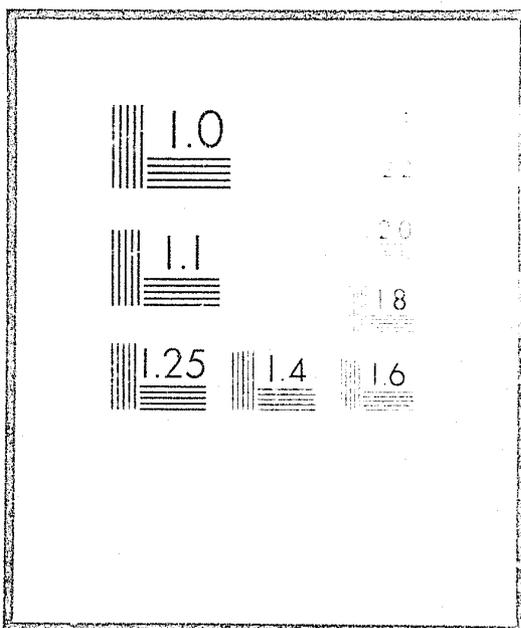


# NCJRS

This microfiche was produced from documents received for inclusion in the NCJRS data base. Since NCJRS cannot exercise control over the physical condition of the documents submitted, the individual frame quality will vary. The resolution chart on this frame may be used to evaluate the document quality.



Microfilming procedures used to create this fiche comply with the standards set forth in 41CFR 101-11.504

Points of view or opinions stated in this document are those of the author(s) and do not represent the official position or policies of the U.S. Department of Justice.

U.S. DEPARTMENT OF JUSTICE  
LAW ENFORCEMENT ASSISTANCE ADMINISTRATION  
NATIONAL CRIMINAL JUSTICE REFERENCE SERVICE  
WASHINGTON, D.C. 20531

Date filmed

7/16/76

# Voice Identification Research

NCJ-000 481



NATIONAL INSTITUTE OF LAW ENFORCEMENT  
AND CRIMINAL JUSTICE

U.S. DEPARTMENT OF JUSTICE  
LAW ENFORCEMENT ASSISTANCE ADMINISTRATION

Voice Identification  
Research

RP 72-1  
FEBRUARY 1972



Submitted by

The Department of  
Michigan State Police  
East Lansing, Michigan

This report was prepared by the Michigan State Police  
Department of Forensic Science, East Lansing, Michigan  
and is being published for the benefit of the law enforcement  
community. It is the property of the Michigan State Police  
Department and should not be distributed outside the  
Department.

U.S. Department of Justice  
Law Enforcement Assistance Administration  
National Institute of Law Enforcement and Criminal Justice

## PREFACE

The criminal or mentally disturbed person who uses his voice to prey on victims for extortion, bomb scares, nuisance, obscene or threatening phone calls is extremely difficult to identify. He is usually left to go his nefarious way until he commits an overt act that causes his arrest. The one clue he leaves is his voice. If it can be established that voice identification is possible and practical, it will be an extremely useful tool in the identification and prosecution of criminals.

Personal identification has many forms and degrees of quality. Because each Human is unique, he can be singled out. This is done routinely through photography, fingerprints and handwriting. The theory has been propounded and supported that a person can be identified through his voice. Many people can identify associates aurally, but this is not sufficiently reliable for forensic purposes. There is evidence that positive voice identification by other means is possible and sufficiently reliable for use in our courts. There have been outspoken dissenters to this hypothesis, particularly in acoustical societies.

Much of the dispute centers around claims made by Mr. Lawrence Kersta. Voice spectrography was developed at Bell Telephone Laboratories by Potter, Kopp and Green. An instrument called the Sound Spectrograph was developed, and was produced commercially as the Kay Sonograph. Mr. Kersta, a former member of the Bell Telephone research staff, dealt for many years with voice spectrograms. He made instrumental refinements on the sound spectrograph so that it would be more adaptable to personal identification.

In 1962, Kersta reported that voice spectrograms could provide a reliable means of identification. His methods were applied in several criminal investigations and subsequently voice identification testimony was presented in court. It became apparent from court decisions and resistance from the scientific field that there was a need for further study to replicate Kersta's work if voice identification was to be established as a scientific method.

Consequently, the Department of Michigan State Police developed a program with the following goals:

- (1) To establish Voice Identification as an aid to Law Enforcement.
- (2) To validate the Kersta method of Voice Identification.
- (3) To develop new methods of talker identification, through speech signals that might complement the voice spectrograph.
- (4) To evaluate the practical application of Voice Identification and prepare an operational manual for law enforcement.

The program was supported by a grant from the U.S. Department of Justice. The Department of Michigan State Police is indebted to the following organizations for their exceptional assistance in carrying out the proposed research.

The Sensory Sciences Research Center, Stanford Research Institute, Menlo Park, California, made an interpretive survey of the literature concerning methods for measuring speaker recognition. Initial research and experimentation in new areas of talker identification was proposed. Their survey relates the state of the art as it appeared at the outset of this project, and it provides the introductory information for this report.

The Audiology and Speech Sciences Department, Michigan State University, East Lansing, Michigan conducted research to identify speakers through voiceprints. (Voiceprints is a copyright term used by the firm established by Kersta, Voiceprint Laboratories, Inc., Somerville, N.J., to describe a particular graphic display made by an instrument that gives continuous display in time of the energy present in frequency bands.) This research is really the core of the first two years of the project.

The School of Criminal Justice, Michigan State University, East Lansing, Michigan, made a feasibility study of talker identification to determine its practical application to Law Enforcement.

The Department of Michigan State Police, East Lansing, Michigan, coordinated these efforts and

also perused the practical application of voiceprint identification to criminal investigation.

Department of Michigan State Police personnel for this project were:

- Captain Raymond H. McConnell, Project Director, September 1968–February 1969.
- Captain Wallace VanStratt, Ass't. Project Director, September 1968–February 1969; Proj-

- ect Director, February 1969–December 1970
- Lieut. Robert Earhart, Assistant Project Director, February 1969–December 1970.
- Det. Sgt. Ernest Nash, Voice Identification Technician.
- Detective Lewis Wilson, Voice Identification Technician.
- Elzora Conley, Secretary.

## TABLE OF CONTENTS

|  | Page |   | Page |
|--|------|---|------|
| Preface  | iii  | IV. Discussion and Conclusions  | 57   |
|  |      | V. Extension of Results from Forensic Models to Real Cases                          | 58   |
| <b>SUMMARY</b>   |      |   |      |
| I. Review of Procedures for Speaker Recognition                                      | 5    | <b>PART 3</b>   |      |
| II. Research of Speaker Identification by the Spectrographic Method                  | 9    | Some Guidelines for the Use of Voiceprint Identification Techniques                 | 61   |
| III. Practical Application of Voice Identification to Criminal Justice Investigation | 15   | I. Introduction   | 65   |
| IV. Training   | 16   | II. Procedure   | 65   |
| V. A Look to the Future  | 17   | III. Training   | 66   |
|  |      | IV. Summary of Experiment   | 67   |
|  |      | V. Interpretation of Results  | 68   |
|  |      | VI. Conclusions and Guidelines  | 68   |
| <b>PART 1</b>  |      |   |      |
| Summary Review of Procedures for Speaker Recognition                                 | 19   | <b>PART 4</b>   |      |
| I. Introduction  | 23   | The Practical Application of Voice Identification in Criminal Investigations        | 61   |
| II. Interspeaker and Intraspeaker Variability  | 23   | I. Introduction   | 75   |
| III. Speaker Recognition by Listening  | 24   | II. Methods   | 75   |
| IV. Speaker Recognition by Visual Comparison of Spectrograms                         | 25   | III. Results  | 76   |
| V. Speaker Recognition by Machine  | 29   | IV. Training of Voiceprint Examiners  | 78   |
| VI. Future Developments in Speaker Recognition                                       | 30   | V. A Look to the Future   | 79   |
| <b>PART 2</b>  |      |   |      |
| Michigan State University Voice Identification Project                               | 35   | <b>APPENDIXES</b>   |      |
| I. Introduction  | 39   | A. Master Tables of Results from Trials   | 81   |
| II. Results From the First Cycle of the Project                                      | 48   | B. An Examination of Conditional Variations for Voice Identification Trials         | 119  |
| III. Results From the Second Cycle of the Project                                    | 52   | C. An Examination of the Nature of Differences Among Voice Identification Examiners | 137  |
|  |      | An Examination of the Types of Errors Made by Examiners                             | 145  |

## SUMMARY

## CONTENTS

|  | Page |
|--|------|
| I. Review of Procedures for Speaker Recognition .....                              | 5    |
| A. Introduction .....  | 5    |
| B. Interspeaker and intraspeaker variability .....                                 | 5    |
| C. Speaker recognition by listening .....  | 6    |
| D. Speaker recognition by visual comparison of spectrograms .....                  | 6    |
| E. Speaker recognition by machine .....  | 8    |
| F. Future developments in speaker recognition .....                                | 8    |
| II. Research of Speaker Identification by the Spectrographic Method .....          | 9    |
| A. Introduction .....  | 9    |
| B. Experimental procedure .....  | 10   |
| C. Results of the first cycle of the project .....                                 | 11   |
| D. Results from the second cycle of the project .....                              | 13   |
| E. Discussion and conclusions .....  | 13   |
| F. Extension of results from forensic models to real cases .....                   | 14   |
| III. Practical Application of Voice Identification to Criminal Investigation ..... | 15   |
| A. Preparation .....   | 15   |
| B. Execution .....   | 16   |
| C. Results .....   | 16   |
| D. Conclusions .....   | 16   |
| IV. Training .....   | 16   |
| V. A Look to the Future .....  | 17   |
| References .....   | 18   |

## I. Review of Procedures for Speaker Recognition

### A. Introduction

When a person speaks he produces a complex acoustic signal that contains various kinds of information. This signal serves primarily to convey a linguistic message. Listeners who are familiar with the language can transcribe or at least repeat what the speaker said. Besides conveying a message the speech signal also reflects some of the anatomy and physiology of the speaker.

There are three general methods of speaker recognition. These are speaker recognition by listening, speaker recognition by visual comparison of spectrograms, and speaker recognition by machine. Speaker recognition by listening is, of course, the method used in everyday life. A possible limitation of this method is that it is entirely subjective. No matter how accurate and reliable listeners may be they are usually unable to describe the criteria upon which their decisions are based and thus they are unable to justify their conclusions in a court of law.

Speaker recognition by visual comparison of spectrograms is considered to be a more objective method. Spectrograms are visual displays of the speech signal. They exhibit graphic features that can be discussed in a fairly objective manner. These features are also interpreted subjectively in arriving at an overall decision. For this reason there has been much interest in a third method, namely, speaker recognition by machine. Although machine decisions are inherently objective, they are, as of now, often less accurate for speaker recognition purposes than comparable human decisions. Current research efforts in speaker recognition by machine are specifically directed toward overcoming this limitation.

All methods of speaker recognition are based on the fact that a given word or phrase tends to be uttered differently by different speakers. There is much variability in the speech signal and some of this variability is undoubtedly related to particular speaker differences. The nature of speaker variabil-

ity is discussed as background material to provide the reader with an understanding of principles of speaker recognition.

### B. Interspeaker and intraspeaker variability

It is well-known that the pronunciation of a given word or phrase tends to vary from speaker to speaker. Acoustical analyses of utterances of several speakers typically reveal many dissimilarities. This effect is referred to as interspeaker (between-speaker) variability. Interspeaker variability in the speech signal can be attributed in part to organic differences in the structure of the vocal mechanism and, in part, to learned differences in the use of the vocal mechanism during speech production. Organic differences may be related to regional, social and cultural factors.

Not so well-known is the fact that a particular speaker rarely utters a given word twice in exactly the same way, even when the utterances are produced in succession. This is referred to as intraspeaker (within-speaker) variability. In generating an utterance a speaker strives to produce appropriate respiratory, laryngeal, and articulatory activity that will lead to understandable speech. But many details of the resulting waveform will change from utterance to utterance depending upon rate of speaking, mood of the speaker, emphasis given to various words, and many other variables.

The success of any method of speaker recognition depends on the degree to which interspeaker variability is greater than intraspeaker variability. Both forms of speaker variability are extremely difficult to quantify, because speaker variability is a reflection of many differences in speech production. It cannot be meaningfully expressed in terms of a single measure. The measurement of speaker variability requires an understanding of how specific differences in speech production are manifested in the speech signal. But such an understanding is not yet available.

### C. Speaker recognition by listening

Several kinds of tests have been devised to study different aspects of speaker recognition by listening. All tests employ the same basic procedure. Speakers drawn from a prescribed population are recorded, while reading selected speech material. The recordings are edited and presented to listeners, and the listeners carry out a recognition task. Each step in this procedure introduces variables that can influence the resulting performance.

The objective of most studies on speaker recognition by listening is, of course, to appraise the likelihood that a listener's judgment might be in error.

McGehee (1937) studied the reliability with which listeners can recognize unfamiliar voices. Groups of listeners participated in two experimental sessions that were separated in time, from one day to five months. The results indicate that the reliability of recognition decreases rapidly as the time interval between sessions is extended beyond two weeks.

The effect of increasing the number of speakers heard during the first sessions was also investigated. When one of two speakers heard during the first session spoke again during a second session two days later, 77 percent of the listeners recognized his voice. When five speakers participated in the first sessions, only 46 percent of the listeners could recognize one of their voices two days later. Vocal disguise was also found to be effective in lowering recognition scores.

These results are illustrative of many of the results reported in the scientific literature. They illustrate the important fact that the speech waveform carries information relevant for distinguishing among talkers. However, the ability of listeners to identify speakers by their voice alone falls far short of 100 percent reliability. The quest for a more reliable means of identifying speakers on the basis of their voices has led to the study of speaker recognition by visual comparison of spectrograms and speaker recognition by machine. These two approaches will be briefly described in the following sections.

### D. Speaker recognition by visual comparison of spectrograms

This method of speaker recognition makes use of an instrument that converts the speech signal

into a visual display. The instrument is called a sound spectrograph and the display it provides is a spectrogram. Spectrograms of different utterances of a given word or phrase are presented to a trained observer who attempts to determine whether some utterances were produced by a common speaker. The sound spectrograph consists of four basic parts: (1) a magnetic recording device, (2) a variable electronic filter, (3) a paper-carrying drum that is coupled to the magnetic recording device, and (4) an electric stylus that marks the paper as the drum rotates. The magnetic recording device is used to record a short sample of speech. The duration of the speech sample corresponds to the time required for one revolution of the drum. Then the speech sample is played repeatedly in order to analyze its spectral contents. For each revolution of the drum, the variable electronic filter passes only a certain band of frequencies, and the energy in the frequency band activates the electric stylus so that a straight line of varying darkness is produced across the paper. The degree of darkness represents the varying amplitude of the speech signal at the specified time within the given frequency band. As the drum revolves, the variable electronic filter moves to higher and higher frequencies, and the electric stylus moves parallel to the axis of the drum. Thus a pattern of closely-spaced lines is generated on the paper. This pattern, which is the spectrogram, has the dimensions of frequency, time, and amplitude.

The spectrogram provides a permanent visual record of a speech signal. Such records may be studied in detail, point for point comparisons may be made among spectrograms, and judgments of similarity may be expressed in quantitative terms. Thus, the spectrogram has obvious appeal in legal applications. It is likely that the full potential of the spectrogram as a tool for achieving speaker recognition has not yet been reached.

The general procedure used in experiments employing the spectrogram as a means of speaker recognition is as follows: speakers are recorded while reading selected words or phrases. Spectrograms are prepared from the recordings. Two or more spectrograms of different utterances of the same words or phrases are presented to trained observers, and the observers carry out a recognition task. As is the case with speaker recognition by listening, each step in this procedure introduces variables that can affect performance, that is, the ability of the observer to

match correctly spectrograms that represent the same speaker.

The fallibility of the observer is a crucial issue in the legal use of this method of speaker recognition (Borders, 1966; Ladefoged and Vanderslice, 1967; McDade, 1968; Bolt et al., 1970). Although a machine (the sound spectrograph) is used to prepare spectrograms, the interpretation of spectrograms is an art rather than a science. In the first trial in which spectrograms were allowed as evidence, the jury could not reach an agreement as to how much weight this evidence should be given (McDade, 1968). The conviction of Edward Lee King was reversed by a Court of Appeals because "The Voiceprint identification process has not reached a sufficient level of scientific certainty to be accepted as identification evidence."

Claims by Kersta and others of the reliability of the Voiceprint for achieving speaker recognition are based largely on the results of unpublished experiments, thus the scientific community cannot appraise the design of these experiments and the validity of the conclusions reached (Ladefoged and Vanderslice, 1967). The published results of one series of experiments (Kersta, 1962b) could not be duplicated by other investigators. Young and Campbell (1967), and also Stevens, Williams, Carbonell, and Woods (1968), obtained much higher error scores than those reported by Kersta (1962a, 1962b). Such disagreements make the publication of detailed descriptions of future experiments extremely desirable and necessary.

In the first experiments concerned with the question of Voiceprint, the observers were required to sort spectrograms into groups that represented different speakers (Kersta, 1962a, 1962b). Later experiments employed the multiple-choice identification test (Kersta, 1962c; Young and Campbell, 1967; Stevens, Williams, Carbonell, and Woods, 1968). There have been no reports of experiments dealing directly with the type of identification task commonly encountered in criminal investigations. Ladefoged and Vanderslice (1967) argued that the reliability of Voiceprint identification in practical cases cannot be predicted from the results of the published studies.

It has been claimed that spectrogram recognition performance is essentially unaffected by the loss of teeth, tonsils, or adenoids, the aging process, and attempts to disguise the voice, such as changing the fundamental frequency, whispering, mimicking an-

other voice, and ventriloquism (Kersta, 1962c, Anon., 1965). However, in the absence of supporting experimental data, these claims cannot be considered established facts.

According to Kersta (1962b), the probability that two speakers have similar enough vocal-tract dimensions and articulation patterns to produce indistinguishable spectrograms is extremely small. This belief, which appears to underlie many experiments, has not been formally translated into a hypothesis that can be tested with a finite population of speakers. There is evidence that two arbitrarily selected speakers can occasionally produce very similar spectrograms (Ladefoged and Vanderslice, 1967).

Stevens, Williams, Carbonell, and Woods (1968) examined the ability of observers to distinguish between familiar and unfamiliar speakers in a 32-item identification-discrimination test. The observer was given eight reference spectrograms that represented eight "familiar" speakers. There were two experimental conditions; either four or 16 of the 32 test spectrograms represented "unfamiliar" speakers who were not represented by the reference spectrograms. Most of the familiar speakers were recognized as such, and they were subsequently correctly identified. Many of the unfamiliar speakers, however, were erroneously recognized as familiar speakers. As a point of comparison, listening tests were conducted using the same speakers and the same test format. Spectrograms were not employed in these tests. A comparison of the two sets of data reveals that there were considerably more acceptances of unfamiliar speakers in the visual tests than in the oral tests. When only four of the 32 test items represented unfamiliar speakers, there were also more false rejections of familiar speakers in the visual tests. Thus, speaker recognition by listening was found, in this study, to be the more accurate method. It must be pointed out that the observers employed by Stevens et al. had very little training. One would expect better performance from highly-trained observers, but this study does demonstrate that speaker recognition by spectrogram matching is neither obvious nor easily achieved.

The above discussion may be summarized as follows: In view of the use of the visual comparison of spectrograms for speaker identification as evidence in courts of law, the fallibility of the observer must be studied further (Bolt et al., 1970). Future

experiments should be carefully designed so as to avoid possible artifacts in the results. A detailed description of the experimental procedure, accompanied by the obtained data, should be published or otherwise be made available to the scientific community. Claims should be clearly differentiated from proven facts. The spectrographic method for speaker identification has obvious potential in various investigative and forensic applications.

### E. Speaker recognition by machine

Two approaches have been used to study the feasibility of speaker recognition by machine. One approach is to have the machine generate and examine amplitude-frequency-time matrices of specific speech samples. The other approach is to have the machine extract speaker-dependent parameters from the speech signal and subject them to a statistical analysis. Each approach has led to a number of recognition techniques.

In the first case, the utterances of specific speech samples are usually processed by a spectrum analyzer that consists of a bank of bandpass filters, rectifiers, and smoothing circuits. The outputs of the analyzer are periodically sampled, and the amplitudes are quantized for further processing by computer. Each utterance is represented in the computer by a data matrix. The rows of the matrix correspond to the frequency bands of the spectrum analyzer, the columns correspond to the temporal locations of the sample spectra, and each matrix cell contains the measured amplitude level. Such a matrix may be thought of as a "digital spectrogram." For each phrase, word, or phoneme used, several matrices representing different utterances by the same speaker are combined to form a single reference matrix for that speaker. A reference matrix is thus constructed for each speaker participating in a recognition experiment. The speaker to be recognized is represented by a test matrix. Depending on the type of recognition to be performed the test matrix is compared with all, or one of the reference matrices. The degree of similarity between the test matrix and each reference matrix is computed, and the results are used to arrive at a decision.

There are two basic recognition tasks, identification and discrimination. In the identification task, several reference matrices are used and it is assumed

that the speaker represented by the test matrix is also represented by one of the reference matrices; thus, the reference matrix that is most similar to the test matrix is expected to identify the speaker represented by the test matrix. In the discrimination task only one reference matrix is used and the speaker represented by the test matrix may or may not be represented by this reference matrix. Decision rules are selected to specify when the test and reference matrices are similar enough to represent the same speaker. A summary of six studies resulted in a range of correct scores from 89 to 95 percent.

Techniques using statistical analyses of speech parameters involve two distinct processes: (1) the extraction from the speech signal of parameters thought to be useful for differentiating among speakers, and (2) the application of decision rules to combinations of parameter values that represent particular speech samples.

Questions regarding the most appropriate speech parameters have generally not been resolved as well as have questions regarding optimal decision rules. Various kinds of parameters have been examined, using both waveform analyses and spectroanalyses of the speech signal. Studies conducted by Clarke and Becker (1969), Hargraves and Starkweather (1963), Smith (1962), Ramishvili (1966), Edie and Sebestyen (1962), and Floyd (1964) have considered many speech parameters and several decision techniques. In general, results have been promising but it is clear that much work remains to be done before automatic recognition techniques attain high reliability.

### F. Future developments in speaker recognition

The previous material describes in general terms the current status of speaker recognition by listeners, by visual examination of spectrograms, and by machine. Here we will comment briefly on the potential of each of these methods.

#### (1) Speaker recognition by listeners

There is little likelihood that much can be done, or should be done, to improve the average individual's ability to recognize speakers by voice. Identification based on the average individual's recognition of voice will undoubtedly remain un-

reliable although in some cases it may be admitted as evidence. Thus, it would appear that the potential of speaker recognition by listeners is quite limited.

#### (2) Speaker recognition by visual examination of spectrograms

It is unlikely that this method has achieved its full potential. There has been too little systematic study of spectrogram features to determine optimal procedures for discriminating among talkers. Whereas the speech spectrograph should prove to be an increasingly valuable tool for investigative purposes it is unlikely that it will ever, under all circumstances, permit positive identification by voice.

## II. Research of Speaker Identification by the Spectrographic Method

### A. Introduction

The method of speaker recognition researched by the Audiology and Speech Sciences Department at Michigan State University is based upon the visual examination and comparison of spectrograms. Speech spectrography was developed at Bell Research Laboratories by Potter et al. (1947). This type of spectrography is accomplished by the use of an instrument called a sound spectrograph, which transforms speech into a visual display, a spectrogram. The spectrogram portrays three main parameters of speech: time (horizontal axis), frequencies (vertical axis) and relative amplitude (degree of darkness of the different spectrographic regions). Each phoneme, word or phrase is correlated with a characteristic spectrographic pattern. The general aspect of patterns corresponding to different utterances of the same word are similar, in such a way that a person specially trained in "reading" spectrograms, could determine with more or less accuracy which words or phrases were portrayed by a particular pattern. However, "interspeaker" and "intraspeaker" variabilities are also portrayed by the spectrographic patterns. In fact, the spectrograms of different utterances of the same word or phrase by the same or by different speakers are never *exactly* alike.

Kersta (1962) claimed that spectrograms of several utterances of the same words by a given speaker

#### (3) Speaker recognition by machine

This method of speaker recognition may prove to be the most promising. Computers are now capable of performing fast and accurate analyses of speech waveforms. Various parameters may be abstracted from the speech waveform and analyzed to determine those features most useful for distinguishing among talkers. Freedom to choose these optimal parameters may enable machine performance to exceed that of listeners or of trained observers using spectrograms as these two latter methods suffer from strict and arbitrary limitations upon processing equipment. To achieve improved or perfect performance the relevant speech parameters must be properly identified and incorporated into the analysis and decision processes of the machine.

always contain more similar spectral features than those produced by different speakers. Kersta concluded, therefore, that speaker identification by visual examination of spectrograms, has to be reliable. According to Kersta, speaker recognition by visual inspection of spectrograms consists of subjectively matching similarities found in pairs of spectrograms from the same person, that are not found in pairs of spectrograms from different persons. The dissimilarities presented by the matched spectrograms are disregarded; they are assumed to be a result of intraspeaker variability. To back his claim, Kersta published the results of experiments he performed at Bell Research Laboratories. In these experiments he observed fewer than 1 percent of wrong identification.

Matching similarities through the combination eye-brain is essentially a subjective method, but the examiner can display objectively these similarities in a court of law to support his subjective conclusions. The possibility of objective displays for legal application has perhaps made this method of speaker recognition quite appealing. It should not be assumed that Kersta's method precludes listening to the known and unknown voices. On the contrary, the examiner who selects samples of the voices to be spectrally compared must listen to the samples in addition to examining the spectrograms visually. The spectrograph commercialized by Kersta under the trade name of "Voiceprint" (Presti, 1967)

has a special playback mechanism, that allows proper selection and continuous listening of samples prior to their being fed into the processing circuits. Since 1962 Kersta has been producing legal testimony on speaker identification by using his method, as well as offering training in this "art" to law enforcement officers.

Several speech scientists (Bolt et al., 1970) have expressed their concern for the legal application of "voiceprinting," prior to proving its accuracy and reliability through controlled experimentation.

Tosi (1967, 1968) evaluated the "voiceprinting" method by analyzing the data derived from 236 experimental trials of identification performed by nine of Kersta's trainees, as well as by participating himself in the training courses given by Kersta at the Voiceprint Laboratories. These trials of identification of one speaker among 50, using five clue words, yielded an error of 6.3 percent. In his report to the Michigan Department of State Police, Tosi could only conclude that the "Kersta method shows promise," suggesting the need for an independent study, one that would include variables not considered by Kersta in his experiments and that would further test "voiceprinting."

Such a study, "Michigan State University Voice Identification Project," was conducted at the Department of Audiology and Speech Sciences of Michigan State University from 1968 to 1970, under a contract with the Michigan Department of State Police, through a grant from the United States Department of Justice.

To prepare for the experimentation, 250 speakers were randomly drafted from a population of approximately 25,000 male students at Michigan State University. The speakers recorded 9 clue words spoken in isolation, in a fixed context and in a random context, and repeated six times in each session. Three different types of transmissions were used.

These recordings were processed through a "Voiceprint" spectrograph using an expanded scale of frequencies from 50Hz to 4,000Hz. Half the spectrograms from the first recording session were considered "known" spectrograms and half were considered as produced by "unknown" speakers. All the spectrograms yielded from the second recording session were assumed to correspond to "unknown" speakers.

Newspaper advertisements were used to announce the opening of examiner positions. Applicants were screened prior to their selection as examiners and

only those who performed successfully on the screening tests were considered for participation in the study. The 29 persons accepted received approximately one month of training prior to the starting of the experimental trials.

As part of the training, each examiner was instructed to consider the following objective points of spectrograms: (a) mean frequencies of vowel formants, (b) formants band-widths, (c) gaps and types of vertical striations, (d) slopes of formants, (e) durations, (f) characteristic patterns of fricatives and interformant energies. As the examiner progressed, he was given increasingly more difficult tasks to perform.

Spectrograms used during the training period were not used during the experiment. Listening to the known and unknown voices was excluded from this study.

The examiners were grouped into three panels according to sex and background. The first panel consisted of women ranging from 17 to 60 years of age, with various levels of education, from high school to four years of college. The second panel included male undergraduates from several departments of Michigan State University. The third panel was formed exclusively from the Criminal Justice Department of the University. Further, each panel was divided into three sub-panels, one of 3 examiners, one of 2 examiners, and one of a single individual. These nine sub-panels performed the same experimental tasks, yielding 9 answers for each different type of trial. Examiners were rotated within the three different sub-panels.

## B. Experimental procedure

The experiment was divided into two cycles. In the first cycle, the examiner had 9 clue words to examine and compare. In the second cycle, he had 6 words. There were 486 different types of tasks involving every possible combination of the variables tested. Each combination was reiterated four times by nine-sub-panels of examiners, using different spectrograms in each reiteration. Therefore, the total number of trials in this experiment was 34,992.

The tasks of the examiners in the *open trials* consisted of deciding whether the "matching" spectrograms were or were not produced by one of the "known" speakers, and if so, which "known" speaker produced them. Three kinds of errors were possible:

(1) Error A: A match did exist but the exam-

iner selected the wrong one. (false identification)

(2) Error B: A match did exist but the examiner failed to recognize it.

(3) Error C: A match did not exist although the examiner selected one (false identification)

In the *closed trials*, the examiners had to decide which "known" speaker produced the matching spectrograms. Since a match always existed, only one kind of error was possible. This error was labeled *Error D*.

Each sub-panel was forced to reach a common decision in each trial. Each member had to indicate his degree of confidence in the forced decision, based on the following scale: 1 = almost uncertain; 2 = fairly uncertain; 3 = fairly certain; 4 = almost certain.

The result of each trial was finally quantified with two expressions, conveying the right or wrong response and the averaged self confidence grade on such a response. All answer sheets were filed in a room protected with an electric alarm system connected with the University Police Headquarters.

Examiners usually followed the same procedure to complete each trial of speaker identification. The steps in this procedure were: (1) comparing the spectrograms of the unknown and known voices by a rather fast scan; (2) discarding those known voices spectrograms that appeared subjectively to the examiner as containing no significant similarities with the unknown voice spectrograms. Usually these steps reduced to a very few the known voices spectrograms to be further examined; (3) continuing the scanning by folding and superimposing each of the remaining known spectrograms on the matching spectrograms. This procedure provided the examiners with a better technique in searching for similarities and reduced even more the number of suspected known spectrograms; (4) If the previous steps did not produce a positive decision, the examiners counted the number of similarities they found between each of the suspected known spectrograms and the matching spectrograms. The known spectrogram which presented more points of similarities was supposed to be chosen as a correct response in the case of closed trials. For open trials the procedure was essentially the same, but complicated by the circumstance that the examiners had to decide between two possible alternatives: "there is not a match," or "there is a match, being speaker *n* the same as the unknown speaker." (5) Subpanel members arrived at a common decision for each trial through discussion. (6) After the decision was

reached, each subpanel member assessed this decision by registering his personal rate of confidence on the common decision. He used the grading scale described earlier and recorded his judgment on the subpanel answer sheet, which were given to the research assistant for tabulation on the master tables.

After completion of each cycle, the results from the master tables were coded in IBM cards and processed through the 3600 CDC computer to calculate error percentages and perform an analysis of variance to test significances and interactions of the different variables involved in the experiment.

## C. Results of the first cycle of the project

The first cycle, using 9 clue words, was completed in approximately 8 months. The results were processed through a CDC 3600 computer. Table 1 presents the pooled percentages of correct responses produced by the examiners under each of the main conditions tested in the project during this first cycle.

No significant statistical difference was detected between the one, two and three utterances of the same clue words or between the three different types of recording transmissions. (see table 1) Other levels of the variables tested showed statistically significant differences.

Another analysis of variance was performed to test the differences in the performance of panels and sub-panels of examiners. No significant difference was detected between panel types, but a significant difference was detected between sub-panel sizes: Sub-panels of three members performed slightly better than the other sub-panels. However, composition as well as size may have been a contributing factor. The staff observed that the best examiners often exerted positive influence on those less motivated.

The grand mean percentage of errors from the 17,496 trials of speaker identification performed during the first cycle of the project was 8.9 percent. This grand mean was composed of 4.3 percent errors of wrong matching or false identification (errors A+C+D), and 4.6 percent of failures to recognize a match when it actually existed (error B).

In order to construct models relevant to the forensic point of view, trials were grouped according to the following characteristics:

has a special playback mechanism, that allows proper selection and continuous listening of samples prior to their being fed into the processing circuits. Since 1962 Kersta has been producing legal testimony on speaker identification by using his method, as well as offering training in this "art" to law enforcement officers.

Several speech scientists (Bolt et al., 1970) have expressed their concern for the legal application of "voiceprinting," prior to proving its accuracy and reliability through controlled experimentation.

Tosi (1967, 1968) evaluated the "voiceprinting" method by analyzing the data derived from 236 experimental trials of identification performed by nine of Kersta's trainees, as well as by participating himself in the training courses given by Kersta at the Voiceprint Laboratories. These trials of identification of one speaker among 50, using five clue words, yielded an error of 6.3 percent. In his report to the Michigan Department of State Police, Tosi could only conclude that the "Kersta method shows promise," suggesting the need for an independent study, one that would include variables not considered by Kersta in his experiments and that would further test "voiceprinting."

Such a study, "Michigan State University Voice Identification Project," was conducted at the Department of Audiology and Speech Sciences of Michigan State University from 1968 to 1970, under a contract with the Michigan Department of State Police, through a grant from the United States Department of Justice.

To prepare for the experimentation, 250 speakers were randomly drafted from a population of approximately 25,000 male students at Michigan State University. The speakers recorded 9 clue words spoken in isolation, in a fixed context and in a random context, and repeated six times in each session. Three different types of transmissions were used.

These recordings were processed through a "Voiceprint" spectrograph using an expanded scale of frequencies from 50Hz to 4,000Hz. Half the spectrograms from the first recording session were considered "known" spectrograms and half were considered as produced by "unknown" speakers. All the spectrograms yielded from the second recording session were assumed to correspond to "unknown" speakers.

Newspaper advertisements were used to announce the opening of examiner positions. Applicants were screened prior to their selection as examiners and

only those who performed successfully on the screening tests were considered for participation in the study. The 29 persons accepted received approximately one month of training prior to the starting of the experimental trials.

As part of the training, each examiner was instructed to consider the following objective points of spectrograms: (a) mean frequencies of vowel formants, (b) formants band-widths, (c) gaps and types of vertical striations, (d) slopes of formants, (e) durations, (f) characteristic patterns of fricatives and interformant energies. As the examiner progressed, he was given increasingly more difficult tasks to perform.

Spectrograms used during the training period were not used during the experiment. Listening to the known and unknown voices was excluded from this study.

The examiners were grouped into three panels according to sex and background. The first panel consisted of women ranging from 17 to 60 years of age, with various levels of education, from high school to four years of college. The second panel included male undergraduates from several departments of Michigan State University. The third panel was formed exclusively from the Criminal Justice Department of the University. Further, each panel was divided into three sub-panels, one of 3 examiners, one of 2 examiners, and one of a single individual. These nine sub-panels performed the same experimental tasks, yielding 9 answers for each different type of trial. Examiners were rotated within the three different sub-panels.

## B. Experimental procedure

The experiment was divided into two cycles. In the first cycle, the examiner had 9 clue words to examine and compare. In the second cycle, he had 6 words. There were 486 different types of tasks involving every possible combination of the variables tested. Each combination was reiterated four times by nine-sub-panels of examiners, using different spectrograms in each reiteration. Therefore, the total number of trials in this experiment was 34,992.

The tasks of the examiners in the *open trials* consisted of deciding whether the "matching" spectrograms were or were not produced by one of the "known" speakers, and if so, which "known" speaker produced them. Three kinds of errors were possible:

(1) Error A: A match did exist but the exam-

iner selected the wrong one. (false identification)

(2) Error B: A match did exist but the examiner failed to recognize it.

(3) Error C: A match did not exist although the examiner selected one (false identification)

In the *closed trials*, the examiners had to decide which "known" speaker produced the matching spectrograms. Since a match always existed, only one kind of error was possible. This error was labeled *Error D*.

Each sub-panel was forced to reach a common decision in each trial. Each member had to indicate his degree of confidence in the forced decision, based on the following scale: 1 = almost uncertain; 2 = fairly uncertain; 3 = fairly certain; 4 = almost certain.

The result of each trial was finally quantified with two expressions, conveying the right or wrong response and the averaged self confidence grade on such a response. All answer sheets were filed in a room protected with an electric alarm system connected with the University Police Headquarters.

Examiners usually followed the same procedure to complete each trial of speaker identification. The steps in this procedure were: (1) comparing the spectrograms of the unknown and known voices by a rather fast scan; (2) discarding those known voices spectrograms that appeared subjectively to the examiner as containing no significant similarities with the unknown voice spectrograms. Usually these steps reduced to a very few the known voices spectrograms to be further examined; (3) continuing the scanning by folding and superimposing each of the remaining known spectrograms on the matching spectrograms. This procedure provided the examiners with a better technique in searching for similarities and reduced even more the number of suspected known spectrograms; (4) If the previous steps did not produce a positive decision, the examiners counted the number of similarities they found between each of the suspected known spectrograms and the matching spectrograms. The known spectrogram which presented more points of similarities was supposed to be chosen as a correct response in the case of closed trials. For open trials the procedure was essentially the same, but complicated by the circumstance that the examiners had to decide between two possible alternatives: "there is not a match," or "there is a match, being speaker *n* the same as the unknown speaker." (5) Subpanel members arrived at a common decision for each trial through discussion. (6) After the decision was

reached, each subpanel member assessed this decision by registering his personal rate of confidence on the common decision. He used the grading scale described earlier and recorded his judgment on the subpanel answer sheet, which were given to the research assistant for tabulation on the master tables.

After completion of each cycle, the results from the master tables were coded in IBM cards and processed through the 3600 CDC computer to calculate error percentages and perform an analysis of variance to test significances and interactions of the different variables involved in the experiment.

