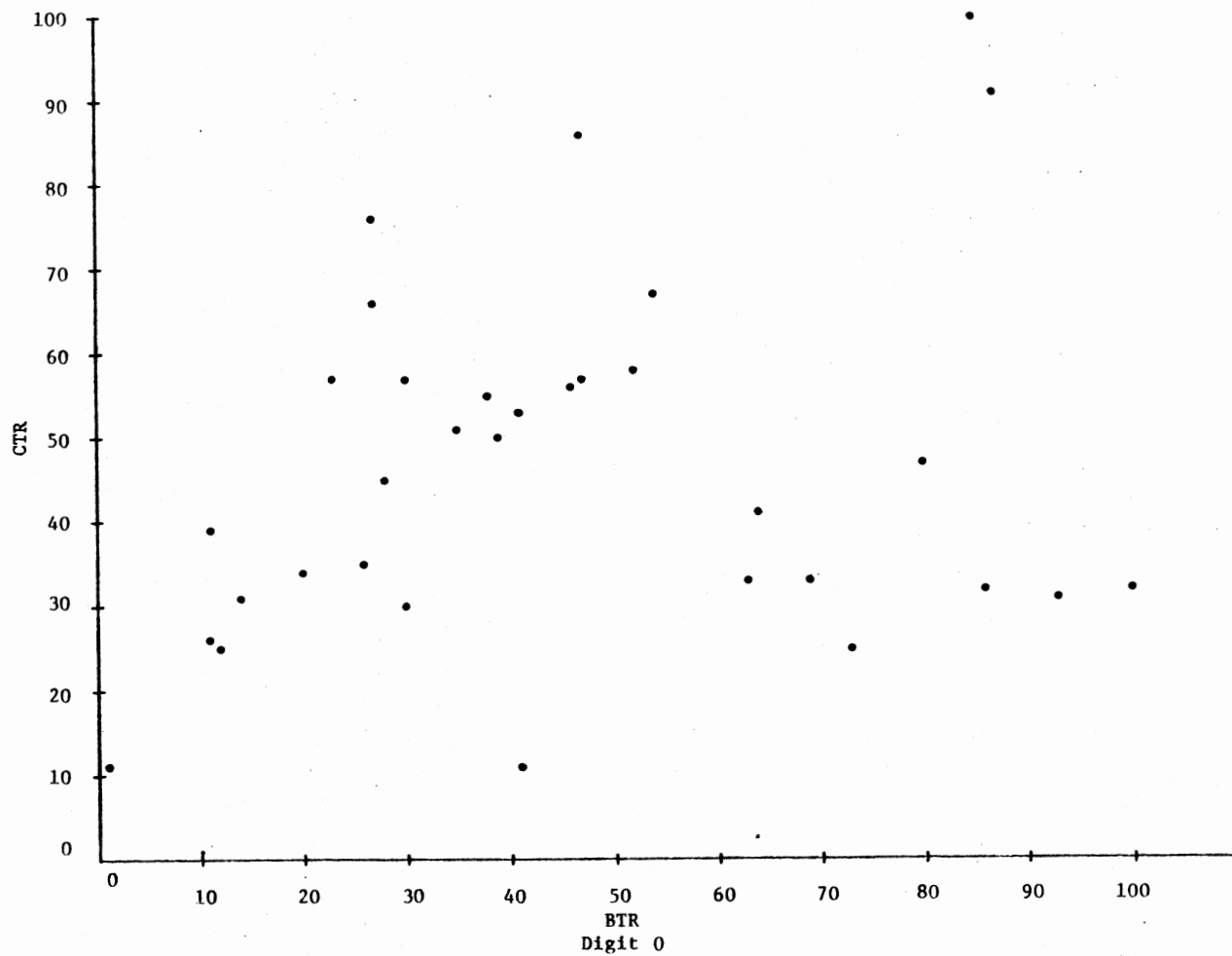


Figure 95. Plots of CTR vs BTR for Digit Zero

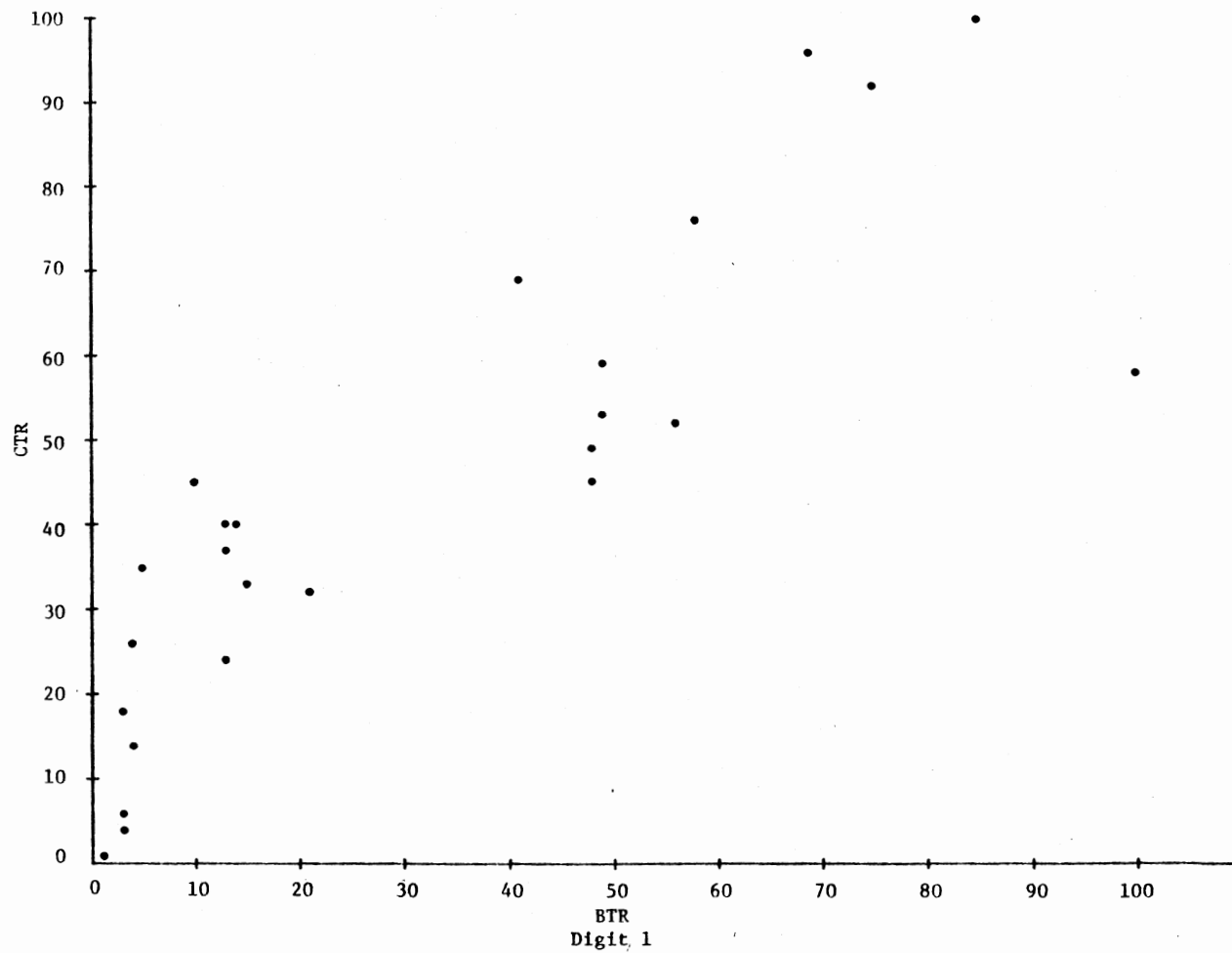


Figure 96. Plots of CTR vs BTR for Digit One

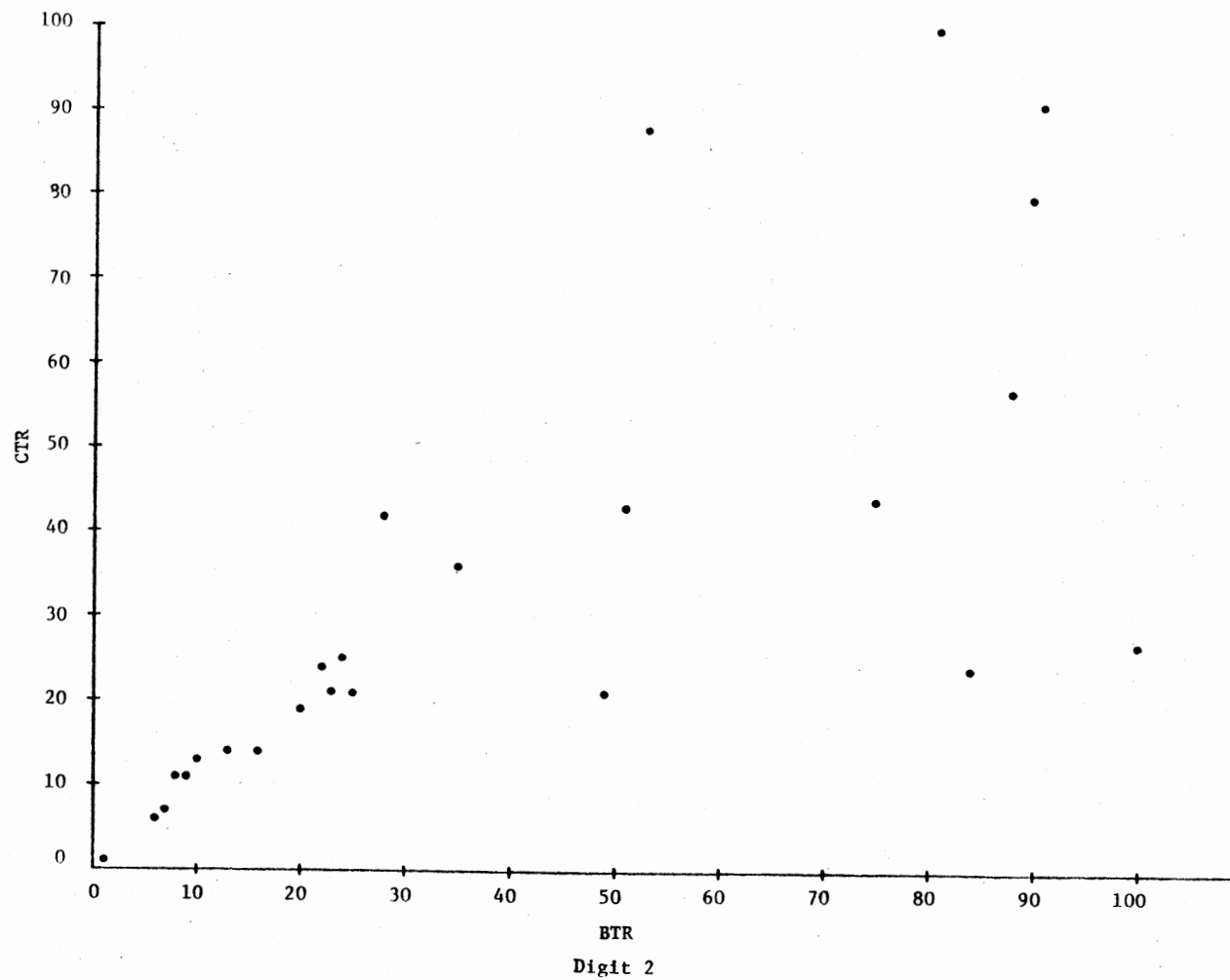


Figure 97. Plots of CTR vs BTR for Digit Two

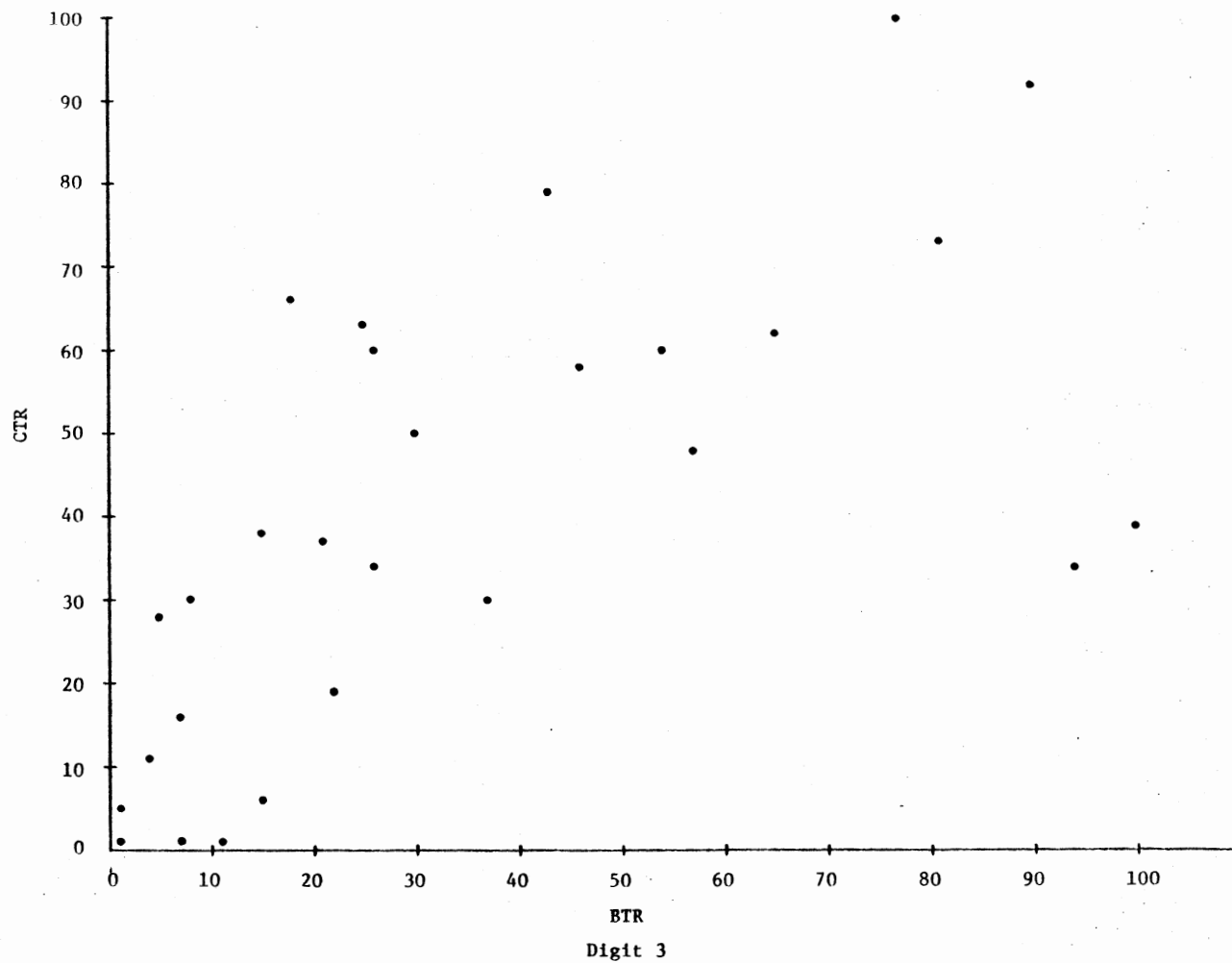


Figure 98. Plots of CTR vs BTR for Digit Three

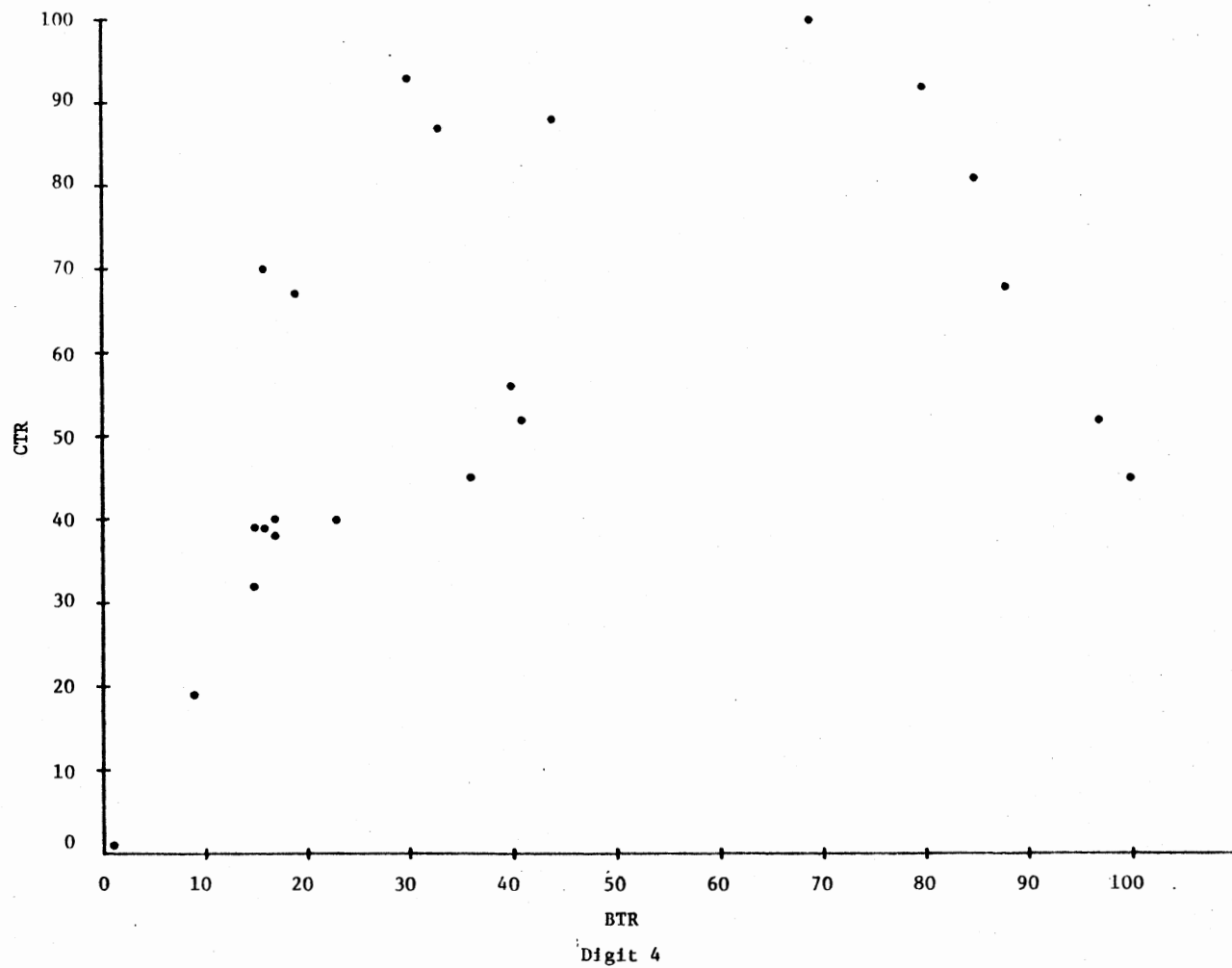


Figure 99. Plots of CTR vs BTR for Digit Four

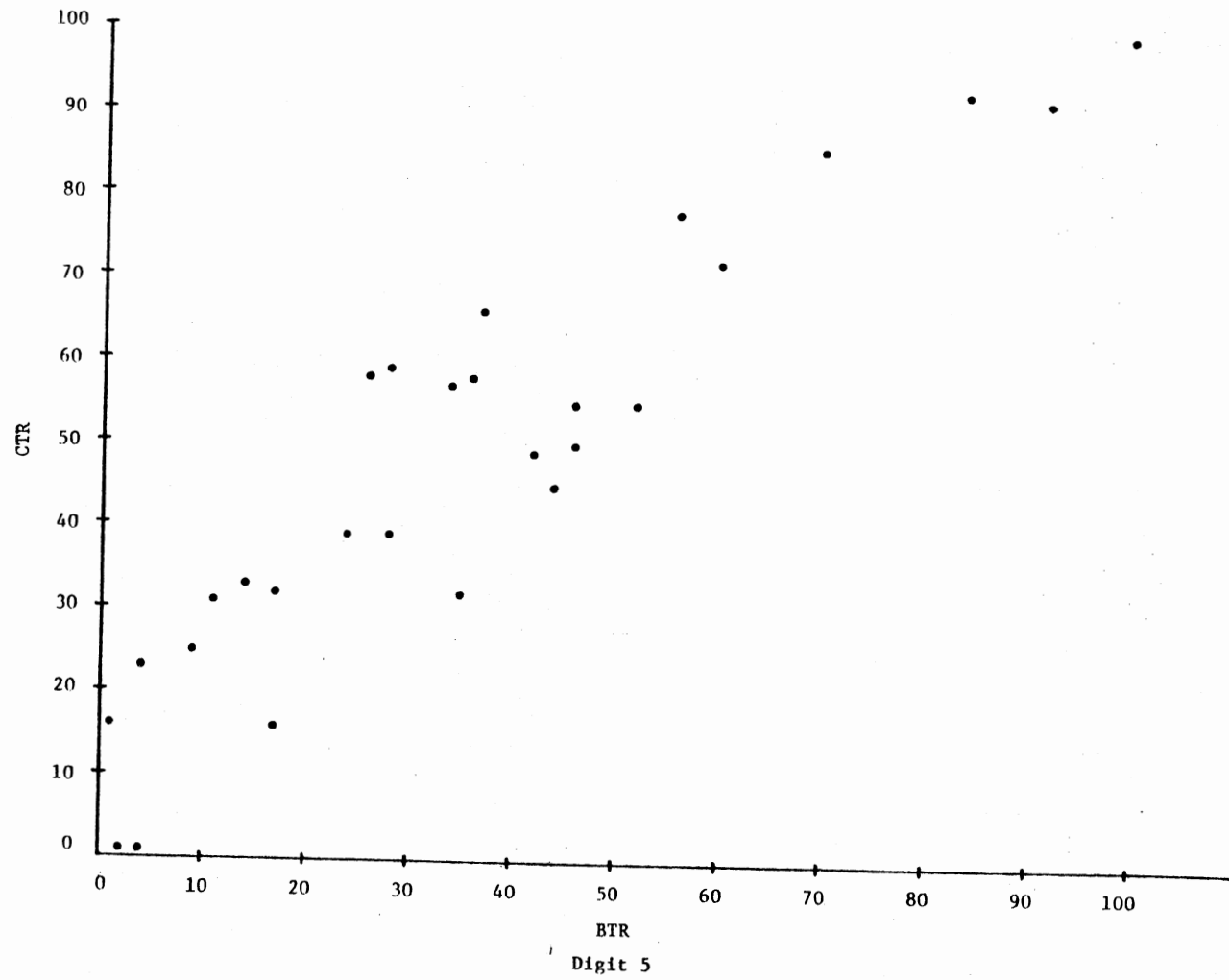


Figure 100. Plots of CTR vs BTR for Digit Five

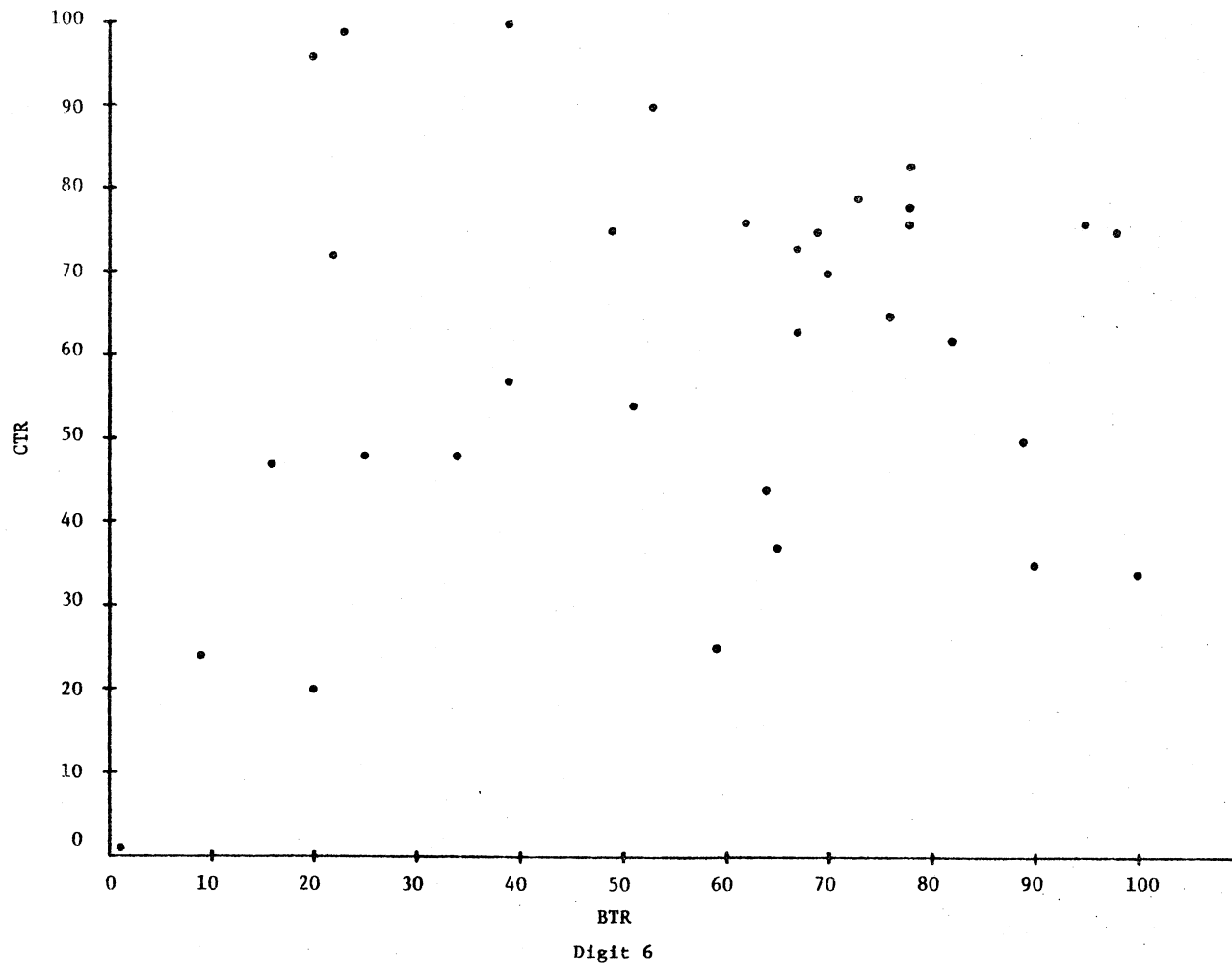


Figure 101. Plots of CTR vs BTR for Digit Six

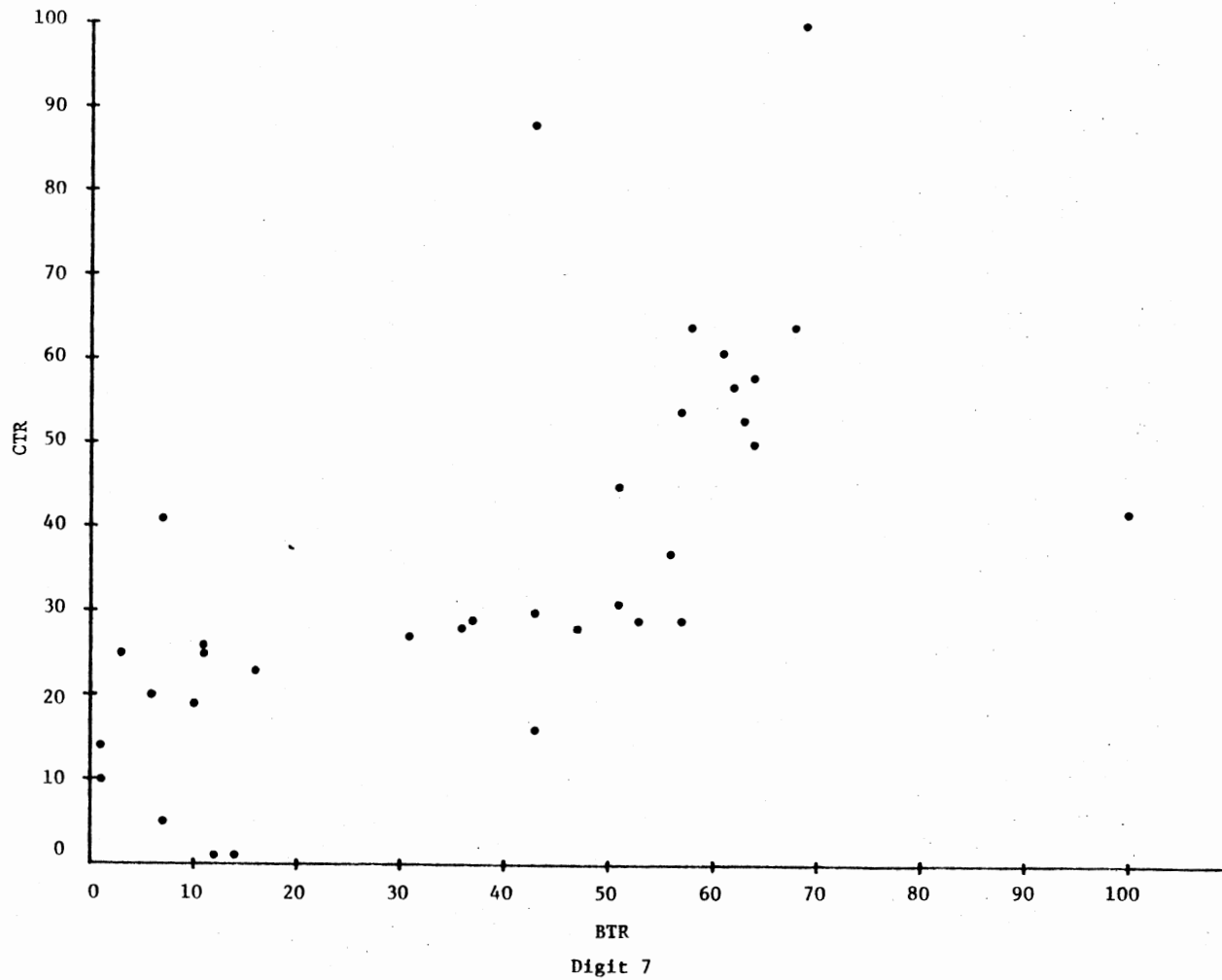


Figure 102. Plots of CTR vs BTR for Digit Seven

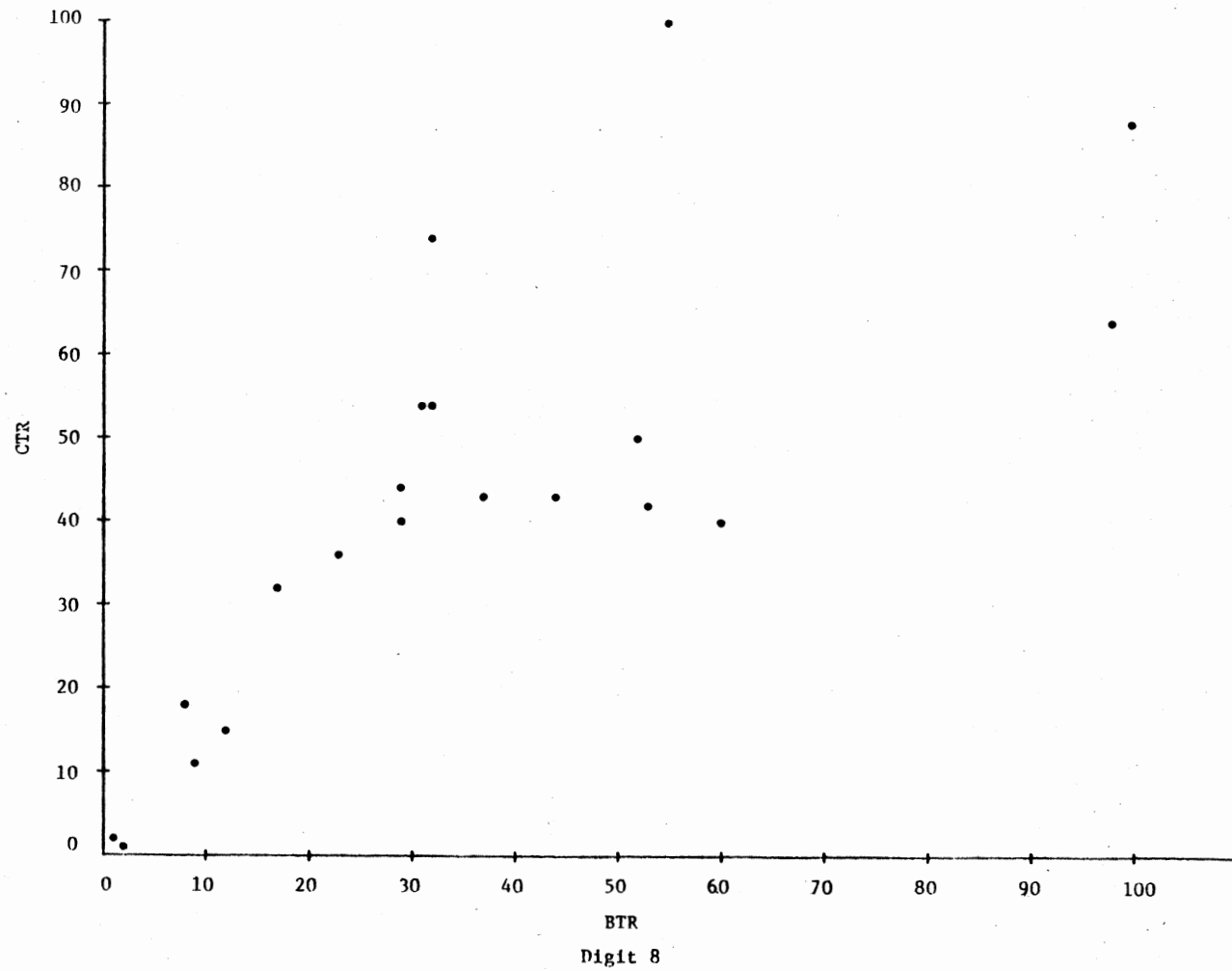


Figure 103. Plots of CTR vs BTR for Digit Eight

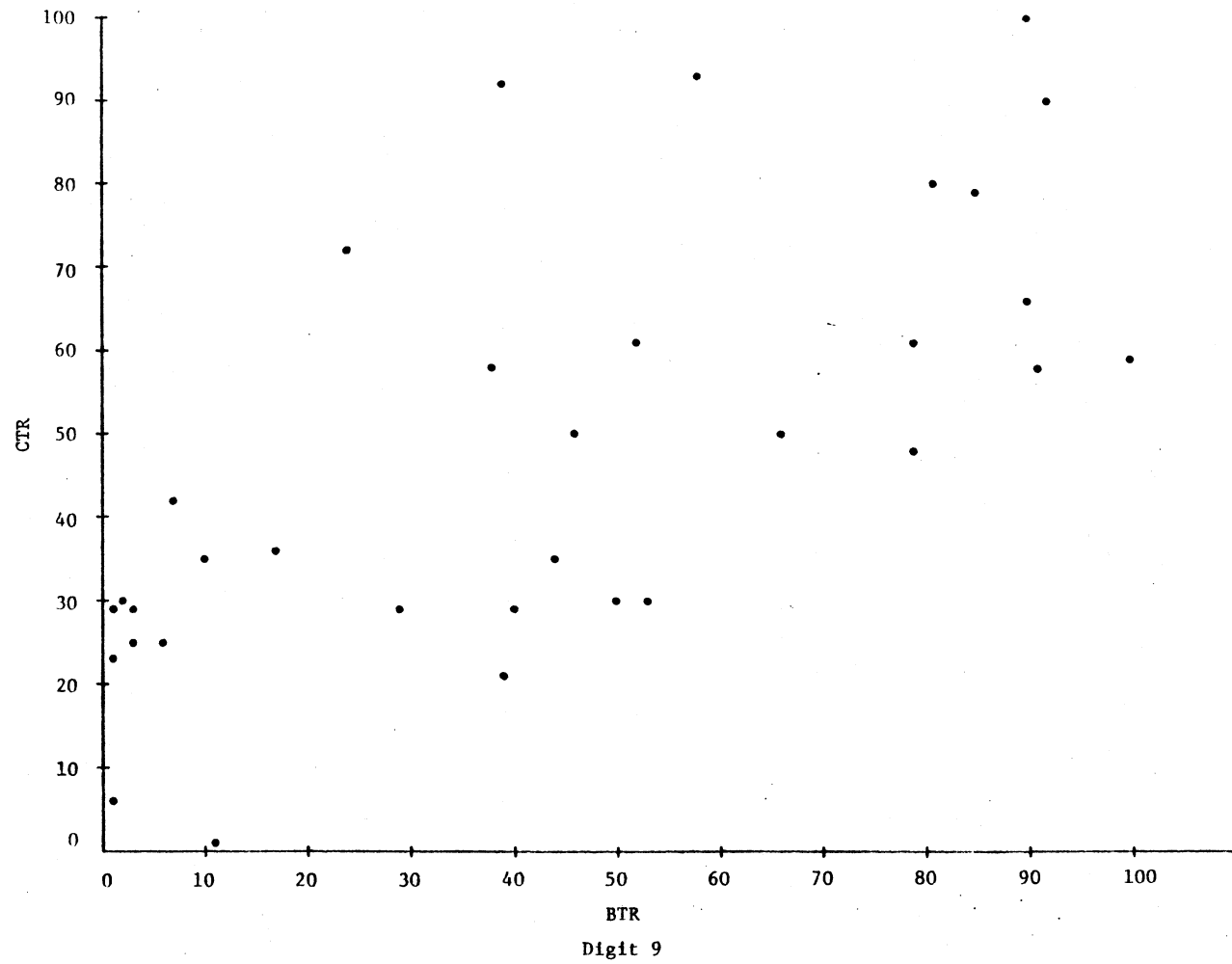


Figure 104. Plots of CTR vs BTR for Digit Nine

REFERENCES

- [1] L. R. Rabiner and R. W. Schafer, "Digital Processing of Speech Signals," Signal Processing Series, R. V. Oppenheim (Series Ed.). Englewood Cliffs: Prentice Hall, Inc. 1978.
- [2] D. R. Reddy, "Speech Recognition by Machines," Proceedings of the IEEE, Vol. 64, No. 4, 501-531, Apr. 1976.
- [3] P. Denes and M. W. Mathews, "Spoken Digit Recognition Using Time-Frequency Pattern Matching," J. Acoust. Soc. Am., Vol. 32, 1450-1455, 1960.
- [4] M. R. Sambur and L. R. Rabiner, "A Speaker-Independent Digit-Recognition System," Bell System Technical Journal, Vol. 54, No. 1, 81-102, Jan. 1975.
- [5] P. W. Ross, "A Limited-Vocabulary Adaptive Speech-Recognition System," Journal of the Audio Engineering Society, Vol. 15, No. 4, 414-418, Oct. 1967.
- [6] H. R. Dixon and H. F. Silverman, "A General Language-Operated Decision Implementation System (GLODIS), Its Application to Continuous Speech Segmentation," IEEE Trans. on Acous. Sp. Sig. Proc., Vol. ASSP-24, No. 2, 137-162, Apr. 1976.
- [7] D. R. Reddy, "Segmentation of Speech Sounds," J. Acoust. Soc. Am., Vol. 40, No. 2, 307-312, March 1966.
- [8] H. Kasuya and H. Wakita, "An Approach to Segmenting Speech into Vowel- and Nonvowel-like Intervals," IEEE Trans. on Acous. Sp. and Sig. Proc., Vol. ASSP 27, No. 4, 319-327, Aug. 1979.
