

Automatic Speaker Recognition Using Vocoded Speech

Stephanie S. Everett

Naval Research Laboratory, Washington, DC

ABSTRACT

Automatic speaker recognition (ASR) offers potential benefit for numerous applications, including identification of users of communication channels such as the telephone and channels using processed or vocoded speech. Currently the listener must subjectively determine whether the person on the other end of the line is who he or she claims to be. However, when the speech is processed, recognition of voices can be very difficult for human listeners. A series of tests was conducted to evaluate the feasibility of automatic speaker recognition with processed or vocoded speech. The analog outputs of six different voice processors were used as input to a real-time ASR system. Recognition accuracy results for the processed speech were 70% to 95% using a 2500 Hz bandwidth input filter, and 75% to 95% using a 4000 Hz input filter. These results indicate that ASR using vocoded speech is definitely feasible, though further research is needed to determine which speech parameters are best suited for use with each voice processor.

INTRODUCTION

Over the last 20 years, the research community has spent a great deal of time and money developing methods of recognizing speakers automatically based on voice input alone. Such systems could be of considerable benefit in numerous applications, including access control for restricted areas and information, communication security, and verification of computer users through terminals accepting voice input.

One logical application of automatic speaker recognition (ASR) to the problem of communication security is the task of verifying the identity of speakers over the telephone or over communication channels using processed or vocoded speech. Currently the listener must subjectively identify the speaker based only on his or her recollection of the person's voice. If the listener has never spoken with the person over this particular type of channel, or has never met him or her before, it can be nearly impossible to know for sure whether the person is who he or she claims to be.

Though some research has been done on ASR using telephone (or telephone-quality) speech, very little is known about ASR using processed or

vocoded speech. All voice processing systems have been carefully designed to maximize the quality and intelligibility of the synthesized speech. However, the analyses performed to allow bandwidth compression or encoding of the speech signal frequently remove or distort certain characteristics of the original speech. It is not known what effects, if any, this processing has on the portions of the speech signal relevant to speaker identity. A recent study showed that over a 2400 bits per second (bps) linear predictive coding (LPC) voice channel people could identify the familiar voices of their coworkers only about 70% of the time [1]. Other studies have shown that automatic speaker recognition systems actually perform better than human listeners under certain conditions [2].

As a preliminary evaluation of the potential for performing ASR with vocoded speech a series of tests was conducted where the outputs from several voice processors were used as input to an existing ASR system.

ASR SYSTEM DESCRIPTION

The ASR system used in this feasibility evaluation was selected to meet the following requirements:

- perform text-independent speaker identification,
- operate in real time,
- recognition procedure completely automatic (i.e., no hand-marking of speech segments, determination of silent periods, etc.),
- capable of providing recognition results using less than 3 seconds of input speech,
- capable of operating on a set of at least 20 speakers, including both males and females, and
- have a reported overall recognition accuracy of at least 90% using high-quality input speech for a set of not less than 10 speakers.

The ASR system is based on LPC analysis of the speech, and uses a multiple parameter recognition algorithm. This involves calculating the mean, variance, and covariance for all the frames in the training utterance. The mean vector for the reflection coefficients is then calculated over the test utterance and compared to the mean vectors in the model using the Mahalanobis distance measure. A more complete description of this system may be found in Reference 3.

11.2.1

Text-independent speaker recognition was required for this series of tests because it was felt that this would give a better indication of how well the voice processors reproduce all sounds for a variety of voices, rather than only those sounds in each individual's code word or phrase as with text-dependent recognition. The task chosen was speaker identification, where the ASR system must choose the speaker's identity from a set of known voices. This is somewhat more difficult than speaker verification, where the ASR system is given an identity claim and need only accept or reject that claim based on a preset distance measure threshold.

TEST DESCRIPTION

To test the performance of the ASR system an audio source tape was generated containing five phonetically balanced sentences from each of 20 speakers (10 males, 10 females) for a total of 100 different utterances. For each test condition the first two sentences from a given speaker were used for model generation and the remaining three sentences were used as test utterances.