## C. Results of the first cycle of the project

The first cycle, using 9 clue words, was completed in approximately 8 months. The results were processed through a CDC 3600 computer. Table 1 presents the pooled percentages of correct responses produced by the examiners under each of the main conditions tested in the project during this first cycle.

No significant statistical difference was detected between the one, two and three utterances of the same clue words or between the three different types of recording transmissions. (see table 1) Other levels of the variables tested showed statistically significant differences.

Another analysis of variance was performed to test the differences in the performance of panels and sub-panels of examiners. No significant difference was detected between panel types, but a significant difference was detected between sub-panel sizes: Sub-panels of three members performed slightly better than the other sub-panels. However, composition as well as size may have been a contributing factor. The staff observed that the best examiners often exerted positive influence on those less motivated.

The grand mean percentage of errors from the 17,496 trials of speaker identification performed during the first cycle of the project was 8.9 percent. This grand mean was composed of 4.3 percent errors of wrong matching or false identification (errors A+C+D), and 4.6 percent of failures to recognize a match when it actually existed (error B).

In order to construct models relevant to the forensic point of view, trials were grouped according to the following characteristics:

TABLE 1.—First Cycle—Results of an Analysis of Variance of the Correct Responses Produced Under Each of the Main Conditions Tested

| Condition   | Pooled percentage of correct responses | Probability of the difference between levels, less than: |
|---|--|--|
| Number of utterances of the same clue word:             |  |  |
| 1 utterance .....                                       | 91.29                                  | n.s.   |
| 2 utterances .....                                      | 90.96                                  | n.s.   |
| 3 utterances .....                                      | 92.49                                  | n.s.   |
| Different types of recording transmission:              |  |  |
| (a) directly into tape recorder .....                   | 92.42                                  | n.s.   |
| (β) through a telephone line in quiet environment ..... | 91.31                                  | n.s.   |
| (γ) through a telephone line in noisy environment ..... | 91.02                                  | n.s.   |
| Context of the clue words spoken:                       |  |  |
| (I) in isolation .....                                  | 95.77                                  | 0.01   |
| (II) in a fixed context .....                           | 92.39                                  | 0.01   |
| (III) in a random context .....                         | 86.59                                  |  |
| Different number of "known" speakers:                   |  |  |
| 10 speakers .....                                       | 93.03                                  |  |
| 20 speakers .....                                       | 91.87                                  | n.s.   |
| 40 speakers .....                                       | 89.58                                  | 0.01   |
| Time-elapsd between recordings:                         |  |  |
| Contemporary matching spectrograms .....                | 95.21                                  | 0.01   |
| Non-contemporary matching spectrograms .....            | 87.95                                  |  |
| Awareness of examiners:                                 |  |  |
| Closed trials .....                                     | 94.48                                  | 0.01   |
| Open trials .....                                       | 90.14                                  |  |

(1) Awareness of the examiners (closed or open trials).

(2) Time-elapsd (contemporary or non-contemporary matching spectrograms).

(3) Context (clue words spoken in isolation, fixed context or random context).

Two groups of trials are especially pertinent to the forensic point of view: Open trials determined by the use of non-contemporary spectrograms and clue words spoken in a fixed or in a random context. In fact, all real cases of forensic speaker identification would include these particular variables. Total error percentages yielded by these two groups were 14.35 percent and 18.26 percent respectively. Approximately one-third of these errors were errors of false identification (errors A+C) and two-thirds

were failures to recognize a match when it actually existed (error B). The percentages of *false identifications* observed in these forensic models were 4.22% and 6.43% for clue words in fixed and in random context, respectively. Percentages of *failures to recognize a match* when it actually existed were 10.13 percent and 11.83 percent respectively.

Another group of trials—closed trials, contemporary matching spectrograms and clue words spoken in isolation—produced findings germane to the goals of the study. Since these were essentially the variables tested by Kersta in 1962 and examined again by Tosi in 1968, their importance to the present study can be seen. The error percentage for this group was 0.51 percent, the minimum/lowest error percentage observed in the project, as expected. However, this group of trials does not

fit any type of forensic model and has no direct application.

The upper limit of the range of errors was found in another group of trials characterized by: open-match types exclusively, non-contemporary matching spectrograms, and clue words spoken in a random context. This group, which also does not fit any type of forensic model, yielded 29.01 percent error. This extreme error percentage was composed of 5.35 percent of false identifications (error A) and 23.66 percent of failures to recognize an existing match (error B).

#### D. Results from the second cycle of the project

The second cycle was undertaken to determine the effect a reduction in the number of clue words would/or would not have on the accuracy of the examiners.

Comments concerning the performance of the examiners is relevant. Only the more motivated persons remained with the project long enough to complete the second cycle; some of those who quit considered the task extremely boring. Many of the less motivated examiners did not perform well. These examiners tended to take an excessive number of rest periods and showed little concern for reaching the best possible decision in each trial, behavior which was viewed as hampering performances.

In overall conditions no significant difference was found between the two cycles. The examiners were correct 91.58% during the nine clue word cycle vs 91.24% during the six clue word cycle.

The improvement observed in the second cycle for the particular group of open trials which used non-contemporary spectrograms of clue words spoken in a random context could be explained on the basis of the learning process the examiners experienced during the first cycle. They assessed the open trials with non-contemporary spectrograms of clue words spoken in a random context as the most difficult tasks that produced the largest percentage of errors. The staff was aware that during the second cycle most of the examiners considered this particular type of trial as a challenge, devoting more time and special attention in the search for the correct answers.

The fact that results of the second cycle did not differ substantially from those of the first cycle must

not be interpreted as meaning that decreasing the number of available clue words from nine to six is not generally significant. The learning process the examiners experienced during the previous eight months devoted to the completion of the first cycle possibly interacted with these results, thus compensating for the fewer number of clue words available.

#### E. Discussion and conclusions

The results from the "Michigan State University Voice Identification Project" suggest that experienced examiners can identify or eliminate one unknown speaker from among as many as 40 known speakers, with little difference in accuracy being evidence in the use of nine or six clue words. The expected percentage of errors made by examiners who are forced to reach a positive decision in every trial of speaker identification they perform, (using *exclusively* visual examination of spectrograms), varies according to the conditions involved in each type of trial.

Closed trials, involving contemporary spectrograms of clue words spoken in isolation, yielded fewer than 1 percent error of false identifications. Since these conditions were essentially the one employed by Kersta, it can be concluded that the present study has confirmed the figures reported by Kersta in 1962. In the 1968 Tosi's evaluation of "Voiceprinting," the error percentage reported for similar type of trials was approximately six percent. This discrepancy can be explained on the basis of individual differences among examiners. In fact, considering the performance of each examiner separately in that evaluation, the range of error percentage was 14 to 0 percent.

The second goal of the present study was to test forensic models that included the following variables:

(a) Random chance that the unknown speaker is or is not among the known ones ("open trials");

(b) non-contemporary spectrograms (spectrograms of the unknown speaker obtained at a different time from the spectrograms of the known speakers);

(c) same sentences uttered by known and unknown speakers ("fixed context" or different sen-

tences including the same clue words, "random context.")

The error observed was approximately 15 percent for fixed context, of which approximately five percent were errors of false identifications (errors A + C) and approximately 10 percent were failures of recognizing a match when it actually existed (error B). For models including "random context," the total error was approximately 18 percent. This percentage was composed of approximately six percent of errors of false identifications and approximately 12 percent of failures of recognizing a match when it existed.

These findings suggest that if an experienced examiner, using only Visual Inspection of Spectrograms for legal purposes of identification and excluding any kind of listening, is forced to reach a positive decision in each case (devoting approximately 15 minutes to complete the task), his expected error range would be 14-18 percent. The probability that his wrong decisions will eliminate a guilty person is 75 percent of the total expected error. The probability that when in error this examiner will accuse an innocent person is 25 percent of the total expected error.

Under the specified conditions, the expected range of false identification is 5-6% and the expected range of the elimination of a guilty person is 10-12%.

Analysis of the ratings in the scale of self confidence used by the examiners in this project showed that approximately 60% of their wrong decisions were graded as "uncertain". This finding suggests that the examiners errors could have been reduced to approximately 40% of the observed figure, were these examiners not forced to reach a positive decision for the trials in which they felt uncertain.

Clearly, the repeated errors apply to experimental trials in which the examiners used visual inspection of spectrograms exclusively, devoting an average of 15 minutes per trial in reaching a forced positive decision. It could be hypothesized that if in addition to visual comparisons of spectrograms the examiners would not have been forced to reach a decision when uncertain, and allowed to listen to the unknown and known voices, the errors might have been further reduced. The experiment performed by Stevens et al. (1968), as well as the opinion of some phoneticians and linguists who feel that speaker recognition by listening is more accurate than by visual comparison of spectrograms,

seem to confirm this hypothesis. A further study including forensic models, similar to the ones used in the present experiment might result in important additional information if trained examiners could both listen and make visual comparisons of spectrograms. Also, the present study should be complemented by the testing of disguised voices and non-contemporary spectrograms obtained from spans of time longer than one month.

#### F. Extension of results from forensic models to real cases

A group of speech scientists (Bolt et al., 1970) have expressed concern about the use of spectrographic evidence in court, before this method has been validated by controlled experimentation. The question arises: assuming that the results from the statistical forensic models studied in the present experiment could be applied toward such a validation, how would the conditions in practical legal cases differ from the conditions in the statistical models? In what way would these real conditions possibly alter the error expectancy disclosed by the models?

Main differences of conditions that could exist between models and real cases are as follows:

1. *Population of known voices.* In the models of the present study the number of known voices varied from 10 to 40, drafted from a closed catalog of 250 speakers, representing a statistical sample of a homogeneous population of 25,000 persons. In forensic cases, the catalog of known voices could theoretically include millions of samples, if the voice spectrogram of the criminal would be compared with filed voice spectrograms of the population of the world, or even the United States of America. Obviously, conclusions derived from an experimental study of a small population of speakers can not be extrapolated to populations of millions of individuals. However, this is not the case in the present practical situations that police must handle. In these cases the catalog of known voices is *open*, true, but *limited* to a few suspected persons. It seems reasonable to assume that the intra and interspeaker variabilities within such a reduced group of suspected persons would not differ substantially from the variabilities that existed within the highly homogeneous group of experimental speakers utilized in the present study. There-

fore, it seems advisable to disregard size of the population of known voices as a differential characteristic that could hamper extrapolation of experimental results from the present study.

2. *Availability of time and responsibility of the examiners.* In the present study the examiners devoted an average of 15 minutes to reach a positive conclusion in each trial. Whether such a conclusion was the right or the wrong one, no effect could take place whatsoever over the examiner or the speaker. In forensic cases, the professional examiner normally may devote all the time necessary to reach a conclusion. He is aware of the consequences that a wrong decision could mean to his professional status as well as the consequences to the speaker whom he might erroneously identify. It seems reasonable to conclude that the differential characteristics between experimental and professional examiners might help to improve the accuracy of the professional examiners.

3. *Type of decisions examiners are urged to reach in each trial.* In the statistical model the examiners were forced to reach a positive conclusion in each trial, even if they were uncertain of the correct response. In real forensic cases, the professional examiner is permitted to make the following alternative decisions:\*

- a. Positive identification.
- b. Positive elimination.
- c. Possibility that the unknown speaker is one of the suspected persons, but more evidence is necessary in order to reach a definite conclusion.
- d. Possibility that the unknown speaker is none

### III. Practical Application of Voice Identification to Criminal Investigation

#### A. Preparation

To determine the usefulness of Voice Identification to Criminal Investigation, certain preparations were necessary. Two men with exemplary records as State Police Officers and Latent Identification Technicians were given Voice Identification train-

\*These are the alternative decisions that Sgt. Nash, head of the Voice Identification Unit of the Michigan Department of State Police is presently making.

of the available suspected persons but more evidence is necessary to reach a definite conclusion.

- c. Unable to reach any conclusion with the available voice samples.

4. *Availability of clues.* In the statistical models of this study, only spectrograms of nine and six clue words were available to the examiners for visual inspection. In real forensic cases the examiner must necessarily listen first to the unknown and known voices while processing the spectrograms for visual comparison. The professional examiner is entitled to request as many samples as he deems necessary to reach a positive conclusion. Combination of methods of voice recognition by listening and by visual inspection of spectrograms can enhance the accuracy of his conclusion. Listening also insures the comparison of the same sounds. Moreover, by using this combination the professional examiner can objectively sustain in court his opinion, by presenting the spectrographic similarities.

5. *Research panelist vs forensic examiner.* A common problem of both the Kersta and Tosi experiments was that of maintaining enthusiasm on the part of the panelist. Incentive was high in the beginning. However, after convincing themselves that reliable determinations could be made, some panelists lost interest. The challenge was gone. There was no punishment or consequence connected with mistakes.

The forensic examiner, on the other hand, is aware of his responsibility to be unbiased, the consequences of a mistake to both individuals and his own profession. The very nature of the work provides incentives to excel.

ing with Mr. Lawrence Kersta at Voiceprint Laboratories in New Jersey. Equipment was purchased. A laboratory space was provided. Experience was gained through continued work and communication with Kersta and Tosi. Police investigators were informed as to how Voice Identification might be useful, what evidence was necessary to conduct an examination, and how best to obtain known and unknown tape recordings. Over 4600 officers received such training with many more reached through radio and television.

## B. Execution

Since the inception of a voice identification service, in 1967, 291 cases have been submitted to the Voice Identification Unit, mostly from Michigan Police and Fire Departments. However, requests from all departments were honored and assistance was rendered to such places as Indianapolis, Indiana; Riverside, California; Orlando, Florida; Los Angeles, California; St. Paul, Minnesota; Ladue, Missouri; Erie, Pennsylvania; Chicago, Illinois; Dade County, Florida; Astoria, Oregon, and South Miami, Florida.

These cases involved 27 different types of crimes, ranging from nuisance telephone calls to murder.

## C. Results

673 voices were examined by the study of 42,432 spectrograms. 105 persons were identified as the unknown or questioned voice on tape recordings. 172 persons were eliminated as the unknown or questioned voice on tape recordings. For various

reasons, a definite opinion could not be rendered concerning the other 396 persons.

It was not always possible to obtain information from the investigating officers that would refute or substantiate the opinions of the voice identification examiners. However, it was reported that in thirty cases, those persons identified by voice identification techniques later made confessions or admissions correlating voice identification opinions. No information was found to prove the wrong person had been identified by voice identification techniques.

## D. Conclusions

Voice Identification by spectrographic analysis has a definite usefulness in the investigation of crime.

Given a sufficient quantity and quality of known and unknown voice recordings to work with, a qualified voice identification examiner, can arrive at opinions that have an accuracy level comparable to other types of subjective examinations now made in Forensic Laboratories.

## IV. Training

The application of voice identification techniques in actual cases pre-supposes the use of examiners who are educated, well trained and experienced.

It is important that the education include an understanding of the speech and hearing process. Although it does not bear directly on the visual comparison of spectrograms, it does provide the examiner with a better understanding of differences that occur within separate utterances of the same word by one speaker. This will help him understand and explain when slight differences exist.

Listening to the recordings is also an important part of the identification process because the examiner must be assured that he is comparing the same sounds. In a training or research project where the examiner is presented with two prepared spectrograms to compare, both could be labeled "the". However, if in actuality one spectrogram was made from the sound "thee" and the other spectrogram was made from the sound "thuh", identification would be impossible. Knowledge of the various

sounds of the spoken word, what sounds are germane to the identification process and how to listen to these sounds is necessary in the proper labeling of the spectrograms to be compared.

Basic training in theory of voice identification, the production of spectrograms and the comparison process are necessary in the early development of the voice identification examiner. However, this formal schooling does not sufficiently prepare an individual to undertake the responsibility of examining voice identification evidence and to give opinions in forensic cases. As in other forensic sciences that are subjective in nature, there must be experience and testing in the comparison of spectrograms until the examiner can demonstrate his ability to unerringly resolve the problems submitted to him. This does not preclude the fact that in some cases he may not be able to arrive at a definite opinion.

It has been demonstrated in the research by the Audiology and Speech Sciences Department of

Michigan State University that voice identification by the visual comparison of spectrograms is possible. The successful use of this method in forensic cases and in court, therefore, will ultimately depend on the reliability of the trained and experienced examiner.

The proper training of examiners is of utmost importance to the successful use of the voiceprint identification technique. To this end, the Michigan State University School of Criminal Justice makes the following recommendations with regard to training and educational requisites:

1. Ideally, the voiceprint identification expert should hold a baccalaureate degree in either speech science or physical science.
2. While it has been demonstrated that acceptable second generation trainees can be recruited from a general population, law enforcement technicians with comparative identification expertise are the recommended trainees.

## V. A Look to the Future

There are other research projects that should be initiated to extend the effectiveness of the voice spectrograph in criminal investigation. This would include experimentation with the identification of disguised voices and non-contemporary recordings. However, this should not deter its use by forensic laboratories or interfere with efforts to present voice identification testimony in court. In this respect, voice identification is no different than other forensic sciences in that there are always new questions to be answered.

Research is planned for speaker recognition by machine. This method could very well become an effective process to substantiate, extend or replace opinions now rendered by voice identification examiners using spectrographic techniques.

The possibility of using the spectrograph to identify sounds other than the human voice should not be overlooked. As an example, let us imagine that

3. In the absence of a baccalaureate degree as suggested above, the following college level courses are strongly urged as a prerequisite to eventual use of the voiceprint identification technique: Phonetics, acoustics (with the accompanying basic physics), speech science, linguistics, audiology and basic electronics.

4. Thorough training in the preparation of tape recordings and voice spectrograms.

5. A carefully supervised training program in voice spectrogram identification until the trainee reaches a 99% level of accuracy in closed trials working with spectrograms made from a homogeneous population.

6. Upon satisfactory completion of a training program similar to what has been outlined above, the trainee should then undergo apprenticeship instruction with an experienced supervisor.

an anonymous bomb threat is received and recorded. The sound of a motor can be distinguished as part of the background noise. If the motor noise, through sound spectrography, can be identified as to type, it might help investigators locate the source of the call. Again let us imagine that a woman calls the police and says she is about to be shot. An explosive sound ends the conversation. The sound spectrograph in this case may be effective in identifying the explosive noise as a firearm, perhaps a rifle rather than a pistol, and of large caliber.

As time passes, investigators and examiners alike will discover new applications of the sound spectrograph as it relates to criminal investigations. It remains now for more agencies and individuals to become involved in developing expertise and gaining experience in order that this relatively new technique can reach its full potential for solving crime.

## REFERENCES

- Anonym, "Voice Print Identification." *Criminalistics* (W. W. Turner, ed.), Aqueduct Books, Rochester (1965).
- Bolt, R. H., F. S. Cooper, E. E. David, Jr., P. B. Denes, J. M. Pickett, and K. N. Stevens, "Speaker Identification by Speech Spectrograms: A Scientists' View of its Reliability for Legal Purposes." *J. Acoust. Soc. Am.* 47, 597-612 (1970).
- Broders, W., "Voiceprint Allowed as Evidence; Ruling Called First of Its Kind." *The New York Times*, Tues., April 12 (1966).
- Clarke, F. R., and R. W. Becker, "Comparison of Techniques for Discriminating Among Talkers." *J. Speech and Hearing Res.* 12, 747-761 (1969).
- Edie, J., and G. S. Sebestyen, "Voice Identification General Criteria." Rept. RADC-TDR-62-278, Rome Air Development Center, Air Force Systems Command, Griffiss AFB (May 1962).
- Floyd, W., "Voice Identification Techniques." Rept. RADC-TDR-64-312, Rome Air Development Center, Research and Technology Division, Air Force Systems Command, Griffiss AFB (Sept. 1964).
- Hargreaves, W. A., and J. A. Starkweather, "Recognition of Speaker Identity." *Language and Speech* 6, 63-67 (1963).
- Kersta, L. G., "Voiceprint Identification." *J. Acoust. Soc. Am.* 34, 725 (A) (1962a).
- Kersta, L. G., "Voiceprint Identification." *Nature* 196, No. 4861, 1253-1257 (1962b).
- Kersta, L. G., "Voiceprint-Identification Infallibility." *J. Acoust. Soc. Am.* 34, 1978 (A) (1962c).
- Ladefoged, P., and R. Vanderslice, "The Voiceprint Mystique." Working Papers in Phonetics 7, University of California, Los Angeles (Nov. 1967).
- McDade, T., "The Voiceprint." *The Criminologist*, No. 7, 52-70 (Feb. 1968).
- McGehee, F., "The Reliability of the Identification of the Human Voice." *J. Gen. Psychol.* 17, 249-271 (1937).
- Presti, A., "High Speed Sound Spectrograph." *Journal of the Acoustical Society of America*, 40, 628-634 (1966).
- Ramishvili, G. S., "Automatic Voice Recognition." *Engineering Cybernetics*, No. 5, 84-90 (Sept.-Oct. 1966).
- Smith, J. E. K., "Decision-Theoretic Speaker Recognizer." *J. Acoust. Soc. Am.* 34, 1988 (A) (1962).
- Stevens, K. N., C. E. Williams, J. R. Carbonell, and B. Woods, "Speaker Authentication and Identification: A Comparison of Spectrographic and Auditory Presentations of Speech Material." *J. Acoust. Soc. Am.* 44, 1596-1607 (1968).
- Tosi, O., Personal communication (1970).
- Young, M. A., and R. A. Campbell, "Effects of Context on Talker Identification." *J. Acoust. Soc. Am.* 42, 1250-1254 (1967).

## PART 1

### Summary Review of Procedures for Speaker Recognition

By

Karl D. Kryter, Ph.D.  
Richard W. Becker, M.S.  
Frank R. Clarke, Ph.D.

Michael H.L. Hecker, B.M.  
Stephen E. Stuntz, M.A.  
Fausto Poza, M.S.

James R. Young, Ph.D.

SENSORY SCIENCES RESEARCH CENTER  
STANFORD RESEARCH INSTITUTE  
Menlo Park, California 94025

This study was performed during the first year of the grant. The first portion of the text is a summary review of a monograph written by Michael H.L. Hecker.

## CONTENTS

|  | Page |
|--|------|
| I. Introduction  | 23   |
| II. Interspeaker and Intraspeaker Variability                | 23   |
| III. Speaker Recognition by Listening                        | 24   |
| IV. Speaker Recognition by Visual Comparison of Spectrograms | 25   |
| V. Speaker Recognition by Machine                            | 29   |
| VI. Future Developments in Speaker Recognition               | 30   |
| A. Speaker recognition by listeners                          | 30   |
| B. Speaker recognition by visual examination of spectrograms | 32   |
| C. Speaker recognition by machine                            | 32   |
| References   | 33   |

## I. Introduction

When a person speaks he produces a complex acoustic signal that contains various kinds of information. This signal serves primarily to convey a linguistic message. Listeners who are familiar with the language can transcribe or at least repeat what the speaker said. Besides conveying a message the speech signal also reflects some of the anatomy and physiology of the speaker. For example, listeners can often determine the speaker's sex, his approximate age, his emotional state, and whether or not he is suffering from an illness (such as the common cold). Of particular interest is the ability of listeners to distinguish among the speech characteristics of different speakers. This ability is the basis of one method of speaker recognition.

There are three general methods of speaker recognition. These are speaker recognition by listening, speaker recognition by comparison of spectrograms, and speaker recognition by machine. Each of these methods is described in greater detail in separate sections of this report. Speaker recognition by listening is, of course, the method used in everyday life. It has been studied for a longer period of time and appears to be more accurate and reliable than either of the other methods as they are now practiced. A possible limitation of this method is that it is entirely subjective. No matter how accurate and reliable listeners may be they are

usually unable to describe the criteria upon which their decisions are based and thus they are unable to justify their conclusions in a court of law.

Speaker recognition by visual comparison of spectrograms is considered to be a more objective method. Spectrograms are visual displays of the speech signal. They exhibit graphic features that can be discussed in a fairly objective manner. But these features are still interpreted subjectively in arriving at an overall decision. For this reason there has been much interest in a third method, namely, speaker recognition by machine. Although machine decisions are inherently objective, they are, as of now, often less accurate for speaker recognition purposes than comparable human decisions. Current research efforts in speaker recognition by machine are specifically directed toward overcoming this limitation.

All methods of speaker recognition are based on the fact that a given word or phrase tends to be uttered differently by different speakers. There is much variability in the speech signal and some of this variability is undoubtedly related to particular speaker differences. The nature of speaker variability is discussed as background material to provide the reader with an understanding of principles of speaker recognition.

## II. Interspeaker and Intraspeaker Variability

It is well-known that the pronunciation of a given word or phrase tends to vary from speaker to speaker. Acoustical analyses of utterances of several speakers typically reveal many dissimilarities. This effect is referred to as interspeaker (between-speaker) variability. Interspeaker variability in the speech signal can be attributed in part to organic differences in the structure of the vocal mechanism and, in part, to learned differences in the use of

the vocal mechanism during speech production. Organic differences may be determined by heredity, sex, and age. Learned differences may be related to regional, social, and cultural factors.

Not so well-known is the fact that a particular speaker rarely utters a given word twice in exactly the same way, even when the utterances are produced in succession. This is referred to as intraspeaker (within-speaker) variability. In generating

an utterance a speaker strives to produce appropriate respiratory, laryngeal, and articulatory activity that will lead to understandable speech. But many details of the resulting waveform will change from utterance to utterance depending upon rate of speaking, mood of the speaker, emphasis given to various words, and many other variables.

The success of any method of speaker recognition depends on the degree to which interspeaker

### III. Speaker Recognition by Listening

Several kinds of tests have been devised to study different aspects of speaker recognition by listening. All tests employ the same basic procedure. Speakers drawn from a prescribed population are recorded, while reading selected speech material. The recordings are edited and presented to listeners, and the listeners carry out a recognition task. Each step in this procedure introduces variables that can influence the resulting performance. These variables include the size and homogeneity of the speaker group, the selection of speech materials, the size and training of the listener group, the mode of presentation of speech material, and the specific task assigned to the listeners. Each of these classes of variables is discussed in some detail by Hecker (1970).

The objective of most studies on speaker recognition by listening is, of course, to appraise the likelihood that a listener's judgment might be in error. In fact one of the first studies of this kind was motivated by a legal question of fallibility that arose in the Lindbergh case of 1935 (McGehee, 1937). Lindbergh claimed that he recognized the voice of the defendant as the voice of his son's kidnapper, heard almost three years earlier. Although Lindbergh's testimony was accepted by the court, the defense argued that such recognition was not entitled to much weight as evidence.

McGehee studied the reliability with which listeners can recognize unfamiliar voices. Groups of listeners participated in two experimental sessions that were separated in time, from one day to five months. During the first session they heard an unfamiliar speaker read a paragraph of text. During the second session they heard the same paragraph read successively by five speakers, including the

variability is greater than intraspeaker variability. Both forms of speaker variability are extremely difficult to quantify, because speaker variability is a reflection of many differences in speech production. It cannot be meaningfully expressed in terms of a single measure. The measurement of speaker variability requires an understanding of how specific differences in speech production are manifested in the speech signal. But such an understanding is not yet available.

speaker from the first session. The ability of the listeners to recognize the speaker whom they heard in the first session was investigated as a function of the time interval between the two sessions. The results, which are shown in Table I indicate that the reliability of recognition decreases rapidly as the time interval between sessions is extended beyond two weeks.

TABLE I.—Percent Correct Recognition of Unfamiliar Male Speakers After Various Intervals of Time (After McGehee, 1937.)

| Days |     |     | Weeks |     |     | Months |     |     |
|------|-----|-----|-------|-----|-----|--------|-----|-----|
| 1    | 2   | 3   | 1     | 2   | 3   | 1      | 2   | 3   |
| 83%  | 83% | 81% | 81%   | 69% | 51% | 57%    | 35% | 13% |

The effect of increasing the number of speakers heard during the first sessions was also investigated. When one of two speakers heard during the first session spoke again during a second session two days later, 77 percent of the listeners recognized his voice. When five speakers participated in the first sessions, only 46 percent of the listeners could recognize one of their voices two days later. Vocal disguise was also found to be effective in lowering recognition scores. In this experiment only one speaker was heard during the first session. He disguised his voice by changing its fundamental frequency. During the second session he used his normal voice. With a time interval of only one day, correct recognition was reduced by 13 percentage points.

These results are illustrative of many of the results reported in the scientific literature. They

illustrate the important fact that the speech waveform carries information relevant for distinguishing among talkers. However, the ability of listeners to identify speakers by their voice alone falls far short of 100 percent reliability. The quest for a more reliable means of identifying speakers on the

### IV. Speaker Recognition by Visual Comparison of Spectrograms

This method of speaker recognition makes use of an instrument that converts the speech signal into a visual display. The instrument is called a sound spectrograph, and the display it provides is a sound spectrogram (or Voiceprint, a trade name owned by Voiceprint Laboratories, Somerville, New Jersey). Spectrograms of different utterances of a given word or phrase are presented to a trained observer who attempts to determine whether some utterances were produced by a common speaker. Because the method has obvious applications in criminology, many studies have been concerned with its reliability as a means of positive identification. The sound spectrograph consists of four basic parts: (1) a magnetic recording device, (2) a variable electronic filter, (3) a paper-carrying drum that is coupled to the magnetic recording device, and (4) an electric stylus that marks the paper as the drum rotates. The magnetic recording device is used to record a short sample of speech. The duration of the speech sample corresponds to the time required for one revolution of the drum. Then the speech sample is played repeatedly in order to analyze its spectral contents. For each revolution of the drum, the variable electronic filter passes only a certain band of frequencies, and the energy in the frequency band activates the electric stylus so that a straight line of varying darkness is produced across the paper. The darkness of the line at any point on the paper indicates how much energy is present in the speech signal at the specified time within the given frequency band. As the drum revolves, the passband of the variable electronic filter moves to higher and higher frequencies, and the electric stylus moves parallel to the axis of the drum. Thus a pattern of closely-spaced lines is generated on the paper. This pattern, which is the spectrogram, has the dimension of frequency, time, and amplitude.

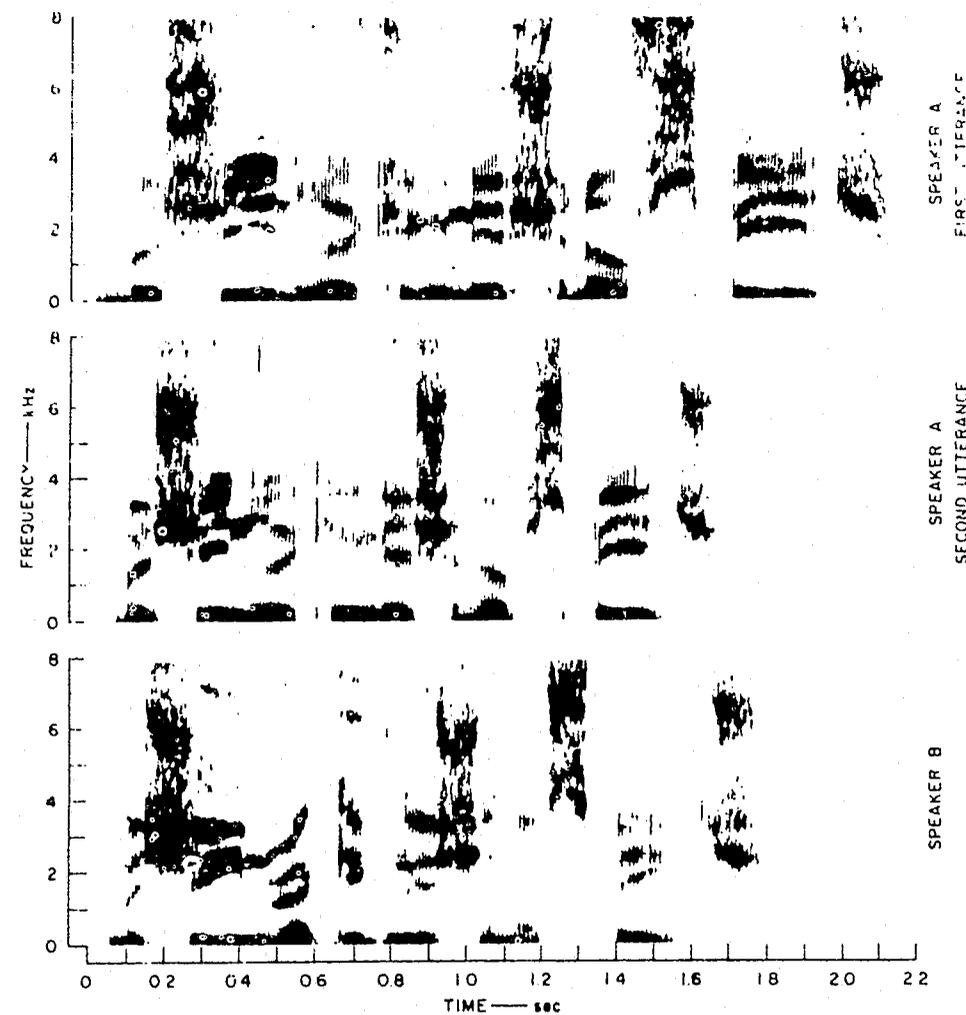
basis of their voices has led to the study of speaker recognition by visual comparison of spectrograms and speaker recognition by machine. These two approaches will be briefly described in the following sections.

Figure 1 shows three spectrograms. Since the spectrograms portray different utterances of the same phrase, each spectral feature of one utterance has a grossly similar counterpart in another utterance. The variability in corresponding spectral features appears to be somewhat greater between the two speakers (interspeaker variability) than between the two utterances by the same speaker (intraspeaker variability).

The spectrogram provides a permanent visual record of a speech signal. Such records may be studied in detail, point for point comparisons may be made among spectrograms, and judgments of similarity may be expressed in quantitative terms. Thus, the spectrogram has obvious appeal in legal applications. It is likely that the full potential of the spectrogram as a tool for achieving speaker recognition has not yet been reached.

However, the sound spectrogram has inherent limitations for speaker recognition applications. The display was designed to show differences among words and phonemes. It was not a purpose of the design to reveal differences between talkers. Thus, no attempt was made to have the device extract parameters from the speech waveform that might optimize speaker recognition performance. Further, a basic characteristic of all spectrum analyzers is that their frequency resolution can be increased only at the expense of temporal resolution and vice versa. The capability of a particular instrument to resolve frequency differences and temporal events is determined primarily by the bandwidth of its analyzing bandpass filter. Although the sound spectrograph contains two bandpass filters with different bandwidths (45 Hz and 300 Hz), the choice of either filter represents a compromise. Those features that might eventually prove to be the most useful ones for differentiating among speakers are not necessarily revealed in either the narrow-band or the wideband spectrogram.

Figure 1.—Sound spectrograms of three utterances of the Phrase "Machine Recognition of Speech" (after Young and Hecker, 1968.)



Because of the limited resolving power of the sound spectrograph, it is possible that spectrograms prepared from slightly different utterances of the same word cannot be differentiated by human observers. While the differences among the utterances would be evident in oscillographic recordings (which describe the utterances most completely) these differences may be obscured in the sound spectrogram. Therefore, when two spectrograms appear to be identical in all respects, it cannot be concluded that they necessarily represent the same speech signal. This limitation can be particularly severe in cases where the speech signals under analysis are distorted, or embedded in noise.

The general procedure used in experiments employing the spectrogram as a means of speaker recognition

is as follows: speakers are recorded while reading selected words or phrases. Spectrograms are prepared from the recordings. Two or more spectrograms of different utterances of the same words or phrases are presented to trained observers, and the observers carry out a recognition task. As is the case with speaker recognition by listening, each step in this procedure introduces variables that can affect performance, that is, the ability of the observer to match correctly spectrograms that represent the same speaker. The most important variables are described in detail by Hecker (1970), and will not be discussed in this report.

The fallibility of the observer is a crucial issue in the legal use of this method of speaker recognition (Borders, 1966; Ladefoged and Vanderslice,

1967; McDade, 1968; Bolt et al, 1970). Although a machine (the sound spectrograph) is used to prepare spectrograms, the interpretation of spectrograms is an art rather than a science. When this fact is pointed out to members of a jury they may be unable to evaluate the reliability of this means of identification. In the first trial in which spectrograms were allowed as evidence, the jury could not reach an agreement as to how much weight this evidence should be given (McDade, 1968). The conviction of Edward Lee King was reversed by a Court of Appeals because "The Voiceprint identification process has not reached a sufficient level of scientific certainty to be accepted as identification evidence in cases where the life or liberty of a defendant may be at stake." (Kennedy, 1968)

The use of the term Voiceprint, and the degree to which the analogy between Voiceprints and fingerprints has been emphasized (Kersta, 1962a, 1962b; Anon., 1965; McDade, 1968) are rather unfortunate. There is an important difference between spectrograms and fingerprints that is too seldom considered. The intraspeaker variability of the speech signal can be substantial. And this variability is, of course, demonstrated in spectrograms that represent a particular speaker. The variability exhibited by the whorls and ridges on a particular person's fingers is essentially zero (Ladefoged and Vanderslice, 1967; Bolt et al, 1970). Any difficulty in matching fingerprints is caused by the fact that fingerprints may be incomplete or smeared. As a means of identification, fingerprints must be regarded as being considerably more foolproof than the spectrograms.

Claims by Kersta and others of the reliability of the Voiceprint for achieving speaker recognition are based largely on the results of unpublished experiments, thus the scientific community cannot appraise the design of these experiments and the validity of the conclusions reached (Ladefoged and Vanderslice, 1967). The published results of one series of experiments (Kersta, 1962b) could not be duplicated by other investigators. Young and Campbell (1967), and also Stevens, Williams, Carbonell, and Woods (1968), obtained much higher error scores than those reported by Kersta (1962a, 1962b). Such disagreements make the publication of detailed descriptions of future experiments extremely desirable and necessary.

In the first experiments concerned with the question of Voiceprint, the observers were required to sort spectrograms into groups that represented dif-

ferent speakers (Kersta, 1962a, 1962b). Later experiments employed the multiple-choice identification test (Kersta, 1962c; Young and Campbell, 1967; Stevens, Williams, Carbonell, and Woods, 1968). There have been no reports of experiments dealing directly with the type of identification task commonly encountered in criminal investigations. Ladefoged and Vanderslice (1967) argued that the reliability of Voiceprint identification in practical cases cannot be predicted from the results of the published studies.