- [9] L. R. Rabiner and M. R. Sambur, "An Algorithm for Determining the Endpoints of Isolated Utterances," Bell System Technical Journal, Vol. 54, No. 2, 297-315, Feb. 1975.
- [10] W. A. Lea, Trends in Speech Recognition. Englewood Cliffs, New Jersey: Prentice Hall, Inc., 1980.
- [11] J. L. Flanagan, Speech Analysis, Synthesis and Perception, Vol. 3. New York: Springer-Verlag, 1972.

- [12] L. R. Rabiner and M. R. Sambur, "Some Preliminary Experiments in the Recognition of Connected Digits," IEEE Trans. on Acous. Sp. and Sig. Proc., Vol. ASSP 24, No. 2, 244-256, Apr. 1976.
- [13] K. H. Davis, R. Biddulph and S. Balshek, "Automatic Recognition of Spoken Digits," J. Acoust. Soc. Am., Vol. 24, No. 6, 637-642, Nov. 1952.
- [14] T. G. von Keller, "An On-Line Recognition System for Spoken Digits," J. Acoust. Soc. Am., Vol. 49, No. 4 (Part 2), 1288-2196, 1971.
- [15] P. N. Sholtz and R. Bakis, "Spoken Digit Recognition Using Vowel-Consonants Segmentation," J. Acoust. Soc. Am., Vol. 34, No. 1, 1-5, Jan. 1962.
- [16] A. Cohen and S. G. Nooteboom, Structure and Process in Speech Perception. New York: Springer-Verlag, 1975.
- [17] J. D. Markel and A. H. Gray, Linear Prediction of Speech. New York: Springer-Verlag, 1966.
- [18] F. D. Minifie, T. J. Hixon and F. Williams, Normal Aspects of Speech, Hearing, and Language. Englewood Cliffs, New Jersey: Prentice-Hall, Inc., 1973.
- [19] W. R. Zemlin, Speech and Hearing Science Anatomy and Physiology. Englewood Cliffs, New Jersey: Prentice-Hall, Inc., 1968.
- [20] C. W. Kanliner and R. West, Phonetics. New York: Harper and Brothers Publishers, 1960.
- [21] G. Fant, Speech Communication, Speech Wave Processing and Transmission, Vol. 1. New York: John Wiley and Sons, 1975.
- [22] C. H. Coker, "A Model of Articulatory Dynamics and Control," Proceedings of the IEEE, Vol. 64, No. 4, 452-459, Apr. 1976.
- [23] C. J. Weinstein, S. S. McCandless, L. F. Mondschein and V. W. Zue, "A System for Acoustic-Phonetic-Analysis of Continuous Speech," IEEE Trans. on Acous. Sp. and Sig. Proc., Vol. ASSP-23, No. 1, 54-67, Feb. 1975.
- [24] K. N. Stevens and A. S. House, "Development of a Quantitative Description of Vowel Articulation," J. Acoust. Soc. Am., Vol. 27, No. 3, 484-492, May 1955.