The tape was played through each of the six different voice processors listed in Table 1. The processed output speech was recorded for use as input to the ASR system. The system was not specifically tailored to any of the voice processors used in this investigation, though an ideal transmission channel (back-to-back) was used in each case.

One other test condition was designed to simulate performance of speaker recognition prior to resynthesis at the receiver of a 2400 bps LPC voice processing system (assuming an ideal trans-

Table 1. Coding methods and transmission rates of voice processing systems used for evaluating the feasibility of performing ASR using vocoded speech.

Coding Method	Data Rate (bps)
Pulse Code Modulation (PCM)	64,000
Continuously Variable Slope Delta (CVSD)	16,000
Residual-Excited Linear Prediction (RELP)	16,000 & 9,600
Adaptive Predictive Coding (APC)	9,600
Linear Predictive Coding (LPC)	2,400

mission channel). This test used the clear text source tape and required a slight modification of the ASR algorithm to include quantization of the LPC parameters. This algorithm uses 20 acoustic features derived from LPC analysis: 10 reflection coefficients and 10 cepstral coefficients. In the standard algorithm the LPC prediction and reflection coefficients are derived from the autocorrelation coefficients using Levinson's recursion. The prediction coefficients are then converted to cepstral coefficients as shown in Figure 1a [4]. For the quantized LPC parameters test condition the ASR system was modified as shown in Figure 1b. The reflection coefficients

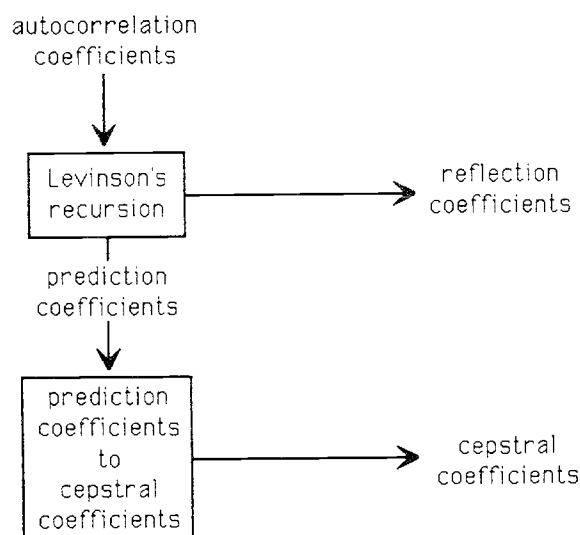


Fig. 1(a)

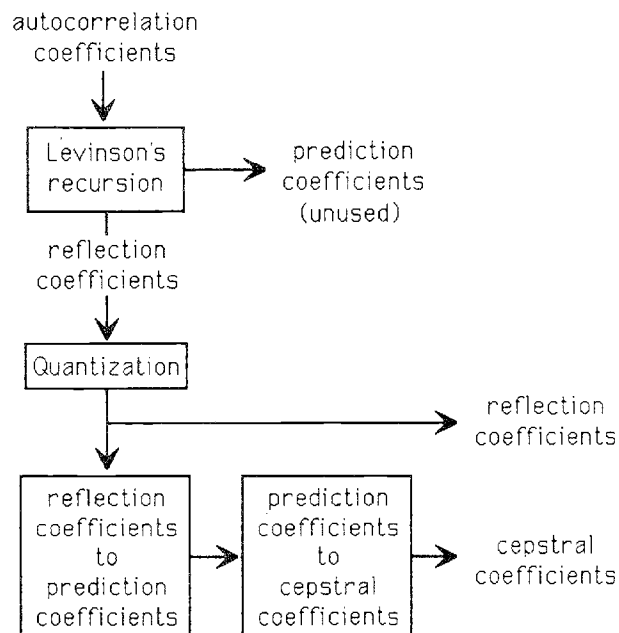
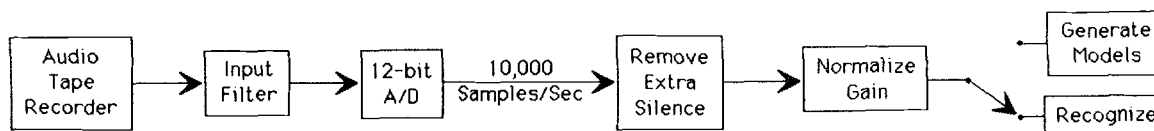


Fig. 1(b)

Figure 1. ASR signal processing procedure. Figure 1(a) shows the usual ASR algorithm; 1(b) shows the modifications used for the quantized LPC parameters test condition [4].