It has been claimed that spectrogram recognition performance is essentially unaffected by the loss of teeth, tonsils, or adenoids, the aging process, and attempts to disguise the voice, such as changing the fundamental frequency, whispering, mimicking another voice, and ventriloquism (Kersta, 1962c, Anon., 1965). However, in the absence of supporting experimental data, these claims cannot be considered established facts. Furthermore, when the speech signal is degraded, as it may well be when transmitted by a typical telephone system, many of the above-mentioned factors can be expected to reduce the reliability of this method.

According to Kersta (1962b), the probability that two speakers have similar enough vocal-tract dimensions and articulation patterns to produce indistinguishable spectrograms is extremely small. This belief, which appears to underlie many experiments, has not been formally translated into a hypothesis that can be tested with a finite population of speakers. There is evidence that two arbitrarily selected speakers can occasionally produce very similar spectrograms (Ladefoged and Vanderslice, 1967). This situation is illustrated in Fig. 2 for the word "you." Findings of this kind suggest that the range of one speaker's pronunciations of a given word (intraspeaker variability) may partially overlap the range of another speaker's pronunciations of the same word, and argue for the use of a large number of different words in making an identification. There is also evidence of considerable similarity among spectrograms representing different members of a family (Kersta, 1965a), and this suggests another source of observer fallibility.

Stevens, Williams, Carbonell, and Woods (1968) examined the ability of observers to distinguish between familiar and unfamiliar speakers in a 32-item identification-discrimination test. The observer was given eight reference spectrograms that represented eight "familiar" speakers. There were two experimental conditions; either four or 16 of the 32 test

Figure 2.—Similar Spectrograms of the word "YOU" uttered by two arbitrarily selected speakers (after Ladefoged and Vanderslice, 1967.)



spectrograms represented "unfamiliar" speakers who were not represented by the reference spectro-

grams. The results of this study are shown in Table II. Most of the familiar speakers were recognized as such, and they were subsequently correctly identified. Many of the unfamiliar speakers, however, were erroneously recognized as familiar speakers, especially when they appeared as often as the familiar speakers. As a point of comparison, listening tests were conducted using the same speakers and the same test format. Spectrograms were not employed in these tests. These data are shown in Table III. A comparison of the two sets of data reveals that there were considerably more acceptances of unfamiliar speakers in the visual tests than in the oral tests. When only four of the 32 test items represented unfamiliar speakers, there were also more false rejections of familiar speakers in the visual tests. Thus, speaker recognition by listening was found, in this study, to be the more accurate method. It must be pointed out that the observers employed by Stevens et al had very little training. One would expect better performance from highly-trained observers, but this study does demonstrate that speaker recognition by spectrogram matching is neither obvious nor easily achieved.

Data based upon carefully controlled experiment

TABLE II.—Percent Correct Recognition of Familiar and Unfamiliar Male Speakers by Visual Comparison of Spectrograms. (Data are shown for two experimental conditions. After Stevens, et al., 1968.)

| 4 of 32 test items by unfamiliar speakers  |               |            |
|--|---------------|------------|
| Speaker                                    | Recognized as |            |
|  | Familiar      | Unfamiliar |
| Familiar                                   | 80            | 20         |
| Unfamiliar                                 | 31            | 69         |
| 16 of 32 test items by unfamiliar speakers |               |            |
| Speaker                                    | Recognized as |            |
|  | Familiar      | Unfamiliar |
| Familiar                                   | 90            | 10         |
| Unfamiliar                                 | 17            | 53         |

TABLE III.—Percent Correct Recognition of Familiar and Unfamiliar Male Speakers by Listening. (Data are shown for two experiment conditions. After Stevens et al., 1968.)

| 1 of 32 test items by unfamiliar speakers  |               |            |
|--|---------------|------------|
| Speaker                                    | Recognized as |            |
|  | Familiar      | Unfamiliar |
| Familiar                                   | 88            | 12         |
| Unfamiliar                                 | 6             | 94         |
| 16 of 32 test items by unfamiliar speakers |               |            |
| Speaker                                    | Recognized as |            |
|  | Familiar      | Unfamiliar |
| Familiar                                   | 92            | 8          |
| Unfamiliar                                 | 8             | 92         |

using well-trained observers will soon be available. In a program sponsored by the National Institute of Law Enforcement and Criminal Justice, U.S. Department of Justice, through the Michigan State Police, scientists at Michigan State University have been examining speaker recognition by visual comparison of spectrograms as a function of several variables including: quality of recordings, context of words used in the identification task, number of speakers in the comparison population, number of words used for identification purposes, and number of samples of each word (Tosi, 1970). These data, which should soon be published, will provide a good determination of the reliability of speaker recognition by the current technique of making visual comparisons of speech spectrograms.

The above discussion may be summarized as fol-

lows: In view of the use of the visual comparison of spectrograms for speaker identification as evidence in courts of law, the fallibility of the observer must be studied further (Bolt et al., 1970). Future experiments should be carefully designed so as to avoid possible artifacts in the results. A detailed description of the experimental procedure, accompanied by the obtained data, should be published or otherwise be made available to the scientific community. Claims should be clearly differentiated from proven facts, and statements establishing an analogy between Voiceprints and fingerprints should be avoided. Although the spectrographic method for speaker identification has obvious potential in various investigative and forensic applications, its reliability as a means of identification has not yet been established.

## V. Speaker Recognition by Machine

Two approaches have been used to study the feasibility of speaker recognition by machine. One approach is to have the machine generate and

examine amplitude-frequency-time matrices of specific speech samples. The other approach is to have the machine extract speaker-dependent parameters

from the speech signal and subject them to a statistical analysis. Each approach has led to a number of recognition techniques.

In the first case, the utterances of specific speech samples are usually processed by a spectrum analyzer that consists of a bank of bandpass filters, rectifiers, and smoothing circuits. The outputs of the analyzer are periodically sampled, and the amplitudes are quantized for further processing by computer. Each utterance is represented in the computer by a data matrix. The rows of the matrix correspond to the frequency bands of the spectrum analyzer, the columns correspond to the temporal locations of the sample spectra, and each matrix cell contains the measured amplitude level. Such a matrix may be thought of as a "digital spectrogram." For each phrase, word, or phoneme used, several matrices representing different utterances by the same speaker are combined to form a single reference matrix for that speaker. A reference matrix is thus constructed for each speaker participating in a recognition experiment. The speaker to be recognized is represented by a test matrix. Depending on the type of recognition to be performed the test matrix is compared with all, or one of the reference matrices. The degree of similarity between the test matrix and each reference matrix is computed, and the results are used to arrive at a decision.

There are two basic recognition tasks, identification and discrimination. In the identification task, several reference matrices are used and it is assumed that the speaker represented by the test matrix is also represented by one of the reference matrices; thus, the reference matrix that is most similar to the test matrix is expected to identify

## VI. Future Developments in Speaker Recognition

The previous material describes in general terms the current status of speaker recognition by listeners, by visual examination of spectrograms, and by machine. Here we will comment briefly on the potential of each of these methods.

### A. Speaker recognition by listeners

There is little likelihood that much can be

the speaker represented by the test matrix. In the discrimination task only one reference matrix is used and the speaker represented by the test matrix may or may not be represented by this reference matrix. Decision rules are selected to specify when the test and reference matrices are similar enough to represent the same speaker. A summary description of six studies is presented in Table IV. For each experimental study this table gives the speech materials used, the configuration of the data matrix, the number of utterances included in the reference and test matrices, the recognition task, the number of speakers involved, and an overall measure of performance. Obtained percent correct scores range from 89 to 95 percent.

Techniques using statistical analyses of speech parameters involve two distinct processes: (1) the extraction from the speech signal of parameters thought to be useful for differentiating among speakers, and (2) the application of decision rules to combinations of parameter values that represent particular speech samples.

Questions regarding the most appropriate speech parameters have generally not been resolved as well as have questions regarding optimal decision rules. Various kinds of parameters have been examined, using both waveform analyses and spectroanalyses of the speech signal. Studies conducted by Clarke and Becker (1969), Hargraves and Starkweather (1963), Smith (1962), Ramishvili (1966), Edie and Sebestyen (1962), and Floyd (1964) have considered many speech parameters and several decision techniques. In general, results have been promising but it is clear that much work remains to be done before automatic recognition techniques attain high reliability.

done, or should be done, to improve the average individual's ability to recognize speakers by voice. Identification based on the average individual's recognition of voice will undoubtedly remain unreliable although in some cases it may be admitted as evidence. Trained linguists, on the other hand, are reported to be very good at recognizing various dialects and the geographical region of origin of speakers. They are sometimes employed in the investigation phase of law enforcement and have been

TABLE IV.—Summary Description of Six Recognition Techniques Using Specific Cue Material

| Study                        | Speech Material | Matrix Configuration |           |               |                | Utterances Incl. |             | Recogn. Task | Speakers | Perform. % |
|------------------------------|-----------------|----------------------|-----------|---------------|----------------|------------------|-------------|--------------|----------|------------|
|                              |                 | Frequency Bands      | Range kHz | Interval msec | Amplitude bits | Ref. Matrix      | Test Matrix |              |          |            |
| Pruzansky (1963)             | 10 Words        | 17                   | 0.2-7.0   | 10            | 10             | 3                | 1           | Ident.       | 10       | 89         |
| Pruzansky and Mathews (1964) | 10 Words        | 17                   | 0.1-10.0  | 10            | 10             | 3                | 1           | Ident.       | 10       | 93         |
| Ramishvili (1965)            | 10 Words        | 7                    | 0.2-10.0  | 50            | 2              | 10               | 1           | Ident.       | 20       | 92         |
| Li et al. (1966)             | 3 Phrases*      | 15                   | 0.3-4.0   | 20            | 12             | 10+              | 1           | Discr.       | 20       | 90         |
| Glenn and Kleiner (1968)     | Conson. [n]     | 25                   | 1.0-3.5   | —             | 6              | 10               | 10          | Ident.       | 30       | 93         |
| Mecker (1967)                | 4 Vowels        | 19                   | 0.2-8.0   | 40+           | +              | 20               | 3           | Discr.       | 11       | 95         |

\*Used only first 500 msec of each utterance.

+Used relative frequencies of occurrence of three spectral slopes.

used as expert witnesses in legal proceedings. It is very possible that some linguists are far superior to the untrained individual in achieving reliable speaker recognition. However, we know of no studies that have directly investigated this possibility, nor do we know of any plans to do so. Thus, it would appear that the potential of speaker recognition by listeners is quite limited.

### B. Speaker recognition by visual examination of spectrograms

It is unlikely that this method has achieved its full potential. There has been too little systematic study of spectrogram features to determine optimal procedures for discriminating among talkers. While the current performance of analyzing machines can undoubtedly be improved upon, the fact remains that the spectrograph was not designed to emphasize features useful for distinguishing among talkers and it discards much information that may be of value for this purpose. Whereas the speech spectrograph should prove to be an increasingly valuable tool for investigative purposes it is unlikely

that it will ever, under all circumstances, permit positive identification by voice.

### C. Speaker recognition by machine

This method of speaker recognition should prove to be the most promising. Computers are now capable of performing fast and accurate analyses of speech waveforms. Various parameters may be abstracted from the speech waveform and analyzed to determine those features most useful for distinguishing among talkers. Freedom to choose these optimal parameters should enable machine performance to exceed that of listeners or of trained observers using spectrograms as these two latter methods suffer from strict and arbitrary limitations upon processing equipment. While it is not scientifically obvious that absolutely positive identification by voice alone will ever be achieved by any method, speaker recognition by machine has the best chance of attaining this goal. To achieve improved or perfect performance the relevant speech parameters must be properly identified and incorporated into the analysis and decision processes of the machine.

## REFERENCES

- Anonym, "Voice Print Identification." *Criminalistics* (W. W. Turner, ed.), Aqueduct Books, Rochester (1965).
- Bolt, R. H., F. S. Cooper, E. E. David, Jr., P. B. Denes, J. M. Pickett, and K. N. Stevens, "Speaker Identification by Speech Spectrograms: A Scientists' View of its Reliability for Legal Purposes." *J. Acoust. Soc. Am.* *17*, 597-612 (1970).
- Broders, W., "Voiceprint Allowed as Evidence; Ruling Called First of Its Kind." *The New York Times*, Tues., April 12 (1966).
- Clarke, F. R., and R. W. Becker, "Comparison of Techniques for Discriminating Among Talkers." *J. Speech and Hearing Res.* *12*, 747-761 (1969).
- Edie, J., and G. S. Sebestyen, "Voice Identification General Criteria." Rept. RADC-TDR-62-278, Rome Air Development Center, Air Force Systems Command, Griffiss AFB (May 1962).
- Floyd, W., "Voice Identification Techniques." Rept. RADC-TDR-64-312, Rome Air Development Center, Research and Technology Division, Air Force Systems Command, Griffiss AFB (Sept. 1964).
- Glenn, J. W., and N. Kleiner, "Speaker Identification Based on Nasal Phonation." *J. Acoust. Soc. Am.* *43*, 368-372 (1968).
- Hargreaves, W. A., and J. A. Starkweather, "Recognition of Speaker Identity." *Language and Speech* *6*, 63-67 (1963).
- Hecker, M. H. L., "Speaker Recognition." Monograph A.S.H.A., (1970) (In press)
- Kennedy, H., "Appeals Court Reverses State's First 'Voiceprint' Conviction." *Los Angeles Times*, Fri., Oct. 11 (1968).
- Kersta, L. G., "Voiceprint Identification." *J. Acoust. Soc. Am.* *34*, 725 (A) (1962a).
- Kersta, L. G., "Voiceprint Identification." *Nature* *196*, No. 4861, 1253-1257 (1962b).
- Kersta, L. G., "Voiceprint-Identification Infallibility." *J. Acoust. Soc. Am.* *34*, 1978 (A) (1962c).
- Ladefoged, P., and R. Vanderslice, "The Voiceprint Mystique." Working Papers in Phonetics *7*, University of California, Los Angeles (Nov. 1967).
- Li, K. P., J. E. Dammann, and W. D. Chapman, "Experimental Studies in Speaker Verification, Using an Adaptive System." *J. Acoust. Soc. Am.* *40*, 966-978 (1966).
- McDade, T., "The Voiceprint." *The Criminologist*, No. 7, 52-70 (Feb. 1968).
- McGehee, F., "The Reliability of the Identification of the Human Voice." *J. Gen. Psychol.* *17*, 249-271 (1937).
- McGehee, F., "An Experimental Study in Voice Recognition." *J. Gen. Psychol.* *31*, 53-65 (1941).
- Meeker, W. F., "Speaker Authentication Techniques." Techn. Rept. ECOM-02526-F, U.S. Army Electronics Command, Ft. Monmouth, N.J. (Dec. 1967).
- Pruzansky, S., "Pattern Matching Procedure for Automatic Talker Recognition." *J. Acoust. Soc. Am.* *35*, 354-358 (1963).
- Pruzansky, S., and M. V. Mathews, "Talker-Recognition Procedure Based on Analysis of Variance." *J. Acoust. Soc. Am.* *36*, 2041-2047 (1964).
- Ramishvili, G. S., "Automatic Recognition of Speaking Persons." Rept. FTD-TT-65-1079, Foreign Technology Division, Air Force Systems Command, Wright-Patterson AFB (Dec. 1965).
- Ramishvili, G. S., "Automatic Voice Recognition." *Engineering Cybernetics*, No. 5, 84-90 (Sept.-Oct. 1966).
- Smith, J. E. K., "Decision-Theoretic Speaker Recognizer." *J. Acoust. Soc. Am.* *34*, 1988 (A) (1962).
- Stevens, K. N., C. E. Williams, J. R. Carbonell, and B. Woods, "Speaker Authentication and Identification: A Comparison of Spectrographic and Auditory Presentations of Speech Material." *J. Acoust. Soc. Am.* *44*, 1596-1607 (1968).
- Tosi, O., Personal communication (1970).
- Young, J. R., and M. H. L. Hecker, "Some Observations on the Problem of Machine Recognition of Speech." Proc. 1968 National Electronics Conf., Chicago (Dec. 1968).
- Young, M. A., and R. A. Campbell, "Effects of Context on Talker Identification." *J. Acoust. Soc. Am.* *42*, 1250-1254 (1967).

## PART 2

### Michigan State University Voice Identification Project

By

Oscar Tosi, Ph.D.

Charles Pedrey, Ph.D.

Herbert Oyer, Ph.D.

Julie Nicol, B.A.

William Lashbrook, Ph.D.

Donald Riggs, A.S.

DEPARTMENT OF AUDIOLOGY AND SPEECH SCIENCES  
MICHIGAN STATE UNIVERSITY  
East Lansing, Michigan

Consultants to the study:

Irwin Pollack, Ph.D., University of Michigan

John W. Black, Ph.D., Ohio State University

Lawrence Kersta, M.S., Voiceprint Laboratories

Mac Steer, Ph.D., Purdue University

Kenneth Stevens, Ph.D., Massachusetts Institute of Technology

## CONTENTS

|   | Page |
|---|------|
| I. Introduction .....   | 39   |
| A. Purpose .....  | 41   |
| B. Variables tested .....   | 41   |
| C. Speakers and spectrograms .....                                    | 41   |
| D. Examiners and training of examiners .....                          | 43   |
| E. Experimental procedure .....                                       | 45   |
| II. Results From the First Cycle of the Project .....                 | 48   |
| III. Results From the Second Cycle of the Project .....               | 52   |
| IV. Discussion and Conclusions .....                                  | 57   |
| V. Extension of Results From Forensic Models to Real Cases .....      | 58   |
| A. Population of known voices .....                                   | 58   |
| B. Availability of time and responsibility of the examiners .....     | 58   |
| C. Type of decisions examiners are urged to reach in each trial ..... | 58   |
| D. Availability of clues .....  | 59   |
| References .....  | 60   |

## I. Introduction

Studies concerned with methods for identifying persons are important because of the legal ramifications and because of the forensic involvements associated with the application of these studies.

Fingerprinting, photographic and antropometric techniques are the most commonly used methods of identification. In some instances mapping the contours of teeth and nasal cavities and/or mapping labial impressions have also proved to be useful means of identification.

In the present era of widely used telephone, radio and tape recorder communication, the voice of an individual is often the only available clue for identification. The problem that speaker recognition poses is essentially different from the problem of fingerprinting or any other type of technique of identification using clues that are "invariant." Indeed, the voice of an individual is far from being invariant as are his fingerprints. Usually no person utters the same word twice with all characteristics being *exactly* the same; laymen, for the most part, are not aware that such differences occur.

Speech scientists refer to these differences as "intraspeaker variability." As yet "intraspeaker variability" is not well understood; nor has it been quantified or correlated with specific acoustical parameters of the speech signal.

In contrast, differences between the same words uttered by different speakers are quite apparent to any listener; such differences are labeled "interspeaker variability." This variability stems mainly from anatomical differences in vocal tracts and from learned differences in the use of the speech mechanism. "Interspeaker variability"—like "intraspeaker variability"—has not, as yet, been quantified or correlated with specific acoustical parameters of the speech signal. Nevertheless, observers—although they may not understand the rules of detection—do detect "interspeaker variability" in much the same manner as they detect differences in the handwritings or photographs of different persons, even to the point of identifying a person through these clues. It is to be noted that handwritings and pho-

tographs as well as speech involve "intraperson" and "interperson" variability.

Hecker (1970), in a critical study of the methods presently available for speaker identification, classifies these methods into three general areas: a) speaker identification by listening; b) speaker identification by machine; and c) speaker identification by visual examination of spectrograms. In essence all three procedures are based on the assumption that "interspeaker variability" is always greater than "intraspeaker variability," regardless of the parameters involved in these variabilities. As yet this assumption has not been proved. Since the parameters responsible for variabilities are not well determined and quantified, at the present time the only way to prove scientifically that "interspeaker variability" is greater than "intraspeaker variability" is by inference. Such an inference can be produced by proper evaluation of empirical data obtained through experiments of speaker identification. An inference thus derived might be affected by data which could be confounded by both effects from speakers and from the method of identification used.

Review of the literature concerning the three methods of identification described by Hecker, reveals certain deficiencies in each method. Speaker identification by listening, one of the methods discussed, is far from being 100 percent accurate. It is an entirely subjective method; an expert witness using only this method would be unable to justify his conclusions in a court of law. Besides, the task of comparing voices purely by listening becomes a difficult one when several speakers are involved. In this case the method necessitates that the examiner relies a great deal on auditory memory since listening dichotically to the known and unknown voices is most inconvenient.

Speaker identification by machine, a second method available, is presently less accurate than any method involving human examiners. Questions regarding the most appropriate speech parameters for machine recognition, as well as questions regarding

optimal decision rules are still far from being answered. This situation is comparable to that of other fields of recognition, handwriting and fingerprinting, for which no recognition machines are presently available. In the future, hard research might bring knowledge to overcome the present limitations of speaker identification by machine. It is difficult to predict just when "totally" reliable machines will become available. But even if a speaker recognition machine were available, the human expert, trained in phonetics, spectrography and related areas, would be required to select the proper samples from the unknown and known voices to feed the machine, to evaluate the machine output, and possibly to check the results by using an alternative method.

The third method of speaker recognition is based upon the visual examination and comparison of spectrograms. Speech spectrography was developed at Bell Research Laboratories by Potter *et al.* (1947). This type of spectrography is accomplished by the use of an instrument called a sound spectrograph, which transforms speech into a visual display, a spectrogram. The spectrogram portrays three main parameters of speech: time (horizontal axis), frequencies (vertical axis) and relative amplitude (degree of darkness of the different spectrographic regions). Each phoneme, word or phrase is correlated with a characteristic spectrographic pattern. The general aspect of patterns corresponding to different utterances of the same word are similar, in such a way that a person specially trained in "reading" spectrograms, could determine with more or less accuracy which words or phrases were portrayed by a particular pattern. However, "interspeaker" and "intraspeaker" variabilities are also portrayed by the spectrographic patterns. In fact, the spectrograms of different utterances of the same word or phrase by the same or by different speakers are never *exactly* alike.

Kersta (1962) claimed that spectrograms of several utterances of the same words by a given speaker always contain more similar spectral features than those produced by different speakers. Kersta concluded, therefore, that speaker identification by visual examination of spectrograms, has to be reliable. According to Kersta, speaker recognition by visual inspection of spectrograms consists of subjectively matching similarities found in pairs of spectrograms from the same person, that are not found in pairs of spectrograms from different persons. The dissimilarities presented by the matched

spectrograms are disregarded; they are assumed to be a result of intraspeaker variability. To back his claim, Kersta published the results of experiments he performed at Bell Research Laboratories. In these experiments he observed fewer than 1 percent of wrong identifications.

Matching similarities through the combination eye-brain is essentially a subjective method, but the examiner can display objectively these similarities in a court of law to support his subjective conclusions. The possibility of objective displays for legal application has perhaps made this method of speaker recognition quite appealing. Unfortunately this characteristic has led to its being labeled "voiceprinting." This word conveys the erroneous impression that speaker recognition by visual examination of spectrograms can be equated with fingerprinting. Nor should it be assumed that Kersta's method precludes listening to the known and unknown voices. On the contrary, the examiner who selects samples of the voices to be spectrally compared must listen to the samples in addition to examining the spectrograms visually. The spectrograph commercialized by Kersta under the trade name of "Voiceprint" (Presti, 1967) has a special playback mechanism, that allows proper selection and continuous listening of samples prior to their being fed into the processing circuits. Since 1962 Kersta has been producing legal testimony on speaker identification by using his method, as well as offering training in this "art" to law enforcement officers.

Several speech scientists (Bolt *et al.*, 1970) have expressed their concern for the legal application of "voiceprinting," prior to proving its accuracy and reliability through controlled experimentation.

Tosi (1967, 1968) evaluated the "voiceprinting" method by analyzing the data derived from 236 experimental trials of identification performed by nine of Kersta's trainees, as well as by participating himself in the training courses given by Kersta at the Voiceprint Laboratories. These trials of identification of one speaker among 50, using five clue words, yielded an error of 6.3 percent. In his report to the Michigan Department of State Police, Tosi could only conclude that the "Kersta method shows promise," suggesting the need for an independent study, one that would include variables not considered by Kersta in his experiments and that would further test "voiceprinting."

Such a study, "Michigan State University Voice Identification Project," was conducted at the De-

partment of Audiology and Speech Sciences of Michigan State University from 1968 to 1970, under a contract with the Michigan Department of State Police, through a grant from the United States Department of Justice.

#### A. Purpose

The two year Voice Identification Project, conducted at the Department of Audiology and Speech Sciences, Michigan State University, from 1968 to 1970 had the twofold purpose of:

1. replicating Kersta's experimental trials of speaker identification by visual examination of spectrograms, and
2. testing other types of trials of speaker identification that included variables most relevant to forensic applications of this method, not reported in the Kersta studies.

#### B. Variables tested (Figure 1)

1. Different number of clue words: nine and six clue words.
2. Different number of utterances or examples: from the same clue word produced by each speaker: one, two, or three utterances.
3. Different types of recording transmissions: ( $\alpha$ ) directly into a tape recorder; ( $\beta$ ) through a telephone line in a quiet environment, and ( $\gamma$ ) through a telephone line in a noisy environment (50 dB<sub>L<sub>p</sub></sub> of white noise measured at the head of the speaker).
4. Context of the clue words used for identification: three types of contexts were tested: (I) clue words spoken in isolation; (II) clue words spoken in a fixed context—same sentences produced by the known and unknown speakers were compared; and (III) clue words spoken in a random context—different sentences containing the clue words were compared. The clue words used in this experiment were: it-is-on-you-and-the-I-to-me. These words were selected because of their high percentage of occurrence in English.
5. Different number of "known" speakers included in each trial of identification: 10, 20, or 40 "known" speakers.
6. Time elapsed between recordings from the same speaker: these recordings were obtained during two different recording sessions, held one month

apart. Spectrograms obtained from the first recording session were divided into two sets. One set was assumed to correspond to a known speaker and labeled "known spectrograms." The speaker number was stamped directly in the "known spectrogram." The other set was assumed to correspond to an unknown speaker and labeled "contemporary matching spectrograms." The speaker number was coded in the matching spectrograms. Spectrograms obtained from the second recording session were assumed to correspond to an unknown speaker and labeled "non-contemporary matching spectrograms." The speaker number was coded in the "non-contemporary matching spectrograms."

7. Awareness of the examiners: two different conditions were tested: a) closed trials: in which the examiners were aware that the "unknown" speaker was among the "known" ones, and b) open trials: in which the examiners were not aware whether or not the "unknown" speaker was among the "known" ones. Although the examiners were given only "open trials," the researcher presented randomly to them two different types of open trials: open trials including the unknown speaker among the known speakers (trials referred to as "open-match") and open trials in which the unknown speaker was not included among the known speakers (trials referred to as "open-no match").

Listening to the unknown and known voices was excluded from this study, based solely on *visual examination* of spectrograms.

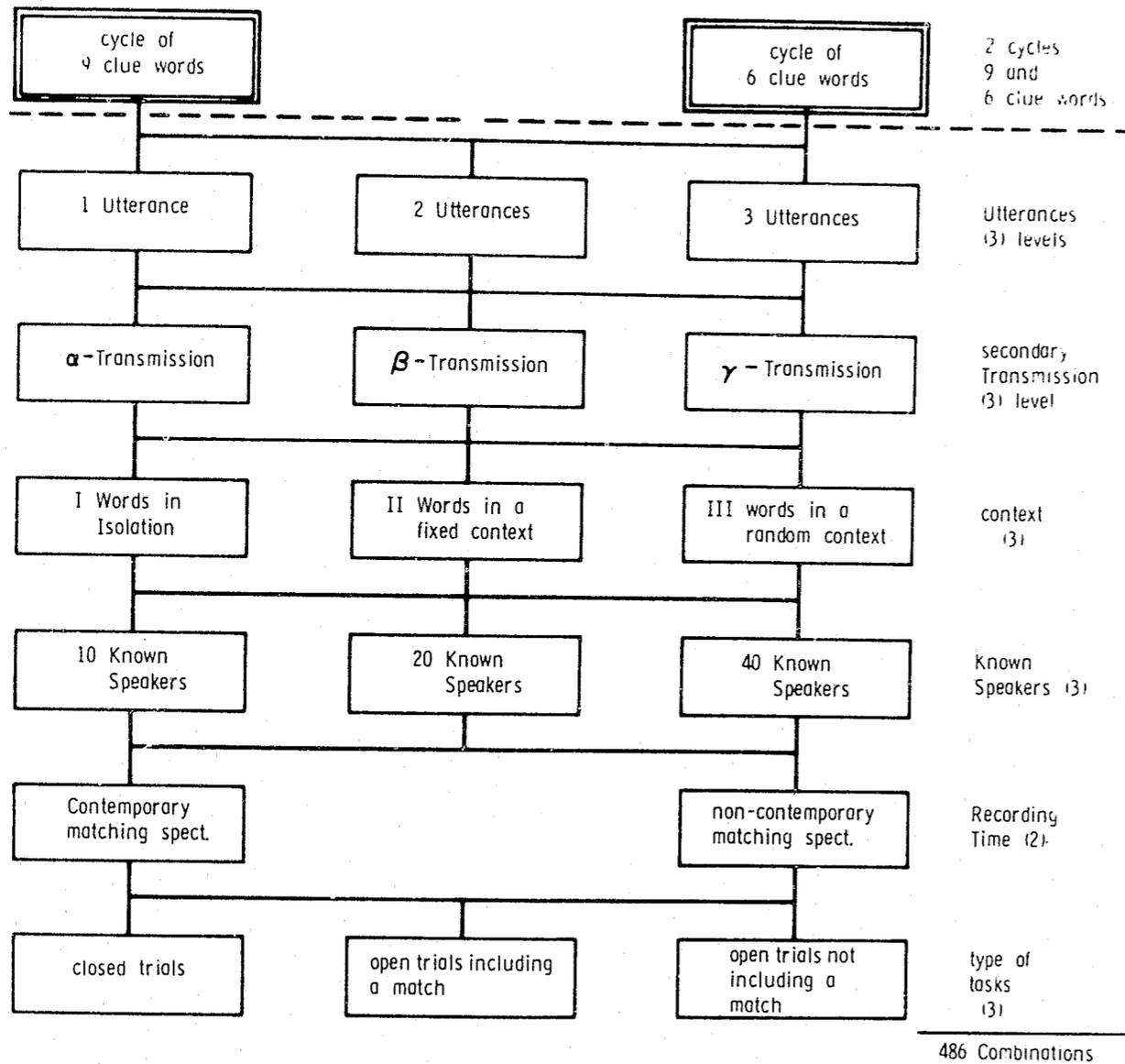
#### C. Speakers and spectrograms

Two hundred fifty speakers, randomly drafted from a population of approximately 25,000 male students at Michigan State University, participated in this experiment. The mean age of these speakers was 19.8 years; the range 17 to 27 years; standard deviation 2.1 years. This population excluded foreign students; all speakers were natives of the United States of America, utilizing general American English dialect with no speech defects.

The speakers recorded nine clue words during two sessions held one month apart. The nine clue words were spoken in isolation, in a fixed context and in a random context provided by six different sentences. The same texts were repeated six times by each speaker in each recording session. The recordings were obtained by using three different

Figure 1.—Organization of the experimental design. For each cycle including 9 or 6 clue words, 486 different types of trials were performed by 9 subpanels of examiners. This process was reiterated 4 times. Total number of responses for each type of trial  $9 \times 4 = 36$ . Total number of responses for each cycle:  $486 \times 36 = 17,496$

## MSU Voice Identification Project



types of transmission:  $\alpha$ ,  $\beta$ , and  $\gamma$  as described on page 41.

These recordings were processed through a "Voiceprint" spectrograph using an expanded scale of frequencies from 50 Hz to 4,000 Hz. The spectrograms yielded by each speaker during the first recording session were divided into two groups. Spectrograms of the first group were assumed to have been produced by "known" speakers. They were designated "known" spectrograms and consequently labeled with the corresponding speaker number. The second group of spectrograms were considered as produced by "unknown" speakers. They were designated "contemporary matching spectrograms;" the corresponding speaker number was therefore coded in these spectrograms. All the spectrograms yielded by the 250 speakers during the second recording session were assumed to correspond to "unknown" speakers. They were designated "non-contemporary matching spectrograms." The speaker numbers were also coded in these spectrograms. All coded numbers were covered with masking tape, so as to be invisible to the examiners.

### D. Examiners and training of examiners

Newspaper advertisements were used to announce the opening of examiner positions. Applicants were screened prior to their selection as examiners. All 29 examiners who participated in this study passed the screening tests. The process of applicant screening included:

1. Lecturing briefly to the applicant on spectrograms and spectrographic data;
2. Showing the applicant pairs of very similar spectrograms of a sentence produced by the same speaker, and pairs of quite different spectrograms of the same sentence produced by different speakers;
3. Asking the applicant to decide by visual inspection which two out of four spectrograms had been produced by the same speaker. Each applicant was subjected to three different trials of this type.
4. Those applicants who performed successfully the above trials were given a final test. The applicant was asked to decide which two, among eleven spectrograms of nine words spoken in isolation, were produced by the same speaker.

Only those applicants who performed successfully on the tests were considered for participation in the study as examiners. Those accepted received

approximately one month of training prior to starting the experimental trials. This training consisted of:

1. lectures on phonetics and spectrography, and a discussion of the variables to be tested in the experiment;

2. performance of closed trials of identification of one speaker among ten, by using words in isolation, contemporary matching spectrograms. These trials were first performed under direct supervision of the researcher; after a few days the examiner was left on his own. He was informed of his mistakes and allowed to compare the right matches with the wrong ones that he had selected. The examiner was instructed to make mainly subjective decisions similar to those made when comparing different photographs of faces or different hand-writings. In the event that he felt uncertain about his judgment of matching several spectrograms of known voices which appeared to be subjectively very similar to the spectrogram of the unknown voice, the examiner was instructed to consider the following objective points of similarity between spectrograms of the unknown and known voices:

- (a) similar mean frequencies of vowel formants,
- (b) formants band-widths, (c) gaps and type of vertical striations, (d) slopes of formants, (e) durations, (f) characteristic patterns of fricatives and interformant energies. According to Kersta these similarities are often present and more numerous in pairs of spectrograms of the same words produced by the same speaker at different times, than in pairs of spectrograms of the same words produced by different speakers.

3. After each examiner performed these closed trials including contemporary matching spectrograms and words in isolation, with a success better than 96 percent, he was given other types of tasks that were increasingly more difficult. These "advanced" tasks included: open trials, non-contemporary matching spectrograms, and words in fixed and random contexts.

After one month of training the actual experiment started. Spectrograms used for training were not used during the experiment. Listening to the unknown and known voices was excluded from this study.

The examiners were grouped into three panels according to sex and background (Figure 2). The first panel consisted of women ranging from 17 to 60 years of age, with various levels of education, from high school up to four years of college. The

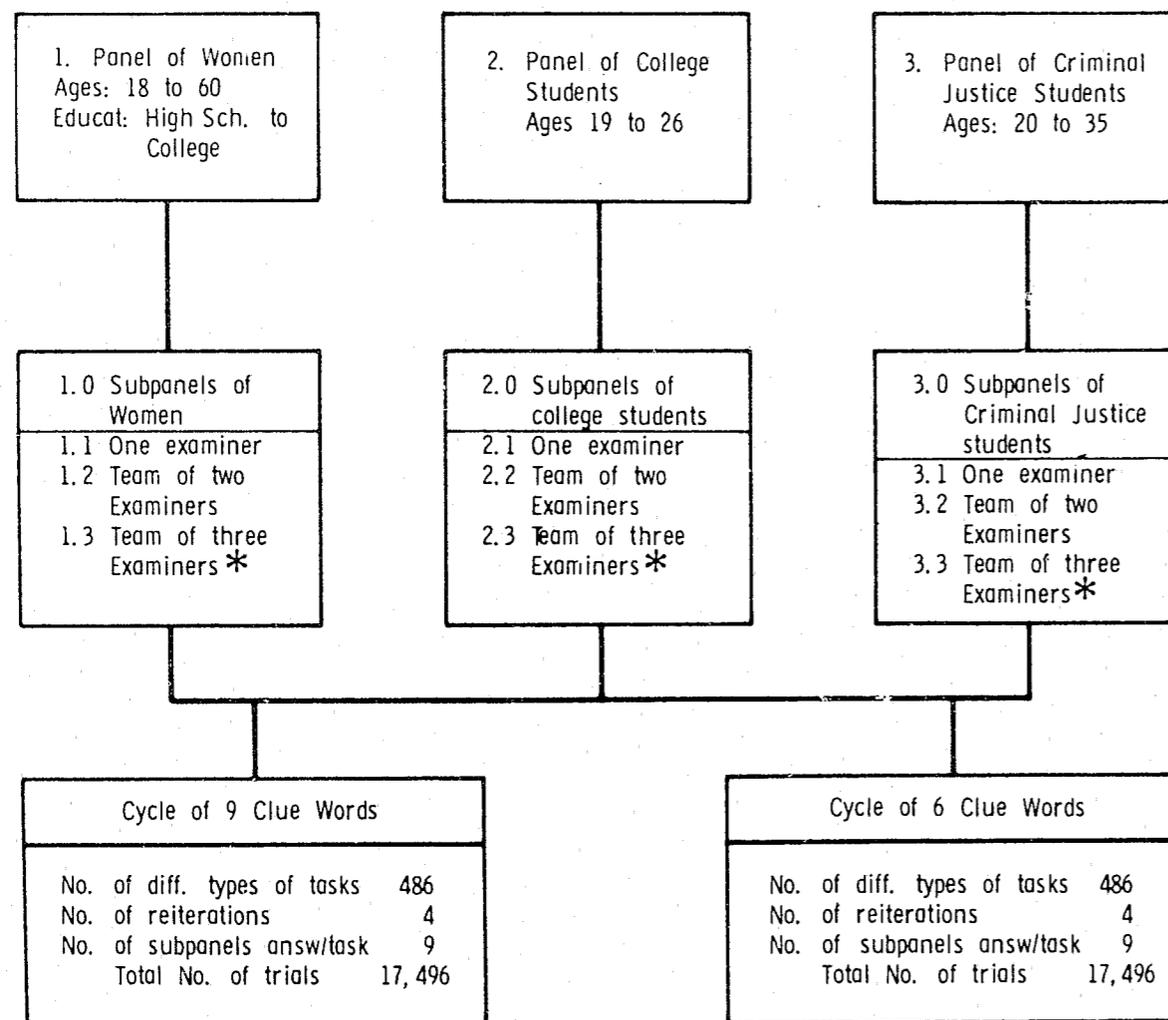
second panel included male undergraduate students from several departments of Michigan State University. The third panel was formed exclusively by students from the Criminal Justice Department of this University. Further, each panel was divided into three subpanels: one formed by a team of three examiners, one by a team of two, and the third subpanel consisted of a single individual. These nine subpanels performed the same experimental tasks, yielding nine answers for each different type of trial.