- [25] R. Schwartz and J. Makhoul, "Where the Phones are: Dealing with Ambiguity in Acoustic-Phonetic Recognition," IEEE Trans. on Acous. Sp. and Sig. Proc., Vol. ASSP-23, No. 1, 50-53, Feb. 1975.

- [26] I. Lehiste and G. E. Peterson, "Transitions, Glides, and Dip-Thongs," J. Acoust. Soc. Am., Vol. 33, No. 3, 268-277, Mar. 1961.
- [27] J. Makhoul, "Linear Prediction: A Tutorial Review," Proceedings of the IEEE, Vol. 63, 561-580, Apr. 1975.
- [28] O. Fujimura, "Analysis of Nasal Consonants," J. Acoust. Soc. Am., Vol. 34, No. 12, 1865-1875, Dec. 1962.
- [29] B. S. Atal and S. L. Hanauer, "Speech Analysis and Synthesis by Linear Prediction of the Speech Wave," J. Acoust. Soc. Am., Vol. 50, No. 2 (Part 2), 637-655, Aug. 1971.
- [30] H. Wakita, "Direct Estimation of the Vocal Tract Shape by Inverse Filtering of Acoustic Speech Waveforms," IEEE Trans. Audio and Electroacoustics, Vol. AU-21, No. 25, 417-427, Oct. 1973.
- [31] R. F. Purton, "Speech Recognition Using Autocorrelation Analysis," IEEE Trans. Audio and Electroacoustics, Vol. AU-16, No. 2, 235-239, June 1968.
- [32] L. R. Rabiner, M. R. Sambur and C. E. Schmidt, "Applications of a Nonlinear Smoothing Algorithm to Speech Processing," IEEE Trans. on Acous. Sp. and Sig. Proc., Vol. ASSP-23, No. 6, 552-557, Dec. 1975.
- [33] W. J. Hess, "Time-Domain, Digital Segmentation of Connected Natural Speech," International Conference on Artificial Intelligence, 4th TBILIST, USSR, 491-498, 1975.
- [34] G. Fant, "The Acoustics of Speech," Proc. Third Intern. Congr. Acoust., 188-201, 1959.
- [35] H. K. Dunn, "The Calculation of Vowel Resonances and an Electrical Vocal Tract," J. Acoust. Soc. Am., Vol. 22, 740-753, 1950.
- [36] B. L. Bowerman and R. T. O'Connell, Time Series and Forecasting. Dunbury Press, A Division of Wadsworth, Inc., 1979.
- [37] H. Lane, "The Motor Theory of Speech Perception: A Critical Review," Psychological Review, Vol. 72, No. 4, 275-309, 1965.
- [38] R. A. Fisher, "The Use of Multiple Measurements in Taxonomic Problems," Ann. Eugen., Vol. 7, 1979-188, 1936.
- [39] H. Suzuki, H. Kasuya and K. Kido, "The Acoustic Parameters for Vowel Recognition Without Distinction of Speakers," Proceedings of the Conf. on Speech Communication and Processing, Cambridge, England, 92-96, Mar. 1967.