11.2.2

Figure 2. Block diagram of the testing configuration.



were quantized according to the DoD standard, then converted first to prediction coefficients and finally to cepstral coefficients. The modified ASR algorithm was used in both training and recognition for this one test condition.

A block diagram of the testing configuration is shown in Figure 2. To perform the tests the audio signal from the tape recorder was passed through an input filter and into a 12 bit analog-to-digital (A/D) converter. The signal was sampled at 10,000 samples per second. For each test condition the two or three sentences of each speaker's training or test set were grouped together and digitized as a single file, giving 40 such files for each condition. Excess silence at the beginning and end of the files was removed to reduce the amount of memory required to store the data. This was done using a simple energy-thresholding algorithm, leaving approximately one

half second of silence as padding outside the detected endpoints to ensure that no speech sounds were removed. In order to use the entire 12 bit range of the digital representation the samples in each file were multiplied by a factor calculated to clip 1% of the samples.

Two series of tests were conducted using input filters of 2500 Hz and 4000 Hz bandwidth, respectively. The results for both sets of tests are discussed in the next section.

TEST RESULTS

Table 2 summarizes the results of the ASR system performance tests described in the preceding section. The results are based upon all of the available speech data. Tests were performed only within conditions, i.e., recogni-

2500 Hz Bandwidth Input Filter								Speaker Number	4000 Hz Bandwidth Input Filter							
Clear	PCM 64000	REL P 16000	CVSD 16000	REL P 9600	APC 9600	LPC 2400	Quant 2400		Clear	PCM 64000	REL P 16000	CVSD 16000	REL P 9600	APC 9600	LPC 2400	Quant 2400
.	10	1
.	.	.	1	.	1	1	.	2	1	.
.	3
5	5	5	6	10	6	5	10	4	.	5	5	6	5	5	5	.
.	13	10	10	5	10	.	.	.	10	.	10	6
.	6
.	7
.	.	.	14	.	14	.	.	8	.	.	.	14	.	14	.	.
.	9
.	10
.	19	.	20	11	14	.	.
14	.	.	14	.	14	14	.	12	.	.	.	14	.	14	14	.
.	.	.	16	13	.	.	20	19	19	.	.	20
.	14
.	15
.	16
.	.	.	19	17
.	18
.	19
.	16	20	19	.	.
2	2	1	6	1	6	4	4	Total Errors	1	1	2	4	3	5	4	2
90	90	95	70	95	70	80	80	Percent Recognition	95	95	90	80	85	75	80	90

Table 2. Summary of results for ASR system performance tests. Speakers 1-10 are males, 11-20 are females. A dot indicates that the speaker was identified correctly under the given condition; a number indicates an identification error. The number entered is the identified speaker, while the row number is the actual speaker. The total number of errors is shown at the bottom of the table for each condition, along with the recognition accuracy measured with a granularity of 5%.

tion tests were always conducted using the same input speech condition as had been used in the generation of the speaker models.

With the 2500 Hz bandwidth input filter recognition accuracy for the clear speech was 90%. Accuracy for the processed speech ranged from 70% to 95%. LPC parameter quantization alone lowered the accuracy to 80%. It is interesting that synthesized speech from the two residual-excited linear predictive (RELP) systems gave better results than the original clear speech, though the difference is not significant. It is conceivable that the speech processing performed by these systems actually improves speaker discrimination, but further investigation would be needed to confirm this hypothesis.