The examiners belonging to each panel were rotated within the three different subpanels of one, two and three members that comprised every panel. This procedure permitted a better observation of effects due exclusively to panel type and subpanel size, minimizing individual effects on the results. However, individual variations among examiner performances were detected and qualitatively evaluated by the staff.

A total of 29 examiners were employed in the project. This total number provided personnel for

Figure 2.—Panels and Subpanels of Examiners

## MSU Voice Identification Project



\* For this cycle the subpanels of 3 members were replaced by subpanels of 2 members

completing the nine subpanels utilized, as well as compensating for the rate of attrition due to resignations. All examiners received the same screening and training as described above. They were paid a flat rate per hour. Toward the end of the first cycle of the project, in which nine clue words were used for the identification trials, the rate of examiners' attrition increased, reducing the crew to only 15 examiners. Since less than four months remained for the termination of the contract, hiring and training new examiners seemed inconvenient at that point. This opinion was reinforced by the observation that the examiners who quit were the less motivated of the group, often complaining of boredom and fatigue. Decision was made to complete the second cycle of the project, using only these 15 more motivated and reliable examiners. Although teams of three members had to be reduced to teams of two members, no significant effect from this alteration was expected.

### E. Experimental procedure

The experiment was divided into two cycles, the first using nine clue words and the second, six clue words. There were 486 different types of tasks involving every possible combination of the variables tested. Each combination was reiterated four times in an unsystematic manner by nine subpanels of examiners, using different spectrograms in each reiteration. Consequently there were 36 answers for each of the 486 different tasks, yielding a total of 17,496 trials per cycle. Therefore, the total number of trials involved in this experiment was 34,992. One-third of the trials were the "closed" type; two-thirds were "open" trials of which 50 percent were "open-match" and 50 percent were "open-no match," randomly presented to the examiners.

The examiners worked three hours daily, completing as many tasks as possible. Examiners were encouraged to take rest periods as needed to avoid fatigue, a condition that might hamper examiner performance.

To perform identification trials the examiners were provided with a set of spectrograms assumed to correspond with known speakers ("known" spectrograms) and spectrograms assumed to correspond to an "unknown" speaker ("matching" spectrograms). These spectrograms were arranged on 27 specially designed tables and shelves. Each table was used for a specific task, and was supplied daily

with different spectrograms. The matching spectrograms (unknown voices spectrograms) were secured on the tables with masking tape covering the coded numbers. The known voices spectrograms were placed on the shelves attached to each table.

Examiners of each subpanel were provided with answer sheets specifying the type of trial arranged on each table. (See Figure 3). Judgments from each trial were recorded by the examiners on these answer sheets. After the completion of each answer sheet, the members of the respective subpanel were informed of their mistakes and given an opportunity to inspect the correct matching spectrograms. This procedure encouraged continuous learning on the part of the examiners.

The tasks of the examiners in the *open trials* consisted of deciding whether the "matching" spectrograms were or were not produced by one of the "known" speakers; and if they were, which "known" speaker produced them. In these open trials three kinds of errors were possible: (See Figure 4).

1) *Error A*: a match did exist but the examiner selected the wrong one. (False Identification)

2) *Error B*: a match did exist but the examiner failed to recognize it.

3) *Error C*: a match did not exist although the examiner selected one. (False Identification)

Errors A and C can be together labeled "errors of wrong matching" or "false identifications."

In the *closed trials* the examiners had to decide which "known" speaker produced the matching spectrograms. In these closed trials, since a match always existed, only one kind of error was possible. This error was labeled *Error D*. (False Identification or Wrong Matching)

Each subpanel was forced to reach a common positive decision in each trial; the decision was arrived at through discussion. In addition each member of the subpanel had to indicate his confidence—or lack of confidence—in this common positive decision. The following scale of self confidence was used for grading: 1 = almost uncertain; 2 = fairly uncertain; 3 = fairly certain; 4 = almost certain. Figure 3 is a copy of the answer sheets used to record decisions and confidence ratings.

Each subpanel used a different answer sheet, upon which up to 27 trials could be recorded. After completion of these 27 trials, the members of each subpanel submitted their answer sheet to the research assistant. She analyzed the recorded answers

# ANSWER SHEET

SUBPANEL # \_\_\_\_\_ NAME(S) \_\_\_\_\_

DATE \_\_\_\_\_

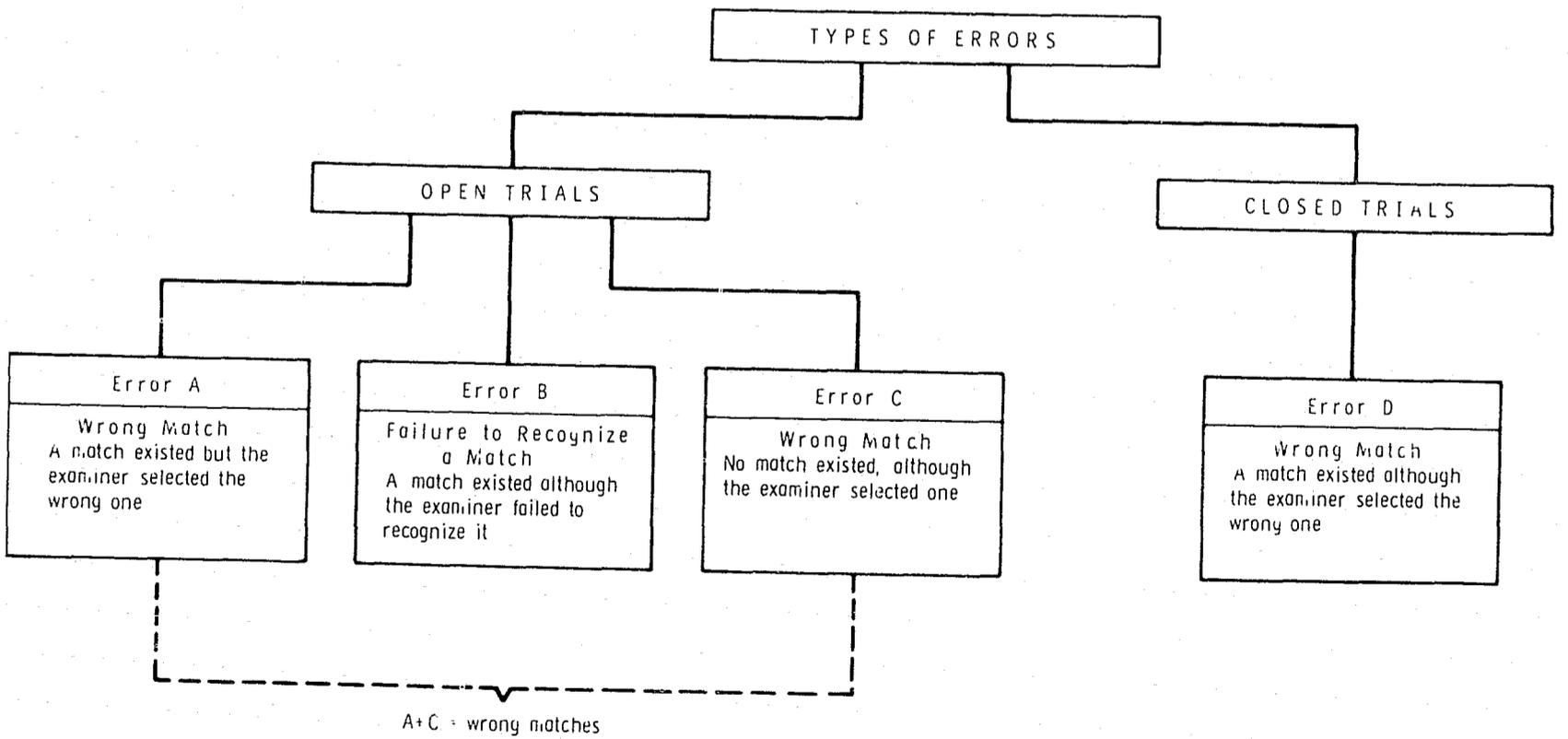
TRIAL # \_\_\_\_\_

W =     U =   
 C =     N =

| TASK | ROM | GRK | SPK | CL.  |  | S/MILAR | SCALE | RATE | M. SPECTR |
|------|-----|-----|-----|------|--|---------|-------|------|-----------|
|      |     |     |     | -OP. |  |         |       |      |           |
| 1    |     | α   | 10  |      |  |         |       |      |           |
| 2    |     | α   | 20  |      |  |         |       |      |           |
| 3    |     | α   | 40  |      |  |         |       |      |           |
| 4    |     | β   | 10  |      |  |         |       |      |           |
| 5    |     | β   | 20  |      |  |         |       |      |           |
| 6    |     | β   | 40  |      |  |         |       |      |           |
| 7    |     | γ   | 10  |      |  |         |       |      |           |
| 8    |     | γ   | 20  |      |  |         |       |      |           |
| 9    |     | γ   | 40  |      |  |         |       |      |           |
| 10   |     | α   | 10  |      |  |         |       |      |           |
| 11   |     | α   | 20  |      |  |         |       |      |           |
| 12   |     | α   | 40  |      |  |         |       |      |           |
| 13   |     | β   | 10  |      |  |         |       |      |           |
| 14   |     | β   | 20  |      |  |         |       |      |           |
| 15   |     | β   | 40  |      |  |         |       |      |           |
| 16   |     | γ   | 10  |      |  |         |       |      |           |
| 17   |     | γ   | 20  |      |  |         |       |      |           |
| 18   |     | γ   | 40  |      |  |         |       |      |           |
| 19   |     | α   | 10  |      |  |         |       |      |           |
| 20   |     | α   | 20  |      |  |         |       |      |           |
| 21   |     | α   | 40  |      |  |         |       |      |           |
| 22   |     | β   | 10  |      |  |         |       |      |           |
| 23   |     | β   | 20  |      |  |         |       |      |           |
| 24   |     | β   | 40  |      |  |         |       |      |           |
| 25   |     | γ   | 10  |      |  |         |       |      |           |
| 26   |     | γ   | 20  |      |  |         |       |      |           |
| 27   |     | γ   | 40  |      |  |         |       |      |           |

Figure 4.—Classification of Errors

## MSU Voice Identification Project



and graded the right responses with the numeral one (1); she graded wrong responses with the numeral zero (0). The numeral zero was followed by the letter A, B, C or D according to the type of error made. The numbers expressing the self confidence of each member of the subpanel on each common response were averaged; so the result of each trial was finally quantified with two expressions, conveying the right or wrong response and the averaged self confidence grade on such a response.

The results from trials, randomly presented to the subpanels, were transcribed and grouped in master tables to facilitate subsequent analysis. (See Appendix A). These master tables were designed to provide the experimenter with a systematic means for arranging in an orderly fashion the randomly administered types of trials. Checks of the correctness of the transcriptions were made daily, and all answer sheets were filed in a room protected with an electronic alarm system connected with the University Police Headquarters.

Examiners usually followed the same procedure to complete each trial of speaker identification. The steps in this procedure were: 1) comparing the spectrograms of the unknown and known voices by a rather fast scan; 2) discarding those known voices spectrograms that appeared subjectively to the examiner as containing no significant similarities with the unknown voice spectrograms. Usually these steps reduced to a very few the known voices spectrograms to be further examined; 3) continuing the scanning by folding and superimposing each of the

remaining known spectrograms on the matching spectrograms. This procedure provided the examiners with a better technique in searching for similarities and reduced even more the number of suspected known spectrograms; 4) If the previous steps did not produce a positive decision, the examiners counted the number of similarities they found between each of the suspected known spectrograms and the matching spectrograms. These similarities were listed on page 45. The known spectrogram which presented more points of similarities was supposed to be chosen as a correct response in the case of closed trials. For open trials the procedure was essentially the same, but complicated by the circumstance that the examiners had to decide between two possible alternatives: "there is not a match," or "there is a match, being speaker *n* the same as the unknown speaker." 5) Subpanel members arrived at a common decision for each trial through discussion. 6) After the decision was reached, each subpanel member assessed this decision by registering his personal rate of confidence on the common decision. He used the grading scale described earlier and recorded his judgment on the subpanel answer sheet, which were given to the research assistant for tabulation on the master tables.

After completion of each cycle, the results from the master tables were coded in IBM cards and processed through the 3600 CDC computer to calculate error percentages and perform an analysis of variance to test significances and interactions of the different variables involved in the experiment. Appendix B includes a complete report of the analysis of variance performed.

## II. Results From the First Cycle of the Project

During the first cycle of the present study, nine clue words were used in all experimental trials of speaker identification. This cycle was completed in approximately eight months. As described earlier, there were 486 different types of trials involving each possible combination of the variables tested. Each combination was reiterated four times by nine subpanels of examiners, using different spectrograms in each reiteration. Therefore, there were 36 answers for each of these 486 different tasks, yielding a total of 17,496 trials of speaker identification. One-third of these trials were the "closed" type;

two-thirds were "open" trials, 50 percent "open-match" and 50 percent "open-no match," randomly presented to the examiners.

The results from this first cycle were processed through a CDC 3600 Computer. Table 1 presents the pooled percentages of correct responses produced by the examiners under each of the main conditions tested in the project during this first cycle. Percentages from the different levels of each main condition level were computed by collapsing all other conditions and pooling responses from the nine examiner subpanels obtained in each of

the four reiterations of the 486 different types of tasks. A detailed report of the analysis of variance performed with the subpanels' responses is included in Appendix B. The results of this analysis of variance are also shown in Table 1. No significant statistical difference was detected between the following levels:

1. one, two or three utterances or examples of the same clue word. Percentage of correct responses from all trials using one utterance of each of the nine clue words, was 91.29 percent; using two utterances, 90.96 percent; and using three utterances, 92.49 percent.

2.  $\alpha$ ,  $\beta$ , and  $\gamma$  types of recording transmissions. Percentage of correct responses from all trials using  $\alpha$  recording transmissions was 92.42 percent; using  $\beta$  recording transmission, 91.31 percent; and using  $\gamma$  recording transmission, 91.02 percent.

Other levels of the variables tested showed statistically significant differences. Another analysis of variance was performed to test the differences in the performance of panels and subpanels of examiners (Appendix C). Table 2 shows the results of this analysis, as well as the percentages of correct responses from each panel and subpanel, pooled over all the common trials of speaker identification they performed. No significant difference was detected between panel types, but a significant difference was detected between subpanel sizes: subpanels of three members performed slightly better than the other subpanels. However, it must be pointed out that size alone may not be responsible for this difference in performance. It is believed that composition of the subpanel may also have been a contributing factor. Examiners rotated freely among the subpanels of their assigned panel. It is quite

TABLE 1.—First Cycle—Results of an Analysis of Variance of the Correct Responses Produced under Each of the Main Conditions Tested

| Condition  | Pooled percentage of correct responses | Probability of the difference between levels, less than: |
|--|--|--|
| Number of utterances of the same clue word:                      |  |  |
| 1 utterance .....  | 91.29                                  | n.s.   |
| 2 utterances .....   | 90.96                                  | n.s.   |
| 3 utterances .....   | 92.49                                  |  |
| Different types of recording transmissions:                      |  |  |
| ( $\alpha$ ) directly into a tape recorder .....                 | 92.42                                  | n.s.   |
| ( $\beta$ ) through a telephone line in quiet environment .....  | 91.31                                  | n.s.   |
| ( $\gamma$ ) through a telephone line in noisy environment ..... | 91.02                                  |  |
| Context of the clue words spoken:                                |  |  |
| (I) in isolation .....   | 95.77                                  | 0.01   |
| (II) in a fixed context .....                                    | 92.39                                  | 0.01   |
| (III) in a random context .....                                  | 86.59                                  |  |
| Different number of "known" speakers:                            |  |  |
| 10 speakers .....  | 93.03                                  | n.s.   |
| 20 speakers .....  | 91.87                                  | 0.01   |
| 40 speakers .....  | 89.58                                  |  |
| Time-elapsd between recordings:                                  |  |  |
| contemporary matching spectrograms .....                         | 95.21                                  | 0.01   |
| non-contemporary matching spectrograms .....                     | 87.95                                  |  |
| Awareness of examiners:  |  |  |
| closed trials .....  | 91.48                                  | 0.01   |
| open trials .....  | 90.14                                  |  |

TABLE 2.—First Cycle—Results of the Analysis of Variance of the Percentages of Correct Responses Produced by Panels and Subpanels of Examiners

|               |                                     | Percentage of correct responses | Probability of the difference between levels |
|---------------|-------------------------------------|---------------------------------|--|
| Panel Type    | Panel 1.0 Women                     | 91.30                           | n.s.   |
|               | Panel 2.0 College students          | 92.13                           | n.s.   |
|               | Panel 3.0 Criminal Justice students | 91.31                           |  |
| Subpanel Size | Subpanels of 1 member               | 91.06                           | n.s.   |
|               | Subpanels of 2 members              | 90.36                           | <0.01  |
|               | Subpanels of 3 members              | 93.31                           |  |

possible, therefore, that each three-member sub-panel had among its members one high quality examiner; there would be less chance of this occurring in the other two subpanels. The staff observed that the best examiners often exerted positive influence on those less motivated.

The grand mean percentage of errors from the 17,496 trials of speaker identification performed during the first cycle of the project, including nine clue words was 8.9 percent. This grand mean was composed of 4.3 percent errors of wrong matching or false identifications (errors A + C + D), and 4.6 percent of failures to recognize a match when it actually existed (error B). These percentages do not offer any kind of specific information because they were pooled from trials involving many different conditions. In order to construct models relevant to the forensic point of view, trials were grouped according to the following characteristics:

1. awareness of the examiners (closed or open trials).
2. time-elapsed (contemporary or non-contemporary matching spectrograms).
3. context (clue words spoken in isolation, in a fixed context or in a random context).

There were 972 closed trials and 1944 open trials (match and no-match) for each of the six possible combinations of time-elapsed and context levels. Errors from each group were counted and percentages were computed with respect to the total number of trials of each group (972 or 1944 respectively). Table 3 shows these 12 percentages. Figure 5 displays graphically the figures from Table 3. The consistent patterns of this graph might constitute an indicator of the examiners' reliability, considering that each point of the graph represents

grouped data obtained from trials examiners performed in an almost random sequence.

Two groups of trials are especially pertinent to the forensic point of view: open trials determined by the use of non-contemporary spectrograms and clue words spoken in a fixed or in a random context. In fact, all real cases of forensic speaker identification would include these particular variables, regardless of the number of known speakers and examiners involved. The total error percentages yielded by these two groups were 14.35 percent and 18.26 percent respectively.

A break-down of these total percentages shows that approximately one-third of the errors of false identifications (errors A + C) and two-thirds were failures to recognize a match when it actually existed (error B). In summary, the percentages of false identifications observed in these forensic models were 4.22 percent and 6.43 percent for clue words in fixed and in random contexts, respectively. Percentages of failures to recognize a match when it actually existed were 10.13 percent and 11.83 percent respectively.

Another group of trials—closed trials, contemporary matching spectrograms and clue words spoken in isolation—produced findings germane to the goals of the study. Since these were essentially the variables tested by Kersta in 1962 and examined again by Tosi in 1968, their importance to the present study can be seen. The error percentage for this group was 0.51 percent, the minimum/lowest error percentage observed in the project, as expected. However, this group of trials does not fit any type of forensic model and has no direct application.

The upper limit of the range of errors was found

Figure 5.—First Cycle.—Pooled percentage errors from 6 groups of 972 closed trials and 6 groups of 1944 open trials.

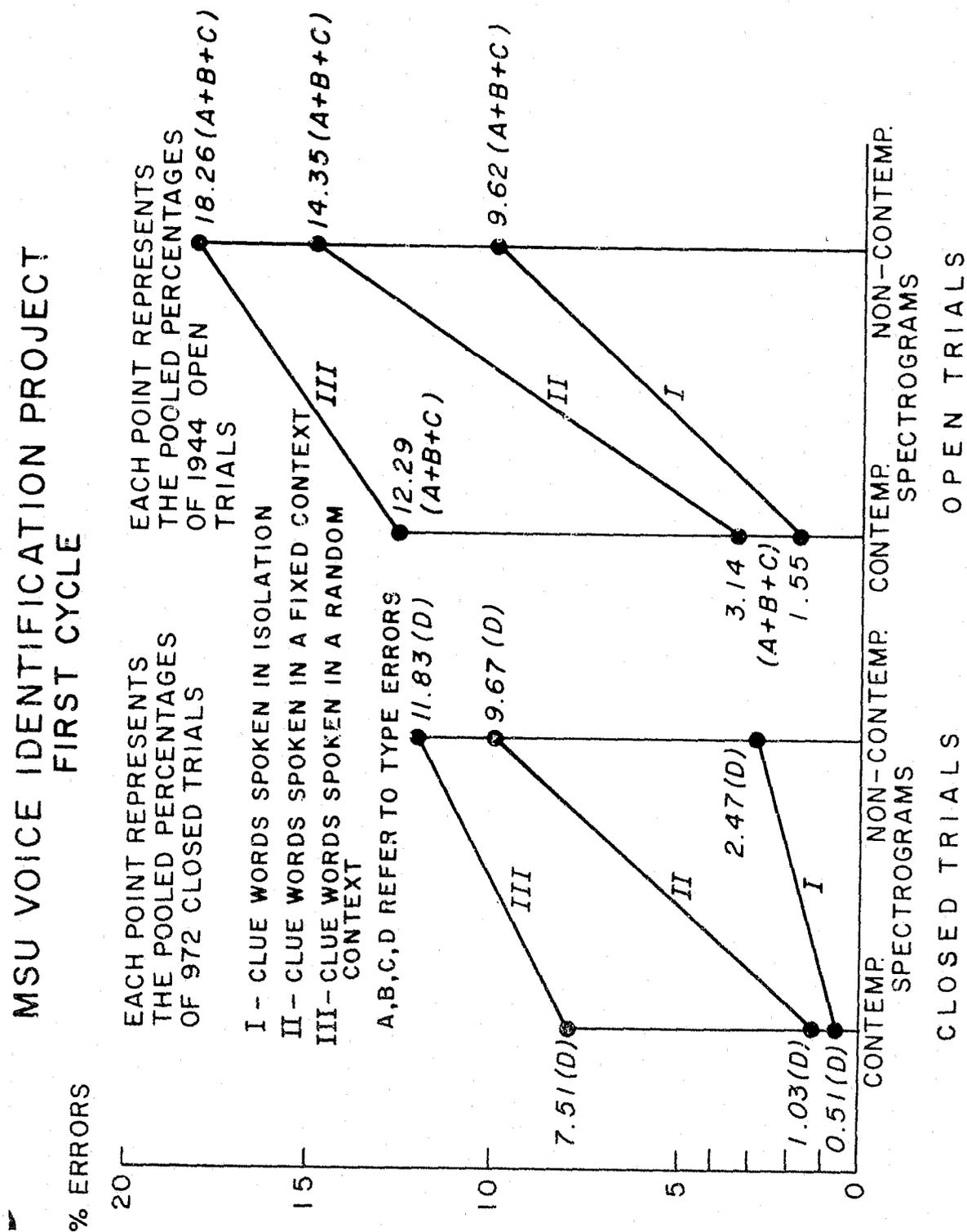


TABLE 3.—First Cycle—Pooled Error Percentages of Six Groups of 972 Closed Trials and Six Groups of 1944 Open Trials of Speaker Identification

| Context                                     | Closed Trials                      |  | Open Trials                             |  |
|---|------------------------------------|--|---|--|
|   | Contemporary Matching Spectrograms | Non-contemporary Matching Spectrograms | Contemporary Matching Spectrograms      | Non-contemporary Matching Spectrograms   |
| (I) Clue words spoken in isolation          | 0.51 (D)                           | 2.47 (D)                               | 1.55 (A+B+C)<br>0.36 (A+C)<br>1.19 (B)  | 9.62 (A+B+C)<br>2.37 (A+C)<br>7.25 (B)   |
| (II) Clue words spoken in a fixed context   | 1.03 (D)                           | 9.67 (D)                               | 3.11 (A+B+C)<br>1.49 (A+C)<br>1.65 (B)  | 14.35 (A+B+C)<br>1.22 (A+C)<br>10.13 (B) |
| (III) Clue words spoken in a random context | 7.51 (D)                           | 11.83 (D)                              | 12.29 (A+B+C)<br>4.01 (A+C)<br>8.28 (B) | 18.26 (A+B+C)<br>6.43 (A+C)<br>11.83 (B) |

Note.—A, B, C, D refers to the type of errors as described on page 15.

in another group of trials characterized by: open-match types exclusively, non-contemporary matching spectrograms, and clue words spoken in a random context. This group, which also does not fit any type of forensic model, yielded 29.01 percent

error. This extreme error percentage was composed of 5.35 percent of false identifications (error A) and 23.66 percent of failures to recognize an existing match (error B).

### III. Results From the Second Cycle of the Project

The second cycle of the project, using six clue words instead of nine, was undertaken after all 17,496 trials of identification for the first cycle were completed. The purpose of this replication was to obtain information concerning the effect a reduction in the number of clue words would/or would not have on the accuracy of the examiners. It should be pointed out that the major findings of the Voice Identification Project were those associated with the first cycle. The single purpose of the second cycle replication was to compare effects from both cycles.

Second cycle results may be somewhat biased since examiners gained experience during the preceding eight months, the time required to complete the first cycle. This experience or learning was in addition

to the training that all examiners received prior to the start of the study.

Comment concerning the performance of the examiners is relevant. Only the more motivated persons remained with the project long enough to complete the second cycle; some of those who quit considered the task extremely boring. The staff observed that many of the less motivated examiners did not perform well. These examiners tended to take an excessive number of rest periods and showed little concern for reaching the best possible decision in each trial, behavior which was viewed as hampering performance.

A procedure that would have eliminated the effect of an increasing experience of the examiners on the second cycle—testing randomly nine and six

clue words over the entire duration of the project—was not feasible. Such a design would have presented two logistical problems. First, the simultaneous use of segmented spectrograms containing six and three clue words would have complicated and hampered the results of the trials using nine clue words spoken in fixed and in random contexts. Second, predicting the time necessary for the completion of both cycles was hazardous. Therefore, since the nine clue words were considered more relevant to forensic models than the six clue words, an effort was made to secure first the complete performance of all trials using nine clue words, and to leave the six clue words trials as a replication.

Main conditions tested in the second cycle which produced significant differences between their various levels were: context, number of speakers, awareness of the examiners (closed and open trials) and time-elapsed (contemporary and non-contemporary spectrograms). These significant differences paralleled the results found for the first cycle in terms of comparisons within each main condition. Table

4 presents the results of an analysis of variance of the percentage errors yielded by each level within each of the main conditions that involved significant differences.

Table 5 presents the percentage of correct responses for the six main conditions tested in both first and second cycles, as well as the results of a statistical test of the differences. The test revealed significant differences between the two cycles, with  $p < .5$  percent, in the following instances: context (words spoken in isolation, 95.77 percent vs 93.83 percent); number of utterances (one utterance, 91.29 percent vs 89.71 percent), yielding the nine clue words cycle a larger percentage of correct responses. In overall conditions no significant difference was found between the two cycles. The examiners were correct 91.58 percent during the nine clue words cycle vs 91.24 percent during the six clue words cycle.

As was done with trials of the first cycle, trials of the second cycle were grouped according to the following characteristics:

TABLE 4.—Second Cycle—Results of an Analysis of Variance of the Correct Responses Produced under Each of the Main Conditions Tested

| Condition   | Pooled percentage of correct responses | Probability of the difference between levels, less than: |
|---|--|--|
| Number of utterances of the same clue word:       |  |  |
| 1 utterance                                       | 89.71                                  | n.s.   |
| 2 utterances                                      | 91.62                                  | n.s.   |
| 3 utterances                                      | 92.39                                  |  |
| Different types of recording transmissions:       |  |  |
| (a) directly into a tape recorder                 | 91.11                                  | n.s.   |
| (b) through a telephone line in quiet environment | 92.20                                  | n.s.   |
| (c) through a telephone line in noisy environment | 91.10                                  |  |
| Context of the clue words spoken                  |  |  |
| (I) in isolation                                  | 93.83                                  | 0.01   |
| (II) in a fixed context                           | 91.68                                  | 0.01   |
| (III) in a random context                         | 88.20                                  |  |
| Different number of "known" speakers              |  |  |
| 10 speakers                                       | 93.79                                  | 0.01   |
| 20 speakers                                       | 90.10                                  | n.s.   |
| 40 speakers                                       | 89.52                                  |  |
| Time-elapsed between recordings:                  |  |  |
| contemporary matching spectrograms                | 95.13                                  | 0.01   |
| non-contemporary matching spectrograms            | 87.35                                  |  |
| Awareness of examiners:                           |  |  |
| closed trials                                     | 91.31                                  | 0.01   |
| open trials                                       | 89.71                                  |  |

TABLE 5.—Summary of a Statistical Test of the Difference between Percentages of Correct Responses from Both First and Second Cycles, under the Main Conditions Tested

| Condition   | 1st Cycle | 2nd Cycle | Difference 1st Cycle-2nd Cycle | Probability of the diff. between 1st & 2nd cycle, less than: |
|---|-----------|-----------|--------------------------------|--|
| Number of utterances of the same clue word:       |           |           |                                |  |
| 1 utterance                                       | 91.29     | 89.71     | 1.58                           | 0.05   |
| 2 utterances                                      | 90.96     | 91.62     | -0.66                          | n.s.   |
| 3 utterances                                      | 92.19     | 92.39     | 0.10                           | n.s.   |
| Different types of recording transmissions:       |           |           |                                |  |
| (a) directly into a tape recorder                 | 92.42     | 91.41     | 1.01                           | n.s.   |
| (b) through a telephone line in quiet environment | 91.31     | 91.20     | 0.11                           | n.s.   |
| (c) through a telephone line in noisy environment | 91.02     | 91.10     | -0.08                          | n.s.   |
| Context of the clue words spoken:                 |           |           |                                |  |
| (I) in isolation                                  | 95.77     | 93.83     | 1.94                           | 0.05   |
| (II) in a fixed context                           | 92.39     | 91.68     | 0.71                           | n.s.   |
| (III) in a random context                         | 86.59     | 88.20     | -1.61                          | n.s.   |
| Different number of "known" speakers:             |           |           |                                |  |
| 10 speakers                                       | 93.30     | 93.79     | -0.49                          | n.s.   |
| 20 speakers                                       | 91.87     | 90.10     | 1.47                           | n.s.   |
| 40 speakers                                       | 89.58     | 89.52     | 0.06                           | n.s.   |
| Time-elapsed between recordings:                  |           |           |                                |  |
| contemporary matching spectrograms                | 95.21     | 95.13     | 0.08                           | n.s.   |
| non-contemporary matching spectrograms            | 87.95     | 87.35     | 0.60                           | n.s.   |
| Awareness of examiners:                           |           |           |                                |  |
| closed trials                                     | 94.18     | 94.31     | 0.17                           | n.s.   |
| open trials                                       | 90.14     | 89.71     | 0.43                           | n.s.   |

1. awareness of the examiners (closed and open trials).
2. time-elapsed (contemporary and non-contemporary spectrograms).
3. context (clue words spoken in isolation, in a fixed context and in a random context).

There were six groups of 972 closed trials each and six groups of 1944 open trials each, as described on page 50. Table 6 presents a comparison of the error percentages obtained from the first and second cycles for each of these 12 groups. Results of a statistical test of the differences between first and second cycle for each group are also shown in Table 6. Figure 6 displays graphically the data from Table 6. The errors from the second cycle were slightly larger in ten of these groups, but no significant differences were detected. The two remaining groups of open trials using non-contemporary spectrograms were significantly different in the first and

second cycle, with probability  $p < 5$  percent. Open trials using non-contemporary spectrograms and words spoken in isolation produced 13.23 percent error in the second cycle vs. 9.62 percent error in the first cycle. Open trials using non-contemporary spectrograms and clue words spoken in a random context produced 14.84 percent error in the second cycle vs. 18.26 percent error in the first cycle. Proportion of false identifications (errors A + C) and failures to recognize an existing match (error B) did not vary much in the two cycles. The improvement observed in the second cycle for the particular group of open trials which used non-contemporary spectrograms of clue words spoken in a random context could be explained on the basis of the learning process the examiners experienced during the first cycle. They assessed the open trials with non-contemporary spectrograms of clue words spoken in a random context as the most difficult tasks

Figure 6.—Second Cycle.—Pooled percentage errors from 6 groups of 972 closed trials and 6 groups of 1944 open trials.

### MSU VOICE IDENTIFICATION PROJECT SECOND CYCLE

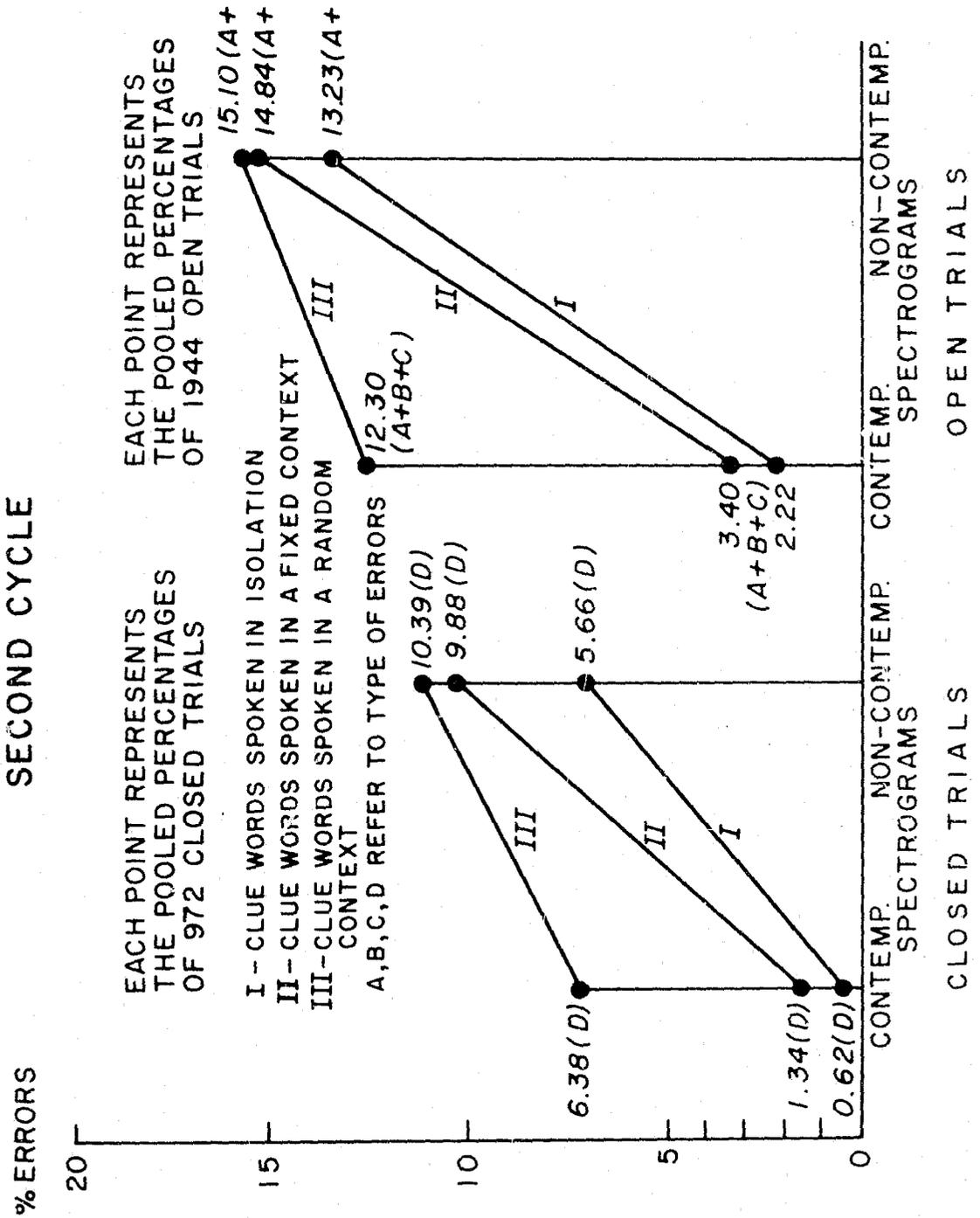


TABLE 6.—Summary of a Statistical Test of the Differences between Error Percentages from 1st and 2nd Cycles, of 12 Groups of Trials of Speaker Identification

| Type of trial condition | Time elapsed condition                 | Context condition                         | Error percentages from 1st cycle trials  | Error percentages from 2nd cycle trials  | Diff in error percentages between 1st and 2nd cycle | Probability of the diff between 1st & 2nd cycle, less than: |
|-------------------------|--|---|--|--|---|---|
| Closed Trials           | contemporary matching spectrograms     | (I) clue words spoken in isolation        | 0.51 (D)                                 | 0.62 (D)                                 | 0.11  | n.s.  |
|                         |  | (II) clue words spoken in fixed context   | 1.03 (D)                                 | 1.34 (D)                                 | 0.31  | n.s.  |
|                         |  | (III) clue words spoken in random context | 7.51 (D)                                 | 6.38 (D)                                 | -1.13   | n.s.  |
|                         | non-contemporary matching spectrograms | (I) clue words spoken in isolation        | 2.47 (D)                                 | 3.66 (D)                                 | 3.19  | n.s.  |
|                         |  | (II) clue words spoken in fixed context   | 9.67 (D)                                 | 9.88 (D)                                 | 0.21  | n.s.  |
|                         |  | (III) clue words spoken in random context | 11.83 (D)                                | 10.39 (D)                                | -1.44   | n.s.  |
| Open Trials             | contemporary matching spectrograms     | (I) clue words spoken in isolation        | 1.55 (A+B+C)<br>0.36 (A+C)<br>1.23 (B)   | 2.22 (A+B+C)<br>0.52 (A+C)<br>1.70 (B)   | 0.67<br>0.16<br>0.47                                | n.s.<br>n.s.<br>n.s.  |
|                         |  | (II) clue words spoken in fixed context   | 3.11 (A+B+C)<br>1.49 (A+C)<br>1.65 (B)   | 3.10 (A+B+C)<br>1.09 (A+C)<br>2.31 (B)   | 0.26<br>0.10<br>0.66                                | n.s.<br>n.s.<br>n.s.  |
|                         |  | (III) clue words spoken in random context | 12.29 (A+B+C)<br>4.01 (A+C)<br>8.28 (B)  | 12.30 (A+B+C)<br>1.96 (A+C)<br>10.31 (B) | -0.01<br>2.05<br>2.06                               | n.s.<br>0.05<br>0.05  |
|                         |  | (I) clue words spoken in isolation        | 9.62 (A+B+C)<br>2.37 (A+C)<br>7.25 (B)   | 13.23 (A+B+C)<br>1.22 (A+C)<br>9.01 (B)  | 3.61<br>1.85<br>1.76                                | 0.05<br>0.05<br>0.05  |
|                         | non-contemporary matching spectrograms | (I) clue words spoken in isolation        | 14.35 (A+B+C)<br>1.22 (A+C)<br>10.13 (B) | 11.81 (A+B+C)<br>1.27 (A+C)<br>12.68 (B) | 1.59<br>0.05<br>2.55                                | n.s.<br>n.s.<br>0.05  |
|                         |  | (II) clue words spoken in random context  | 18.26 (A+B+C)<br>6.43 (A+C)<br>11.83 (B) | 15.10 (A+B+C)<br>4.81 (A+C)<br>10.29 (B) | -3.16<br>1.62<br>1.54                               | 0.05<br>0.05<br>0.05  |

Note: A, B, C, and D refer to the type of error.

that produced the largest percentage of errors. The staff was aware that during the second cycle most of the examiners considered this particular type of trial as a challenge, devoting more time and special attention in the search for the correct answers. Besides, at this point of their training, the examiners were able to consider the extra clues offered by common phonemes included in the non-clue words which completed the random contexts.