- [40] B. S. Atal, "Linear Prediction of Speech--Recent Advances with Applications to Speech Analysis," Reprinted from: Speech Recognition. New York: Academic Press, Inc., 221-230, 1975.
- [41] J. D. Markel and A. Gray, "On Autocorrelation Equations as Applied to Speech Analysis," IEEE Trans. on Audio and Electroacoustics, Vol. AU-21, No. 2, 69-79, Apr. 1973.
- [42] F. L. Wightman and D. M. Green, "The Perception of Pitch," American Scientist, Vol. 62, No. 2, 208-215, Mar.-Apr. 1974.
- [43] J. B. Thomas and R. J. Niederjohn, "Enhancement of Speech Intelligibility at High Noise Levels by Filtering and Clipping," Journal of the Audio Engineering Society, Vol. 16, No. 4, 412-415, Oct. 1968.
- [44] M. J. Hunt, J. S. Bridle and J. N. Holmes, "Interactive Digital Inverse Filtering and Its Relation to Linear Prediction Methods," International Conference on Acoustics, Speech and Signal Processing, 78CH1285-6 ASSP, 15-18, 1978.
- [45] C. J. Weinstein, S. S. McCandless, L. F. Mondschein and V. W. Zue, "A System for Acoustic-Phonetic Analysis of Continuous Speech," IEEE Trans. on Acoust., Sp. and Sig. Proc., Vol. ASSP-23, No. 1, 54-67, Feb. 1975.
- [46] B. T. Oshika, V. W. Zue, R. V. Weeks, H. N. Neu, and J. Aurbach, "The Role of Phonological Rules in Speech Understanding Research," IEEE Trans. on Acoust. Sp. and Sig. Proc., Vol. ASSP-23, No. 1, 104-112, Feb. 1975.
- [47] H. F. Silverman and N. R. Dixon, "A Parametrically Controlled Spectral Analysis System for Speech," IEEE Trans. on Acoust., Sp. and Sig. Proc., Vol. ASSP-22, No. 5, 362-381, Oct. 1974.
- [48] N. R. Dixon and H. F. Silverman, "A General Language-Operated Decision Implementation System (GLODIS): Its Application to Continuous-Speech Segmentation," IEEE Trans. on Acoust., Sp. and Sig. Proc., Vol. ASSP-24, No. 2, 137-162, Apr. 1976.
- [49] H. Dudley and S. Balashek, "Automatic Recognition of Phonetic Patterns in Speech," J. Acous. Soc. Am., Vol. 30, No. 8, 721-732, Aug. 1958.
- [50] P. Ladefoged and D. E. Broadbent, "Information Conveyed by Vowels," J. Acous. Soc. Am., Vol. 29, No. 1, 98-104, Jan. 1957.
- [51] G. W. Hughes and M. Halle, "Spectral Properties of Fricative Consonants," J. Acous. Soc. Am., Vol. 28, No. 2, 303-310, Mar. 1956.

- [52] D. R. Reddy and P. J. Vicens, "A Procedure for the Segmentation of Connected Speech," Journal of the Audio Engineering Society, Vol. 16, No. 4, 404-411, Oct. 1968.
- [53] A. S. House, "On Vowel Duration in English," J. Acous. Soc. Am., Vol. 33, No. 9, 1174-1178, Sept. 1961.
- [54] B. S. Atal and L. R. Rabiner, "A Pattern Recognition Approach to Voiced-Unvoiced-Silence Classification with Applications to Speech Recognition," IEEE Trans. on Acoust., Sp. and Sig. Proc., Vol. ASSP-24, No. 3, 201-212, June 1976.
- [55] G. E. Peterson and H. L. Barney, "Control Methods Used in a Study of the Vowels," J. Acous. Soc. Am., Vol. 24, No. 2, 175-184, Mar. 1952.
- [56] R. Viswanathan and J. Markhoul, "Quantization Properties of Transmission Parameters in Linear Predictive Systems," IEEE Trans. on Acoust., Sp. and Sig. Proc., Vol. ASSP-23, No. 3, 309-321, June 1975.
- [57] J. D. Markel, "Digital Inverse Filtering--A New Tool for Formant Trajectory Estimation," IEEE Transactions on Audio and Electroacoustics, Vol. AU-20, No. 2, 129-137, June 1972.
- [58] A. E. Rosenberg and M. R. Sambur, "New Techniques for Automatic Speaker Verification," IEEE Trans. on Acoust., Sp. and Sig. Proc., Vol. ASSP-23, No. 2, 1969-175, Apr. 1975.
- [59] S. S. McCandless, "An Algorithm for Automatic Formant Extraction Using Linear Prediction Spectra," IEEE Trans. on Acoust., Sp. and Sig. Proc., Vol. ASSP-22, No. 2, Apr. 1974.
- [60] R. W. Schafer, "A Survey of Digital Speech Processing Techniques," IEEE Transactions on Audio and Electroacoustics, Vol. AU-20, No. 1, 28-35, Mar. 1972.
- [61] J. Markhoul, "Spectral Linear Prediction: Properties and Applications," IEEE Trans. on Acoust., Sp. and Sig. Proc., Vol. ASSP-23, No. 3, 283-296, June 1975.
- [62] L. Molho, "Automatic Acoustic-Phonetic Analysis of Fricatives and Plosives," IEEE International Conference on Acoustics, Speech and Signal Processing, 182-189, 1976.
- [63] S. Brooks and F. Fallside, "A Technique for Converting the Linear Prediction Areas Model of Speech to a Simple Articulatory Model," IEEE International Conference on Acoustics, Speech, and Signal Processing, 71-78, 1976.

- [64] W. Bezdel and H. J. Chandler, "Results of an Analysis and Recognition of Vowels by Computer Using Zero-Crossing Data," Proceedings of the IEEE, Vol. 112, No. 11, 2060-2066, 1965.
- [65] W. Bezdel and J. S. Bridle, "Speech Recognition Using Zero-Crossing Measurements and Sequence Information," Proceedings of the IEEE, Vol. 116, No. 4, 617-622, 1969.
- [66] I. B. Thomas and R. J. Niederjohn, "A Preliminary Analysis Technique for Speech Sound Classification," Proceedings of the National Electronics Conference, Vol. 25, 685-690, 1969.
- [67] P. F. Abboud, N. A. Bezirgan, W. M. Erwin, M. A. Khouri, E. N. McCarus, and R. M. Rammuny, Introduction to Modern Standard Arabic Pronunciation and Writing. Ann Arbor, Mich.: University of Michigan, Department of Eastern Studies, 1968.
- [68] A. E. Rosenberg, "Effect of Glottal Pulse Shape on the Quality of Natural Vowels," J. Acoust. Soc. Am., Vol. 49, No. 2 (Part 2), 583-590, 1971.
- [69] M. Abdul-Rauf, "Arabic for English Speaking Students," Published by The Islamic Center, Manufactured by McGregor and Warner, Inc., Library of Congress Catalog Number 76-59592, 4th. Ed., 1979.
- [70] A. Y. Ali, "The Holy Quran," Published by American Trust Publishers for The Muslim Students' Association of the United State and Canada, Library of Congress Catalog Number 77-78098, ISBN No. 0-89259-006-8, pp. xvii, June 1977.

APPENDIX A

DATA ACQUISITION

Data Acquisition

It would, of course, be more desirable to develop a digit recognition system for digits spoken by randomly selected untrained speakers rather than for a trained speaker or group of trained speakers. In pursuit of this objective, seven speakers were selected, two American males, two American females and three international males. Each person spoke the digits zero through nine in American English. Also several combinations of three connected digits were spoken by the seven speakers. In addition, three sets of digits zero through nine were spoken in a sound proof chamber. Comparably, five Arabic speaking people were randomly selected and they spoke the digits zero through nine in Arabic.

The analog-to-digital conversion system used to prepare the data for computer processing is shown in Figure 105. A low-pass second order filter with a gain of 2 and 4 KHz cut-off frequency is used as shown in Figure 106. An 8 KHz sampling frequency is used. Furthermore, a 10 bit quantizer is used with an analog signal range of ± 10.24 volt. Therefore 20 mV/Bit resolution was obtained.

The speakers were instructed to depress a button as they began their speech. This caused a pulse to be entered on tape. The data entry procedure is not begun until this, triggering pulse, is detected in order to activate the system. Following detection of the triggering pulse, the voice data entry computer routine is activated, but the A/D conversion does not begin until the speech energy exceeds a preset threshold level. When this level is exceeded, the voice

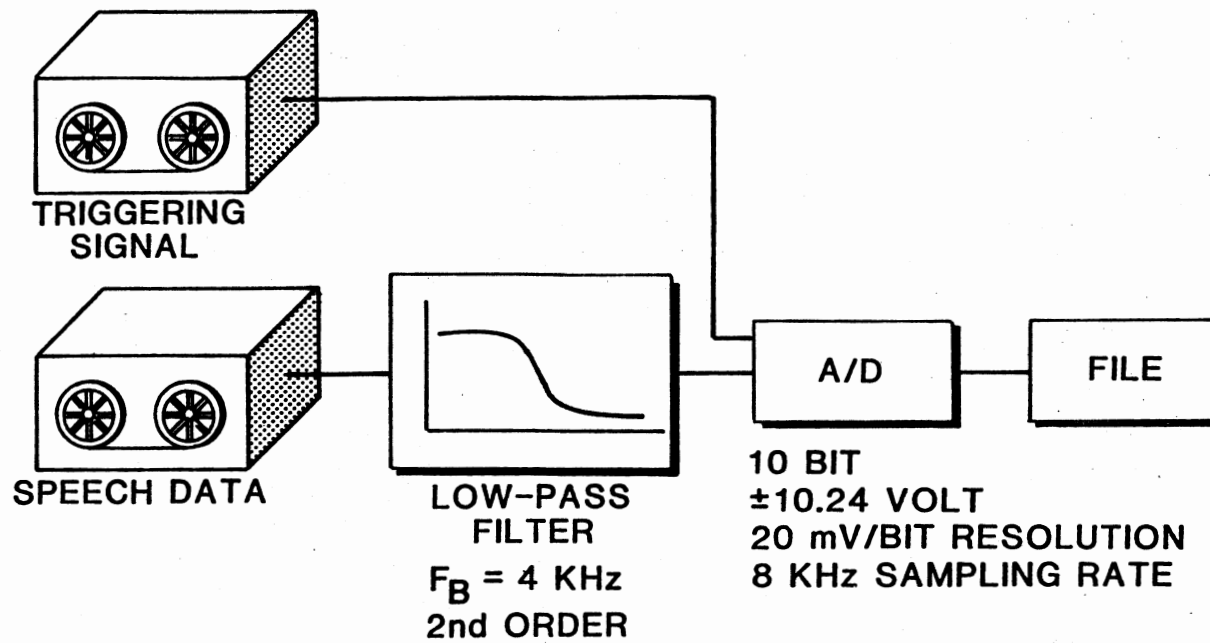
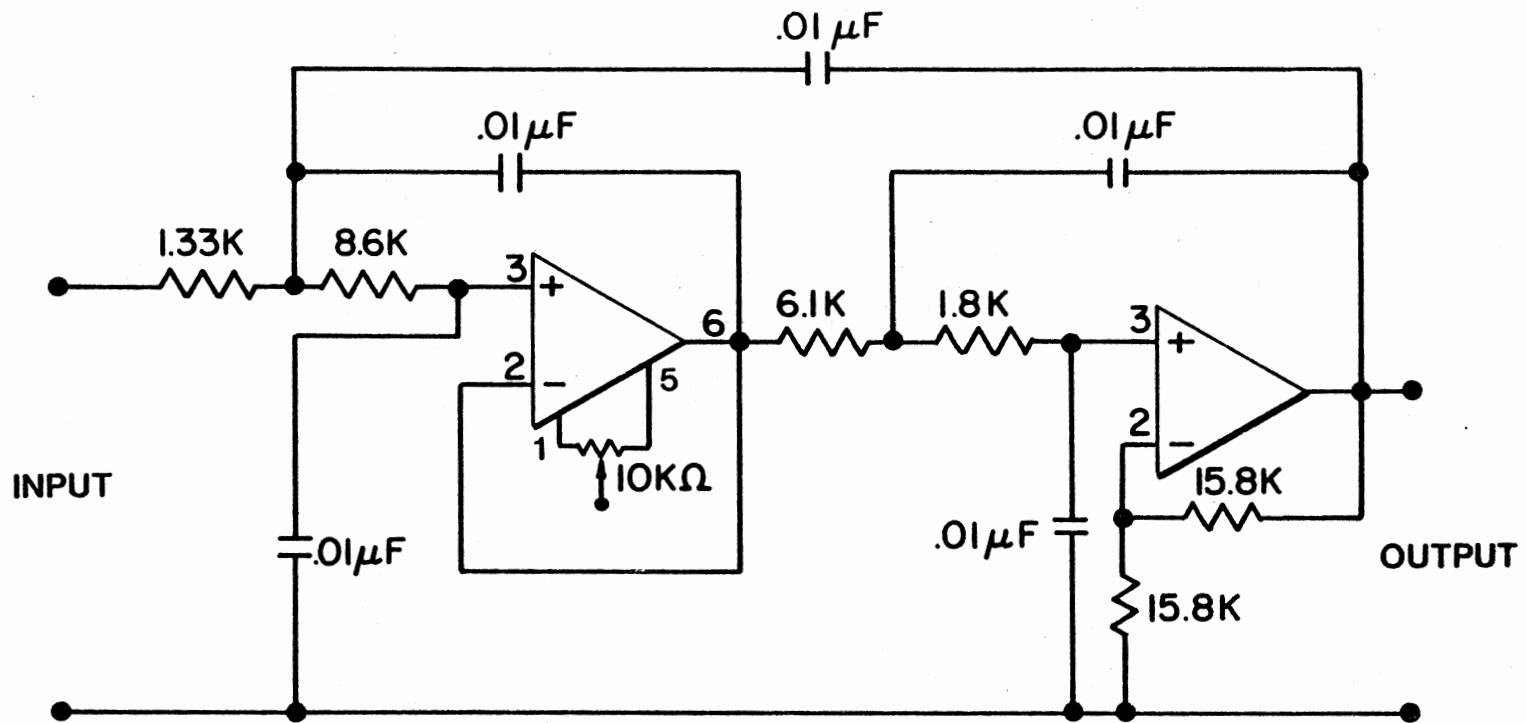


Figure 105. Data Acquisition



Pin **4** V_{CC-}

Pin **7** V_{CC+}

Pin **1** and **5** input for offset voltage (Null Circuit)

Figure 106. Second Order Low-Pass Filter

APPENDIX B

COMPUTER SUBROUTINES

Subroutine FRMS

This subroutine computes the RMS energy using 128 data points per frame. The details are explained in Chapter III. A flow chart is given in Figure 107.