With the 4000 Hz bandwidth filter the range of scores was roughly the same, but recognition accuracy for individual conditions changed by as much as 10%. Accuracy improved or remained the same for all vocoders except the RELP processors which both degraded somewhat. Accuracy for clear speech and the pulse code modulated (PCM) processor improved to 95%, and accuracy using the quantized LPC parameters improved to 90%. It should be noted that Speaker 4 was recognized correctly under only two conditions. There is no readily apparent explanation for this, however it was verified that there were no procedural errors and that there was no unusual imbalance in the phonetic content of the training and testing sentences for this person.

It is curious that recognition accuracy for the RELP processors should degrade slightly with the wider bandwidth filter, while accuracy for all the other vocoders improved or remained the same. One possible explanation for this is the type of signal processing performed in these systems. These processors synthesize an accurate representation of the low frequency portion, or baseband, of the speech signal and approximate the high frequency portion using information in the baseband residual. *For this reason the high frequency spectrum of the synthesized speech often differs significantly from that of the original speech, particularly for voiced sounds.* This inaccuracy is not noticeable to human listeners, but could have a marked effect on an ASR algorithm that weights all portions of the speech spectrum equally. The speaker identification information is therefore contained primarily in the baseband for these coders, and the amount of additional information in the high frequencies is outweighed by the inaccuracy of the spectral representation. With this in mind, it is not surprising that these processors do not benefit from the wider bandwidth. The 16,000 bps RELP degraded less than the 9600 bps because it has a wider baseband.

The relatively poor performance of the 16000 bps CVSD and the 9600 bps APC vocoders can be attributed primarily to the wideband quantization noise generated by these algorithms. Though not overly distracting to human listeners, this noise corrupts the speech signal enough that the subsequent LPC analysis is unable to extract the necessary speaker identification characteristics.

This does not necessarily indicate that these vocoders are unsuitable for ASR. It does suggest, however, that a non-LPC based approach might be better suited for use with voice processing systems of this sort.

High bit rate channels such as the 64,000 bps PCM system tested here probably generate speech of sufficient quality to be used with any type of ASR algorithm.

SUMMARY

Automatic speaker recognition offers great potential benefit for a variety of applications, including communication security. The ability to perform ASR using processed or vocoded speech would allow verification of communication channel users and access control for the channel itself.

A series of tests was conducted to evaluate the potential for performing ASR using processed input speech. In these tests the analog output of six different voice processors was used as input to an existing real-time ASR system. The vocoders used a variety of processing algorithms and had data transmission rates ranging from 2400 to 64,000 bps.

The results of the tests indicate that ASR using processed or vocoded speech is definitely feasible. However, further research is required to determine exactly which parameters, analysis methods and distance measures produce optimum ASR results for a given voice processing system or type of system.

Despite the many years of research in speech acoustics, relatively little is known about the speaker-specific characteristics of voices, or about which cues human listeners use in identifying speakers. Further investigation in these areas would also lead to increased recognition accuracy for ASR systems.

ACKNOWLEDGEMENTS

The author wishes to acknowledge the contributions to this study of Alan Higgins and Joe Naylor of ITT Defense Communications Division in San Diego who performed the testing and provided numerous insights into this problem.

REFERENCES

- [1] A. Schmidt-Nielsen and K. R. Stern, "Identification as a Function of Familiarity of Known Voices Talking Over an Unprocessed Channel and an LPC Voice Processor," NRL Memorandum Report 5382, July 1984.
- [2] A. E. Rosenberg, "Listener Performance in Speaker Verification Tasks," IEEE Trans. Audio Electroacoustics, Vol. AU-21, pp. 221-225, 1973.
- [3] E. H. Wrench, Jr., "A Realtime Implementation of a Text Independent Speaker Recognition System," Proceedings IEEE 1981 ICASSP, pp. 193-196, March 1981.
- [4] A. Higgins and J. Naylor, "Final Report on Contract N00014-84-C-2130," ITT Defense Comm. Div., San Diego, CA, July 1984.