In summary, the fact that results of the second

#### IV. Discussion and Conclusions

The results from the "Michigan State University Voice Identification Project" suggest that experienced examiners can identify or eliminate one unknown speaker from among as many as 40 known speakers, with little difference in accuracy being evidence in the use of nine or six clue words. The expected percentage of errors made by examiners who are forced to reach a positive decision in every trial of speaker identification they perform, (using exclusively visual examination of spectrograms), varies according to the conditions involved in each type of trial.

Closed trials, involving contemporary spectrograms of clue words spoken in isolation, yielded fewer than 1 percent error of false identifications. Since these conditions were essentially the ones employed by Kersta, it can be concluded that the present study has confirmed the figures reported by Kersta in 1962. In the 1968 Tosi's evaluation of "Voiceprinting," the error percentage reported for similar type of trials was approximately six percent. This discrepancy can be explained on the basis of individual differences among examiners. In fact, considering the performance of each examiner separately in that evaluation, the range of error percentage was 14 to 0 percent.

The second goal of the present study was to test forensic models that included the following variables:

- (a) random chance that the unknown speaker is or is not among the known ones ("open trials");
- (b) non-contemporary spectrograms (spectrograms of the unknown speaker obtained at a different time from the spectrograms of the known speakers);

- (c) same sentences uttered by known and unknown speakers ("fixed context" or different sen-

cycle did not differ substantially from those of the first cycle must not be solely interpreted as meaning that decreasing the number of available clue words from nine to six is not generally significant. The learning process the examiners experienced during the previous eight months devoted to the completion of the first cycle possibly interacted with the results of the second cycle, thus compensating for the fewer number of clue words available.

tences including the same clue words "random context.")

The error observed was approximately 15 percent for fixed context, of which approximately five percent were errors of false identifications (errors A + C) and approximately 10 percent were failures of recognizing a match when it actually existed (error B). For models including "random context," the total error was approximately 18 percent. This percentage was composed of approximately six percent of errors of false identifications and approximately 12 percent of failures of recognizing a match when it existed.

These findings suggest that if an experienced examiner, using only Visual Inspection of Spectrograms for legal purposes of identification and excluding any kind of listening, is forced to reach a positive decision in each case (devoting approximately 15 minutes to complete the task), his expected error range would be 14-18 percent. The probability that his wrong decisions will eliminate a guilty person is 75 percent of the total expected error. The probability that when in error this examiner will accuse an innocent person is 25 percent of the total expected error.

In summary: under the specified conditions the expected range of false identifications is 5-6 percent and the expected range of elimination of a guilty person is 10-12 percent.

Analysis of the ratings in the scale of self confidence used by the examiners in this project showed that approximately 60 percent of their wrong decisions were graded as "uncertain," with numbers 1 and 2. This finding suggests that the examiners' errors could have been reduced to approximately 40 percent of the observed figures, were these ex-

examiners not forced to reach a positive decision for the trials in which they felt uncertain.

Clearly, the reported errors apply to experimental trials in which the examiners used visual inspection of spectrograms exclusively, devoting an average of 15 minutes per trial in reaching a forced positive decision. It could be hypothesized that if in addition to visual comparisons of spectrograms the examiners would not have been forced to reach a decision when uncertain, and allowed to listen to the unknown and known voices, the errors might have been further reduced. The experiment per-

## V. Extension of Results From Forensic Models to Real Cases

A group of speech scientists (Bolt *et al.*, 1970) have expressed concern about the use of spectrographic evidence in court, before this method has been validated by controlled experimentation. The question arises: assuming that the results from the statistical forensic models studied in the present experiment could be applied toward such a validation, how would the conditions in practical legal cases differ from the conditions in the statistical models? In what way would these real conditions possibly alter the error expectancy disclosed by the models?

Main differences of conditions that could exist between models and real cases are as follows:

A. *Population of known voices.* In the models of the present study the number of known voices varied from 10 to 40, drafted from a closed catalog of 250 speakers, representing a statistical sample of a homogeneous population of 25,000 persons. In forensic cases, the catalog of known voices could theoretically include millions of samples, if the voice spectrogram of the criminal would be compared with filed voice spectrograms of the population of the world, or even the United States of America. Obviously, conclusions derived from an experimental study of a small population of speakers can not be extrapolated to populations of millions of individuals. However, this is not the case in the present practical situations that police must handle. In these cases the catalog of known voices is *open*, true, but *limited* to a few suspected persons. It seems reasonable to assume that the intra and interspeaker variabilities within such a reduced group of suspected persons would not differ substantially from the variabilities that existed

formed by Stevens *et al.* (1968), as well as the opinion of some phoneticians and linguists who feel that speaker recognition by listening is more accurate than by visual comparison of spectrograms, seem to confirm this hypothesis. A further study including forensic models, similar to the ones used in the present experiment might result in important additional information if trained examiners could both listen and make visual comparisons of spectrograms. Also, the present study should be complemented by the testing of disguised voices and non-contemporary spectrograms obtained from spans of time longer than one month.

within the highly homogeneous group of experimental speakers utilized in the present study. Therefore, it seems advisable to disregard size of the population of known voices as a differential characteristic that could hamper extrapolation of experimental results from the present study.

B. *Availability of time and responsibility of the examiners.* In the present study the examiners devoted an average of 15 minutes to reach a positive conclusion in each trial. Whether such a conclusion was the right or the wrong one, no effect could take place whatsoever over the examiner or the speaker. In forensic cases, the professional examiner normally may devote all the necessary time to reach a conclusion. He is aware of the consequences that a wrong decision could mean to his professional status as well as the consequences to the speaker whom he might erroneously identify. It seems reasonable to conclude that the differential characteristics between experimental and professional examiners might help to improve the accuracy of the professional examiners.

C. *Type of decisions examiners are urged to reach in each trial.* In the statistical model the examiners were forced to reach a positive conclusion in each trial, even if they were uncertain of the correct response. In real forensic cases, the professional examiner is permitted to make the following alternative decisions:<sup>(1)</sup>

- (a) Positive identification.
- (b) Positive elimination.

<sup>(1)</sup> These are the alternative decisions that Sgt. Nash, head of the Voice Identification section of the Michigan Department of State Police, is presently making.

(c) Possibility that the unknown speaker is one of the suspected persons, but more evidence is necessary in order to reach a positive identification.

(d) Possibility that the unknown speaker is none of the available suspected persons but more evidence is necessary to reach a positive elimination.

(e) Unable to reach any conclusion with the available voice samples.

These possibilities of alternative decisions could confer an extremely high reliability to the positive identifications or eliminations. The following information released by Sgt. Nash, of the Voice Identification Unit of the Michigan Department of State Police is cited as an illustration: From a total of 673 voice examinations, a positive identification was reached in 88 instances. Later on, most of the accused persons admitted culpability or were convicted by evidence other than that produced by their voices. In 172 cases, the conclusion was: "positive elimination." "Possibility of identification" or "elimination" was the conclusion in 31 other examinations. Finally, in 382 cases, the examiner concluded that he was "unable to reach any conclusion due to the lack of and/or poor voice samples."

D. *Availability of clues.* In the statistical models of this study, only spectrograms of nine and six clue words were available to the examiners for visual inspection. In real forensic cases the examiner must necessarily listen first to the unknown and known

voices while processing the spectrograms for visual comparison. The professional examiner is entitled to request as many samples as he deems necessary to reach a positive conclusion. Combination of methods of voice recognition by listening and by visual inspection of spectrograms can enhance the accuracy of his conclusion. Moreover, by using this combination the professional examiner can objectively sustain in court his opinion, by presenting the spectrographic similarities.

In conclusion, it is the opinion of the writer—based on his experience obtained through the performance of the present study and the observation of the practical work in the field done in the Voice Identification Unit of the Michigan Department of State Police—that the Federal Department of Justice should encourage the training of Voice Identifier Experts, who must be properly tested and certified prior to being recognized by the United States Courts as expert witnesses in the field.

Qualified personnel, the expert witnesses in the field, will continue to provide valuable service even if a satisfactory voice recognition machine is developed in the future. With a recognition machine available, the trained personnel would be demanded to prepare the necessary samples to feed the machine, to evaluate the results and to check the results of the machine by an alternative method, for instance the spectrographic one.

## REFERENCES

- Bolt, R. *et al.*, "Speaker Identification by Speech Spectrograms: A Scientist's View of its Reliability for Legal Purposes," *Journal of the Acoustical Society of America*, 47, 597 (1970).
- Hecker, M., "Speaker Recognition: An Interpretive Survey of the Literature," *ASHA Monographs*, no. 16 (1971).
- Kersta, L., "Voiceprint Identification," *Nature*, 196, No. 4861, 1253-1257 (1962).
- Potter, R., *et al.*, *Visible Speech*, Dover, New York (1966).
- Presti, A., "High Speed Sound Spectrograph," *Journal of the Acoustical Society of America*, 40, 628-634 (1966).
- Stevens, K., C. Williams, J. Carbonell and B. Woods "Speaker Authentication and Identification: A Comparison of Spectrographic and Auditory Presentations of Speech Material," *Journal of the Acoustical Society of America* 44, 1596-1607 (1968).
- Tosi, Oscar, "Evaluation of the Voiceprinting Method," Report to the Michigan Department of State Police (1967).
- Tosi, Oscar, "Speaker Identification through Acoustic Spectrography," *Comptes Rendus* of the XIV International Congress of Logopedics and Phoniatics, Paris, 138-141 (1968).

## PART 3

### Some Guidelines for the Use of Voiceprint Identification Techniques

By

Ralph F. Turner, M. S.  
Vern Rich, M. S.

Clarence Romig, M. S.  
James Hennessy, M. S.

SCHOOL OF CRIMINAL JUSTICE  
MICHIGAN STATE UNIVERSITY  
East Lansing, Michigan

## *CONTENTS*

|                                      | Page |
|--------------------------------------|------|
| I. Introduction .....                | 65   |
| II. Procedure .....                  | 65   |
| III. Training .....                  | 66   |
| IV. Summary of Experiment .....      | 67   |
| V. Interpretation of Results .....   | 68   |
| VI. Conclusions and Guidelines ..... | 68   |

## I. Introduction

Two experiments were conducted by the Michigan State University School of Criminal Justice. A complete account of this work is found in James Hennessy's paper, "An Analysis of Voiceprint Identification." This is an unpublished Master's Degree thesis on file in the Michigan State University Library.

As a result of this work and in accordance with the original project agreement, some guidelines are set forth which are believed to be important for law enforcement agencies who may be considering use of the voiceprint identification technique. These guidelines are based on observations of Dr. Tosi's experiment and the School of Criminal Justice experiment.

Personal identification through the use of voice spectrograms was introduced to the criminal justice arena in 1962 by Mr. Lawrence Kersta of the Bell Telephone Laboratories. The technique was used in limited, experimental fashion by law enforcement agencies, and since 1966 the record indicates that various courts have been asked to deal with this technique as an evidentiary problem. In addition to the legal questions pertaining to admissibility of this new kind of evidence, critics of the technique, primarily those having expertise in speech science, challenged the accuracy and validity of the method. It is not surprising, therefore, to find law enforcement practitioners, attorneys and judges confronted with a familiar problem. The one new element, however, is the relative speed with which this technique has achieved public interest, and, if not carefully controlled, may become a premature tool in the administration of justice's collection of scientific techniques.

The period of less than a decade is a relatively

short time for the Anglo-American system of justice to adopt a technique that, in the final analysis, may have an important bearing on the future life or liberty of an individual involved in the adjudication of a problem of law. The judicial acceptance of the technique of personal identification by means of fingerprint pattern recognition and comparison required several decades of testing, challenging, and research before it acquired its present status. The use of the polygraph in criminal investigation dates back to the early 1930's. Today the introduction of the results of such examinations are heard in a court of law only under the most unusual conditions and after many safeguards have been imposed. Breath testing to determine blood alcohol levels, a technique which also came on the scene in the middle 1930's, was quickly adopted by law enforcement agencies. The appellate record is again replete with accounts challenging the validity of this method. There are also numerous examples of gross mismanagement of this technique, resulting in miscarriages of justice in favor of both the guilty and the innocent.

These brief references to the evolution of technical methods used in law enforcement procedures and the subsequent introduction of the results in a court of law should clearly indicate the path which law enforcement practitioners must follow if the new technique of voiceprint identification is to become a useful method, serving the highest goals of the justice system. To this end the Michigan State University School of Criminal Justice has attempted to set forth some preliminary guidelines for consideration by law enforcement agencies contemplating the use of the voiceprint identification technique. These guidelines are based on available evidence at this time.

## II. Procedure

Two projects were conducted by Michigan State University during the period 1968-70. Both were

supported by a grant from the Law Enforcement Assistance Administration to the Michigan State

Police, who, in turn, contracted with (a) the Department of Audiology and Speech Science and (b) the School of Criminal Justice. The major project was carried out under the direction of Dr. Oscar Tosi of the speech science program at Michigan State University. The results of that work,

### III. Training

It has been said earlier that this report is to suggest some practical guidelines for law enforcement agencies contemplating the use of voiceprint identification techniques. If law enforcement practitioners are willing to acknowledge errors of the past, with particular reference to polygraph and breath testing techniques, and if they are seriously interested in avoiding a repetition of these errors, one observation becomes abundantly clear from the Michigan State University work. A proper training program is essential to the successful use of this technique. If this seems obvious or unimaginative to the reader, let him be reminded that the appellate court record is replete with far too many examples of scientific techniques being performed by law enforcement employees with inadequate training and education. Given the current image of the criminal justice system in the United States, responsible practitioners must insure the technical excellence and capabilities of those entrusted with the interpretation of this kind of evidence.

The original Kersta experiment<sup>1</sup> utilized high school students in large scale spectrogram identification problems. They were subjected to a training program devised by Kersta. Each of the eight identifiers, 16 to 17 year old high school girls, was given one week of training in voiceprint reading and detection of unique clues to be found in voiceprints. Kersta's work reports a total error of 1% in words taken from context.<sup>2</sup> Tosi, and his colleague Nash, of the Michigan State Police, received training in voice spectrogram recognition and identification from Kersta and are identified (for this report) as first generation trainees. Tosi, in his experiment, trained his group of identifiers in a manner similar to Kersta; however, there were certain modifications.

hereinafter referred to as the Tosi Report (and included in the Michigan State Police comprehensive report), provide the basis for some of the recommendations made in this report. Conclusions drawn from work done by the School of Criminal Justice are also included in this report.

Dr. Tosi's training consisted of:<sup>3</sup>

"(a) Lectures on phonetics, spectrography and a discussion of the variables to be tested in the experiment;

(b) Performance on closed trials of identification of one speaker among ten, by using words in isolation, contemporary matching spectrograms. These trials were first performed under direct supervision of the researcher; after a few days the examiner was left on his own. He was informed of his mistakes and allowed to compare the right matches with the wrong ones that he had selected. The examiner was instructed to make mainly subjective decisions, similar to those made when comparing photographs or handwritings. In the event that he felt uncertain about this judgment of spectrograms which appeared to be subjectively very similar, the examiner was instructed to consider the following objective points of similarity: (1) similar mean frequencies of vowel formants; (2) formants band-widths; (3) gaps and type of vertical striations; (4) slopes of formants; (5) characteristic patterns of fricatives and interformant energies. These similarities are often present in pairs of spectrograms of the same words produced by the same speaker at different times.

(c) After each examiner performed these closed trials, including contemporary matching spectrograms and words in isolation, with a success better than 96%, he was given other types of tasks that were increasingly more difficult. These "advanced" tasks included: open trials, non-contemporary matching spectrograms, and words in fixed and random contexts.

After one month of training the actual experi-

ment started. Spectrograms used for training were not used during the experiment."

It is noted that Tosi holds doctorate degrees in both physics and speech science. Nash is an experienced fingerprint identification expert. Tosi's second generation trainees included women ranging in age from 17 to 60 and a group of male undergraduate students. In one closed trial test involving 972 trials, 0.9% errors were reported. This result is similar to that obtained by Kersta and is mentioned merely to indicate that Tosi was able to replicate Kersta's early work.

### IV. Summary of Experiment

Following is a summary of two experiments conducted by the School of Criminal Justice. A full account is found in "An Analysis of Voiceprint Identification" by James J. Hennessy.<sup>4</sup> (1) Tape recordings were made in a dormitory reading room and lobby using portable equipment. The speakers included 12 male and 8 female graduate students. Six of the original speakers were rerecorded a week later. Spectrograms were prepared using the same instrument employed in the Tosi experiment. The identification trials were arranged to include thirty tasks. An analysis of the results showed that the two identifiers had an average of 70% accuracy in their identifications.

(2) A teller's window in the Cashier's Office of Michigan State University served as the site of the field recordings. Three tape recorders, two Wollensaks and one Uher, were placed in the teller's booth. The microphones were placed on a cardboard stand in the middle of the top of the counter area, facing directly outward at an upward angle of 45 degrees. The positions of the microphones were randomly changed four times to prevent any one microphone from being on one side or in the exact middle all the time. The recorders were run at a speed of 7½ inches per second. The recording levels were kept at a constant level.

The Uher was taken to a stand in an alcove of

The School of Criminal Justice project used two male identifiers. One was a senior in the School of Criminal Justice and a criminalist major; the other was a general law enforcement major. The first identifier chosen had already received training and had worked in the Audiology and Speech Sciences Department's voiceprint project. The second had received no training. He had seven years' experience in law enforcement<sup>5</sup> and was introduced to the voiceprint identification technique by attending some of Dr. Tosi's lectures and learning "matching techniques" from his colleagues.

the Cashier's Office. The Uher recorded the laboratory type samples of the known spectrograms.

There was a considerable amount of noise from typewriters, change machines, adding machines, and other people talking, entering, and leaving the Cashier's Office.

A Bruel and Kjaer Precision Sound Level Meter was utilized to measure the exact sound levels at the teller's window and in the alcove.<sup>6</sup> Readings were taken at two different times. The ranges of the sound levels for the teller's window were 64 to 80 db's (C scale) for the stand in the alcove, the range was 51 to 68 db's (C scale). Thirty-four db's is the sound level in a library; 54 db's is the sound level in a typical business office; 65 db's is the sound level of average conversational speech; 74 db's is the sound level of average street traffic; 88 db's is the sound level of the inside of a bus; and 94 db's is the average sound level inside a New York subway train.<sup>7</sup> As can be seen from these figures, the sound level of the alcove was substantially lower than the rather noisy level of the teller's window.

Standard Scotch Brand, 5 inch reel, magnetic tape (190 series) was used to record the speech of the volunteers.

It had been decided that an equal number of males and females, plus a random number of extra speakers, in case a recording was not good, should be obtained. In addition, permission had been obtained from the director of the Audiology and

<sup>1</sup> Lawrence G. Kersta, "Voiceprint Identification," *Nature*, 196 (Dec. 29, 1962), 1253-57.

<sup>2</sup> Kersta, "Voice Print Identification," *op. cit.*, 1253-57.

<sup>3</sup> *An Experiment on Voice Identification by Visual Inspection of Spectrograms*. Paper read by Dr. Tosi at meeting of Acoustical Society of America, Houston, Texas, November, 1970.

<sup>4</sup> Hennessy, James J. "An Analysis of Voiceprint Identification." Unpublished Master's Degree thesis. Michigan State University Library, 1970.

<sup>5</sup> Hennessy, "An Analysis of Voiceprint Identification," *op. cit.*

<sup>6</sup> *Instructions and Application of the Precision Sound Level Meter* (Naerum, Denmark: Bruel and Kjaer, 1965).

<sup>7</sup> *Ibid.*, p. 4

Speech Sciences Department's research project to state that the Audiology and Speech Sciences Department was the department conducting the research. A sign was posted on the stand in the alcove and on the side of the teller's window being used, stating that this was the voiceprint research project.

The two sentences used were:

Please give me my money; I want it.

My money is on the counter; please give it to me.

The procedures for the recordings were as follows: One technician was stationed in the teller's window. Another technician was stationed near the door of the Cashier's Office. As a possible speaker either came in or went out, the technician approached him and asked if he would like to participate in a research project. If the person accepted, a brief statement of the research goals and procedures was made. He was then led to the teller's window where the technician there recorded his name, age, nationality and state and local address, if he were or were not an American citizen. The

technician also recorded the sex of the speaker and gave him a speaker number. The speaker's instructions were given to him. He was to read, twice, the two sentences. No attempt was made to position him in front of the microphones. The speaker stated his speaker number and then the two sentences. When this was done, he was led to the alcove where he again stated his speaker number and repeated the two sentences twice. At this station, however, the speaker himself held the microphone approximately eight to twelve inches from his lips. The range of the distance of the speaker from the microphones in the teller's window was approximately 16 to 28 inches. Some speakers leaned on the counter, others stood back from it.

Eighty-four speakers in all were recorded, 42 males and 42 females. Most of the speakers were natives of lower Michigan. Almost all were undergraduates of 18 to 20 years of age.

The spectrograms used were open, contemporary, fixed context, and from one sex at a time. The average accuracy of identification by the two identifiers in this experiment was 59%.

## V. Interpretation of Results

An interpretation of the above described results suggest a negative experiment, clearly not compatible with those reported by Kersta and Tosi. The original goal, however, of the School of Criminal Justice involvement in the over-all voiceprint experiment was to produce, if possible, some guidelines which might be useful to law enforcement agencies. Viewed in this context, the experiment has produced some results which are quite important at this stage of development in the use of voiceprint identification techniques by criminal justice agencies. Reference has been made earlier to the difficulties that have plagued the widespread and ready acceptance of polygraph and breath testing techniques. While early users of these techniques were warned of the dangers of premature

use of the method without adequate testing and preparation, seldom have negative results been interpreted in a positive fashion to forestall avoidable errors.

While there may be some exceptions to the following general pattern, history and experience indicate that many forensic experts have acquired their expertise by the apprenticeship method. Fortunately, the majority of these learning methods have turned out favorably for the criminal justice system. The unfortunate experiences, however, constitute a blemish on the record which cannot be tolerated by a society whose foundation is based on a rule of law. It is for these reasons that the experiment should be viewed in a positive fashion.

## VI. Conclusions and Guidelines

In an effort to provide some guidelines for an interpretation of the present status of voiceprint

identification as it may be used in the criminal justice system, the following observations are made

1. Voiceprint identification techniques as they may be used by law enforcement agencies are relatively new. Eight years is a very short time to move from the introduction of a technique to expecting judicial acceptance of results obtained by this method.

2. The Tosi experiment, which is acknowledged to be a carefully controlled and important replication experiment, indicates the following results:

(a) In a closed trial experiment the percentage of error was 0.51%. Other experiments produced errors up to 29.1% depending on the conditions of the trials.

(b) In his "Extension of Results from Forensic Models to Real Cases," Tosi states that given the circumstances under which an actual case is investigated, a properly trained voice spectrogram identification expert can expect to achieve the same level of accuracy, i.e., 1% error.

3. The Tosi experiment indicates that second generation trainees can produce an acceptable level (1% error) of accuracy in their work.

4. The Criminal Justice experiment indicates that second generation trainees following an apprenticeship method of study, doing work under uncontrolled conditions, and not using equipment (other than the spectrograph) such as was used in the Tosi experiment, did not achieve acceptable (70% and 59% accuracy) results.

5. Proper training of identifiers is of utmost importance to the successful use of the voiceprint identification technique.

6. To this end, the following recommendations are made with regard to training and education requisites:

(a) Ideally, the voiceprint identification expert should hold a baccalaureate degree in either speech science or physical science. Forensic science laboratories today generally require baccalaureate degrees as a minimum educational prerequisite.

(b) While it has been demonstrated that acceptable second generation trainees can be recruited from a general population, law enforcement technicians with comparative identification expertise may be the preferred source of recruiting trainees.

(c) In the absence of a baccalaureate degree as suggested above, the following college level courses are strongly urged as a prerequisite to eventual use of the voiceprint identification technique: phonetics, acoustics (with the accompanying basic physics instruction), speech science, linguistics, audiology and basic electronics.

(d) Thorough training in the preparation of tape recordings and voice spectrograms is essential. The Tosi experiment demonstrates that proficiency in these techniques can be transmitted to trainees.

(e) A carefully supervised training program in voice spectrogram identification until the trainee reaches a 99% level of accuracy in closed trials working with spectrograms made from a homogeneous population is the ultimate goal.

(f) Upon satisfactory completion of a training program similar to what has been outlined above, the trainee should then undergo apprenticeship instruction with an experienced supervisor. This training period will utilize actual case evidence and the supervisor will indicate when he feels the student is qualified to render opinions based on his own observations.

## PART 4

### The Practical Application of Voice Identification in Criminal Investigations

Department of Michigan State Police  
East Lansing, Michigan

## CONTENTS

|  | Page |
|--|------|
| I. Introduction .....  | 75   |
| II. Methods .....  | 75   |
| III. Results .....   | 76   |
| A. Equipment recommendations .....                             | 76   |
| B. Educational programs for law enforcement .....              | 76   |
| C. Central voice identification file .....                     | 77   |
| D. Application of the voiceprint technique in real cases ..... | 77   |
| IV. Training of Voiceprint Examiners .....                     | 78   |
| V. A Look to the Future .....                                  | 79   |

## I. Introduction

In 1966, the Department of Michigan State Police became aware of a new identification technique that could be an aid to law enforcement. The communications media related work being performed by Lawrence Kersta using a sound spectrograph to identify recorded voices.

Voice spectrography was developed at Bell Company Laboratories by Potter, Kopp and Green. Mr. Lawrence Kersta, a former member of the Bell Telephone Company research staff, dealt for many years with voice spectrograms. He became interested in finding out whether speaker identification was possible and reliable on the basis of this type of spectrogram. Mr. Kersta reported in a convention of the Acoustical Society of America, in 1962, that after performing controlled experimentation, he concluded that voice spectrograms could be used as a reliable means of identification.

Mr. Kersta claimed that he accumulated evidence to support his conclusions by using a panel of twelve high school girls, whom he trained in voiceprint matching. They identified speakers among different sized speaker-utterance matrices taken from a population of 123 speakers. According to Kersta,

this panel made 99.75 percent correct identification of the speakers.

In 1966 Mr. Lawrence Kersta was contacted by the Michigan Department of State Police concerning the effectiveness of his system. Subsequently their Latent Fingerprint Technician, Detective Ernest Nash, met with Mr. Kersta at Voiceprint Laboratories to discuss the feasibility and adaptability of Voice Identification as an aid to law enforcement.

As a result of this meeting, Det. Ernest Nash and Det. Lewis Wilson received training at Voiceprint Laboratories in 1967. Both officers were experienced and expert Latent Identification technicians. Because of a desire for impartial consultation, Dr. Oscar Tosi of Michigan State University, who has as credentials a Doctorate in both Physics and Speech Science, was contracted to accompany the Technicians. At the conclusion of the training course, Dr. Tosi submitted a report which indicated the Kersta system to be significant. He suggested, however, that there was a need for further scientific study to replicate Mr. Kersta's work to further establish Voice Identification as a scientific method.

## II. Methods

As a result of this study, and to further implement the practical application of voice identification using a sound spectrograph, the Department of Michigan State Police took the following action:

1. Purchased equipment.
2. Initiated a program to educate regional law enforcement officers in the collection, preservation and applications of voice identification evidence.
3. Instituted a centralized Voice Identification file

by collecting voice recordings of known individuals.

4. Increased the experience and extended the expertise of technicians Nash and Wilson by actual case work and through continued association with Tosi and Kersta.

After three years of concentrated experience in the application of Voice Identification in criminal cases, of which two years were a part of this federally funded program, the following results are reported:

### III. Results

#### A. Equipment recommendations

The equipment needed for the collection, preservation and preparation of voice identification evidence is not extensive. However, experience gained during this research supplies some guidelines that should be helpful to any agency contemplating a voice identification program.

The most important piece of equipment is a sound spectrograph. Although there are other makes available, the Voiceprint Sound Spectrograph was specifically developed for voice identification and in our experience provides the most satisfactory results. There is more than one model voiceprint Spectrograph. All models perform equally well for voice identification, including some models developed to analyze sounds other than voice and used principally for medical research.

It is also necessary to have a device capable of recording speaker utterances. Such a recorder should be of the magnetic tape type, although experience indicates that recordings made on dictaphone belts can be used. The quality of the recording device is not usually critical, but must have the capacity to record the frequencies necessary for intelligible speech. However, evidence tapes are often made under less than ideal conditions and it seems reasonable to conclude that good reliable instrumentation will increase the likelihood of obtaining usable recordings. The tape used should be polyester backed and at least 1.0 mil thick. If a cassette recorder is employed, the user is cautioned against a cassette that records one hour on each side. This tape is so thin that it is too likely to break, foul the winding mechanism and fail to record the desired information.

Most of the criminal cases involving voice identification are the result of telephone conversations. Therefore it is of prime importance to have the ability to make quality recordings from the telephone.

The telephone pickup should be of the inductive type devised in such a manner that the recording will be made from the back of the ear piece. One model fits over the ear piece and thus eliminates the possibility of dislodging the pickup while handling the phone. An important advantage of the inductive type is that it does not require the ma-

nipulation of the telephone wires or recording equipment.

A variety of patch cords and other connectors are necessary so that recordings can be made from various model recorders and for making duplicate recordings.

A well equipped laboratory will have a bandpass filter. Some recordings have extraneous noise that interferes with the analysis. In many instances, this unwanted noise can be filtered out with a bandpass filter without damaging the information available from the speech signal. Although not critical to the operation, this equipment will make it possible to render a definite opinion in a greater percentage of cases.

A better service is rendered by a laboratory that has at its disposal several makes and models of tape recorders. This is especially important if examinations are being conducted for many different agencies where there is no control over the types of recorders used for obtaining the evidence recordings. Many small police agencies do not have recording equipment and it is important for them to be able to obtain such equipment on a temporary basis when confronted with an investigation requiring this ability.

A sound proof room should be available for making recordings in the laboratory and for listening to the sounds being analyzed. Because of the nature of some cases, obscene language has been recorded. A sound proof room will allow the study of these tapes without subjecting other employees to the content.

#### B. Educational programs for law enforcement

Before law enforcement agencies could be expected to submit evidence for voice identification, it was necessary to inform them as to how such evidence might be useful, what evidence was necessary in order to conduct an examination, and how best to obtain known and unknown tape recordings. It is estimated that over 4600 police officers received direct information in a classroom setting from technicians Nash and Wilson. At the same time, other citizens were familiarized with the voiceprint technique through service club appearances, radio and television.

As this type of information was completely new, some very basic procedures were disseminated. The following investigative hints were found to be practical:

1. The questioned and known voices should be recorded on the same tape, utilizing the same recorder, whenever possible.

2. Any instrument that will record on 1/4" tape can be used. However, a poor quality recording may interfere with identification or elimination.

3. The tape should be 1/4" with at least 1.0 mil of polyester or mylar backing. Tape with acetate backing should not be used.

4. A new or bulk erased tape should be employed for each case.

5. Tapes recorded at speeds slower than 17/8 i.p.s. do not usually contain sufficient frequency response for positive identification by voiceprints.

6. Enough tape should be available to record all anticipated conversation.

7. At the beginning of the evidence tape, pertinent data such as the date, time, location, telephone number, case number should be recorded.

8. If the victim is to record the incoming call, instructions should be given on recorder operation and elimination of background noise.

9. Telephone companies can be helpful in identifying the telephone number of the anonymous caller. In most instances, they can be prepared to identify the next incoming call within five minutes after being called to assist.

10. When recording the known voice, use a prepared text that contains the same words and phrases as the questioned recording.

11. On several occasions, officers obtained good recordings of the unknown voice but made poor recordings of the known voices. This was caused by the improper placement of the microphone and the failure to eliminate background noise.

As part of a program of instructions given at the Second Annual Criminal Advocacy Institute, technician Nash participated in a teaching program with Practising Law Institute of New York City. Instructions were held in New York City, N.Y.; Las Vegas, Nevada; Miami Beach, Florida and Dallas, Texas. The Institute was attended by prosecutors, trial lawyers and judges from throughout the United States. The style of presentation of voice identification information was through a moot trial conducted by experienced lawyers and judges.

Forensic scientists were informed about voice identification through speaking engagements at the

annual meeting of the American Academy of Forensic Sciences and the semi-annual meeting of Law Enforcement, Science and Technology. In addition, many people traveled to East Lansing to view first hand the work being performed in voice identification.

#### C. Central voice identification file

As experience was gained in the examination of voiceprints, it became apparent that plans to classify and file actual voiceprints were not practical. The voiceprint did not adapt readily to a classification system. A more practical method at this early stage of development was to store samples of the voice offenders on master tapes. The speaker is identified by name on the recording. The location of the voice on the recording is noted on the name card. At this writing, there has been little reference made to this file of voices. However, it will become more useful as the file grows. It is anticipated that the computer will eventually solve the problem of storing and retrieving voice identification information. Part of the research by Stanford Research Institute will be concerned with this possibility.

#### D. Application of the voiceprint technique in real cases

Since the inception of a voice identification program by the Department of Michigan State Police in 1967, 291 cases have been submitted to the Voice Identification Unit, mostly from Michigan police and fire departments. However, requests from all departments were honored and examinations were conducted for agencies in Indianapolis, Indiana; Riverside, California; Orlando, Florida; Los Angeles, California; St. Paul, Minnesota; Ladue, Missouri; Erie, Pennsylvania; Chicago, Illinois; Dade County, Florida; Astoria, Oregon; and South Miami, Florida. To indicate the wide variety of crimes that voice identification can become a part of, the 27 types of complaints received during this time are listed below:

| Type of Crime           | Number of cases |
|-------------------------|-----------------|
| Obscene telephone calls | 94              |
| False fire alarms       | 46              |
| Bomb Scares             | 28              |
| Threats                 | 26              |

| Type of Crime                         | Number of cases |
|---------------------------------------|-----------------|
| Nuisance telephone calls              | 25              |
| Extortion and blackmail               | 18              |
| Murder                                | 11              |
| Breaking and Entering                 | 8               |
| Kidnapping                            | 5               |
| Robbery                               | 5               |
| Rape                                  | 3               |
| Abortion                              | 2               |
| Attempted murder                      | 2               |
| Arson                                 | 2               |
| Bribery                               | 2               |
| Accosting and Soliciting              | 2               |
| Larceny                               | 2               |
| Fraud                                 | 1               |
| Weapon violation                      | 1               |
| Harrassing                            | 1               |
| Radical act                           | 1               |
| Impersonating a Police Officer        | 1               |
| False Report of a crime               | 1               |
| Annoying                              | 1               |
| Abduction                             | 1               |
| Gross Indecency                       | 1               |
| Intelligence Informant Identification | 1               |

673 voices were examined by the study of 42,432 spectrograms. 105 persons were identified as the

#### IV. Training of Voiceprint Examiners

The application of voice identification techniques in actual cases pre-supposes the use of examiners who are educated, well trained and experienced.

It is important that the education include an understanding of the speech and hearing process. Although it does not bear directly on the visual comparison of spectrograms, it does provide the examiner with a better understanding of differences that occur within separate utterances of the same word by one speaker. This will help him understand and explain when slight differences exist.

Listening to the recordings is also an important part of the identification process because the examiner must be assured that he is comparing the same sounds. In a training or research project where the examiner is presented with two prepared spectrograms to compare, both could be labeled "the". However, if in actuality one spectrogram was made from the sound "thee" and the other spectrogram was made from the sound "thuh", identification would be impossible. Knowledge of

unknown or questioned voice on tape recordings. 172 persons were eliminated as the unknown or questioned voice on tape recordings. For various reasons, a definite opinion could not be rendered concerning the other 396 persons.

It was not always possible to obtain information from the investigating officers that would refute or substantiate the opinions of the voice identification examiners. However, it was reported that in thirty cases, those persons identified by voice identification techniques later made confessions or admissions correlating voice identification opinions. No information was found to prove the wrong person had been identified by voice identification techniques.