Subroutine FBTR

In this routine the area functions from the reflection coefficients are computed. The BTR, CTR and FTR and SFBR are then estimated as discussed in Chapter III. A flow chart is given in Figure 108.

Subroutine Smooth

This routine performs 3 and 5 point median smoothing of a discretized input signal. A 3-point Hanning window smoother is used in conjunction with the median smoothers. A detailed discussion is given in Chapter III.

Subroutine Quantz

This routine performs the linear quantization of the data in an array. The resulting values are in the range 1-100 as explained in Chapter IV.

Subroutine FDIP

The maxima and the minima in the RMS energy contours are located and then passed to the dip-classification routine. A flow chart is given in Figure 109.

NPF = 128

Function: Find RMS energy
128 data point per frame

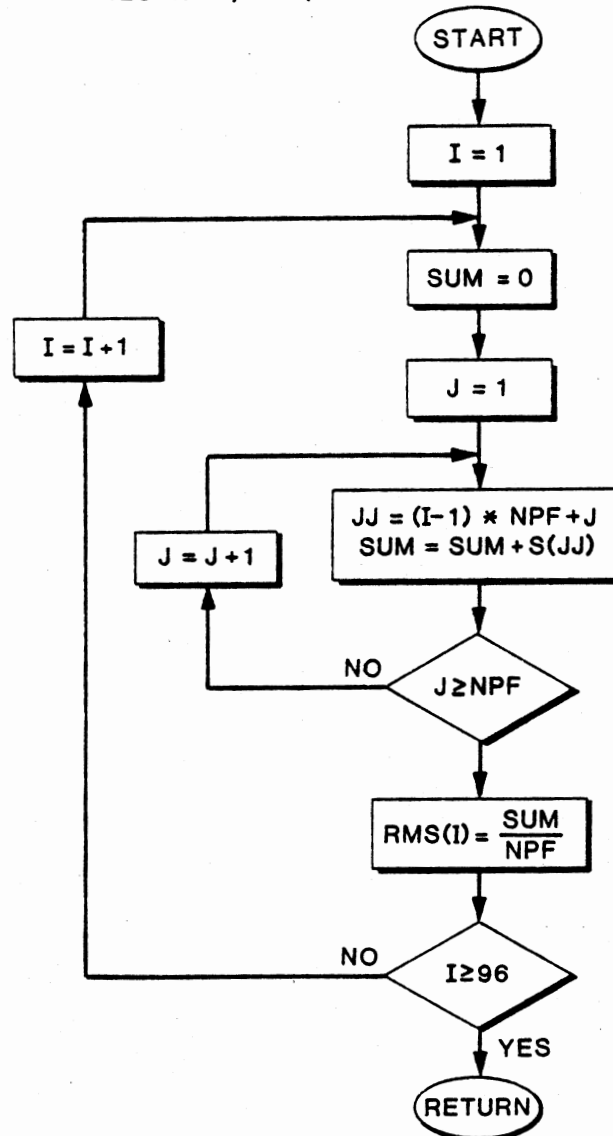


Figure 107. Subroutine FRMS

Function: Find the BTR and SFBR from the reflection coefficients

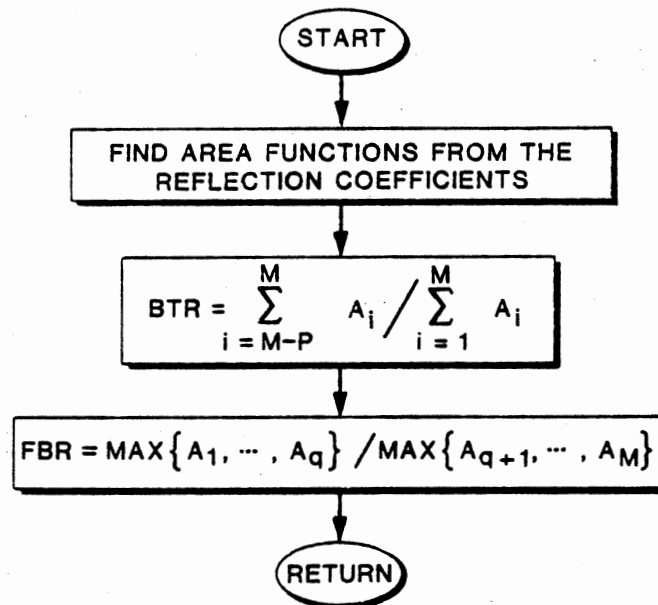


Figure 108. Subroutine FBTR

Function: Find the maximums and minimums
in the RMS energy

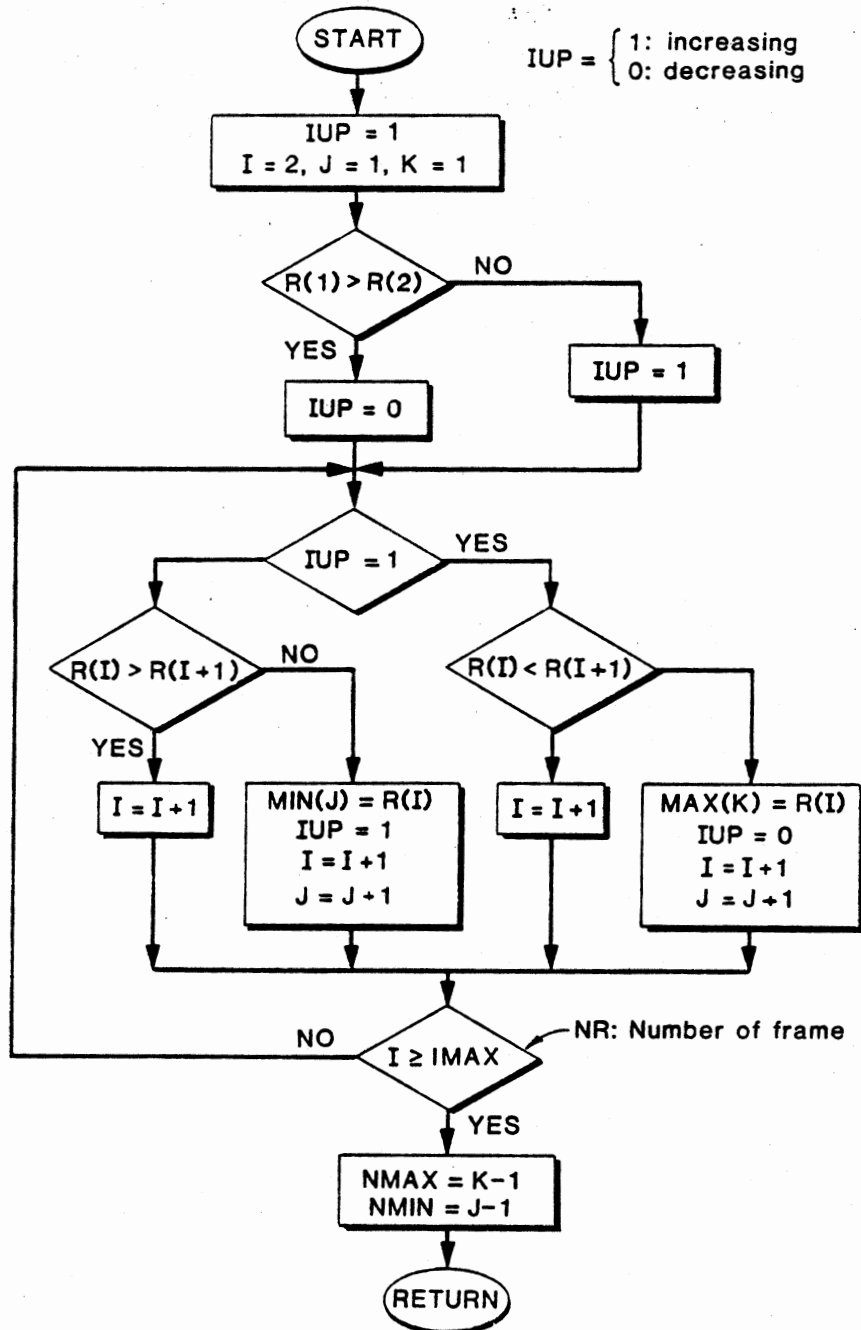


Figure 109. Subroutine FDIP

Subroutine FCSD

It finds the statistical distribution of $V_i^+ - V_j^-$ and $V_{i+1}^+ - V_j^-$, which will be used later to determine the slicing functions SL2 and SL2, as explained in Chapter III.

Subroutine Class

This routine classifies and identifies the three types of dips according to the slicing functions and the sign of Z_1 , Z_2 , and Z_3 as discussed in Chapter III. A flow chart is given in Figure 110.

Subroutine DWIND

This subroutine windows the data prior to low-frequency pre-emphasis. It uses a 150 point Hamming window.

Subroutine Auto [17]

This routine is for implementing the autocorrelation method of linear prediction, and uses 128 points.

Subroutine DIGIT

This routine recognizes the spoken digit using decisions based on phoneme sequence as discussed in Chapters II, IV, and V.

Subroutine Spect

This routine finds the ratio of the energy in the high frequency to that in the low frequency.

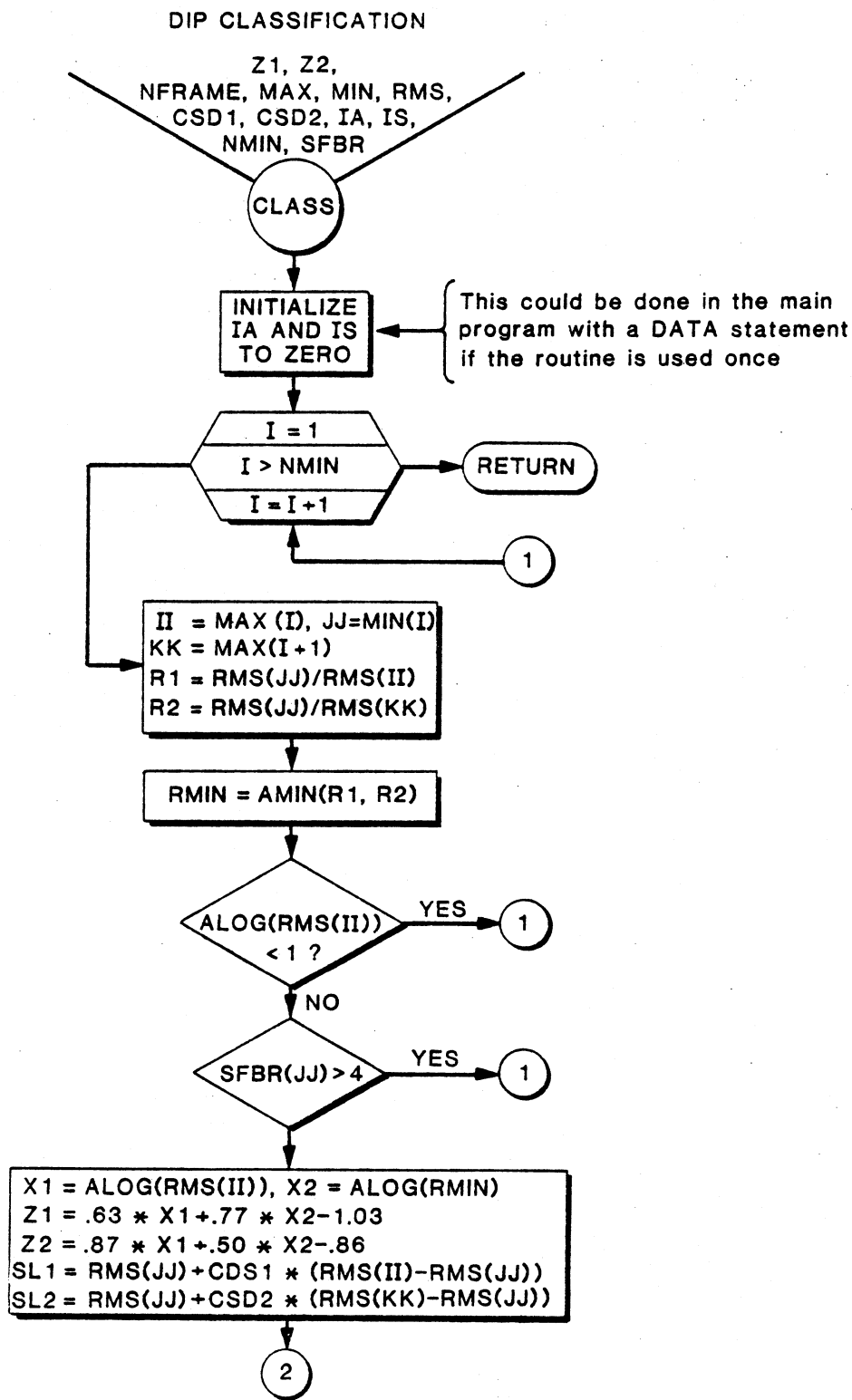


Figure 110. Subroutine Class

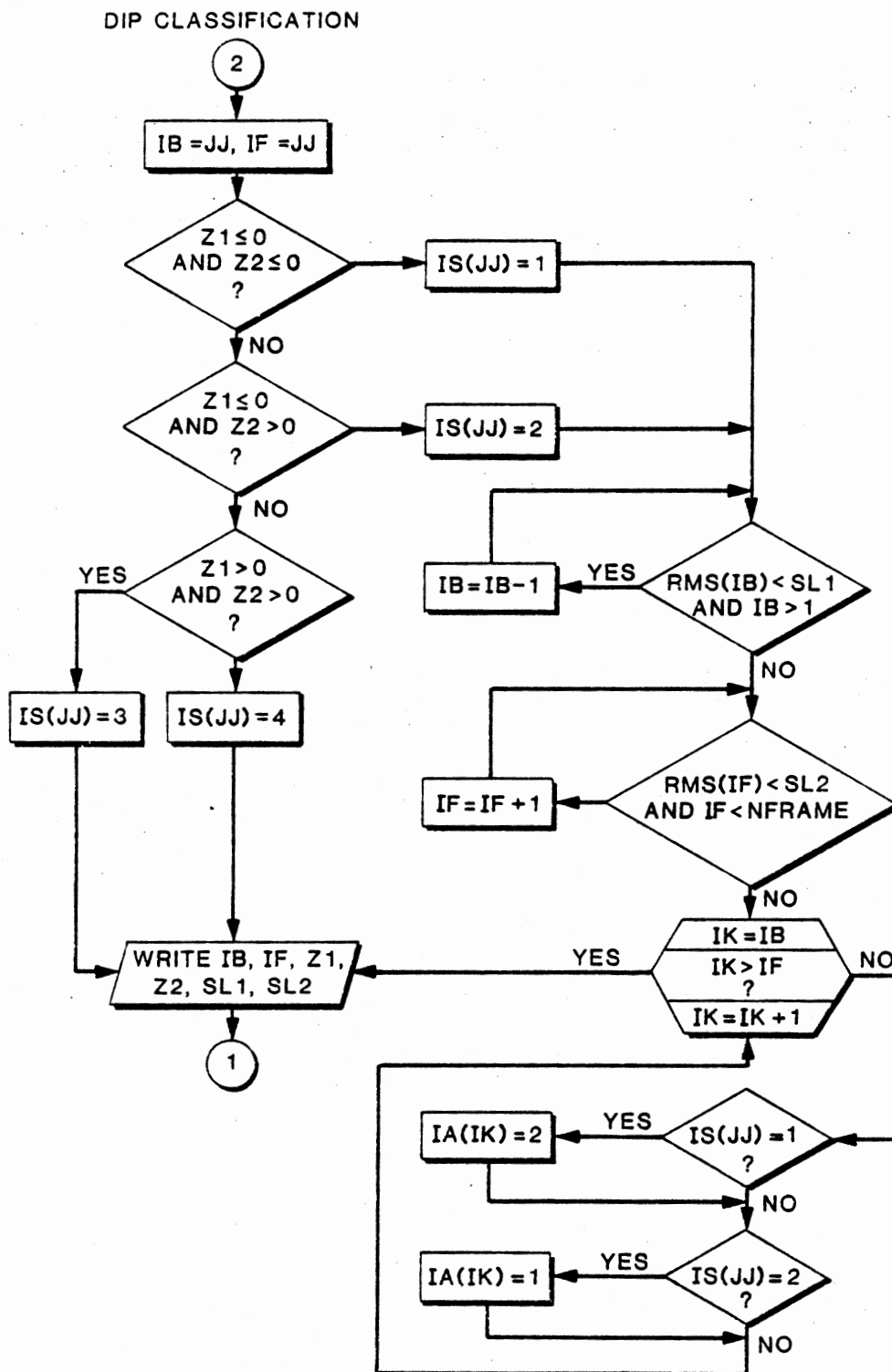


Figure 110. (Continued)

Subroutine UGETIO

This routine plots the RMS energy, BTR, CTR, FTR and the decision for vowel, vowel-like and non-vowel intervals.

```

C
C
C.....MAIN PROGRAM.....
C
C
C      INTEGER EV
C      DIMENSION IX(10000)
C      DIMENSION INoise(400)
C      DIMENSION RMS(400),SRMS(400),ORMS(400)
C      DIMENSION B(256),RANGE(4),EVLVL(7)
C      DIMENSION BTR(400) , SFBR(400) , S(15000) , S1(256) , SISM(512)
C      DIMENSION IS(400)
C      DIMENSION IA(400) , MAX(100) , MIN(100)
C      DIMENSION X(400)
C      DIMENSION KC(21)
C      DIMENSION IIRA(400) , EV(7)
C      INTEGER V,VL,ATYPE,BTYPE,STYPE,HTYPE,STYPE1,STYPE2,FINTYP
C      DATA LP/0/,NV,V,VL,NVL,ELNK/'NV',V',VL',NVL',',',/
C      DATA RANGE/16.,13.,7.,1./
C
C
C
C.....READ THE INPUT DATA.....
C      READ(5,100) ISTART,ISTCP
C
C      READ(3,1) IFILE,INUMR
100  FORMAT(2I6)
C      *WRITE(LP,111) IFILE,INUMR
111  FORMAT('1',10X,'FILE -- ',14,10X,'NUMBER OF POINTS',16)
C      READ(8,2)(IX(K),K=1,ISTCP)
2    FORMAT (20I4)
C      IDIF = ( ISTCP - ISTART + 1 )
C      NF=IDIF/NPS
C      IF (NF*NPS.NE.IDIF) NF=NF+1
C      NDIM=NF*NPS+ISTART
C      IK=ISTCP+1
C      DO 3 I=IK,NDIM
C          IX(I)=0
C
C
C
C      Y1 = 0
C      DO 19 I=1,512
19  Y1 = FLUAT(IX(I)) + Y1
C      DCLEV = Y1 / 512
C
C      WRITE (LP,*) DCLEV
C
C
C.....SCALE THE GIVEN DATA.....
C      J=0
C      DO 6 I=ISTART,NDIM
C          J=J+1
C          S(J)=(FLUAT(IX(I))+.3)*SCALE
C      CONTINUE
C
C      CALL FRMS (NF,NPS,S,RMS)
C
C      WRITE(LP,193)(RMS(J),J=1,NF)
193  FORMAT('1'////11X,'UNSMOOTHED RMS VALUES -- '//
C          (11X,SE20.7))
C      CALL PLOTFR(NF,RMS)
C      CALL SMOOTH (RMS,SRMS,NF)
C      CALL PLOTFR(NF,SRMS)
C      WRITE(LP,197)(SRMS(J),J=1,NF)

```

```

0000010
0000020
0000030
0000040
0000050
0000060
0000070
0000080
0000090
0000100
0000110
0000120
0000130
0000140
0000150
0000160
0000170
0000180
0000190
0000200
0000210
0000220
0000230
0000240
0000250
0000260
0000270
0000280
0000290
0000300
0000310
0000320
0000330
0000340
0000350
0000360
0000370
0000380
0000390
0000400
0000410
0000420
0000430
0000440
0000450
0000460
0000470
0000480
0000490
0000500
0000510
0000520
0000530
0000540
0000550
0000560
0000570
0000580
0000590
0000600
0000610
0000620
0000630
0000640
0000650
0000660
0000670

```

```

197      WRITE(LP,197)(SRMS(J),J=1,NF)
      FORMAT('1'///11X,'SMOOTHED RMS VALUES --'//
            (11X,5E20.7))
      CALL QUANTZ(SRMS,QRMS,NF,RMAX,RMIN)
      CALL PLUTER(NF,QRMS)
      WRITE(LP,199)(QRMS(J),J=1,NF)
199      FORMAT('1'///11X,'QUANTIZED RMS VALUES --'//
            (11X,5E20.7))
C
C
C      WRITE (LP,7)
7      FORMAT (1HL)
C
C
C      CALL FUIP(QRMS,MAX,MIN,NMAX,NMIN,NF)
      CALL FCSD(NMAX,NMIN,RMS,MAX,MIN,CSD1,CSD2)
C
C      WRITE (0,7)
C
C      WRITE(0,1010) NMAX
      WRITE(0,1020) NMIN
      WRITE(0,1030)(MAX(J),J=1,NMAX)
      WRITE(0,1040)(MIN(J),J=1,NMIN)
      WRITE(0,1050) CSD1,CSD2
1010  FORMAT(' NMAX= ',(1X,10I10))
1020  FORMAT(' NMIN= ',(1X,10I10))
1030  FORMAT(' MAX= ',/(1X,10I10))
1040  FORMAT(' MIN= ',/(1X,10I10))
1050  FORMAT(' CSD1= ',1PE14.6,/, ' CSD2= ',1PE14.6)
C
C
      NPW=200
      DO 5 I=1,NF
      DO 4 J=1,NPW
          K = J + (I-1)*NPW
          S1(J) = S(K)
4      CONTINUE
      SAVE=S(K)
      M1 = 14
      IPT = IPT
C
      CALL DWINDF (S1,NPW,10,1,0,XPRE)
      XPRE=SAVE
      CALL SLP (NPS,S1,M1,IF1,K1,K2,ALPHA,RC)
      CALL FBTR(RC,14,3,4,ANS1,ANS2)
C
      BTR(1)=ANS1
      SFBR(1)=ANS2
C
      CALL SMOOTH(S1,S1SM,NPS)
      CALL SPECT(S1SM,HLR)
C      HLR(1) = HLR
C      CONTINUE
C      WRITE (0,7)
C      WRITE(0,17) HLR
17      FORMAT(1HU,10X,'HLR --',F10.3)
C
      WRITE(0,1060) (BTR(KL),KL=1,NF)
1060  FORMAT(' BTR= ',/(1X,1PE14.6))
      WRITE (0,7)
      WRITE(0,1070) (SFBR(KL),KL=1,NF)
1070  FORMAT(' SFBR= ',/(1X,1PE14.6))
      WRITE (0,7)
      CALL PLUTER(NF,BTR)
      CALL PLUTER(NF,SFBR)

```

```

C0000670
00000680
00000690
00C00700
00000710
00000720
00000730
00000740
00000750
00000760
00000770
00000780
00000790
00000800
00000810
00000820
00000830
00000840
00000850
00000860
00000870
00000880
00000890
00C00900
00000910
00000920
00000930
00000940
00000950
00000960
00000970
00000980
00000990
00001000
00001010
00001020
00001030
00001040
00001050
00001060
00001070
00001080
00001090
00001100
00001110
00001120
00001130
00001140
00001150
00001160
00001170
00C01180
00001190
00C01200
00001210
00001220
00001230
00001240
00001250
00001260
00001270
00C01280
00001290
00001300
00001310
00001320
00001330

```



```

CALL PLOTTER(NF,SFBR)
GO TO 99
C
CSD1 = 0.30
CSD2 = 0.17
CALL CLASS (NF,MAX,MIN,QRMS,CSD1,CSD2,IA,IS,NMIN,SFER)
L
C
C PRELIMINARY DECISIONS
C
C PRINT THE BTR, SFBR, AND HLR VALUES AND THEIR TYPE
C
C WRITE (LP,41)
490 41 FORMAT('1'////' ',T21,'** PRELIMINARY DECISIONS **'////' ',
5 T13,'FRAME',T27,'RMS',T37,'IA',T43,'ATYPE',T57,'BTR',T63,
3 'BTYPE',T77,'SFBR',T83,'STYPE',T97,'HLR',T103,'HTYPE')
C
DO 1600 J=1,NF
L
C
C ATYPE=BLNK
C IF (IA(J) .EQ. 2) ATYPE=NV
C IF (IA(J) .EQ. 1) ATYPE=VL
C IF (IA(J) .EQ. 0) ATYPE=V
C
C BTYPE=BLNK
500 C IF (BTR(J) .GT. 0.15) BTYPE=NV
C IF (BTR(J) .LE. 0.15) BTYPE=V
C
C IF (SFBR(J) .GT. (-2.5)) STYPE=NV
C IF (SFBR(J) .LE. (-2.5)) STYPE=V
C IF (SFBR(J) .GT. 0.) STYPE=VL
C
C HTYPE=BLNK
C IF (HLRA(J) .LE. .2) HTYPE=NVL
C IF (HLRA(J) .GT. .2) HTYPE=VL
L
C
C WRITE (LP,43) J,SRMS(J),IA(J),ATYPE,BTR(J),BTYPE,
43 SFBR(J),STYPE,HLRA(J),HTYPE
3 FORMAT(' ',T11,I5,T21,F10.3,T36,I3,T44,A3,T51,F10.3,
3 T64,A3,T71,F10.3,T84,A3,T91,F10.3,T104,A3)
C
C 1600 CONTINUE
C
C FIRST SECONDARY DECISIONS
C PRELIMINARY DECISION BASED ON SFBR AND HLR
C
C WRITE(LP,47)
47 FORMAT('1'////' ',T13,'FRAME',T27,'SFBR',T41,'HLR',T49,'STYPE1')
C
C DO 1800 J=1,NF
C STYPE1=NVL
C IF (SFBR(J) .LE. (-2.5) .OR. HLRA(J) .GT. .2) STYPE1=VL
C WRITE(LP,51) J,SFBR(J),HLRA(J),STYPE1
51 FORMAT(' ',T11,I5,T21,F10.2,T36,F10.2,T49,A3)
1800 CONTINUE
C
C WRITE (6,7)
C DO 18J1 J=1,NF
C STYPE1=VL
C IF (SFBR(J) .GT. (-2.5) .OR. HLRA(J) .LE. .2) STYPE1=NVL
C WRITE(LP,51) J,SFBR(J),HLRA(J),STYPE1
1801 CONTINUE
C
C SECONDARY DECISION
C
C WRITE(LP,59)
C FORMAT('1'////' ',T21,'** SECONDARY DECISION **'////
00001330
00001340
00001350
00001360
00001370
00001380
00001390
00001400
00001410
00001420
00001430
00001440
00001450
00001460
00001470
00001480
00001490
00001500
00001510
00001520
00001530
00001540
00001550
00001560
00001570
00001580
00001590
00001600
00001610
00001620
00001630
00001640
00001650
00001660
00001670
00001680
00001690
00001700
00001710
00001720
00001730
00001740
00001750
00001760
00001770
00001780
00001790
00001800
00001810
00001820
00001830
00001840
00001850
00001860
00001870
00001880
00001890
00001900
00001910
00001920
00001930
00001940
00001950
00001960
00001970
00001980
00001990

```

```

59      FORMAT('1'////' ',T21,'*** SECONDARY DECISION ***'//
          ' ',T13,'FRAME',T28,'IA',T43,'BTR',T57,'STYPE2'//)
C
      DO 1900 K=1,NF
          STYPE2=V
          IF (IA(K).EQ.2 .OR. BTR(K).GT. .15) STYPE2=NV
          WRITE(LP,61)K,IA(K),BTR(K),STYPE2
61      FORMAT(' ',T11,IS,T25,17,T35,F10.3,T53,A3)
1900    CONTINUE
C
          WRITE (6,7)
          DO 1901 K=1,NF
              STYPE2=NVL
              IF (IA(K).EQ.1 .OR. BTR(K).LE. .15) STYPE2=VL
              WRITE(LP,61)K,IA(K),BTR(K),STYPE2
1901    CONTINUE
C
          WRITE (6,7)
          DO 1902 K=1,NF
              STYPE2=NV
              IF (IA(K).EQ.0 .OR. BTR(K).LE. .15) STYPE2=V
              WRITE(LP,61)K,IA(K),BTR(K),STYPE2
1902    CONTINUE
C
          FINAL DECISION
C
          WRITE(LP,63)
63      FORMAT('1'////' ',T21,'*** FINAL DECISION ***'////' ',T13,
          'FRAME',T27,'RMS',T35,'FINTYP'//)
C
          DO 2100 K=1,NF
              FINTYP=VL
              IF (IA(K).EQ.2 .OR. BTR(K).GT. .15 .OR. SFBR(K).GT.(-2.5)
          .OR. HLRA(K).LE. .2) FINTYP=NVL
              WRITE(LP,67)K,SRMS(K),FINTYP
67      FORMAT(' ',T11,IS,T21,F10.2,T36,A3)
2100    CONTINUE
C
              WRITE (LP,7)
              DO 2101 K=1,NF
                  FINTYP=NVL
                  IF (IA(K).EQ.1 .OR. BTR(K).LE. .15 .OR. SFBR(K).LE.(-2.5)
          .OR. HLRA(K).GT. .2) FINTYP=VL
                  WRITE(LP,67)K,SRMS(K),FINTYP
2101    CONTINUE
C
              WRITE (LP,7)
              DO 2102 K=1,NF
                  FINTYP=V
                  IF (IA(K).EQ.0 .OR. BTR(K).LE. .15 .OR. SFBR(K).LE.(-2.5)
          .OR. HLRA(K).GT. .2) FINTYP=V
                  WRITE(LP,67)K,SRMS(K),FINTYP
2102    CONTINUE
C
          SUM=0.
          DO 6000 J=1,NF
C
              SUM=SUM+SRMS(J)
6000    CONTINUE
C
          AVE=SUM/FLOAT(NF)
          WRITE(LP,205)SUM,NF,AVE,SMEAN,SCALE
00002190
00002000
00002010
00002020
00002030
00002040
00002050
00002060
00002070
00002080
00002090
00002100
00002110
00002120
00002130
00002140
00002150
00002160
00002170
00002180
00002190
00002200
00002210
00002220
00002230
00002240
00002250
00002260
00002270
00002280
00002290
00002300
00002310
00002320
00002330
00002340
00002350
00002360
00002370
00002380
00002390
00002400
00002410
00002420
00002430
00002440
00002450
00002460
00002470
00002480
00002490
00002500
00002510
00002520
00002530
00002540
00002550
00002560
00002570
00002580
00002590
00002600
00002610
00002620
00002630
00002640
00002650