From these experiences, it is concluded that Voice Identification by spectrographic analysis has a definite usefulness in the investigation of crime.

Given a sufficient quantity and quality of known and unknown voice recordings to work with, a qualified identification examiner can arrive at opinions that have an accuracy level comparable to other types of subjective examinations now made in Forensic Laboratories.

the various sounds of the spoken word, what sounds are germane to the identification process and how to listen to these sounds is necessary in the proper labeling of the spectrograms to be compared.

Basic training in theory of voice identification, the production of spectrograms and the comparison process are necessary in the early development of the voice identification examiner. However, this formal schooling does not sufficiently prepare an individual to undertake the responsibility of examining voice identification evidence and to give opinions in forensic cases. As in other forensic sciences that are subjective in nature, there must be experience and testing in the comparison of spectrograms until the examiner can demonstrate his ability to unerringly resolve the problems submitted to him. This does not preclude the fact that in some cases he may not be able to arrive at a definite opinion.

It has been demonstrated in the research by the Audiology and Speech Sciences Department of Michigan State University that voice identification by the visual comparison of spectrograms is possible.

The successful use of this method in forensic cases and in court, therefore, will ultimately depend on

the reliability of the trained and experienced examiner.

#### V. A Look to the Future

There are other research projects that should be initiated to extend the effectiveness of the voice spectrograph in criminal investigation. This would include experimentation with the identification of disguised voices and non-contemporary recordings. However, this should not deter its use by forensic laboratories or interfere with efforts to present voice identification testimony in court. In this respect, voice identification is no different than other forensic sciences in that there are always new questions to be answered.

Research is planned for speaker recognition by machine. This method could very well become an effective process to substantiate, extend or replace opinions now rendered by voice identification examiners using spectrographic techniques.

The possibility of using the spectrograph to identify sounds other than the human voice should not be overlooked. As an example, let us imagine

that an anonymous bomb threat is received and recorded. The sound of a motor can be distinguished as part of the background noise. If the motor noise, through sound spectrography, can be identified as to type, it might help investigators locate the source of the call. Again let us imagine that a woman calls the police and says she is about to be shot. An explosive sound ends the conversation. The sound spectrograph in this case may be effective in identifying the explosive noise as a firearm, perhaps a rifle rather than a pistol, and of large caliber.

As time passes, investigators and examiners alike will discover new applications of the sound spectrograph as it relates to criminal investigations. It remains now for more agencies and individuals to become involved in developing expertise and gaining experience in order that this relatively new technique can reach its full potential for solving crime.

## APPENDIX A

### Master Tables of Results from Trials

The word "set" used in these tables refers to a complete sequence of 486 combinations of different levels which define each type of task of speaker identification. The trials com-

pleting the four sets of each cycle of the experiment were performed in an unsystematic manner.

WORDS <sub>3</sub><sub>6</sub><sub>9</sub>    UTTERANCES <sub>1</sub><sub>2</sub><sub>3</sub>     CLOSED     OPEN/MATCH     OPEN/NO MATCH

|        |   | CONTEMPORARY |    |    |             |    |    |             |    |    | NON-CONTEMPORARY |    |    |             |    |    |             |    |    |
|--------|---|--------------|----|----|-------------|----|----|-------------|----|----|------------------|----|----|-------------|----|----|-------------|----|----|
|        |   | 10 SPEAKERS  |    |    | 20 SPEAKERS |    |    | 40 SPEAKERS |    |    | 10 SPEAKERS      |    |    | 20 SPEAKERS |    |    | 40 SPEAKERS |    |    |
| PANELS |   | P1           | P2 | P3 | P1          | P2 | P3 | P1          | P2 | P3 | P1               | P2 | P3 | P1          | P2 | P3 | P1          | P2 | P3 |
| I      | α | S1           | 00 | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  |
|        | β | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 00 | 1  | 1           | 00 | 1  |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  |
|        | γ | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  |
| II     | α | S1           | 1  | 1  | 1           | 1  | 00 | 1           | 00 | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 00 | 1           | 1  | 00 |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 00 | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  |
|        | β | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 00 | 1  | 1           | 1  | 1  |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  |
|        | γ | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 00          | 1  | 1  | 00          | 00 | 00 |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 00 | 1  | 1           | 1  | 1  |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 00 | 1  |
| III    | α | S1           | 00 | 1  | 1           | 1  | 00 | 1           | 00 | 1  | 1                | 1  | 1  | 1           | 00 | 1  | 1           | 1  | 1  |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 00          | 1  | 1  | 1                | 1  | 1  | 1           | 00 | 00 | 00          | 1  | 00 |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 00 | 1  | 00          | 00 | 1  |
|        | β | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 00 | 00 | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 00 |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 00 | 00 | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  |
|        | γ | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 00 | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 00 | 1           | 1  | 1  | 1           | 1  | 1  |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 00 | 00 | 1                | 1  | 1  | 1           | 00 | 1  | 1           | 1  | 1  |

SET I

WORDS <sub>3</sub><sub>6</sub><sub>9</sub>    UTTERANCES <sub>1</sub><sub>2</sub><sub>3</sub>     CLOSED     OPEN/MATCH     OPEN/NO MATCH

|        |   | CONTEMPORARY |    |    |             |    |    |             |    |    | NON-CONTEMPORARY |    |    |             |    |    |             |    |    |
|--------|---|--------------|----|----|-------------|----|----|-------------|----|----|------------------|----|----|-------------|----|----|-------------|----|----|
|        |   | 10 SPEAKERS  |    |    | 20 SPEAKERS |    |    | 40 SPEAKERS |    |    | 10 SPEAKERS      |    |    | 20 SPEAKERS |    |    | 40 SPEAKERS |    |    |
| PANELS |   | P1           | P2 | P3 | P1          | P2 | P3 | P1          | P2 | P3 | P1               | P2 | P3 | P1          | P2 | P3 | P1          | P2 | P3 |
| I      | α | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 00 | 00 | 1           | 1  | 1  | 1           | 1  | 1  |
|        |   | S3           | 1  | 1  | 00          | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 00          | 1  | 1  | 1           | 1  | 1  |
|        | β | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 00 |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 00 | 00 | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 00 | 1  |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 00          | 1  | 1  | 1           | 1  | 00 |
|        | γ | S1           | 1  | 1  | 1           | 1  | 00 | 00          | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 00          | 1  | 1  | 1           | 1  | 1  |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 00          | 1  | 1  |
| II     | α | S1           | 1  | 1  | 1           | 00 | 1  | 1           | 1  | 1  | 1                | 00 | 00 | 00          | 1  | 1  | 1           | 1  | 1  |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 00 | 00 | 1           | 1  | 1  | 1           | 00 | 00 |
|        |   | S3           | 1  | 1  | 1           | 00 | 1  | 1           | 1  | 1  | 1                | 00 | 1  | 1           | 1  | 1  | 1           | 1  | 00 |
|        | β | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 00 | 00 |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 00 | 1  |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 00 | 00 |
|        | γ | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 00 | 1  | 1           | 1  | 1  | 1           | 00 | 1  |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 00 | 00 | 1           | 1  | 1  | 1           | 1  | 00 |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 00          | 1  | 1  | 1           | 1  | 1  |
| III    | α | S1           | 1  | 1  | 1           | 1  | 00 | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  |
|        |   | S2           | 1  | 1  | 1           | 1  | 00 | 1           | 00 | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  |
|        |   | S3           | 1  | 1  | 00          | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  |
|        | β | S1           | 1  | 1  | 1           | 00 | 00 | 1           | 1  | 1  | 00               | 1  | 1  | 1           | 1  | 1  | 1           | 00 | 00 |
|        |   | S2           | 1  | 1  | 1           | 1  | 00 | 00          | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 00 | 1  |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 00 |
|        | γ | S1           | 1  | 00 | 1           | 1  | 1  | 1           | 1  | 00 | 00               | 00 | 00 | 1           | 1  | 1  | 1           | 00 | 1  |
|        |   | S2           | 00 | 1  | 1           | 1  | 1  | 00          | 00 | 00 | 00               | 00 | 1  | 1           | 1  | 1  | 1           | 00 | 1  |
|        |   | S3           | 1  | 1  | 00          | 00 | 1  | 1           | 1  | 1  | 00               | 1  | 1  | 1           | 1  | 1  | 1           | 00 | 00 |

SET I



WORDS <sub>3</sub><sub>6</sub><sub>9</sub> UTTERANCES <sub>1</sub><sub>2</sub><sub>3</sub>  CLOSED  OPEN/MATCH  OPEN/NO MATCH

|        |   | CONTEMPORARY |    |    |             |    |    |             |    |    | NON-CONTEMPORARY |    |    |             |    |    |             |    |    |   |   |
|--------|---|--------------|----|----|-------------|----|----|-------------|----|----|------------------|----|----|-------------|----|----|-------------|----|----|---|---|
|        |   | 10 SPEAKERS  |    |    | 20 SPEAKERS |    |    | 40 SPEAKERS |    |    | 10 SPEAKERS      |    |    | 20 SPEAKERS |    |    | 40 SPEAKERS |    |    |   |   |
| PANELS |   | P1           | P2 | P3 | P1          | P2 | P3 | P1          | P2 | P3 | P1               | P2 | P3 | P1          | P2 | P3 | P1          | P2 | P3 |   |   |
| I      | α | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1 |   |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1 | 1 |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1 | 1 |
|        | β | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1 | 1 |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1 | 1 |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1 | 1 |
|        | γ | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1 | 1 |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1 | 1 |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1 | 1 |
| II     | α | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1 |   |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1 | 1 |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1 | 1 |
|        | β | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1 | 1 |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1 | 1 |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1 | 1 |
|        | γ | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1 | 1 |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1 | 1 |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1 | 1 |
| III    | α | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1 |   |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1 | 1 |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1 | 1 |
|        | β | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1 | 1 |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1 | 1 |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1 | 1 |
|        | γ | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1 | 1 |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1 | 1 |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1 | 1 |

SET III

WORDS <sub>3</sub><sub>6</sub><sub>9</sub> UTTERANCES <sub>1</sub><sub>2</sub><sub>3</sub>  CLOSED  OPEN/MATCH  OPEN/NO MATCH

|        |   | CONTEMPORARY |    |    |             |    |    |             |    |    | NON-CONTEMPORARY |    |    |             |    |    |             |    |    |    |    |
|--------|---|--------------|----|----|-------------|----|----|-------------|----|----|------------------|----|----|-------------|----|----|-------------|----|----|----|----|
|        |   | 10 SPEAKERS  |    |    | 20 SPEAKERS |    |    | 40 SPEAKERS |    |    | 10 SPEAKERS      |    |    | 20 SPEAKERS |    |    | 40 SPEAKERS |    |    |    |    |
| PANELS |   | P1           | P2 | P3 | P1          | P2 | P3 | P1          | P2 | P3 | P1               | P2 | P3 | P1          | P2 | P3 | P1          | P2 | P3 |    |    |
| I      | α | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  |    |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  | 1  |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  | 1  |
|        | β | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  | 1  |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  | 1  |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  | 1  |
|        | γ | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  | 1  |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  | 1  |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  | 1  |
| II     | α | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  |    |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  | 1  |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  | 1  |
|        | β | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  | 1  |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  | 1  |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  | 1  |
|        | γ | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  | 1  |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  | 1  |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  | 1  |
| III    | α | S1           | OB | 1  | OB          | 1  | 1  | OB          | 1  | 1  | OB               | 1  | 1  | OB          | 1  | 1  | OB          | 1  | 1  | 1  | 1  |
|        |   | S2           | 1  | 1  | OB          | 1  | 1  | OB          | OB | 1  | OB               | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  | 1  |
|        |   | S3           | OB | 1  | 1           | OB | 1  | 1           | OB | 1  | OB               | OB | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  | 1  |
|        | β | S1           | 1  | 1  | 1           | 1  | 1  | 1           | OB | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  | 1  |
|        |   | S2           | 1  | 1  | 1           | OB | OB | OB          | OB | OB | 1                | 1  | 1  | OB          | 1  | 1  | OB          | 1  | 1  | 1  | 1  |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | OB | 1  | 1                | 1  | 1  | OB          | OB | OB | OB          | OB | OB | OB | OB |
|        | γ | S1           | 1  | OB | 1           | OB | OB | OB          | 1  | OB | OB               | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  | 1  |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | OB | OB | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  | 1  |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | OB | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  | 1  |

SET III

CONTINUED

1 OF 2

VOICE IDENTIFICATION RESEARCH

NCJ 00481

ANON

72

147D



































WORDS <sub>3</sub> <sub>6</sub> <sub>9</sub>    UTTERANCES <sub>1</sub> <sub>2</sub> <sub>3</sub>     CLOSED     OPEN/MATCH     OPEN/NO MATCH

|        |   | CONTEMPORARY |    |    |             |    |    |             |    |    | NON-CONTEMPORARY |    |    |             |    |    |             |    |    |    |
|--------|---|--------------|----|----|-------------|----|----|-------------|----|----|------------------|----|----|-------------|----|----|-------------|----|----|----|
|        |   | 10 SPEAKERS  |    |    | 20 SPEAKERS |    |    | 40 SPEAKERS |    |    | 10 SPEAKERS      |    |    | 20 SPEAKERS |    |    | 40 SPEAKERS |    |    |    |
| PANELS |   | P1           | P2 | P3 | P1          | P2 | P3 | P1          | P2 | P3 | P1               | P2 | P3 | P1          | P2 | P3 | P1          | P2 | P3 |    |
| I      | α | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | OB               | OB | 1  | 1           | 1  | 1  | 1           | OB | 1  |    |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | OB |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | OB | 1  | OB          | 1  | 1  | 1           | OB | 1  | 1  |
|        | β | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | OB          | 1  | 1  | 1           | OB | 1  | OB |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | OB |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | OB | 1           | 1  | 1  | 1           | OB | 1  | OB |
|        | γ | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | OB |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | OB | 1  | 1           | 1  | 1  | 1           | 1  | OB | OB |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | OB | OB | OB |
| II     | α | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | OB | OB          | 1  | OB |    |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | OB          | 1  | 1  | 1                | 1  | 1  | 1           | OB | 1  | 1           | OB | OB |    |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | OB          | OB | 1  | OB |
|        | β | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | OB | 1           | OB | 1  | 1           | 1  | 1  | 1  |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | OB | 1           | 1  | 1  | 1  |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | OB          | 1  | 1  | 1  |
|        | γ | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | OB | OB | OB |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | OB | 1  | OB |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | OB | 1           | 1  | 1  | 1           | OB | OB | 1  |
| III    | α | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | OB               | OB | 1  | OB          | 1  | 1  | OB          | 1  | OB |    |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | OB | 1           | 1  | 1  | 1           | 1  | OB |    |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | OB | 1  | OB               | 1  | 1  | OB          | 1  | 1  | 1           | 1  | 1  | 1  |
|        | β | S1           | OB | OB | OB          | 1  | 1  | OB          | OB | 1  | 1                | 1  | OB | 1           | 1  | 1  | 1           | 1  | OB | 1  |
|        |   | S2           | OB | OB | OB          | 1  | 1  | 1           | 1  | OB | OB               | 1  | 1  | 1           | 1  | 1  | 1           | 1  | OB | 1  |
|        |   | S3           | OB | OB | OB          | OB | 1  | OB          | OB | OB | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  |
|        | γ | S1           | 1  | OB | 1           | 1  | 1  | 1           | 1  | 1  | OB               | OB | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | OB | 1           | 1  | 1  | 1  |
|        |   | S3           | 1  | OB | 1           | 1  | 1  | 1           | OB | 1  | OB               | 1  | 1  | OB          | 1  | 1  | 1           | 1  | 1  | 1  |

SET IV

# APPENDIX B

## An Examination of Conditional Variations for Voice Identification Trials

Statistical Report No. 1

By

William B. Lashbrook, Ph.D.

VOICE IDENTIFICATION PROJECT

Michigan State University

WORDS <sub>3</sub> <sub>6</sub> <sub>9</sub>    UTTERANCES <sub>1</sub> <sub>2</sub> <sub>3</sub>     CLOSED     OPEN/MATCH     OPEN/NO MATCH

|        |   | CONTEMPORARY |    |    |             |    |    |             |    |    | NON-CONTEMPORARY |    |    |             |    |    |             |    |    |    |
|--------|---|--------------|----|----|-------------|----|----|-------------|----|----|------------------|----|----|-------------|----|----|-------------|----|----|----|
|        |   | 10 SPEAKERS  |    |    | 20 SPEAKERS |    |    | 40 SPEAKERS |    |    | 10 SPEAKERS      |    |    | 20 SPEAKERS |    |    | 40 SPEAKERS |    |    |    |
| PANELS |   | P1           | P2 | P3 | P1          | P2 | P3 | P1          | P2 | P3 | P1               | P2 | P3 | P1          | P2 | P3 | P1          | P2 | P3 |    |
| I      | α | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | OC               | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  |    |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | OC | OC | 1           | 1  | 1  | 1           | 1  | 1  | 1  |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  |
|        | β | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | OC          | 1  | 1  | OC          | 1  | 1  | 1  |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  |
|        | γ | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | OC | 1  | 1  |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | OC | 1  |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | OC |
| II     | α | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  |    |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  |    |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | OC | 1           | OC | 1  | 1           | 1  | 1  | 1  |
|        | β | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | OC | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  |
|        | γ | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  |
|        |   | S2           | 1  | OC | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  |
| III    | α | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  |    |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | OC | 1  | 1           | 1  | 1  |    |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  |
|        | β | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | OC | 1  | 1  |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  |
|        | γ | S1           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  |
|        |   | S2           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | OC | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1  |
|        |   | S3           | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | 1                | 1  | 1  | 1           | 1  | 1  | 1           | 1  | 1  | OC |

SET IV

## Introduction

The purpose of this report is to examine from the statistical point of view the nature of the differences produced in terms of the number of correct speaker identifications that can be attributed to variations in conditions under which data were obtained. Data from the project were scores ranging from zero to nine and representing the correct responses of nine subpanels of examiners, under each condition of the experiment. There were six conditional variables under examination for the voice identification project:

### A. Identification Trials Involving Variations in Transmission—Three Levels:

- a) Spectrograms taken from speakers recorded directly from a microphone.
- β) Spectrograms taken from speakers recorded from a telephone line through a microphone.
- γ) Spectrograms taken from speakers recorded from a telephone line through a microphone with background noise in the transmission.

### B. Identification Trials Involving Variations in Context—Three Levels:

- I) Spectrograms of clue words spoken in isolation.
- II) Spectrograms of clue words spoken in a fixed context.
- III) Spectrograms of clue words spoken in a random context.

### C. Number of "Known" Speakers Involved in an Identification Trial—Three Levels:

- 1) Identification Trials involving 10 speakers.
- 2) Identification Trials involving 20 speakers.
- 3) Identification Trials involving 40 speakers.

### D. Identification Trials Involving System Awareness/Non-Awareness of a Match—Three Levels:

- 1) An identification trial in which there was always a match and the identifiers were aware of the fact. (*Closed-Match*).
- 2) An identification trial in which there was a match and the identifiers were unaware of the fact. (*Open-Match*).

- 3) An identification trial in which there was no match and the identifiers were unaware of the fact. (*Open-No Match*).

### E. Identification Trials Involving Time-Elapsed Variations—Two Levels:

- 1) *Contemporary* "matching" spectrograms assumed to correspond to "unknown" speakers, taken at the same recording session that the "known" speakers' spectrograms were obtained.
- 2) *Non-contemporary* "matching" spectrograms assumed to correspond to "unknown" speakers, taken at a second recording session one month after the first, in which the "known" speakers' spectrograms were obtained.

### F. Identification Trials Involving Variations in the Number of Utterances or Examples of the Same Clue Words—Three Levels:

- 1) *One utterance* of the same words.
- 2) *Two utterances* of the same words.
- 3) *Three utterances* of the same words.

### G. Replication of All Identification Trials.—There were four replications (G) of the project identification trials.

It should be recalled that conditions A, B, and C constitute the Latin matrix aspects of the design under which data for the project were collected. That is, it took each of the nine subpanels (classed by type and size) 27 identification trials in order to complete a matrix. Eighteen matrices were then completed under conditions D, E, and F, producing a total of 486 different types of trials. This entire procedure was replicated (G) four times. Thus, the total number of trials that composed each of the two cycles of the project was 17,496. The different types of trials were performed randomly. Again data of the project were scores from zero to nine representing the correct responses for all nine examiners subpanels under each combination of the given conditions of the project. For the purpose of the statistical analysis that follows, differences within examiners subpanels were cancelled in favor of a more direct test of the effects of conditions on the identification trials themselves. On page 141

of this statistical study a detailed analysis of possible differences between examiners subpanels is reported.

### Statistical design

The basic approach to the statistical analysis involved in the project was to define the replication factor as the number of observations (four) of the nine subpanels under each of the combinations of conditional variables A through F. Thus, it should be readily apparent that the model would have to account for the fact that conditional variables D, E, and F involved repeated measures. In general, the model used for analysis approximated six-way analysis of variance with repeated measures on the last three variables. Variations of conditional variables A, B, and C constituted the parameters of the Latin matrix used in the statistical design of the project. Each of the nine subpanels performed 27 identification trials to complete all possible combinations of these three variable levels. The over-all error effect was divided into eight components, from  $E_1$  to  $E_8$ . That is, a test of the main effects of conditional variables A, B, and C and their interactions would have one unique error factor ( $E_1$ ). The interaction of these variables with D, E, and F (the repeated measure factors) as well as the main effects of D, E, and F would each have a unique error factor ( $E_2$ ). Finally all two and three way interactions of D, E, and F with A, B, and C would each have appropriate error factors ( $E_3, E_4, E_5, E_6, E_7,$  and  $E_8$ ) depending upon the number of repeated measures involved. Table 1 represents the statistical design used to analyze the data statistically.

TABLE 1.—Statistical Design for the Analysis of Voice Identification Data

| Source of variance   | Degree of freedom | Calculated DF |
|--|-------------------|---------------|
| A  | a-1               | 2             |
| B  | b-1               | 2             |
| C  | c-1               | 2             |
| AB   | (a-1) (b-1)       | 4             |
| AC   | (a-1) (c-1)       | 4             |
| BC   | (b-1) (c-1)       | 4             |
| ABC  | (a-1) (b-1) (c-1) | 8             |
| $E_1 = \text{Error}_1 = \text{AG} + \text{BG} + \text{CG} + \text{ABG} + \text{ACG} + \text{BCG} + \text{ABCG} + \text{G}$ | abc (g-1)         | 81            |

Table 1—Continued

| Source of variance   | Degree of freedom       | Calculated DF |
|--|-------------------------|---------------|
| D  | d-1                     | 2             |
| AD   | (a-1) (d-1)             | 4             |
| BD   | (b-1) (d-1)             | 4             |
| CD   | (c-1) (d-1)             | 4             |
| ABD  | (a-1) (b-1) (d-1)       | 8             |
| ACD  | (a-1) (c-1) (d-1)       | 8             |
| BCD  | (b-1) (c-1) (d-1)       | 8             |
| ABCD   | (a-1) (b-1) (c-1) (d-1) | 16            |
| $E_2 = \text{Error}_2 = \text{DG} + \text{ADG} + \text{BDG} + \text{CDG} + \text{ABDG} + \text{ACDG} + \text{BCDG} + \text{ABCDG}$ | abc (d-1) (g-1)         | 162           |
| E  | e-1                     | 1             |
| AE   | (a-1) (e-1)             | 2             |
| BE   | (b-1) (e-1)             | 2             |
| CE   | (c-1) (e-1)             | 2             |
| ABE  | (a-1) (b-1) (e-1)       | 4             |
| ACE  | (a-1) (c-1) (e-1)       | 4             |
| BCE  | (b-1) (c-1) (e-1)       | 4             |
| ABCE   | (a-1) (b-1) (c-1) (e-1) | 8             |
| $E_3 = \text{Error}_3 = \text{EG} + \text{AEG} + \text{BEG} + \text{CEG} + \text{ABEG} + \text{ACEG} + \text{BCEG} + \text{ABCEG}$ | abc (c-1) (g-1)         | 81            |
| F  | f-1                     | 2             |
| AF   | (a-1) (f-1)             | 4             |
| BF   | (b-1) (f-1)             | 4             |
| CF   | (c-1) (f-1)             | 4             |
| ABF  | (a-1) (b-1) (f-1)       | 8             |
| ACF  | (a-1) (c-1) (f-1)       | 8             |
| BCF  | (b-1) (c-1) (f-1)       | 8             |
| ABCF   | (a-1) (b-1) (c-1) (f-1) | 16            |
| $E_4 = \text{Error}_4 = \text{FG} + \text{AFG} + \text{BFG} + \text{CFG} + \text{ABFG} + \text{ACFG} + \text{BCFG} + \text{ABCFG}$ | abc (f-1) (g-1)         | 116           |
| DE   | (d-1) (e-1)             | 2             |
| ADE  | (a-1) (d-1) (e-1)       | 4             |
| BDE  | (b-1) (d-1) (e-1)       | 4             |
| CDE  | (c-1) (d-1) (e-1)       | 4             |
| ABDE   | (a-1) (b-1) (d-1) (e-1) | 8             |
| ACDE   | (a-1) (c-1) (d-1) (e-1) | 8             |
| BCDE   | (b-1) (c-1) (d-1) (e-1) | 8             |

Table 1—Continued

| Source of variance   | Degree of freedom             | Calculated DF |
|--|-------------------------------|---------------|
| ABCDE  | (a-1) (b-1) (c-1) (d-1) (e-1) | 16            |
| $E_5 = \text{Error}_5 = \text{DEG} + \text{ADEG} + \text{BDEG} + \text{CDEG} + \text{ABDEG} + \text{ACDEG} + \text{BCDEG} + \text{ABCDEG}$ | abc (d-1) (e-1) (g-1)         | 162           |
| DF   | (d-1) (f-1)                   | 4             |
| ADF  | (a-1) (d-1) (f-1)             | 8             |
| BDF  | (b-1) (d-1) (f-1)             | 8             |
| CDF  | (c-1) (d-1) (f-1)             | 16            |
| ABDF   | (a-1) (b-1) (d-1) (f-1)       | 16            |
| ACDF   | (a-1) (c-1) (d-1) (f-1)       | 16            |
| BCDF   | (b-1) (c-1) (d-1) (f-1)       | 16            |
| ABCDF  | (a-1) (b-1) (c-1) (d-1) (f-1) | 32            |
| $E_6 = \text{Error}_6 = \text{DFG} + \text{ADFG} + \text{BDFG} + \text{CDFG} + \text{ABDFG} + \text{ACDFG} + \text{BCDFG} + \text{ABCDFG}$ | abc (d-1) (f-1) (g-1)         | 324           |
| EF   | (e-1) (f-1)                   | 2             |
| AEF  | (a-1) (e-1) (f-1)             | 4             |
| BEF  | (b-1) (e-1) (f-1)             | 4             |
| CEF  | (c-1) (e-1) (f-1)             | 4             |
| ABEF   | (a-1) (b-1) (e-1) (f-1)       | 8             |
| ACEF   | (a-1) (c-1) (e-1) (f-1)       | 8             |
| BCEF   | (b-1) (c-1) (e-1) (f-1)       | 8             |
| ABCEF  | (a-1) (b-1) (c-1) (e-1) (f-1) | 16            |
| $E_7 = \text{Error}_7 = \text{EFG} + \text{AEFG} + \text{BEFG} + \text{CEFG} + \text{ABEFG} + \text{ACEFG} + \text{BCEFG} + \text{ABCEFG}$ | abc (e-1) (f-1) (g-1)         | 162           |
| DEF  | (d-1) (e-1) (f-1)             | 4             |
| ADEF   | (a-1) (d-1) (e-1) (f-1)       | 8             |
| BDEF   | (b-1) (d-1) (e-1) (f-1)       | 8             |

Table 1—Continued

| Source of variance   | Degree of freedom                   | Calculated DF |
|--|-------------------------------------|---------------|
| CDEF   | (c-1) (d-1) (e-1) (f-1)             | 8             |
| ABDEF  | (a-1) (b-1) (d-1) (e-1) (f-1)       | 16            |
| ACDEF  | (a-1) (c-1) (d-1) (e-1) (f-1)       | 16            |
| BCDEF  | (b-1) (c-1) (d-1) (e-1) (f-1)       | 16            |
| ABCDEF   | (a-1) (b-1) (c-1) (d-1) (e-1) (f-1) | 32            |
| $E_8 = \text{Error}_8 = \text{DEFG} + \text{ADEFG} + \text{BDEFG} + \text{CDEFG} + \text{ABDEFG} + \text{ACDEFG} + \text{BCDEFG} + \text{ABCDEFG}$ | abc (d-1) (e-1) (f-1) (g-1)         | 324           |
| Total  | abcdefg-1                           | 1,943         |

It will be noted that the between observation (G) effect is relegated to the first of the eight error terms. This procedure is consistent with repeated measure designs. It does not, however, represent a direct test of whether or not there was a significant difference between observations. In order to accomplish the later step the mean sum effects of (G) should be compared with the total sum of squares for (G) interactions with all the remaining conditional variables combinations.

### Results

Table 2 represents the results of the statistical design described in the preceding section, from data of the first cycle of the project. Data used in the analysis was transformed via the square root transformation recommended by Winer (1962) for this type of statistical design:

$$x^1 = \sqrt{x} + \sqrt{x+1}$$

where  $x^1 =$  transformed score

$$x = \text{score}$$

TABLE 2.—Analysis of Variance for Voice Identification Data from the First Cycle

| Source of variance  | df  | Mean square | F-ratio  | Sig. |
|---------------------|-----|-------------|----------|------|
| A-Transmission      | 2   | 0.53474     | 2.1973   |      |
| B-Context           | 2   | 16.05957    | 65.9898  | **   |
| C-No. of Speakers   | 2   | 2.46235     | 10.1180  | **   |
| AB                  | 4   | 0.04041     | 0.1661   |      |
| AC                  | 4   | 0.30651     | 1.2595   |      |
| BC                  | 4   | 0.58995     | 2.4242   |      |
| ABC                 | 8   | 0.15984     | 0.6568   |      |
| Error <sub>1</sub>  | 81  | 0.24336     |          |      |
| D-System Awareness  | 2   | 34.13054    | 114.3056 | **   |
| AD                  | 4   | 0.67458     | 2.2592   |      |
| BD                  | 4   | 2.75975     | 9.2426   | **   |
| CD                  | 4   | 0.54314     | 1.8190   |      |
| ABD                 | 8   | 0.05684     | 0.1904   |      |
| ACD                 | 8   | 0.27265     | 0.9131   |      |
| BCD                 | 8   | 0.07094     | 0.2376   |      |
| ABCD                | 16  | 0.29506     | 0.9882   |      |
| Error <sub>2</sub>  | 162 | 0.29859     |          |      |
| E-Time-Elapsed      | 1   | 31.96546    | 115.4176 | **   |
| AE                  | 2   | 0.97337     | 3.5145   | *    |
| BE                  | 2   | 1.38976     | 5.0180   | **   |
| CE                  | 2   | 0.08860     | 0.3199   |      |
| ABE                 | 4   | 0.15437     | 1.5574   |      |
| ACE                 | 4   | 0.16366     | 0.5909   |      |
| BCE                 | 4   | 0.02578     | 0.0931   |      |
| ABCE                | 8   | 0.08424     | 0.3042   |      |
| Error <sub>3</sub>  | 81  | 0.27695     |          |      |
| F-No. of Utterances | 2   | 0.53388     | 2.2303   |      |
| AF                  | 4   | 0.28702     | 1.1991   |      |
| BF                  | 4   | 0.41008     | 1.7131   | *    |
| CF                  | 4   | 0.66034     | 2.7586   | *    |
| ABF                 | 8   | 0.27043     | 1.1297   |      |
| ACF                 | 8   | 0.37140     | 1.5512   |      |
| BCF                 | 8   | 0.40998     | 1.7127   |      |
| ABCF                | 16  | 0.36536     | 1.5263   |      |
| Error <sub>4</sub>  | 162 | 0.23938     |          |      |
| DE                  | 2   | 7.48119     | 27.0616  | **   |
| ADE                 | 4   | 0.38760     | 1.4021   |      |
| BDE                 | 4   | 0.46482     | 1.6814   |      |
| CDE                 | 4   | 0.06748     | 0.2441   |      |
| ABDE                | 8   | 0.09163     | 0.3315   |      |
| ACDE                | 8   | 0.14154     | 0.5120   |      |
| BCDE                | 8   | 0.26667     | 0.9646   |      |
| ABCDE               | 16  | 0.11377     | 0.4115   |      |
| Error <sub>5</sub>  | 162 | 0.27645     |          |      |
| DF                  | 4   | 0.60533     | 2.4170   | *    |
| ADF                 | 8   | 0.14450     | 0.5770   |      |
| BDF                 | 8   | 0.31999     | 1.2777   |      |
| CDF                 | 8   | 0.20000     | 0.7896   |      |
| ABDF                | 16  | 0.16333     | 0.6522   |      |
| ACDF                | 16  | 0.23981     | 0.9575   |      |
| BCDF                | 16  | 0.19119     | 0.7634   |      |
| ABCDF               | 32  | 0.22085     | 0.8818   |      |
| Error <sub>6</sub>  | 324 | 0.25045     |          |      |
| EF                  | 2   | 1.72674     | 6.3654   | **   |
| AEF                 | 4   | 0.30705     | 1.1319   |      |
| BEF                 | 4   | 0.47424     | 1.7482   |      |

Table 2—Continued

| Source of variance | df   | Mean square | F-ratio | Sig. |
|--------------------|------|-------------|---------|------|
| CEF                | 4    | 0.21082     | 0.7772  |      |
| ABEF               | 8    | 0.33509     | 1.2353  |      |
| ACEF               | 8    | 0.29169     | 1.0753  |      |
| BCEF               | 8    | 0.44239     | 1.6308  |      |
| ABCEF              | 16   | 0.29522     | 1.0883  |      |
| Error <sub>7</sub> | 162  | 0.27127     |         |      |
| DEF                | 4    | 0.67292     | 2.5871  | *    |
| BDEF               | 8    | 0.36589     | 1.4067  |      |
| CDEF               | 8    | 0.22376     | 0.8603  |      |
| ABDEF              | 16   | 0.22201     | 0.8535  |      |
| ACDEF              | 16   | 0.28403     | 1.0920  |      |
| BCDEF              | 16   | 0.22861     | 0.8789  |      |
| ABCDEF             | 32   | 0.35122     | 1.3503  |      |
| Error <sub>8</sub> | 324  | 0.32425     | 1.2504  |      |
| Total              | 1943 | 0.26011     |         |      |

\*p. ≤ 0.05.  
\*\*p. ≤ 0.01.

Results of the main effects

From the results cited in Table 2 it can be seen that the conditional variable "type of transmission" had no significant effect (p. > 0.05) on the ability of the voice identification panelists to make correct identification. In actuality, there was a slight increase in the number of correct identifications for spectrograms made directly from a microphone (92.42 percent) as compared to spectrograms involving the telephone (91.31 percent) and those involving a telephone with background noise (91.02 percent) but these differences were not found to be statistically significant.

Table 2 does reveal a significant difference (p. < 0.01) between variations of the "context" variable. The results show that when the spectrograms were of words spoken in isolation the correct percentage of identification was 95.77 percent. For words spoken in a fixed context the percentage dropped to 92.39 percent and for words spoken in a random context the percentage correct dipped to 86.59 percent. Using Duncan's Multiple Range technique for making individual comparisons it was found that the random context was significantly lower (p. ≤ 0.01) than either the fixed or isolated contexts. Further that words spoken in isolation have a significantly higher (p. ≤ 0.01) number of correct responses than words spoken in either a fixed or random context. (See Table 3).

TABLE 3.—Differences Between Number of Correct Responses for Voice Identification of Clue Words in Specified Contexts

| Context   | Means <sup>1</sup> | Isolation | Fixed  | Random |
|-----------|--------------------|-----------|--------|--------|
|           |                    | 6.0246    | 5.9059 | 5.7127 |
| Isolation | 6.0246             | ....      | **     | **     |
| Fixed     | 5.9059             | **        | ....   | **     |
| Random    | 5.7127             | **        | **     | ....   |

\*\*p. ≤ 0.01

<sup>1</sup>The means in Table 3 are for the transformed raw data; the same data used in the analysis of variance reported in Table 2.

The results indicate a significant difference (p. < 0.01) in terms of correct responses when the numbers of speakers involved in the identification trials were varied. When the identification trials involved matching the spectrograms of 10 speakers the percentage of correct responses was 93.03 percent. When the trials involved 20 speakers the success ratio was 91.87 percent. For trials involving 40 speakers the percentage correct was 89.58 percent. Further analysis revealed that the 40 speaker trials were significantly more difficult (p. ≤ 0.01) than trials involving the spectrograms of 10 or 20 speakers. There was, however, no significant difference (p. > 0.05) between the 10 and 20 speaker trials. (See Table 4).

TABLE 4.—Differences Between Number of Correct Responses for Tasks Involving 10, 20 and 40 Speakers

| Number of Speakers | Means <sup>1</sup> | 10     | 20     | 40     |
|--------------------|--------------------|--------|--------|--------|
|                    |                    | 5.9367 | 5.8916 | 5.8148 |
| 10                 | 5.9367             | ....   | ....   | **     |
| 20                 | 5.8916             | ....   | ....   | **     |
| 30                 | 5.8148             | **     | **     | ....   |

\*\*p. ≤ 0.01

<sup>1</sup>The means in Table 4 are for the transformed raw data; the same data used in the analysis for variance reported in Table 2.

It will be noted from Table 2 that a significant difference (p. < 0.01) was found to be attributable to the system awareness variable. That is, if the trials involved a match and the examiners were aware that a match was there (94.48 percent correct); or if the trials involved a match but the examiners were unaware that the match existed (84.23

percent correct); or if the task did not involve a match but the examiners were unaware of this fact (96.04 percent correct). Further analysis revealed that the system (open-match) where there was a match but the identifiers had no knowledge of the fact resulted in significantly lower numbers of correct responses (p. ≤ 0.01) than the remaining two systems (closed-match and open-no match) which did not differ (p. > 0.05). (See Table 5).