```



```

SUBROUTINE PLOTER(NFRAME,ARRAY)
DIMENSION ANRAY(1)
DIMENSION X(400)
DIMENSION B(400,1)
INTEGER IMAG4(5151),ITITLE(144),ICHAR(10)
REAL RANGE(4)
DATA ICHAR(1)/1H0/,RANGE/4*0.0/

CALL UGETIJ(1,NIN,NOUT)

READ (NIN,1001) (ITITLE(I),I=1,144)
1001 FORMAT (72A1)

DO 21 I=1 , NFRAME
B(I,1) = ARRAY(I)
X(I) = FLOAT(I)
21 CONTINUE

INC=1
MM=1
IY = NFRAME
IOPT=1

CALL USPLT(X,B,IY,NFRAME,MM, INC, ITITLE,RANGE, ICHAR, IOPT, IMAG4, IER)

RETURN
END

SUBROUTINE FDP(RMS,MAX,MIN,NMAX,NMIN,NR)
DIMENSION RMS(1),MAX(100),MIN(100)

IUP=1
I=2
J=1
K=1
DO 165 MAA=1,100

MAX(MAA) = (NR + 1)
MIN(MAA) = (NR + 1)

165 CONTINUE
IF(RMS(1) .GT. RMS(2)) GO TO 99
IUP=1
GO TO 100
99 IUP=0
MAX(K)=1
K=K+1
17 CONTINUE
100 IF (IUP .EQ. 1) GO TO 19

IF (RMS(1) .GE. RMS(I+1)) GO TO 18
MIN(J)=1
IUP=1
J=J+1
18 GO TO 20
19 IF(RMS(1) .LE. RMS(I+1)) GO TO 20
MAX(K)=1
IUP=0
K=K+1
20 I=I+1

IF (I .LT. NR) GO TO 17

```

```

00003310
00003320
00003330
00003340
00003350
00003360
00003370
00003380
00003390
00003400
00003410
00003420
00003430
00003440
00003450
00003460
00003470
00003480
00003490
00003500
00003510
00003520
00003530
00003540
00003550
00003560
00003570
00003580
00003590
00003600
00003610
00003620
00003630
00003640
00003650
00003660
00003670
00003680
00003690
00003700
00003710
00003720
00003730
00003740
00003750
00003760
00003770
00003780
00003790
00003800
00003810
00003820
00003830
00003840
00003850
00003860
00003870
00003880
00003890
00003900
00003910
00003920
00003930
00003940
00003950

```

```

NMAX=K-1
NMIN=J-1
RETURN
END
C
C
C
C.....FIND THE STATISTICAL DISTRIBUTION OF VI-VJ.....
C
SUBROUTINE FCSD(NMAX,NP IN,RMS,MAX,MIN,CSD1,CSD2)
DIMENSION RMS(1),MAX(1),MIN(1)
SUM=0.0
C
DO 244 I=1,NMAX
LL=MAX(I)
NN=MIN(I)
SUM=SUM+(RMS(LL)-RMS(NN))
244 CONTINUE
CSD1=SUM/FLD(1,NMAX)
C
SUM=0.0
NMA=NMAX-1
DO 25 I=1,NMA
LL=MAX(I+1)
NN=MIN(I)
SUM=SUM+(RMS(LL)-RMS(NN))
25 CONTINUE
CSD2=SUM/FLD(1,NMA)
RETURN
END
C
C.....DIP CLASSIFICATION.....
C
SUBROUTINE CLASS(NFRAME,MAX,MIN,RMS,CSD1,CSD2,IA,IS,NMIN,SFBR)
DIMENSION MAX(1),MIN(1),RMS(1),IA(384),IS(384)
DIMENSION SFBR(1)
C
DO 22 I=1,NFRAME
IA(I)=0
IS(I)=0
22 CONTINUE
C
DO 24 I=1,NMIN
II=MAX(I)
JJ=MIN(I)
KK=MAX(I+1)
R1=RMS(JJ)/RMS(II)
R2=RMS(JJ)/RMS(KK)
C
RMIN=AMIN1(R1,R2)
190 IF(ALOG(RMS(II)) .LT. 1.0 .OR. SFBR(JJ) .GT. 4.0) GO TO 24
C
X1=ALOG(RMS(II))
X2=ALOG(RMIN)
C
Z1=0.6J*X1+0.77*X2-1.03
Z2=0.87*X1+0.50*X2-0.86
C
SL1=RMS(JJ)+CSD1*(RMS(II)-RMS(JJ))
SL2=RMS(JJ)+CSD2*(RMS(KK)-RMS(JJ))
C
C
IB=JJ
IF=JJ
C
... DIP ...

```

```

00003950
00003960
00003970
00003980
00003990
00004000
00004010
00004020
00004030
00004040
00004050
00004060
00004070
00004080
00004090
00004100
00004110
00004120
00004130
00004140
00004150
00004160
00004170
00004180
00004190
00004200
00004210
00004220
00004230
00004240
00004250
00004260
00004270
00004280
00004290
00004300
00004310
00004320
00004330
00004340
00004350
00004360
00004370
00004380
00004390
00004400
00004410
00004420
00004430
00004440
00004450
00004460
00004470
00004480
00004490
00004500
00004510
00004520
00004530
00004540
00004550
00004560
00004570
00004580
00004590
00004600
00004610

```



```

      NMIN=J-1
      RETURN
C
      DATA ICHY/10/
      HAM(I) = .54 - .46*CCS((I-1)*CONST)
      PI = 3.14159265
      CONST = 2*PI/N
C
      IF (IHAM.NE.ICHY) GO TO 50
C
      HIGH FREQUENCY PREEPHASIS AND HAMMING WINDOW
C
      IF (PRE.LT.0.) GO TO 30
      B = XPRE
      DO 20 I=1,N
      A = Y(I)
      Y(I) = (A-PRE*B)*HAM(I)
20  B = A
      GO TO 100
C
      LOW FREQUENCY PREEPHASIS AND HAMMING WINDOW
C
30  B = 0.
      DO 40 I=1,N
      A = Y(I) - PRE*B
      Y(I) = A*HAM(I)
40  B = A
      GO TO 100
C
      HIGH FREQUENCY PREEPHASIS AND NO WINDOWING
C
50  IF (PRE.LT.0.) GO TO 80
      B = XPRE
      DO 70 I=1,N
      A = Y(I)
      Y(I) = A - PRE*B
70  B = A
      GO TO 100
C
      LOW FREQUENCY PREEPHASIS AND NO WINDOWING
C
80  B = 0.
      DO 90 I=1,N
      A = Y(I) - PRE*B
      Y(I) = A
90  B = A
100 RETURN
      END
C
      SUBROUTINE SLP (N,X,N,IF,K1,K2,ALPHA,RC)
C
      DIMENSION X(1) , Y(512) , RC(1)
      DIMENSION A(21) , B(21) , R(21)
C
      PI = 3.141592653
      I = 2**IP
      DO 10 J=1,I
10  Y(J) = 0.
      NP = N + I
      DO 20 J=NP,I
20  X(J) = 0.
C
      CALL FFT (X,Y,IP)
C
      L = K2 - K1
      LP = L + 1
      DO 30 J=1,LP
      JK = J + K1

```

```

00002200
00005270
00005280
00005290
00005300
00005310
00005320
00005330
00005340
00005350
00005360
00005370
00005380
00005390
00005400
00005410
00005420
00005430
00005440
00005450
00005460
00005470
00005480
00005490
00005500
00005510
00005520
00005530
00005540
00005550
00005560
00005570
00005580
00005590
00005600
00005610
00005620
00005630
00005640
00005650
00005660
00005670
00005680
00005690
00005700
00005710
00005720
00005730
00005740
00005750
00005760
00005770
00005780
00005790
00005800
00005810
00005820
00005830
00005840
00005850
00005860
00005870
00005880
00005890
00005900
00005910
00005920
00005930

```

```

      JK = J + KI
C
C  J0 X(J) = X(JK)*X(JK) + Y(JK)*Y(JK)
C
      MP = M + 1
      RL = X(LP)
      DO 50 J=1,MP
      R(J) = X(1) + RL
      RL = -RL
      LM = L - 1
      DO 40 K=1,LM
      ARG = (P)*K*(J-1) / L
40  R(J) = 2.*X(K+1)*COS(ARG) + R(J)
50  R(J) = 0.5*R(J) / L
      ALPHA = R(1)
      IF (K(1).NE.0.0) RC(1)=-R(2)/R(1)
      A(1) = 1.
      A(2) = RC(1)
      ALPHA = ALPHA - RC(1)*RC(1)*ALPHA
      DO 90 MINC=2,M
      MF = MINC - 1
      DO 60 J=1,MINC
      JB = MINC - J + 1
60  B(J) = A(JB)
      MF = MF + 1
      S = 0.
      DO 70 JP=1,MF
      MIP = MF-JP+2
70  S = S + R(MIP)*A(JP)
      IF (ALPHA.NE.0.0) RC(M)=-S/ALPHA
      DO 80 JP=2,MF
80  A(JP) = A(JP) + RC(MF)*E(JP-1)
      A(MF+1) = RC(MF)
      ALPHA = ALPHA - RC(MF)*RC(MF)*ALPHA
90  CONTINUE
C
C
C  IF (X(J).NE.0.0) PRINT *,*POWER SPECTRUM= ',X
C
      RETURN
      END
C
C  SUBROUTINE FFT (X,Y,L)
      DIMENSION X(1) , Y(1)
C
C      RADIX = 2FFT
C
      NP = 2**L
      LMX = NP
      SCL = 6.283185383/NP
      DO 20 LD=1,L
      LIX = LMX
      LMX = LMX/2
      ARG = 0.
      DO 10 LM=1,LMX
      C = COS(ARG)
      S = SIN(ARG)
C
      ARG = ARG + SCL
      DO 10 LI=LIX,NP,LIX
      J1 = LI - LIX + LM
      J2 = J1 + LMX
      T1 = X(J1) - X(J2)
      T2 = Y(J1) - Y(J2)
      X(J1) = X(J1) + X(J2)
      Y(J1) = Y(J1) + Y(J2)

```

```

00005930
00005940
00005950
00005960
00005970
00005980
00005990
00006000
00006010
00006020
00006030
00006040
00006050
00006060
00006070
00006080
00006090
00006100
00006110
00006130
00006140
00006150
00006160
00006170
00006180
00006190
00006200
00006210
00006220
00006230
00006240
00006250
00006260
00006270
00006280
00006290
00006300
00006310
00006320
00006330
00006340
00006350
00006360
00006370
00006380
00006390
00006400
00006410
00006420
00006430
00006440
00006450
00006460
00006470
00006480
00006490
00006500
00006510
00006520
00006530
00006540
00006550
00006560
00006570
00006580
00006590
00006600

```



```

X(J2) = C*T1*T2
10 Y(J2) = C*T2 - S*T1
20 SCL = 2.*SCL
C
C      BIT REVERSAL
C
      J = 1
      NV2 = NP/2
      NPM1 = NP - 1
      DO 50 I=1,NPM1
      IF (1.GE.J) GO TO 30
      T1 = X(J)
      T2 = Y(J)
      X(J) = X(I)
      Y(J) = Y(I)
      X(I) = T1
      Y(I) = T2
30 K = NV2
40 IF (K.GE.J) GO TO 50
      J = J - K
      K = K/2
      GO TO 40
50 J = J + K
      RETURN
      END
C
C      SUBROUTINE FBTR(RC,M,IP,IQ,ANS1,ANS2)
C
      DIMENSION AREA(20),RC(1)
      AREA(M+1)=1.0
      DO 10 J=1,M
          I = M-J+1
          AREA(I) = (1.+RC(I))*AREA(I+1)/(1.-RC(I))
10 CONTINUE
C
C.....FIND BIR.....
C
      SUM=0.
      SUM1=0.
      MP=M-IP
      DO 11 I=MP,M
          SUM1=SUM1+AREA(I)
11 CONTINUE
      DO 12 I=1,M
          SUM=SUM+AREA(I)
12 CONTINUE
C
      ANS1=SUM1/SUM
C
C.....FIND FBR.....
C
      TEMP1 = AREA(1)
      DO 13 I=1,IQ
          IF (AREA(I) .GT. TEMP1) TEMP1=AREA(I)
13 CONTINUE
      IR=IQ+1
      IM=M-IR
      TEMP2=AREA(IR)
      DO 14 I=1,IM
          J=I+IR
          IF(AREA(J) .GT. TEMP2) TEMP2=AREA(J)
14 CONTINUE
C
      FBR = TEMP1/TEMP2
C

```

```

00006610
00006620
00006630
00006640
00006650
00006660
00006670
00006680
00006690
00006700
00006710
00006720
00006730
00006740
00006750
00006760
00006770
00006780
00006790
00006800
00006810
00006820
00006830
00006840
00006850
00006860
00006870
00006880
00006890
00006900
00006910
00006920
00006930
00006940
00006950
00006960
00006970
00006980
00006990
00007000
00007010
00007020
00007030
00007040
00007050
00007060
00007070
00007080
00007090
00007100
00007110
00007120
00007130
00007140
00007150
00007160
00007170
00007180
00007190
00007200
00007210
00007220
00007230
00007240
00007250
00007260