TABLE 5.—Differences Between Number of Correct Responses for System-Awareness Conditions

| System-Awareness | Means <sup>1</sup> | Open-Match | Closed-Match | Open-No Match |
|------------------|--------------------|------------|--------------|---------------|
|                  |                    | 5.9865     | 5.6178       | 6.0389        |
| Open-Match       | 5.9865             | ....       | **           | ....          |
| Closed-Match     | 5.6178             | **         | ....         | **            |
| Open-No Match    | 6.0389             | ....       | **           | ....          |

\*\*p. ≤ 0.01

<sup>1</sup>The means in Table 5 are for the transformed raw data; the same data used in the analysis of variance reported in Table 2.

Analysis revealed a significant difference (p. < 0.01) in success ratios for spectrograms made at the two recording sessions. Here it should be recalled that *time-elapsed* differences involved whether or not the spectrogram was contemporary with the spectrogram in the identification task or non-contemporary (made at a second recording session). When the spectrogram to be matched was contemporary to the spectrograms used in the identification tasks the percentage correct was 95.21 percent. When the spectrogram to be matched was non-contemporary to the spectrograms used in the identification task the percentage of correct responses was 87.95 percent.

No significant difference (p. > 0.05) was found that could be attributed to the number of utterances of the same words as used in the identification trials. In actuality the percentage of correct responses for only one utterance of the words was 91.29 percent, for two utterances 90.96 percent and for three utterances 92.49 percent.

Results of the two-way interactions

Table 2 reveals a significant interaction effect (p. < 0.01) between the conditional variables "con-

TABLE 2.—Analysis of Variance for Voice Identification Data from the First Cycle

| Source of variance  | df  | Mean square | F-ratio  | Sig. |
|---------------------|-----|-------------|----------|------|
| A-Transmission      | 2   | 0.53474     | 2.1973   |      |
| B-Context           | 2   | 16.05957    | 65.9898  | **   |
| C-No. of Speakers   | 2   | 2.46235     | 10.1180  | **   |
| AB                  | 4   | 0.04041     | 0.1661   |      |
| AC                  | 4   | 0.30651     | 1.2595   |      |
| BC                  | 4   | 0.58995     | 2.4242   |      |
| ABC                 | 8   | 0.15984     | 0.6568   |      |
| Error <sub>1</sub>  | 81  | 0.24336     |          |      |
| D-System Awareness  | 2   | 34.13054    | 114.3056 | **   |
| AD                  | 4   | 0.67458     | 2.2592   |      |
| BD                  | 4   | 2.75975     | 9.2426   | **   |
| CD                  | 4   | 0.54314     | 1.8190   |      |
| ABD                 | 8   | 0.05684     | 0.1904   |      |
| ACD                 | 8   | 0.27265     | 0.9131   |      |
| BCD                 | 8   | 0.07094     | 0.2376   |      |
| ABCD                | 16  | 0.29506     | 0.9882   |      |
| Error <sub>2</sub>  | 162 | 0.29859     |          |      |
| E-Time-Elapsed      | 1   | 31.96546    | 115.4176 | **   |
| AE                  | 2   | 0.97337     | 3.5145   | *    |
| BE                  | 2   | 1.38976     | 5.0180   | **   |
| CE                  | 2   | 0.08860     | 0.3199   |      |
| ABE                 | 4   | 0.15437     | 1.5574   |      |
| ACE                 | 4   | 0.16366     | 0.5909   |      |
| BCE                 | 4   | 0.02578     | 0.0931   |      |
| ABCE                | 8   | 0.08424     | 0.3042   |      |
| Error <sub>3</sub>  | 81  | 0.27695     |          |      |
| F-No. of Utterances | 2   | 0.53388     | 2.2303   |      |
| AF                  | 4   | 0.28702     | 1.1991   |      |
| BF                  | 4   | 0.41008     | 1.7131   | *    |
| CF                  | 4   | 0.66034     | 2.7586   | *    |
| ABF                 | 8   | 0.27043     | 1.1297   |      |
| ACF                 | 8   | 0.37140     | 1.5512   |      |
| BCF                 | 8   | 0.40998     | 1.7127   |      |
| ABCF                | 16  | 0.36536     | 1.5263   |      |
| Error <sub>4</sub>  | 162 | 0.23938     |          |      |
| DE                  | 2   | 7.48119     | 27.0616  | **   |
| ADE                 | 4   | 0.38760     | 1.4021   |      |
| BDE                 | 4   | 0.46482     | 1.6814   |      |
| CDE                 | 4   | 0.06748     | 0.2441   |      |
| ABDE                | 8   | 0.09163     | 0.3315   |      |
| ACDE                | 8   | 0.14154     | 0.5120   |      |
| BCDE                | 8   | 0.26667     | 0.9646   |      |
| ABCDE               | 16  | 0.11377     | 0.4115   |      |
| Error <sub>5</sub>  | 162 | 0.27645     |          |      |
| DF                  | 4   | 0.60533     | 2.4170   | *    |
| ADF                 | 8   | 0.14450     | 0.5770   |      |
| BDF                 | 8   | 0.31999     | 1.2777   |      |
| CDF                 | 8   | 0.20000     | 0.7896   |      |
| ABDF                | 16  | 0.16333     | 0.6522   |      |
| ACDF                | 16  | 0.23981     | 0.9575   |      |
| BCDF                | 16  | 0.19119     | 0.7634   |      |
| ABCDF               | 32  | 0.22085     | 0.8818   |      |
| Error <sub>6</sub>  | 324 | 0.25045     |          |      |
| EF                  | 2   | 1.72674     | 6.3654   | **   |
| AEF                 | 4   | 0.30705     | 1.1319   |      |
| BEF                 | 4   | 0.47424     | 1.7482   |      |

Table 2—Continued

| Source of variance | df   | Mean square | F-ratio | Sig. |
|--------------------|------|-------------|---------|------|
| CEF                | 4    | 0.21082     | 0.7772  |      |
| ABEF               | 8    | 0.33509     | 1.2353  |      |
| ACEF               | 8    | 0.29169     | 1.0753  |      |
| BCEF               | 8    | 0.44239     | 1.6308  |      |
| ABCEF              | 16   | 0.29522     | 1.0883  |      |
| Error <sub>7</sub> | 162  | 0.27127     |         |      |
| DEF                | 4    | 0.67292     | 2.5871  | *    |
| ADEF               | 8    | 0.36589     | 1.4067  |      |
| BDEF               | 8    | 0.22376     | 0.8603  |      |
| CDEF               | 8    | 0.22201     | 0.8535  |      |
| ABDEF              | 16   | 0.28403     | 1.0920  |      |
| ACDEF              | 16   | 0.22861     | 0.8789  |      |
| BCDEF              | 16   | 0.35122     | 1.3503  |      |
| ABCDEF             | 32   | 0.32425     | 1.2504  |      |
| Error <sub>8</sub> | 324  | 0.26011     |         |      |
| Total              | 1943 |             |         |      |

\*p. ≤ 0.05.  
\*\*p. ≤ 0.01.

Results of the main effects

From the results cited in Table 2 it can be seen that the conditional variable "type of transmission" had no significant effect (p. > 0.05) on the ability of the voice identification panelists to make correct identification. In actuality, there was a slight increase in the number of correct identifications for spectrograms made directly from a microphone (92.42 percent) as compared to spectrograms involving the telephone (91.31 percent) and those involving a telephone with background noise (91.02 percent) but these differences were not found to be statistically significant.

Table 2 does reveal a significant difference (p. < 0.01) between variations of the "context" variable. The results show that when the spectrograms were of words spoken in isolation the correct percentage of identification was 95.77 percent. For words spoken in a fixed context the percentage dropped to 92.39 percent and for words spoken in a random context the percentage correct dipped to 86.59 percent. Using Duncan's Multiple Range technique for making individual comparisons it was found that the random context was significantly lower (p. ≤ 0.01) than either the fixed or isolated contexts. Further that words spoken in isolation have a significantly higher (p. ≤ 0.01) number of correct responses than words spoken in either a fixed or random context. (See Table 3).

TABLE 3.—Differences Between Number of Correct Responses for Voice Identification of Clue Words in Specified Contexts

| Context   | Means <sup>1</sup> | Isolation | Fixed  | Random |
|-----------|--------------------|-----------|--------|--------|
|           |                    | 6.0246    | 5.9059 | 5.7127 |
| Isolation | 6.0246             | ....      | **     | **     |
| Fixed     | 5.9059             | **        | ....   | **     |
| Random    | 5.7127             | **        | **     | ....   |

\*\*p. ≤ 0.01

<sup>1</sup>The means in Table 3 are for the transformed raw data; the same data used in the analysis of variance reported in Table 2.

The results indicate a significant difference (p. < 0.01) in terms of correct responses when the numbers of speakers involved in the identification trials were varied. When the identification trials involved matching the spectrograms of 10 speakers the percentage of correct responses was 93.03 percent. When the trials involved 20 speakers the success ratio was 91.87 percent. For trials involving 40 speakers the percentage correct was 89.58 percent. Further analysis revealed that the 40 speaker trials were significantly more difficult (p. ≤ 0.01) than trials involving the spectrograms of 10 or 20 speakers. There was, however, no significant difference (p. > 0.05) between the 10 and 20 speaker trials. (See Table 4).

TABLE 4.—Differences Between Number of Correct Responses for Tasks Involving 10, 20 and 40 Speakers

| Number of Speakers | Means <sup>1</sup> | 10     | 20     | 40     |
|--------------------|--------------------|--------|--------|--------|
|                    |                    | 5.9367 | 5.8916 | 5.8148 |
| 10                 | 5.9367             | ....   | ....   | **     |
| 20                 | 5.8916             | ....   | ....   | **     |
| 30                 | 5.8148             | **     | **     | ....   |

\*\*p. ≤ 0.01

<sup>1</sup>The means in Table 4 are for the transformed raw data; the same data used in the analysis for variance reported in Table 2.

It will be noted from Table 2 that a significant difference (p. < 0.01) was found to be attributable to the system awareness variable. That is, if the trials involved a match and the examiners were aware that a match was there (94.48 percent correct); or if the trials involved a match but the examiners were unaware that the match existed (84.23

percent correct); or if the task did not involve a match but the examiners were unaware of this fact (96.04 percent correct). Further analysis revealed that the system (open-match) where there was a match but the identifiers had no knowledge of the fact resulted in significantly lower numbers of correct responses (p. ≤ 0.01) than the remaining two systems (closed-match and open-no match) which did not differ (p. > 0.05). (See Table 5).

TABLE 5.—Differences Between Number of Correct Responses for System-Awareness Conditions

| System-Awareness | Means <sup>1</sup> | Open-Match | Closed-Match | Open-No Match |
|------------------|--------------------|------------|--------------|---------------|
|                  |                    | 5.9865     | 5.6178       | 6.0389        |
| Open-Match       | 5.9865             | ....       | **           | ....          |
| Closed-Match     | 5.6178             | **         | ....         | **            |
| Open-No Match    | 6.0389             | ....       | **           | ....          |

\*\*p. ≤ 0.01

<sup>1</sup>The means in Table 5 are for the transformed raw data; the same data used in the analysis of variance reported in Table 2.

Analysis revealed a significant difference (p. < 0.01) in success ratios for spectrograms made at the two recording sessions. Here it should be recalled that *time-elapsed* differences involved whether or not the spectrogram was contemporary with the spectrogram in the identification task or non-contemporary (made at a second recording session). When the spectrogram to be matched was contemporary to the spectrograms used in the identification tasks the percentage correct was 95.21 percent. When the spectrogram to be matched was non-contemporary to the spectrograms used in the identification task the percentage of correct responses was 87.95 percent.

No significant difference (p. > 0.05) was found that could be attributed to the number of utterances of the same words as used in the identification trials. In actuality the percentage of correct responses for only one utterance of the words was 91.29 percent, for two utterances 90.96 percent and for three utterances 92.49 percent.

Results of the two-way interactions

Table 2 reveals a significant interaction effect (p. < 0.01) between the conditional variables "con-

text" and "system-awareness" (BD). Table 6 represents the percentage of correct responses made by the identifiers under all combinations of these two variables.

TABLE 6.—Percentage of Correct Identifications Under Conditions of Context and Speaker-Awareness

| Context   | System Awareness |              |                 |
|-----------|------------------|--------------|-----------------|
|           | Open-Match       | Closed-Match | Closed-No Match |
| Isolation | 98.46            | 90.79        | 98.05           |
| Fixed     | 94.65            | 86.47        | 96.04           |
| Random    | 90.33            | 75.41        | 94.03           |

The results cited in Table 6 conform to expectations.

At all context levels the closed-match system resulted in a lower number of correct identifica-

TABLE 7.—Differences Between Number of Correct Responses Under Conditions of Context by System Awareness

| Context/System | Means <sup>1</sup> | 1/1    | 1/2    | 1/3    | 2/1    | 2/2    | 2/3    | 3/1    | 3/2    | 3/3    |
|----------------|--------------------|--------|--------|--------|--------|--------|--------|--------|--------|--------|
|                |                    | 6.1155 | 5.8564 | 6.1019 | 5.9919 | 5.6881 | 6.0376 | 5.8521 | 5.3089 | 5.9771 |
| 1/1            | 6.1155             | ....   | **     | ....   | *      | **     | ....   | **     | **     | *      |
| 1/2            | 5.8564             | **     | ....   | **     | *      | **     | ....   | **     | **     | *      |
| 1/3            | 6.1019             | ....   | **     | ....   | *      | **     | ....   | **     | **     | *      |
| 2/1            | 5.9919             | *      | *      | *      | ....   | **     | ....   | **     | **     | **     |
| 2/2            | 5.6881             | **     | **     | **     | **     | ....   | ....   | **     | **     | ....   |
| 2/3            | 6.0376             | ....   | **     | ....   | *      | **     | ....   | **     | **     | *      |
| 3/1            | 5.8521             | **     | ....   | **     | ....   | **     | ....   | **     | ....   | **     |
| 3/2            | 5.3089             | **     | **     | **     | **     | **     | **     | **     | ....   | ....   |
| 3/3            | 5.9771             | *      | *      | *      | ....   | **     | ....   | *      | **     | ....   |

\*p. ≤ 0.05.  
\*\*p. ≤ 0.01.

<sup>1</sup>The means in Table 7 are for the transferred raw data; the same data used in the analysis of variance reported in Table 2.

fication task. Table 8 represents the percentage of correct responses made by the examiners under all combinations of these two variables.

As expected under all levels of transmission, tasks involving matching a non-contemporary spectrogram had a lower number of correct responses than tasks involving a contemporary match. The Duncan test found all these differences to be significant (p. ≤ 0.05). What is interesting to note is that under the microphone only transmission level there was a higher number of correct responses on non-

tions. At all levels of system-awareness the isolated words resulted in a higher number of correct identifications (p. ≤ 0.05) followed in turn by the fixed and random contexts. Employing the Duncan Multiple Range Technique on the transformed raw data it was found that under the random context the Open-Match system had a significantly lower (p. ≤ 0.01) number of correct responses than the combination "random context open-no match." This was the only context level where such a significant difference was found. In addition it was found that for the open-no match system the fixed context did not differ significantly (p. > 0.05) from the random context where as in all our systems such a pattern was significant (p. ≤ 0.05). (See Table 7).

A significant interaction (p. < 0.05) was also found between combinations involving "Type of Transmission" and "Time-elapsed" variations in the spectrograms to be matched in a particular identi-

contemporary tasks than for the same tasks at the other two transmission levels. This difference was found to be significant (p. ≤ 0.05). (See Table 9).

Table 2 reveals a significant interaction of the conditional variables of "context" and "time-elapsed" variations. Table 10 represents the percentage of correct responses made by the examiners under all combinations of these variables.

It can be seen from Table 10 that at all three contexts a greater number of correct responses were attained for contemporary tasks than for non-

TABLE 8.—Percentage of Correct Identifications Under Conditions of Transmission and Time-Elapsed

| Type of Transmission           | Non-Con-temporary |           |
|--------------------------------|-------------------|-----------|
|                                | Contemporary      | temporary |
| (a) Microphone Only            | 94.86             | 89.99     |
| (β) Telephone-Microphone       | 95.68             | 86.93     |
| (γ) Telephone-Microphone-Noise | 95.10             | 86.93     |

contemporary tasks. These differences were found via the Duncan technique for individual comparisons to be statistically significant (p. ≤ 0.01). Also, for both the contemporary and non-contemporary tasks the results indicate that the isolated context was superior to the fixed context which in turn was superior to the random context (p. ≤ 0.05). However, in the case of the contemporary tasks the difference between the isolated and fixed context was not found to be statistically significant (p. ≤ 0.05). (See Table 11).

TABLE 9.—Differences Between Number of Correct Responses Under Conditions of Transmission and Time-elapsed

| Transmission/Time-elapsed | Means <sup>1</sup> | 1/1    | 1/2    | 2/1    | 2/2    | 3/1    | 3/2    |
|---------------------------|--------------------|--------|--------|--------|--------|--------|--------|
|                           |                    | 5.9981 | 5.8301 | 6.0232 | 5.7111 | 6.0066 | 5.7173 |
| 1/1                       | 5.9981             | ....   | *      | ....   | *      | ....   | *      |
| 1/2                       | 5.8301             | *      | ....   | *      | *      | *      | *      |
| 2/1                       | 6.0232             | ....   | *      | ....   | *      | ....   | *      |
| 2/2                       | 5.7111             | *      | *      | *      | ....   | *      | ....   |
| 3/1                       | 6.0066             | ....   | *      | ....   | *      | ....   | *      |
| 3/2                       | 5.7173             | *      | *      | *      | ....   | *      | ....   |

\*p. ≤ 0.05.

<sup>1</sup>The means in Table 9 are for the transformed raw data; the same data used in the analysis of variance reported in Table 2.

TABLE 10.—Percentage of Correct Identifications Under Conditions of Context and Time-Elapsed

| Context   | Non-Con-temporary |           |
|-----------|-------------------|-----------|
|           | Contemporary      | temporary |
| Isolation | 98.77             | 92.76     |
| Fixed     | 97.56             | 87.21     |
| Random    | 89.30             | 83.88     |

A significant interaction (p. < 0.05) was found between "number of speakers involved in the identification task" and the number of "utterances" of the words used in the identification task. Table 12 represents the percentage of correct responses made by the identifiers under all combinations of these variables.

It would appear that only when 10 speakers are involved in the identification trials does the num-

TABLE 11.—Differences Between Number of Correct Responses for Context and Time-Elapsed Conditions

| Context/Time-elapsed       | Means <sup>1</sup> | Iso/C  | Iso/NC | Fixed/C | Fixed/NC | Random/C | Random/N |
|----------------------------|--------------------|--------|--------|---------|----------|----------|----------|
|                            |                    | 6.1249 | 5.9243 | 6.0876  | 5.7242   | 5.8155   | 5.6099   |
| Isolation/Contemporary     | 6.1249             | ....   | **     | ....    | **       | **       | **       |
| Isolation/Non-Contemporary | 5.9243             | **     | ....   | **      | *        | *        | **       |
| Fixed/Contemporary         | 6.0876             | ....   | **     | ....    | **       | **       | **       |
| Fixed/Non-Contemporary     | 5.7242             | **     | **     | **      | *        | *        | **       |
| Random/Contemporary        | 5.8155             | **     | *      | **      | ....     | ....     | **       |
| Random/Non-Contemporary    | 5.6099             | **     | **     | **      | **       | **       | ....     |

\*p. ≤ 0.05.  
\*\*p. ≤ 0.01.

<sup>1</sup>The means in Table 11 are for the transformed raw data; the same data used in the analysis of variance reported in Table 2.

**TABLE 12.—Percentage of Correct Identification Under Conditions of Number of Speakers/Utterances**

| Number of Speakers | 1 Utterance | 2 Utterances | 3 Utterances |
|--------------------|-------------|--------------|--------------|
| 10                 | 94.29       | 91.56        | 94.03        |
| 20                 | 90.12       | 92.95        | 92.54        |
| 40                 | 89.46       | 88.37        | 90.89        |

ber of utterances have an effect on the number of correct identifications. Even here the effect is difficult to explain since two utterances of ten speakers differs significantly ( $p. \leq 0.05$ ) (See Table 13) from one and three utterances of ten speakers, but the effect is to reduce rather than increase accuracy. At the 20 and 40 speaker levels the effects of number of utterances was not significant ( $p. > 0.05$ ).

**TABLE 13.—Differences Between Number of Correct Responses for Trials Involving Varying Number of Speakers/Number of Utterances**

| No. Speakers/<br>No. Utterances | Means <sup>1</sup> | 10/1  | 10/2  | 10/3  | 20/1  | 20/2  | 20/3  | 40/1  | 40/2  | 40/3  |
|---------------------------------|--------------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 10/1                            | 5.9765             | ..... | ..... | ..... | ..... | ..... | ..... | ..... | ..... | ..... |
| 10/2                            | 5.8671             | ..... | ..... | ..... | ..... | ..... | ..... | ..... | ..... | ..... |
| 10/3                            | 5.9666             | ..... | ..... | ..... | ..... | ..... | ..... | ..... | ..... | ..... |
| 20/1                            | 5.8298             | ..... | ..... | ..... | ..... | ..... | ..... | ..... | ..... | ..... |
| 20/2                            | 5.9320             | ..... | ..... | ..... | ..... | ..... | ..... | ..... | ..... | ..... |
| 20/3                            | 5.9131             | ..... | ..... | ..... | ..... | ..... | ..... | ..... | ..... | ..... |
| 40/1                            | 5.8181             | ..... | ..... | ..... | ..... | ..... | ..... | ..... | ..... | ..... |
| 40/2                            | 5.7689             | ..... | ..... | ..... | ..... | ..... | ..... | ..... | ..... | ..... |
| 40/3                            | 5.8574             | ..... | ..... | ..... | ..... | ..... | ..... | ..... | ..... | ..... |

\* $p. \leq 0.05$

<sup>1</sup> The means in Table 13 are for the transformed raw data; the same data used in the analysis of variance reported in Table 2.

The data show generally that as the number of speakers involved in the identification trials increases accuracy decreases. However, the difference at the one utterance level is not always significant and at the two utterance level the only difference is between tasks involving 40 speakers and 10 speakers. It would appear that the conditional variables, "number of speakers" and "number of utterances" tend to confound one another. Probably, the logistics of such tasks contribute as much to the variance of correct responses as any combination of conditional effects.

A significant interaction ( $p. < 0.01$ ) between the conditional variables "system-awareness" and "time-elapsed" variations was also reflected in Table 2. Table 14 represents the percentage of correct responses made by the examiners under all combinations of these variables.

Further analysis using the Duncan Technique on the transformed raw scores (See Table 15), found that for all conditions of system-awareness the trials involving matching contemporary spectrograms yielded a significantly higher ( $p. \leq 0.05$ ) number of correct responses than for tasks involving non-

**TABLE 14.—Percentage of Correct Identifications Under Conditions of System-Awareness and Time-Elapsed**

| System-Awareness | Contemporary | Non-Contemporary |
|------------------|--------------|------------------|
| Closed-Match     | 96.95        | 92.01            |
| Open-Match       | 91.02        | 77.44            |
| Open-No Match    | 97.67        | 94.41            |

contemporary spectrograms. It was also found that for both contemporary and non-contemporary tasks the open-match conditions yielded a significantly lower ( $p. \leq 0.01$ ) number of correct responses than the other conditions of system-awareness which did not differ significantly. One point is quite apparent, however, when the examiners were in an open-match system trying to make match non-contemporary spectrograms the percentage of correct responses was disproportionately low.

Table 2 shows a significant interaction ( $p. < 0.05$ ) between conditions of "System-awareness" and "Number of utterances." Table 16 represents the percentage of correct responses made by the examiners under all the conditions of these variables.

**TABLE 15.—Differences Between Number of Correct Responses for Conditions of System-Awareness and Time-Elapsed**

| Awareness/Time-elapsed         | Means <sup>1</sup> | CI/C   | CI/NC  | Op/M-C | Op/M-NC | Op/NM-C | Op/NM-NC |
|--------------------------------|--------------------|--------|--------|--------|---------|---------|----------|
|                                |                    | 6.0674 | 5.9056 | 5.8690 | 5.3666  | 6.0914  | 5.9863   |
| Closed/Contemporary            | 6.0674             | .....  | **     | **     | **      | .....   | *        |
| Closed/Non-Contemporary        | 5.9056             | **     | .....  | .....  | **      | .....   | .....    |
| Open/Match-Contemporary        | 5.8690             | **     | .....  | .....  | **      | .....   | .....    |
| Open/Match-Non-Contemporary    | 5.3666             | **     | **     | **     | .....   | **      | **       |
| Open/No Match-Contemporary     | 6.0914             | .....  | **     | **     | **      | .....   | *        |
| Open/No Match-Non-Contemporary | 5.9863             | *      | .....  | **     | **      | *       | .....    |

\* $p. \leq 0.05$ .

\*\* $p. \leq 0.01$ .

<sup>1</sup> The means in Table 15 are for the transformed raw data; the same data used in the analysis of variance reported in Table 2.

**TABLE 16.—Percentage of Correct Identifications Under Varying Conditions of System-Awareness and Number of Utterances**

| System-Awareness | 1 Utterance | 2 Utterances | 3 Utterances |
|------------------|-------------|--------------|--------------|
| Closed-Match     | 93.26       | 95.68        | 94.50        |
| Open-Match       | 84.47       | 82.25        | 85.96        |
| Open-No Match    | 96.14       | 94.95        | 97.02        |

The results cited in Table 16 when coupled with the individual comparisons contained in Table 17 show that only at the open-match level was there a significant difference ( $p. \leq 0.05$ ) between number of utterances. Two utterances yielded significantly lower ( $p. \leq 0.05$ ) responses than the open-match system at three utterances. The results also indicate that, as was the case with previous interactions involving systems, the open-match condition yielded

**TABLE 17.—Differences Between Number of Correct Responses for Conditions of System-Awareness and Number of Utterances**

| System/Utterances | Means <sup>1</sup> | 1/1    | 1/2    | 1/3    | 2/1    | 2/2    | 2/3    | 3/1    | 3/2    | 3/3    |
|-------------------|--------------------|--------|--------|--------|--------|--------|--------|--------|--------|--------|
|                   |                    | 5.9481 | 6.0269 | 5.9847 | 5.6323 | 5.5382 | 5.6828 | 6.0441 | 6.0029 | 6.0696 |
| 1/1               | 5.9481             | .....  | .....  | .....  | *      | *      | *      | .....  | .....  | *      |
| 1/2               | 6.0269             | .....  | .....  | .....  | *      | *      | *      | .....  | .....  | .....  |
| 1/3               | 5.9847             | .....  | .....  | .....  | *      | *      | *      | .....  | .....  | .....  |
| 2/1               | 5.6323             | *      | *      | *      | .....  | .....  | .....  | *      | *      | *      |
| 2/2               | 5.5382             | *      | *      | *      | .....  | .....  | .....  | *      | *      | *      |
| 2/3               | 5.6828             | *      | *      | *      | .....  | .....  | .....  | *      | *      | *      |
| 3/1               | 6.0441             | .....  | .....  | .....  | *      | *      | *      | .....  | .....  | .....  |
| 3/2               | 6.0029             | .....  | .....  | .....  | *      | *      | *      | .....  | .....  | .....  |
| 3/3               | 6.0696             | *      | .....  | .....  | *      | *      | *      | .....  | .....  | .....  |

\* $p. \leq 0.05$ .

<sup>1</sup> The means in Table 17 are for the transformed raw data; the same data used in the analysis of variance reported in Table 2.

a significantly lower ( $p. \leq 0.05$ ) number of correct responses than the other two systems (which did not differ  $p. > 0.05$ ) regardless of the number of utterances involved in the identification tasks.

The last significant two-way interaction reported

in Table 2 was between the conditional variables of "Time-elapsed" and "Number of utterances." Table 18 represents the percentage of correct identifications under the combination of these two conditional variables.

**TABLE 18.—Percentage of Correct Responses for Conditions of Time-Elapsed and Number of Utterances**

| Time-Elapsed     | 1 Utterance | 2 Utterances | 3 Utterances |
|------------------|-------------|--------------|--------------|
| Contemporary     | 93.28       | 95.92        | 96.43        |
| Non-Contemporary | 89.30       | 86.00        | 88.55        |

When individual comparisons were made on the transformed raw data it was found that regardless of the number of utterances involved in a trial, when the spectrogram to the matched was non-

contemporary to remaining prints, there were significantly lower scores ( $p. \leq 0.01$ ) than when the matching print was contemporary. With respect to the number of utterances, individual comparisons revealed a significantly lower ( $p. \leq 0.05$ ) number of correct responses for identification involving one utterance than those involving two and three utterances for trials involving a contemporary match, whereas for trials involving a non-contemporary match, identifications involving two utterances yielded significantly lower ( $p. \leq 0.05$ ) scores. (See Table 19).

**TABLE 19.—Differences Between Number of Correct Responses for Time-Elapsed Conditions and Number of Utterances**

| Time-elapsed/<br>Utterances | Means <sup>1</sup> | Cont/1 | Cont/2 | Cont/3 | N-Cont/1 | N-Cont/2 | N-Cont/3 |
|-----------------------------|--------------------|--------|--------|--------|----------|----------|----------|
|                             |                    | 5.9469 | 6.0296 | 6.0514 | 5.8028   | 5.6824   | 5.7733   |
| Contemporary/1              | 5.9469             | .....  | •      | •      | **       | **       | **       |
| Contemporary/2              | 6.0296             | •      | .....  | .....  | **       | **       | **       |
| Contemporary/3              | 6.0514             | •      | .....  | .....  | **       | **       | **       |
| Non-Contemporary/1          | 5.8028             | **     | **     | **     | **       | **       | .....    |
| Non-Contemporary/2          | 5.6824             | **     | **     | **     | .....    | .....    | •        |
| Non-Contemporary/3          | 5.7733             | **     | **     | **     | •        | •        | .....    |

\* $p. \leq 0.05$ .

\*\* $p. \leq 0.01$ .

<sup>1</sup>The means in Table 19 are for the transformed raw data; the same data used in the analysis of variance reported in Table 2.

**Trends toward two-way interactions**

There were two two-way interactions that approached significance ( $0.05 < p. < 0.10$ ) which deserve mention. The first involved the conditional variables of "context" and "number of speakers." Table 20 represents the percentage of correct identifications under the combinations of these two variables.

While the results are not statistically significant, it does appear that the most important differential and detrimental effects of "number of speakers" appear in the random context.

The second two-way interaction that approaches significant ( $0.05 < p. < 0.10$ ) involved the conditions of "transmission" and "system-awareness." Table 21 represents the percentage of correct identifications under the combination of these two variables.

The results seem to indicate clearly that type of

**TABLE 20.—Percentage of Correct Responses Under Conditions of Context and Number of Speakers**

| Context       | Number of Speakers |       |       |
|---------------|--------------------|-------|-------|
|               | 10                 | 20    | 40    |
| (I) Isolation | 96.40              | 95.94 | 94.96 |
| (II) Fixed    | 93.31              | 92.64 | 91.20 |
| (III) Random  | 90.18              | 87.04 | 82.56 |

**TABLE 21.—Percentage of Correct Responses Under Conditions of Transmission and System-Awareness**

| Transmission                   | System-Awareness |            |               |
|--------------------------------|------------------|------------|---------------|
|                                | Closed Match     | Open-Match | Open-No Match |
| (a) Microphone                 | 94.24            | 86.99      | 96.04         |
| (β) Telephone-Microphone       | 94.60            | 83.44      | 95.89         |
| (γ) Telephone-Microphone-Noise | 94.60            | 82.25      | 96.19         |

transmission does little to deter spectrogram identification except in the open-match system.

All other two-way interactions reported in Table 2 had a probability of occurrence by chance of more than 10 percent and thus were not considered indicative of trends in the data.

**Significant three-way interactions**

Only one three-way interaction was found to be significant ( $p. < 0.05$ ) as a result of the analysis reported in Table 2. That one involved the conditional variables "System-Awareness," "Time-Elapsed" and "Number of Utterances" (DEF). Table 22 represents the percentage of correct identifications under the combinations of these variables.

The results indicate that for the contemporary identification trials the most significant variable operating was that of system-awareness. The only deviation from this encompassing statement was the fact that for contemporary spectrograms under the open-match system, three utterances yielded significantly higher ( $p. \leq 0.05$ ) scores than one utterance. For the non-contemporary spectrograms the results are less clear. An emerging pattern shows

**TABLE 22.—Percentage of Correct Responses Under Conditions of System-Awareness, Time-Elapsed, and Number of Utterances**

| System-Awareness | Contemporary |       |       | Non-Contemporary |       |       |
|------------------|--------------|-------|-------|------------------|-------|-------|
|                  | 1            | 2     | 3     | 1                | 2     | 3     |
| Closed-Match     | 95.17        | 98.05 | 97.63 | 91.36            | 93.31 | 91.36 |
| Open-Match       | 87.86        | 91.26 | 93.93 | 81.07            | 73.25 | 77.98 |
| Open-No Match    | 96.81        | 98.46 | 97.74 | 95.47            | 91.46 | 96.30 |

that the open-match system yielded lower scores, however, increasing the number of utterances involved in a task seemed to compound the problem. It was also the case that the open-no match system yielded higher scores at three utterances than both the closed-match and the open-match systems. The pattern is clear that for tasks involving non-contemporary matches, regardless of the system or the number of utterances, the scores were significantly lower ( $p. \leq 0.05$ ) than for trials involving contemporary matches. The only exception to this statement is that under the open-no match system with three utterances there was no significant difference ( $p. > 0.05$ ) between contemporary and non-contemporary trials. (See Table 23).

**TABLE 23.—Differences Between Number of Correct Responses for System-Awareness, Time-Elapsed and Number of Utterances Conditions**

| System/<br>Time-<br>Elapsed/<br>Utterances | Means <sup>1</sup> | 1/1/1  | 1/1/2  | 1/1/3  | 1/2/1  | 1/2/2  | 1/2/3  | 2/1/1  | 2/1/2  | 2/1/3  |
|--|--------------------|--------|--------|--------|--------|--------|--------|--------|--------|--------|
|  |                    | 6.0102 | 6.1020 | 6.0902 | 5.8860 | 5.9518 | 5.8792 | 5.7652 | 5.8710 | 5.9707 |
| 1/1/1                                      | 6.0102             | .....  | .....  | .....  | .....  | .....  | .....  | .....  | .....  | .....  |
| 1/1/2                                      | 6.1020             | .....  | .....  | .....  | •      | .....  | .....  | •      | .....  | .....  |
| 1/1/3                                      | 6.0902             | .....  | .....  | .....  | •      | .....  | .....  | •      | .....  | .....  |
| 1/2/1                                      | 5.8860             | .....  | •      | •      | .....  | .....  | .....  | .....  | .....  | .....  |
| 1/2/2                                      | 5.9518             | .....  | .....  | .....  | .....  | .....  | .....  | .....  | .....  | .....  |
| 1/2/3                                      | 5.8792             | .....  | •      | •      | .....  | .....  | .....  | .....  | .....  | .....  |
| 2/1/1                                      | 5.7652             | .....  | •      | •      | .....  | .....  | .....  | .....  | .....  | •      |
| 2/1/2                                      | 5.8710             | .....  | •      | •      | .....  | .....  | .....  | .....  | .....  | .....  |
| 2/1/3                                      | 5.9707             | .....  | .....  | .....  | .....  | .....  | .....  | .....  | .....  | .....  |
| 2/2/1                                      | 5.4994             | •      | •      | •      | •      | •      | •      | •      | •      | •      |
| 2/2/2                                      | 5.2053             | •      | •      | •      | •      | •      | •      | •      | •      | •      |
| 2/2/3                                      | 5.3949             | •      | •      | •      | •      | •      | •      | •      | •      | •      |
| 3/1/1                                      | 6.0653             | .....  | .....  | .....  | .....  | .....  | .....  | .....  | .....  | .....  |
| 3/1/2                                      | 6.1157             | .....  | .....  | .....  | •      | .....  | .....  | •      | .....  | .....  |
| 3/1/3                                      | 6.0933             | .....  | .....  | .....  | .....  | .....  | .....  | .....  | .....  | .....  |
| 3/2/1                                      | 6.0229             | .....  | .....  | .....  | .....  | .....  | .....  | .....  | .....  | .....  |
| 3/2/2                                      | 5.8900             | .....  | •      | •      | .....  | .....  | .....  | .....  | .....  | .....  |
| 3/2/3                                      | 6.0459             | .....  | .....  | .....  | •      | .....  | .....  | •      | .....  | .....  |

TABLE 23.—Continued

| System/<br>Time-<br>Elapsed/<br>Utterances | Means <sup>1</sup> | 2/2/1 | 2/2/2  | 2/2/3 | 3/1/1 | 3/1/2 | 3/1/3 | 3/2/1 | 3/2/2 | 3/2/3 |
|--|--------------------|-------|--------|-------|-------|-------|-------|-------|-------|-------|
|  |                    | 1/1/1 | 6.0102 | •     | •     | •     | •     | •     | •     | •     |
| 1/1/2                                      | 6.1020             | •     | •      | •     | •     | •     | •     | •     | •     | •     |
| 1/1/3                                      | 6.0902             | •     | •      | •     | •     | •     | •     | •     | •     | •     |
| 1/2/1                                      | 5.8860             | •     | •      | •     | •     | •     | •     | •     | •     | •     |
| 1/2/2                                      | 5.9518             | •     | •      | •     | •     | •     | •     | •     | •     | •     |
| 1/2/3                                      | 5.8792             | •     | •      | •     | •     | •     | •     | •     | •     | •     |
| 2/1/1                                      | 5.7652             | •     | •      | •     | •     | •     | •     | •     | •     | •     |
| 2/1/2                                      | 5.8710             | •     | •      | •     | •     | •     | •     | •     | •     | •     |
| 2/1/3                                      | 5.9707             | •     | •      | •     | •     | •     | •     | •     | •     | •     |
| 2/2/1                                      | 5.4494             | •     | •      | •     | •     | •     | •     | •     | •     | •     |
| 2/2/2                                      | 5.2053             | •     | •      | •     | •     | •     | •     | •     | •     | •     |
| 2/2/3                                      | 5.3949             | •     | •      | •     | •     | •     | •     | •     | •     | •     |
| 3/1/1                                      | 6.0653             | •     | •      | •     | •     | •     | •     | •     | •     | •     |
| 3/1/2                                      | 6.1157             | •     | •      | •     | •     | •     | •     | •     | •     | •     |
| 3/1/3                                      | 6.0933             | •     | •      | •     | •     | •     | •     | •     | •     | •     |
| 3/2/1                                      | 6.0229             | •     | •      | •     | •     | •     | •     | •     | •     | •     |
| 3/2/2                                      | 5.8900             | •     | •      | •     | •     | •     | •     | •     | •     | •     |
| 3/2/3                                      | 6.0459             | •     | •      | •     | •     | •     | •     | •     | •     | •     |

\*p. ≤ 0.05.