```

```

C
C      IF(RC(1) .EQ. 0.) GO TO 19
C      SGN=RC(1)/ABS(RC(1))
C      GO TO 190
19      SGN = 1.0
190   ANS2=SGN*FBR
C
C      RETURN
C      END
C
C
C      SUBROUTINE SMOOTH (X,Y,NVAL)
C
C      THIS ROUTINE PERFORMS 3- AND 5-POINT MEDIAN SMOOTHING OF A
C      DISCRETIZED INPUT SIGNAL. A 3-POINT HANNING WINDOW SMCCTHER
C      IS USED IN CONJUNCTION WITH THE MEDIAN SMOOTHERS.
C
C      ARGUMENTS —
C      X - ARRAY CONTAINING SIGNAL REPRESENTATION
C      NVAL - NUMBER OF VALUES STORED IN X
C      Y - SMCCTHERD ARRAY
C
C      DIMENSION X(NVAL),Y(NVAL),ZMED(5),Z(402),W(400)
C      INIEGER PUS
C
C      PERFORM THE 3-POINT MEDIAN SMOOTHING
C
C      Y(NVAL)=X(NVAL)
C      Y(1)=X(1)
C      MAX=NVAL-1
C
C      DO 1000 J=2,MAX
C          X1=X(J-1)
C          X2=X(J)
C          X3=X(J+1)
C
C          FIND THE MEDIAN OF THE VALUES ABCLT X(J)
C
C          IF(X1.LT.X2) GO TO 200
C
C          IF(X3.LE.X2) GO TO 230
C
C          IF(X3.GE.X1) GO TO 240
C
C          GO TO 220
C
C          IF(X3.LE.X1) GO TO 240
C
C          IF(X3.GE.X2) GO TO 230
C
C          X3 IS THE MEDIAN
C          Y(J)=X3
C          GO TO 1000
C
C          X2 IS THE MEDIAN
C          Y(J)=X2
C          GO TO 1000
C
C          X1 IS THE MEDIAN
C          Y(J)=X1
200
230
240

```

```

00007260
00007270
00007280
00007290
00007300
00007310
00007320
00007330
00007340
00007350
00007360
00007370
00007380
00007390
00007400
00007410
00007420
00007430
00007440
00007450
00007460
00007470
00007480
00007490
00007500
00007510
00007520
00007530
00007540
00007550
00007560
00007570
00007580
00007590
00007600
00007610
00007620
00007630
00007640
00007650
00007660
00007670
00007680
00007690
00007700
00007710
00007720
00007730
00007740
00007750
00007760
00007770
00007780
00007790
00007800
00007810
00007820
00007830
00007840
00007850
00007860
00007870
00007880
00007890
00007900
00007910
00007920
00007930

```

240		Y(J)=X1	00007930
C			00007940
1000	CONTINUE		00007950
C			00007960
C			00007970
C			00007980
C	PERFORM THE 3-POINT FANNING SMOOTHING		00007990
C			00008000
C			00008010
C	Y2=Y(1)		00008020
C	Y3=Y2		00008030
C			00008040
C	DC 2000 J=1,MAX		00008050
C			00008060
C	K=J+1		00008070
C	Y1=Y2		00008080
C	Y2=Y3		00008090
C	Y3=Y(K)		00008100
C			00008110
C	Y(J)=.25*(Y1+Y3)+.5*Y2		00008120
C			00008130
C	PERFORM THE SUBTRACTION OF X AND Y SIGNALS		00008140
C			00008150
C	Z(J+1)=X(J)-Y(J)		00008160
2000	CONTINUE		00008170
C			00008180
C	Y(NVAL)=.25*(Y2+Y3)+.5*Y3		00008190
C	Z(NVAL+1)=X(NVAL)-Y(NVAL)		00008200
C			00008210
C			00008220
C			00008230
C	PERFORM A 5-POINT MEDIAN SMOOTHING ON THE X-Y DIFFERENCE		00008240
C			00008250
C			00008260
C	Z(NVAL+2)=Z(NVAL+1)		00008270
C	Z(1)=Z(2)		00008280
C			00008290
C	DU 3300 J=2,MAX		00008300
C			00008310
C	K=J+1		00008320
C	ZMED(1)=Z(K-2)		00008330
C	ZMED(2)=Z(K-1)		00008340
C	ZMED(3)=Z(K)		00008350
C	ZMED(4)=Z(K+1)		00008360
C	ZMED(5)=Z(K+2)		00008370
C			00008380
C			00008390
C			00008400
C			00008410
C			00008420
C			00008430
3150	DC 3150 N=M,5		00008440
C	IF(ZMED(POS).GT.ZMED(N)) POS=N		00008450
C			00008460
C	IF(POS.EQ.L) GO TO 3200		00008470
C			00008480
C	TEMP=ZMED(L)		00008490
C	ZMED(L)=ZMED(POS)		00008500
C	ZMED(POS)=TEMP		00008510
3200	CONTINUE		00008520
C			00008530
C	W(J)=ZMED(3)		00008540
3300	CONTINUE		00008550
C			00008560
C	W(1)=Z(2)		00008570
C	W(NVAL)=Z(NVAL+1)		00008580
C			00008590


```

          PTR=PTR+1
          GO TO 100
C
C
C
200 STATE=-STATE
C
C
C          DISTINGUISH BETWEEN 5 AND 9
C
C          AV=(BTR(1)+BTR(2)+BTR(3))/3.
C          IF(AV.GT. .15) STATE=9
C
C
C          RETURN
C          END

```

```

00009900
00009910
00009920
00009930
00009940
00009950
00009960
00009970
00009980
00009990
00010000
00010010
00010020
00010030
00010040

```

```

SUBROUTINE PLOTFR(NFRAME,ARRAY,ITITLE)
DIMENSION ARRAY(1),ITITLE(1)
DIMENSION X(400)
DIMENSION B(400,1)
INTERFACE IMAG4(5151), ICHAR(10)
REAL RANGE(4)
DATA ICHAR(1)/IHO/, RANGE/4*0./
C
C
C          CALL DGETIG(I,NIN,NCUT)
C
C
C
C          DO 21 I=1, NFRAME
C          B(I,1) = ARRAY(I)
C          X(I) = FLOAT(I)
21 CONTINUE
C
C          INC=1
C          MM=1
C          IY = NFRAME
C          IOPT=1
C
C          CALL USPLT(X,B,IY,NFRAME,MM,INC,ITITLE,RANGE, ICHAR, IOPT, IMAG4, IEF)
C
C          RETURN
C          END

```

```

00000820
00000830
00000840
00000850
00000860
00000870
00000880
00000890
00000900
00000910
00000920
00000930
00000940
00000950
00000960
00000970
00000980
00000990
00001000
00001010
00001020
00001030
00001040
00001050
00001060

```

```

DIMENSION IX(25000),S(25000),SI(512),HAM(512),RC(20)
C DIMENSION ITITLE(144),JTITLE(144),KTITLE(144)
DIMENSION ARRAY(10,40),BARRAY(10,40)
DIMENSION FTR(400),CTR(400),RTR(400),WORK(400),DMT(400)
C READ(5,105) ITITLE,JTITLE,KTITLE
105 FORMAT(72A1)
REAC(5,106) PRE
106 FORMAT(F6.2)
WRITE(6,107) PRE
107 FORMAT (//,10X,'SECOND ORDER PRE-EMPHASIS (GAIN = ',F6.2,')')
NPS = 128
SCALE=102
PI=3.141592
CONST=PI/75.0
DO 9 I=1,150
HAM(I)=.54-.46*CCS((I-1)*CONST)
3 CONTINUE
IDIGT=-1
C.....READ THE INPUT DATA.....
IFL=2
ISTART=1
1 READ(IFL,101) IJUNK,ISTOP
101 FORMAT(2I6)
IDIGT=IDIGT+1
WRITE(6,102) IDIGT,ISTOP
102 FORMAT(10X,'DIGIT -- ',14,10X,'NUMBER OF POINTS',I6)
REAC(IFL,104) (IX(K),K=1,ISTOP)
104 FORMAT (20I4)
IDIF = ( ISTOP - ISTART + 1 )
NF=IDIF/NPS
IF (NF*NPS.NE.IDIF) NF=NF+1
NDIM=NF*NPS+ISTART
IK=ISTOP+1
DO 3 I=IK,NDIM
IX(I)=0
J=1
DO 6 I=ISTART,NDIM
J=J+1
S(J)=FLOAT(IX(I))*SCALE
3 CONTINUE
NPW=150
NPS=128
DO 5 I=1,NF
DO 4 J=1,NPW
K=J+(I-1)*NPS
S1(J)=S(K)
4 CONTINUE
SAVE=S(1)
M1=14
I0=9
K1=0
K2=128
CALL DWINDF(S1,NPW,10,PRE,YFILE,HAY)
DO 10 IM=1,128
S1(IM)=S1(IM+11)
10 CONTINUE
XPRES=SAVE

```

```

CALL SLP(NPS,S1,M1,IP1,K1,K2,ALPHA,RC)
CALL AUTC(NPS,S1,M1,A,ALPHA,PC)
CALL FPTR(RC,14,4,10,ANS1,ANS2,ANS3)
PTR(1)=ANS1
PTR(2)=ANS2
PTR(3)=ANS3
CONTINUE
CALL SMOOTH(FTR,WORK,NF)
CALL QUANTZ(WORK,QMT,NF,XMAX,XMIN)
CALL PLOTFR(NF,QMT,ITITLE)
CALL SMOOTH(CTR,WORK,NF)
CALL QUANTZ(WORK,QMT,NF,XMAX,XMIN)
CALL PLOTFR(NF,QMT,JTITLE)
CALL SMOOTH(BTP,WORK,NF)
CALL QUANTZ(WORK,QMT,NF,XMAX,XMIN)
CALL PLOTFR(NF,QMT,KTITLE)
ARRAY(10IGT+1,1)=FLCAT(NF)
DO 70 I=1,NF
70 ARRAY(10IGT+1,I+1)=CMT(I)
IF(10IGT.LT.3) GO TO 1
WRITE(1,150) ARRAY
150 FORMAT(A0E12.4)
STOP
END

```

```

SUBROUTINE AUTC(N,X,M1,A,ALPHA,RC)
DIMENSION X(1),A(1),RC(1)
DIMENSION R(21),R(21)
MF=M1+1
DO 100 K=1,MF
F(K)=0.0
NK=N-K+1
DO 100 NF=1,NK
NPK=NP+K-1
100 F(K)=R(K)+X(NP)*X(NPK)
ALPHA=R(1)
FC(1)=-R(2)/R(1)
A(1)=1.0
A(2)=R(1)
ALPHA=ALPHA-RC(1)*FC(1)*ALPHA
MF=M1
DO 400 MINC=2,MF
M=MINC-1
DO 200 J=1,MINC
JB=MINC-J+1
200 R(J)=A(JB)
M=M+1
S=0.0
DO 300 IF=1,M
MIF=M-IF+2
300 S=S+F(MIF)*A(IF)
FC(M)=-S/ALPHA
DO 350 IP=2,M
350 A(IF)=A(IF)+RC(M)*E(IP-1)
A(M+1)=FC(M)
ALPHA=ALPHA-RC(M)*FC(M)*ALPHA
400 CONTINUE
FETURN
END

```



```

DIMENSION ARRAY(10,40),BRRAY(10,40),TABLE(10,10)
COMMON ARRAY,BRRAY
C READ(IFL,100) ARRAY
IFL=10
DO 5 IP=1,5
READ(IFL,100) ARRAY
100 FORMAT(40E12.4)
WRITE(6,160) IP
160 FORMAT(1H1,10X,'SET NUMBER : ',14)
WRITE(6,150)
DO 1 ITER=1,10
C TABLE(ITER,ITER)=1.0
NEXT=ITER
IF(NEXT.GT.10) GC TO 2
DO 2 JTER=NEXT,10
CALL CORR(ITER,JTER,COEFF)
TABLE(ITER,JTER)=COEFF
TABLE(JTER,ITER)=CCEFF
2 CONTINUE
1 CONTINUE
WRITE(6,200) TABLE
200 FORMAT(1H0,10F10.3)
150 FORMAT(40X,'RMS CORRELATION ',/,3X,
$ ' ZERO ONE TWO THREE FOUR FIVE
$ ' SIX SEVEN EIGHT NINE')
IFL=IFL+1
9 CONTINUE
STOP
END

SUBROUTINE CORR(ITER,JTER,COEFF)
DIMENSION ARRAY(10,40),BRRAY(10,40)
DIMENSION X(40),Y(40)
COMMON ARRAY,BRRAY
N1=INT(ARRAY(ITER,1))
N2=INT(ARRAY(JTER,1))
N=MAX0(N1,N2)
DO 10 I=1,N
X(I)=ARRAY(ITER,I+1)
Y(I)=ARRAY(JTER,I+1)
10 CONTINUE
XVAR=0.0
YVAR=0.0
XMEAN=0.0
YMEAN=0.0
XYVAR=0.0
DO 20 I=1,N
XMEAN=XMEAN+X(I)
YMEAN=YMEAN+Y(I)
20 CONTINUE
XMEAN=XMEAN/FLOAT(N)
YMEAN=YMEAN/FLOAT(N)
DO 30 I=1,N
RX=X(I)-XMEAN
RY=Y(I)-YMEAN
XVAR=XVAR+RX**2
YVAR=YVAR+RY**2
XYVAR=XYVAR+RX*RY
30 CONTINUE
VAR=XVAR*YVAR
CCEFF=XYVAR/SQRT(VAR)
RETLRN
END

```

VITA²

Ahmed Mahmoud Milyani

Candidate for the Degree of

Doctor of Philosophy

Thesis: A ROBUST PHONETIC DIGIT RECOGNITION SYSTEM FOR DIGITS SPOKEN
IN AMERICAN ENGLISH AND ARABIC

Major Field: Electrical Engineering

Biographical:

Personal Data: Born in Mecca, Saudi Arabia, January 28, 1940, the
son of Mr. Mahmoud M. Milyani and Mrs. Nakeia Y. Zawawi.

Education: Graduated from Ali Abdul-Latif High School, Cairo,
Egypt, in August, 1958; received the H.N.C. degree in Elec-
trical and Electronic Engineering from Cambridge College of
Arts and Technology, Cambridge, England, in August, 1969;
attended Part II courses for the Council of Engineering
Institution Fellowship at Cambridge College of Arts and
Technology, England, from January, 1970 to July, 1970.
received the Master of Science degree in Electrical Engineer-
ing from Oklahoma State University, Stillwater, Oklahoma,
in May, 1976; completed requirements for the Doctor of Phi-
losophy degree at Oklahoma State University, Stillwater,
Oklahoma, in July, 1981.

Professional Experience: Electronics and test engineer, PYE TVT
Ltd., Cambridge, England, September, 1967 to January, 1970;
Electronics and Communication Technician, RAYTHEON Service
Company, Saudi Arabia Branch, Jeddah, Saudi Arabia, September,
1970 to June, 1971; Instructor (full-time faculty member),
University of Petroleum and Minerals, Dhahran, Saudi Arabia,
July, 1971 to June, 1974.

Professional Organizations: Full Member of ETA KAPPA NU Associa-
tion, and The Institute of Electrical and Electronics Engi-
neers.