<sup>1</sup> The means in Table 23 are for the transformed raw data; the same data used in the analysis of variance reported in Table 2.

Trends toward three-way interactions

There was one three-way interaction that approached significance (0.05 < p. < 0.10) which appears to deserve mention. It involved the conditional variables "Context," "Number of Speakers" and "Number of Utterances." Table 24 represents the percentages of correct identifications under the combinations of these variables.

The magnitude of the percentages contained in Table 24 shows rather clearly the effects of the "Context" variable. For every combination, the words spoken in isolation yielded the largest percentage of correct identifications; there was a moderate drop in terms of correct responses for words spoken in a fixed context; and a substantial lowering of the number of correct responses for words

spoken in a random context. Table 24 also shows that for the "Isolation" context there appears very little variance across the combinations of number of speakers and number of utterances. For the "Fixed" and "Random" contexts there is a considerable amount of variance in terms of the cells.

Less clear are the effects of the variable number of speakers in interaction with context and number of utterances. A trend supports the general conclusion that the examiners were not as accurate for trials involving forty speakers as they were when the tasks involved ten or twenty speakers. This statement seems particularly true for words spoken in the "Fixed" and "Random" context.

As has been previously observed, there appears to be no consistent pattern for the effects of the variable number of utterances.

TABLE 24.—Percentage of Correct Responses Under Conditions of Context, Number of Speakers and Number of Utterances

| Context       | 10 Speakers |       |       | 20 Speakers |       |       | 40 Speakers |       |       |
|---------------|-------------|-------|-------|-------------|-------|-------|-------------|-------|-------|
|               | 1 utt       | 2 utt | 3 utt | 1 utt       | 2 utt | 3 utt | 1 utt       | 2 utt | 3 utt |
| (I) Isolation | 96.61       | 96.76 | 95.83 | 95.99       | 95.06 | 96.76 | 94.91       | 95.99 | 93.98 |
| (II) Fixed    | 95.06       | 90.59 | 94.29 | 91.51       | 94.60 | 91.82 | 91.82       | 87.50 | 94.29 |
| (III) Random  | 91.20       | 87.35 | 91.98 | 82.87       | 89.20 | 89.04 | 81.64       | 81.64 | 84.41 |

Trends toward four-way interactions

No significant four-way interactions were found; however, there was one that approached significance (0.05 < p. < 0.10). It involved the conditional

variables: "Transmission," "Context," "Number of Speakers" and "Number of Utterances." Tables 25, 26, and 27 represent the percentages of correct identifications under all combinations of these variables.

TABLE 25.—Percentage of Correct Responses (α) Transmission, Context, Number of Speakers and Utterances

| Context       | 10 Speakers |       |       | 20 Speakers |       |       | 40 Speakers |       |       |
|---------------|-------------|-------|-------|-------------|-------|-------|-------------|-------|-------|
|               | 1 utt       | 2 utt | 3 utt | 1 utt       | 2 utt | 3 utt | 1 utt       | 2 utt | 3 utt |
| (I) Isolation | 96.30       | 94.91 | 94.44 | 96.76       | 95.83 | 98.61 | 97.22       | 98.15 | 95.83 |
| (II) Fixed    | 95.37       | 91.67 | 93.52 | 89.82       | 96.30 | 96.76 | 89.82       | 88.89 | 95.83 |
| (III) Random  | 87.04       | 88.43 | 95.37 | 82.87       | 94.44 | 93.06 | 87.04       | 77.78 | 83.33 |

TABLE 26.—Percentage of Correct Responses (β) Transmission, Context, Number of Speakers and Utterances

| Context       | 10 Speakers |       |       | 20 Speakers |       |       | 40 Speakers |       |       |
|---------------|-------------|-------|-------|-------------|-------|-------|-------------|-------|-------|
|               | 1 utt       | 2 utt | 3 utt | 1 utt       | 2 utt | 3 utt | 1 utt       | 2 utt | 3 utt |
| (I) Isolation | 99.07       | 97.69 | 97.22 | 96.30       | 96.30 | 94.44 | 91.67       | 94.91 | 92.59 |
| (II) Fixed    | 100.00      | 90.28 | 93.52 | 92.13       | 91.67 | 91.20 | 93.52       | 83.80 | 95.83 |
| (III) Random  | 94.91       | 84.72 | 90.28 | 83.33       | 88.89 | 82.87 | 78.24       | 83.33 | 86.57 |

TABLE 27.—Percentage of Correct Responses (γ) Transmission, Context, Number of Speakers and Utterances

| Context       | 10 Speakers |       |       | 20 Speakers |       |       | 40 Speakers |       |       |
|---------------|-------------|-------|-------|-------------|-------|-------|-------------|-------|-------|
|               | 1 utt       | 2 utt | 3 utt | 1 utt       | 2 utt | 3 utt | 1 utt       | 2 utt | 3 utt |
| (I) Isolation | 94.44       | 97.69 | 95.83 | 94.91       | 93.06 | 97.22 | 95.83       | 94.91 | 93.52 |
| (II) Fixed    | 89.82       | 89.82 | 95.83 | 92.59       | 95.83 | 87.50 | 92.13       | 89.82 | 91.20 |
| (III) Random  | 91.67       | 88.89 | 90.28 | 82.41       | 84.26 | 91.20 | 79.63       | 83.30 | 83.33 |

The tables cited above do not reflect a consistent pattern that can be associated with differing levels of transmission represented in the project.

Tables 25, 26, and 27 again support the point that the random context, regardless of the associated variables, yielded lower numbers of correct responses than did the fixed and isolation contexts. It is true, also, that the "Isolation" context represents the least amount of variations for the conditions represented by the other three conditional variables.

The impeding effects of 40 speakers, as opposed to 10 or 20, seems most marked when words were spoken in the random context regardless of the

transmission level or the number of utterances involved in the tasks.

The variables "Number of Speakers" and "Number of Utterances" tend to confound each other in such a manner as to distribute their effects in interaction with "Transmission" and "Context" in a chaotic fashion making interpretation difficult. The interaction of these two variables for combinations above their minimum levels (10 for number of speakers and 1 for number of utterances) probably represents as much a logistical problem for the examiners as a problem of identification.

This section of the report has been confined to an examination of the significant results of the

basic statistical analysis of the spectrogram data. Percentages have been used to describe the results in a readable manner; however, it should again be emphasized that all statistical procedures were performed on the transformed raw data. Many checks were made on the distribution of the raw and transformed data in order to determine if the assumptions of the statistical design were met. These tests allowed for the results herein contained.

## Conclusions

The conclusions of this report will reference those questions asked in the original voice identification project proposal which are relevant to this phase of the total project.

1. *Are the spectrograms of the same words uttered by a speaker on different occasions similar enough to be identified?*

By referencing those conditions of the voice identification experiment which involved matching a specified spectrogram which was non-contemporary to the spectrogram to be matched of the same speaker speaking the same words, it was found that the identifiers made correct responses 84.72 percent of the time. These results combine with conditions of system-awareness which involved a match (92.01 percent for the closed-match; 77.44 percent for the open-match). The answer to question one would appear to be yes, however, it should be noted that when the spectrogram to be matched was contemporary to the matching spectrogram the percentage correct was 93.98 percent (96.95 percent for the closed-match and 91.02 percent for the open-match). These differences were all found to be significant.

2. *Is there a limitation in the time-elapsd among spectrograms taken of the same speaker at different occasions but speaking the same words?*

There was a significant difference in the ability of the identifiers when engaging in tasks involving contemporary spectrograms (95.21 percent correct) vs. trials involving non-contemporary spectrograms (87.95 percent correct). The time span represented in the project represented the lapse of one month between contemporary and non-contemporary spectrograms.

This variable was held constant for the entire project. The most reasonable conclusion would appear to be that: a one month time lapse among spectrograms taken of the same speaker speaking

the same words will produce significant differences. Further, that these differences will tend to impede voice identification. The determination of limitations on the differences in time lapses must await further research.

3. *Are the spectrograms of the same speaker of the same words spoken on different occasions sufficiently different from the spectrograms yielded by any other speaker?*

This question can only be answered indirectly and by examining the number of errors made by the identifiers under certain conditions of the experiment. If we allow the similarity of spectrograms among two or more speakers would produce more incorrect identifications than dissimilarity, then we can examine at least one aspect of the question. In the closed-match system the percentage of false identification for trials involving contemporary spectrograms was 3.05 percent. For tasks in the closed-match system the percentage of false identifications was 7.99 percent. In the open-match system the percentage of error for trials involving contemporary spectrograms was 8.99 percent while for non-contemporary spectrograms the percentage was 22.57 percent. In the open-no match system; for contemporary trials the percentage of error was 2.33 percent and for the non-contemporary trials 5.59 percent. At all system levels the non-contemporary trials produced significantly more incorrect identifications than the contemporary trials. Part of the explanation of these differences may lie in varying degrees of similarity among the spectrograms.

4. *Does the number of utterances of the same word used for voice identification alter the proportion of correct identifications? If so, in what proportion?*

For the over all project there appeared no significant differences among the percentages of correct identification that could be solely attributed to trials involving one, two or three utterances of the same words. The actual percentages were as follows: one utterance 91.29 percent; two utterances 90.96 percent; three utterances 92.49 percent. There were some significant interactions between the variable "Number of Utterances" and other conditions of the project, however, there emerged no significant pattern to the differences. A slight trend was observed among the interactions for two utterances of the same words to produce lower numbers of correct identifications than one or three utterances. The best that can be said for the variance in

number of utterances is that it exhibits many of the characteristics of a confounding variable when associated with voice identification.

5. *Does the number of speakers to be compared with the unknown one alter the proportion of correct identifications? If so, in what proportion?*

The results indicated a significant difference in terms of correct identifications when the number of speakers involved in the trials were varied from 10, to 20 to 40. In general, as the number increased the percentage of errors also increased. It was observed that for 10 speaker trials the percentage of identifications was 93.30 percent; for 20 speaker trials 91.87 percent and for 40 speaker trials 89.58 percent.

The analysis of variance revealed that the only significant difference was between the 10 and 40 speaker trials. In terms of significant interactions between number of speakers and other variables under examination in the project this same general pattern was observed.

6. *Does the percentage of correct responses of the identifiers change if the spectrograms of the speaker to be identified are among the spectrograms of the known or not?*

The best answer to this question comes from examining the nature of the difference in terms of correct identification for trials conducted under the open-match system vs. those under the open-no match system. In these two instances the identifiers had no knowledge as to whether or not the spectrogram of the speaker to be identified was among those spectrograms involved in the trial or not. Analysis revealed a significant difference between these two systems. Then a match could be made, the examiners were correct only 84.23 percent of the time. When a match could not be made, the examiners were correct 96.04 percent of the time. It is important to note that for the open-match system a correct response involved making a match while for the open-no match system a correct response was represented by a claim that no match existed. In all instances of significant interaction this difference was maintained. It is interesting to examine the effect that awareness could produce on examiners. Under this system (closed-match) the percentage correct was 94.48 percent which differed from the open-match condition, but was not significant when compared to the open-no match system.

The conclusion for this particular aspect of the study is that if no awareness of the possibility of

a match exists within a trial (on the part of the examiners) it makes a great deal of difference.

7. *Does the percentage of correct responses obtained from trained examiners change with changes in environmental conditions and contexts of the uttered clue words for identifications?*

This question had to be answered in two parts. For the project overall environmental conditions were equated to the three types of transmission under which the spectrograms were made. Analysis revealed no significant effect directly attributable to variations in transmission. The actual percentages correct were as follows: directly into a tape recorder ( $\alpha$  transmission): only 92.42 percent; through a telephone line in a quiet environment ( $\beta$  transmission): 91.31 percent and through a telephone line in a noisy environment ( $\gamma$  transmission): 91.02 percent.

The second part of question seven refers to the three context levels represented in the project (I-clue words spoken in isolation, II-clue words spoken in a fixed context, III-clue words spoken in a random context). A significant main effect was observed for context. When the spectrograms were of words spoken in isolation the percentage of correct responses was 95.77 percent. For spectrograms of words spoken in a fixed context the percentage correct dropped to 92.39 percent. For spectrograms of words spoken at random the correct percentage was 86.59 percent. All these differences were found to be significant.

No significant interactions were found involving both levels of transmission and those of context. Nor did these two variables (in combinations) interact to any significant degree with the others under investigation.

The conclusion for question seven is that there is no reason to believe that variations in the environmental conditions under which spectrograms are made will alter to any significant degree the percentage of correct responses for trained examiners. However, variations in the contexts in which the words used for identification purposes appear will have a significant effect on the percentage of responses yielded by trained examiners.

8. *Is a trained person able to recognize whether or not spectrograms of the same word were produced by the same speaker?*

This is, of course, the major question posed by the original voice identification project. The question in its most limited sense asks if examiners are able to make correct matches under varying con-

ditions. It will be recalled that the project involved both identification trials (where match did exist) and elimination trials (where matches did not exist). For identification trials trained examiners were found to make a correct match 89.35 percent of the time. The trained examiners made a false match for the identification trials 4.13 percent of the time and said that no match could be found 6.52 percent of the time. It is interesting to note that when no match could be made (non-identifi-

cation trials) the rate of false identification was 3.96 percent which was found not to be significantly different from the rate of false matching for the identification trials.

In general, the evidence is clear in favor of a trained examiner being able to recognize spectrograms of the same words produced by the same speakers. Further, that when errors are committed, a trained examiner is more apt to claim elimination than to say that a match involves the wrong speaker.

## APPENDIX C

### An Examination of the Nature of Differences Among Voice Identification Examiners

Statistical Report No. 2

By

William B. Lashbrook, Ph.D.

VOICE IDENTIFICATION PROJECT

Michigan State University

## Introduction

The purpose of this report is to examine a subset of data stemming from the Michigan State University voice identification project. The chief area of concern will be possible differences that may occur among the voice identification examiners that can be explained in terms of two variables: *panel type* and *subpanel size*. *Panel type* refers to three nominal classification populations from which the project examiners were drawn. *Subpanel size* refers to three nominal conditions (number of persons involved in a particular identification task) under which data were collected for each panel type. This report is concerned only with identification tasks involving the nine clue words: *it, is, on, you, and, the, I, to, me*.

### The selection of voice identification examiners

The parameters of the project required that the persons used as examiners of spectrograms be drawn from three populations. These populations were as follows:

- 1.0 Females with at least a high school education.
- 2.0 Male (non-police administration majors) MSU students.
- 3.0 MSU students with majors in the Police Administration Department.

Initially 18 persons were hired for the project from each of the three populations. All persons received training in voice identification prior to the collection of project data. Data from the first cycle of the project were collected over an approximate 8 month period. While there was an attrition rate, examination indicated that it was unsystematic and that there was no significant difference ( $p > 0.05$ ) between panel types with respect to it. The logistics of the project allowed three periods of data collection to be considered in the first cycle. These periods involved the completion of all identification traces for one, two, and three utterances of nine clue words. The attri-

tion rate was defined as the number of identifiers of a given type dropping out of the project within a particular data collecting period.

Table 1 represents the attrition rates for the project.

TABLE 1.—Panel Attrition Rate for the Voice Identification Project—First Cycle\*

| Data collection period | Panel Type |    |     |
|------------------------|------------|----|-----|
|                        | I          | II | III |
| 1 .....                | 2          | 0  | 1   |
| 2 .....                | 0          | 2  | 1   |
| 3 .....                | 0          | 3  | 2   |

\*Exact probability = 0.3030.

### The placement of voice identification examiners into subpanels

For the purpose of the project a subpanel of examiners was defined in terms of the interaction of the three panel (population) types with the size of the panel completing a matrix. The three panel types were arranged into subpanels composed of one, two, and three examiners each. Assignment to subpanels was done in an unsystematic manner from one identification task set to another. It was assumed with respect to the placement of the voice identification examiners into subpanels that, within a particular panel type the examiners were interchangeable. In order to check this assumption an attempt was made to determine the reliabilities of the identifiers by type for each data collection period of the project. The fact that there was an attrition rate as noted in Table 1 made it reasonable to determine the reliabilities for each data collection period rather than to combine all the project data in order to determine panel reliability. Table 2 represents the reliability estimates obtained, using an analysis of variance approach to reliability determination (Winer, 1962), for the examiners by panel type for each different number of utterances. The data for the determination of the reliability

estimates involved a transformation of the scores for each of examiner panels for each of 72 replications of the Greco-Latin Square matrix used in the statistical design.

TABLE 2.—Reliability Estimates for Voice Identification Panels By Type

| Data Collection Period (Number of Utterances) | Panel Type |       |       |
|---|------------|-------|-------|
|   | 1.0        | 2.0   | 3.0   |
| 1   | 0.566*     | 0.557 | 0.601 |
| 2   | 0.784      | 0.790 | 0.837 |
| 3   | 0.674      | 0.793 | 0.748 |

\* $r = 0.05, \geq 0.286$ .

Table 3 represents the reliability estimates obtained for the identification panels by size for each data collection period.

TABLE 3.—Reliability Estimates for Voice Identification Subpanels By Size\*

| Data Collection Period (Number of Utterances) | Subpanel Size |             |             |
|---|---------------|-------------|-------------|
|   | 0.1 member    | 0.2 members | 0.3 members |
| 1   | 0.528*        | 0.628       | 0.555       |
| 2   | 0.827         | 0.739       | 0.812       |
| 3   | 0.676         | 0.767       | 0.759       |

\* $r = 0.05, \geq 0.268$ .

For the most part, the reliability estimates reported in Tables 2 and 3 support the assumption of interchangeability of raters. There is no evidence to suggest that this assumption was violated by the fact that there was attrition with respect to the examiners.

#### Statistical analysis

The data pertaining to the panels were analyzed via the use of  $3 \times 3 \times 3$  analysis of variance model with repeated measures on two of the three factors. Factor A was equated to the panel types. This factor did not involve repeated measures. Factors B and D were equated to subpanel size and data collection period respectively and did involve repeated measures. Factor C was a replication factor: 72 replica-

tions of the Greco-Latin matrix under which panel data was collected for each of the three data collection periods related to number of utterances. This yielded a total of 216 scores (0-9) for each of the 9 panels under study. The scores were transformed using a square root transformation (see footnote 1). Table 4 represents the basic model under which the data was analyzed (Winer, 1966).

TABLE 4.— $3 \times 3 \times 3$  ANOVA Design with Repeated Measures on Factors B and D\*

| Source of Variation                          | df                 |
|--|--------------------|
| <i>Between Subjects</i> ..... $na-1$         |                    |
| A—Type of panel                              | a-1                |
| Subj. w. groups                              | a(n-1)             |
| <i>Within Subjects</i> ..... $Na$ ( $bd-1$ ) |                    |
| B—Size Subpanel                              | b-1                |
| AB   | (a-1) (b-1)        |
| B x subj. w. groups                          | a(n-1) (b-1)       |
| D—1, 2, 3, utterances                        | d-1                |
| AD   | (a-1) (d-1)        |
| D x subj.                                    | a(n-1) (d-1)       |
| BD   | (b-1) (d-1)        |
| ABD  | (b-1) (d-1)        |
| BD x subj. w. groups                         | (a-1) (b-1) (d-1)  |
|  | a(n-1) (b-1) (d-1) |

\*Assumes A, B, and D as fixed factors.

It is important to note that the "subjects within groups" notation referred to in Table 4 represents the three subpanel sizes within the three panel types. In combination these variables define the nine identification subpanels (N) used in the Voice Identification Project.

#### Results

The study herein reported represents data stemming from nine subpanels attempting 17,496 voice-print identifications under many conditions. Over all the examiners were correct 16,023 times for a percentage of 91.58. The stated purpose of this study was to examine the differences among the examiners that could be explained in terms of panel type and subpanel size. It is important to note that the variables under consideration pertain to the examiners and not to the conditions under

which the identifications were made. This latter analysis constitutes another aspect of the total Voice Identification Project.

Table 5 represents the percentages correct for the nine subpanels classified by type over the three data collection periods (Number of Utterances).

TABLE 5.—Percentage of Correct Identifications for Panels by Type

| Number of Utterances Used | Panel Type |       |       |
|---------------------------|------------|-------|-------|
|                           | 1.0        | 2.0   | 3.0   |
| 1                         | 90.12      | 92.80 | 90.95 |
| 2                         | 90.74      | 92.13 | 90.02 |
| 3                         | 93.06      | 91.46 | 92.95 |

Table 6 represents the percentages correct for the nine subpanels classified by size over the three data collection periods.

TABLE 7.—Percentage of Correct Identifications for Panels By Type and Size

| Number of Utterances | Subpanel Size |       |       |       |       |       |       |       |       |
|----------------------|---------------|-------|-------|-------|-------|-------|-------|-------|-------|
|                      | 1-1           | 1-2   | 1-3   | 2-1   | 2-2   | 2-3   | 3-1   | 3-2   | 3-3   |
| 1                    | 91.36         | 87.65 | 91.36 | 91.20 | 90.28 | 96.91 | 91.05 | 92.28 | 89.51 |
| 2                    | 92.44         | 85.65 | 94.14 | 91.98 | 93.21 | 91.20 | 88.12 | 91.36 | 90.59 |
| 3                    | 92.44         | 90.74 | 95.99 | 89.82 | 90.43 | 94.14 | 91.20 | 91.67 | 95.98 |

TABLE 8.— $3 \times 3 \times 3$  ANOVA Results with Repeated Measures on Factors B and D

| Source of Variation      | df  | Mean Square | F ratio   |
|--------------------------|-----|-------------|-----------|
| <i>Between Subjects:</i> |     |             |           |
| A—Panel Type             | 2   | 0.12597     | 0.2218    |
| Subj. w. groups          | 213 | 0.56786     |           |
| <i>Within Subjects:</i>  |     |             |           |
| B—Subpanel Size          | 2   | 1.61030     | **19.0070 |
| AB                       | 4   | 0.67308     | **7.9447  |
| B x subj. w. groups      | 426 | 0.08472     |           |
| D—Number of Utterances   | 2   | 0.40739     | *3.0883   |
| AD                       | 4   | 0.36904     | *2.8066   |
| D x subj. w. groups      | 426 | 0.13149     |           |
| BD                       | 4   | 0.18245     | *2.6659   |
| ABD                      | 8   | 0.36500     | **5.7331  |
| BD x subj. w. groups     | 352 | 0.06844     |           |

\* $p. F \leq 0.05$ .

\*\* $p. F \leq 0.01$ .

TABLE 6.—Percentage of Correct Identifications for Subpanels by Size

| Number of Utterances Used | Subpanel Size |           |           |
|---------------------------|---------------|-----------|-----------|
|                           | 1 member      | 2 members | 3 members |
| 1                         | 91.20         | 90.07     | 92.59     |
| 2                         | 90.84         | 90.07     | 91.98     |
| 3                         | 91.15         | 90.95     | 95.37     |
|                           | 91.06         | 90.36     | 93.31     |

Table 7 represents the percentages correct for the nine subpanels classified by type and size. It should be remembered that the  $3 \times 3$  panel-size combination define the nine subpanels used in the study.

In order to determine possible differences in the panels attributable to either type or size a  $3 \times 3 \times 3$  repeated measures analysis of variance was run. This analysis was in accord with that discussed in

an earlier section of this report. Data used for this analysis were the transformed raw score for each panel for the 72 replications of the matrix for each of three different numbers of utterances used.

The results summarized in Tables 5 through 8 tend to support the following conclusion:

(1) No significant differences could be found between the examiners panel types attributable to the populations from which the examiners were drawn.

(2) There was a significant difference between the subpanel types. Further analysis (using Duncan's Multiple Range technique) revealed that the three member subpanels had a significantly more ( $p. < 0.01$ ) correct identifications than the single or two member subpanels. See Table 9.

(3) The nine subpanels involved in the project were not equal with respect to the number of correct identifications. While this finding would be

expected, given the fact that the three member panels did better than the one and two member panels, further analysis revealed that subpanel 1.2 with two members had significantly less ( $p < 0.05$ ) correct identifications than any other combination of identifiers. Further analysis also revealed that while the subpanel 3.3 with three members did better than one and two member subpanels of the same type it did not do as well as the other three member panels. See Table 10.

(4) There was a significant difference ( $p < 0.05$ ) in terms of the number of correct identifications between the number of utterances examiners used. Analysis revealed that the number of correct identifications for one utterance (91.29 percent) and two (90.96 percent) were less than for three utterances (92.49 percent), but that the only significant difference ( $p < 0.05$ ) was between utterances three and two.

TABLE 9.—Differences Between Subpanel Sizes via Multiple Range Technique

| Panel Size | Means <sup>1</sup> | Subpanel Size |           |           |
|------------|--------------------|---------------|-----------|-----------|
|            |                    | 1 member      | 2 members | 3 members |
| 1 member   | 5.8835             | 5.8835        | 5.8581    | 5.9543    |
| 2 members  | 5.8581             | ..            | ..        | ..        |
| 3 members  | 5.9543             | **            | **        | ..        |

\*\* $p \leq 0.01$ .

<sup>1</sup>The means in Table 9 are for the transformed raw data; the same data used in the analysis of variance reported in Table 8.

(5) Significant instructions were found ( $p < 0.05$ ) between subpanels by type and size for the three different numbers of utterances used. Further examination of these findings revealed that:

(a) Panel Types 1 and 3 did not achieve the same degree of accuracy as panel Type 2 until the two and three utterances respectively were used. Once obtained, however, there was no significant difference between panels by type.

(b) The two member subpanels consistently across the three different number of utterances used were not as accurate as the three member subpanels. The one member subpanels were not as consistently different from the three member subpanels as the two member subpanels.

(c) As was indicated previously there was a significant difference among the identification pan-

els and this finding extended across the different numbers of utterances used. This variation between panels was to be expected, however, further analysis revealed that there was no consistent pattern to the variations and that result number (3) remains the most cogent thing to be asserted about the individual panels.

### Discussion

The purpose of this section of the report is to address some of the questions posed in the original Voice Identification proposal.

*Question 1: After examiners have been trained, i.e., their learning curves have reached a ceiling, or a relative plateau, what is the percentage of correct responses that can be obtained?*

The evidence indicates that a trained Voice Identification examiner can be expected to make correct identification about 92 percent of the time. While this rate will vary between conditions under which data are gathered, a strong case ( $t = 165.32$ ,  $p < 0.0001$ ) can be made for trained personnel being able to, with a high degree of accuracy, make voice identifications.

*Question 2: What category of persons is most suitable for training as a Voice Identification examiner according to sex, age and background?*

This question cannot be answered in its entirety. Based on the three populations from which the Voice Identification examiner were drawn, there were no significant differences between examiners.

*Question 3: Do Voice Identification examiners perform better working alone or in a team?*

In general, the examiners, when placed in teams (subpanels of two and three members) did slightly better (91.84) than examiners working individually (91.07). This finding, though consistent, was not statistically significant. It was found, however, that a subpanel of three members was significantly better (93.31) than subpanels of two members (90.37) or individuals working singly (91.07). Probability less than 0.05.

*Question 4: What would be the most efficient size of examiners team?*

The term efficiency makes this question difficult to answer. If accuracy is really the issue, then three examiner teams should be recommended. If availability is the issue then one examiner working alone would appear to be as accurate as two-member teams.

TABLE 10.—Differences Between Panels by Type and By Size via Multiple Range Technique

| Subpanel Size | Means <sup>1</sup> | Type 1   |           |           | Type 2   |           |           | Type 3   |           |           |
|---------------|--------------------|----------|-----------|-----------|----------|-----------|-----------|----------|-----------|-----------|
|               |                    | 1 member | 2 members | 3 members | 1 member | 2 members | 3 members | 1 member | 2 members | 3 members |
| Panel Type 1: |                    |          |           |           |          |           |           |          |           |           |
| 1 member      | 5.9177             | ..       | **        | ..        | ..       | ..        | ..        | ..       | ..        | ..        |
| 2 members     | 5.7805             | **       | ..        | ..        | **       | ..        | **        | ..       | **        | ..        |
| 3 members     | 5.9720             | ..       | **        | ..        | ..       | **        | ..        | **       | ..        | ..        |
| Panel Type 2: |                    |          |           |           |          |           |           |          |           |           |
| 1 member      | 5.8765             | ..       | ..        | ..        | ..       | ..        | ..        | ..       | ..        | ..        |
| 2 members     | 5.8902             | ..       | **        | ..        | ..       | **        | ..        | ..       | **        | ..        |
| 3 members     | 5.9775             | ..       | **        | ..        | ..       | **        | ..        | ..       | **        | ..        |
| Panel Type 3: |                    |          |           |           |          |           |           |          |           |           |
| 1 member      | 5.8562             | ..       | **        | ..        | ..       | **        | ..        | ..       | **        | ..        |
| 2 members     | 5.9038             | ..       | **        | ..        | ..       | **        | ..        | ..       | **        | ..        |
| 3 members     | 5.9134             | ..       | **        | ..        | ..       | **        | ..        | ..       | **        | ..        |

\*\* $p \leq 0.01$ .

\* $p \leq 0.05$ .

<sup>1</sup>Means in Table 10 are for the transformed raw data; the same data used in the analysis of variance reported in Table 8.

TABLE 11.—Differences Between Number of Utterances Examined via Multiple Range Technique

| Data Collection Period (Number of Utterances) | Means <sup>1</sup> |              |              |
|---|--------------------|--------------|--------------|
|   | 1 utterance        | 2 utterances | 3 utterances |
| 1   | 5.8916             | 5.8779       | 5.9265       |
| 2   | 5.8916             | 5.8779       | 5.9265       |
| 3   | 5.8916             | 5.8779       | 5.9265       |

\*p. ≤ 0.05.

<sup>1</sup>The means in Table 11 are for the transformed raw data; the same data used in the analysis of variance reported in Table 8.

## An Examination of the Types of Errors Made by Examiners

By

William B. Lashbrook, Ph.D.

Statistical Report No. 3

The purpose of this report is to present a detailed analysis of the type of errors committed by the examiner panels involved in the project. The data for this report are frequencies of errors of a particular type. It is important to note that the classification of errors by type is nominal and that it does not represent the same level of measurement as data referring to the number of correct responses. There were four types of error considered:

- Type A—A match existed, but the examiners made an incorrect match. Open tasks.
- Type B—A match existed, but the examiners failed to make a match. Open tasks.
- Type C—No match existed, but the examiners said that a match existed. Open tasks.
- Type D—A match existed, but the examiners made an incorrect match. Closed tasks.

The difference between Type A and D errors depended upon the knowledge of the task processed by examiners. For Type A errors the examiners had no knowledge as to whether or not the task involved a match. For Type D errors the examiners knew that possible match did exist within the task.

### The distribution of errors

Because of the nature of the data (frequency) a decision was made to use the total errors of a particular type as a basis for examining differences between examiners attributable to panel type, sub-panel size and number of utterances. The total was assumed to represent an ordinal level of measurement (identification panels could be ranked according to the frequency with which they committed errors of a particular type). Statistical analysis involved a three way, distribution free analysis

of variance technique. A separate analysis was run for data from each of the error types.

Table 1 represents the total number of errors for each voice examiner panel by type, subpanel size and number of utterances. It will be recalled that the total number of errors for the examiners (9 wds.) was 1473 out of a possible 17,496. This ratio reduces itself to a percentage of error of 8.42. Additional analysis revealed that of the total number of errors committed:

1. 8.96% were of Type A
2. 51.60% were of Type B
3. 15.68% were of Type C
4. 23.76% were of Type D

TABLE 1.—Frequency of Errors by Type

| Error Type | Number of Utterances | Panels |     |     |     |     |     |     |     |     |
|------------|----------------------|--------|-----|-----|-----|-----|-----|-----|-----|-----|
|            |                      | 1-1    | 1-2 | 1-3 | 2-1 | 2-2 | 2-3 | 3-1 | 3-2 | 3-3 |
| A          | 1                    | 1      | 8   | 2   | 4   | 9   | 2   | 9   | 5   | 5   |
|            | 2                    | 8      | 15  | 4   | 6   | 1   | 6   | 12  | 7   | 10  |
|            | 3                    | 2      | 4   | 0   | 0   | 1   | 1   | 6   | 4   | 0   |
| B          | 1                    | 28     | 31  | 27  | 35  | 21  | 12  | 26  | 25  | 36  |
|            | 2                    | 23     | 38  | 22  | 33  | 30  | 26  | 26  | 33  | 33  |
|            | 3                    | 28     | 36  | 13  | 40  | 47  | 25  | 26  | 25  | 15  |
| C          | 1                    | 5      | 13  | 14  | 7   | 10  | 3   | 12  | 4   | 7   |
|            | 2                    | 11     | 16  | 7   | 9   | 5   | 17  | 17  | 6   | 10  |
|            | 3                    | 8      | 2   | 7   | 13  | 5   | 1   | 8   | 10  | 4   |
| D          | 1                    | 22     | 28  | 13  | 11  | 23  | 3   | 11  | 16  | 20  |
|            | 2                    | 7      | 24  | 5   | 4   | 8   | 8   | 22  | 10  | 8   |
|            | 3                    | 11     | 18  | 6   | 13  | 9   | 11  | 17  | 15  | 7   |

Table 2 represents an analysis of variance for the total frequencies of Type A errors.

TABLE 2.—Three Way AOV for Frequency of Type A Errors

| Source of Variation    | df | x <sup>2</sup> | Sig. |
|------------------------|----|----------------|------|
| A—Panel Type           | 2  | 4.14205        | n.s. |
| B—Subpanel Size        | 2  | 2.30114        | n.s. |
| C—Number of Utterances | 2  | 5.98295        | n.s. |
| AB                     | 4  | 3.22159        | n.s. |
| BC                     | 4  | 4.14205        | n.s. |
| AC                     | 4  | 2.30114        | n.s. |
| ABC                    | 8  | 13.80682       | n.s. |

The results indicate no significant difference between examiners attributable to panel type, sub-panel size or number of utterances used. The effect due to a number of utterances approaches significant ( $p. = 0.0502$ ) and is explainable in terms of a smaller number of Type A errors committed in the three utterance tasks.

Analysis of type B errors

Table 3 represents an analysis of variance for the total frequencies of Type B errors.

TABLE 3.—Three Way AOV for Frequency of Type B Errors

| Source of Variation    | df | x <sup>2</sup> | Sig. |
|------------------------|----|----------------|------|
| A—Panel Type           | 2  | 2.22527        | n.s. |
| B—Subpanel Size        | 2  | 2.22527        | n.s. |
| C—Number of Utterances | 2  | 0.44505        | n.s. |
| AB                     | 4  | 9.34615        | n.s. |
| BC                     | 4  | 5.78571        | n.s. |
| AC                     | 4  | 5.78571        | n.s. |
| ABC                    | 8  | 22.25275       | **   |

\*\*p. x<sup>2</sup>. ≤ .01.

The results indicate no significant difference attributable to direct variations of panel type, size or number of utterances. There was a significant interaction (ABC) which merely supports the position that the 9 subpanels differed among themselves as to the rate of commitment of Type B errors over the three utterances, but that there was no consistent patterns to these differences.

The interaction between panel type and sub-panel size approached significance ( $p. = 0.0534$ ). Most of this difference seems to be accountable in terms of the fact that with respect to Type B errors

panels of one, two, and three members showed more variations within panel Types 1 and 2 than did panels of the Type 3.

Analysis of type C errors

Table 4 represents an analysis of variance for the total frequencies of Type C errors.

TABLE 4.—Three Way AOV for Frequency of Type C Errors

| Source of Variation    | df | x <sup>2</sup> | Sig. |
|------------------------|----|----------------|------|
| A—Panel Type           | 2  | 0.44505        | n.s. |
| B—Subpanel Size        | 2  | 4.00549        | n.s. |
| C—Number of Utterances | 2  | 1.33516        | n.s. |
| AB                     | 4  | 1.33516        | n.s. |
| BC                     | 4  | 5.78571        | n.s. |
| AC                     | 4  | 1.33516        | n.s. |
| ABC                    | 8  | 21.36264       | **   |

\*\*p. ≤ .01.

The results indicate no significant difference attributable to direct variation of panel type, sub-panel size or number of utterances. There was a significant interaction (ABC) which supports the assertion that the individual subpanels differed among themselves as to the rate of commitment of Type C errors over the three utterances.

Analysis of type D errors

Table 5 represents an analysis of variance for the total frequencies of Type D errors.

TABLE 5.—Three Way AOV for Frequency of Type D Errors

| Source of Variation    | df | x <sup>2</sup> | Sig. |
|------------------------|----|----------------|------|
| A—Panel Type           | 2  | 1.35000        | n.s. |
| B—Panel Size           | 2  | 4.05000        | n.s. |
| C—Number of Utterances | 2  | 8.55000        | *    |
| AB                     | 4  | 2.25000        | n.s. |
| BC                     | 4  | 0.45000        | n.s. |
| AC                     | 4  | 0.45000        | n.s. |
| ABC                    | 8  | 12.6000        | n.s. |

\*p. ≤ .05.

The result indicate a significant difference between utterances for Type D errors. The result can be best explained by the fact that there were significantly more Type D errors for one utterance than for two utterances, ( $T = 8, p. ≤ .05$ ).

Discussion

The analysis of the type of errors committed by the Voice Identification panels was performed as a further check on the nature of possible differences between panels due to their type or size. In no case were such differences found. Variations between panels and between number of utterances were to be expected. Possible explanations of the variations when found to be significant appear to be an artifact of the project.

**END**