

# De-Identification of Speech

**Bhaskar Adavi**

Carnegie Mellon University  
Pittsburgh, PA 15213  
*badavi@andrew.cmu.edu*

**Bharat Kandoi**

Carnegie Mellon University  
Pittsburgh, PA 15213  
*bkandoi@andrew.cmu.edu*

**Prashant Malani**

Carnegie Mellon University  
Pittsburgh, PA 15213  
*malani@cmu.edu*

**Sahil Sapre**

Carnegie Mellon University  
Pittsburgh, PA 15213  
*sahilsapre@cmu.edu*

## Abstract

In this paper, we describe an efficient method of de-identification of speech such that the transformation from the source speech is furthest away from the source features, yet fully intelligible. We have designed a speaker ID system that is 91.8% accurate in identifying 20 utterances spoken by 30 speakers - 23 standard American newsreaders, 5 speakers from the CMU Arctic database, and 2 native Indian speakers. We then de-identify these voices using voice conversion such that the speaker ID system trained to correctly identify these speakers gets confused when presented with the de-identified voices as input.

## 1 Introduction

Telephonic dialogues are a popular means of information retrieval. However, they are also a hidden channel for the invasion of the end-users' privacy. Companies could maintain records of what the users said, without their consent or knowledge. For example, recorded conversations between doctors and patients could potentially reveal confidential information about the patients. One way of protecting people from this invasion of privacy is to de-identify the speaker's voice such that the speech still sounds natural and fully intelligible, yet does not reveal information about the identity of the speaker.

In order to truly test the efficacy of any de-identification mechanism, an accurate speaker ID system is required. Such a system should, once trained to correctly identify a set of speakers, be able to identify any new voice input from these speakers with high accuracy. Such a system should also maintain its accuracy when the set of speakers have similar voice characteristics, like pitch and accent.

The organization of this paper is as follows. In section 2, a brief summary of the various feature extraction methods employed in both the speaker ID system as well as voice transformation is discussed. In section 3, the working principle behind the speaker ID system is explained, along with the test results from simulations. In section 4, the Voice Conversion technique is presented, and in section 5 the method used to select the 'best' transformation for de-Identification is outlined. In section 6, a summary for the results of the de-identification algorithm are presented and section 7 contains concluding remarks.

42

## 43 **2 Feature Extraction**

44

### 45 **2.1 Mel-frequency cepstral coefficients (MFCC)**

46 The mel-frequency cepstrum (MFC) is a representation of the short term power spectrum of  
47 a sound, based on a linear cosine transform of a log power spectrum on a nonlinear mel scale  
48 of frequency. Here, the frequency bands are equally spaced on the mel scale, which better  
49 approximates the human auditory response than linearly-spaced bands.

50

### 51 **2.2 Fundamental Frequencies (F0)**

52 Frequency measure that denotes the number of times the vocal folds open and close per  
53 second. The typical male F0 value is 120 Hz, while the value for females is 180 Hz

54

## 55 **3 Speaker Identification System Design**

56 The basic methodology behind the Speaker ID system is to create statistical models that  
57 most closely approximate the voices to be tested, and then compare inputs against these  
58 models to find the model that most closely resembles them. We have improved upon a basic  
59 MFCC-based speaker ID system [1], adding weights for clustering of MFCCs and for F0s.

60

### 61 **3.1 Generation of Gaussian Mixture Models (GMMs)**

62 From each input file that is used to train the system for a particular speaker, the MFCCs and  
63 F0s are used to create GMMs that best approximate these parameters. Using an Expectation  
64 Maximization algorithm, GMM parameters such as means, variances and weights are  
65 obtained which best represent the speakers voice, for each audio file.

66

### 67 **3.2 Clustering of MFCCs**

68 To further improve the accuracy of the speaker ID system, clustering is performed on the  
69 sets of MFCCs for each of the training audio files for a speaker using the k-means method.  
70 The mean of the centroids obtained from each file is computed and stored. The clustering  
71 method used was developed by Esfandiar Zavarehei.

72 The GMM models for the MFCCs and F0s, along with the results of clustering, form the  
73 code book used for testing.

74

### 75 **3.3 Training Phase**

76 To train the speaker ID, we use the spectral parameters of each of the training input files,  
77 and obtain a log likelihood estimate of how close a particular parameter is to each of the  
78 GMMs constructed for that parameter. This process is performed for both MFCCs and F0s.  
79 The mean and standard deviation of the log likelihood estimates is then measured for both  
80 MFCC and F0 values.

81 This process can be repeated for multiple speakers to create a bank of speaker models to test  
82 against.

83

### 84 **3.4 Testing Phase**

85 To test an arbitrary input file for a particular speaker, the extracted MFCCs and F0s of the  
86 test file are sent into the system. The MFCCs are compared against a speakers GMMs  
87 (calculated in section 3.1) and the mean is taken of the log likelihoods that results for each  
88 GMM. A score is then assigned to the average value based on how many standard deviations  
89 it is away from the mean calculated during training. The standard deviation used here is the  
90 one obtained during training. This process is repeated for the F0 values of the test input.

91 The MFCCs of the test input are then clustered to find a centroid. The distance of this

92 centroid is then measured from the mean of the centroids for the model, which was  
93 calculated in section 3.2.

94 Assuming the MFCCs, F0s and centroidal distances are orthogonal to each other, a  
95 Euclidean distance measure is obtained for each speaker model. The system chooses the  
96 model with the least score as its conclusion about the identity of the test voice.

97

#### 98 **4 Voice Transformation**

99 The methodology used to de-identify speech is to transform the input voice such that it is as  
100 far as possible from the original voice, and successfully confuses the speaker ID system into  
101 mistaking it for another.

102 Subsequently, the Festvox transformation tools developed by the speech group at the  
103 Language Technologies Institute, Carnegie Mellon University, are used to perform voice  
104 transformations. This software creates mappings to convert the input voice to a specified  
105 output voice based on the joint probabilities of the two voices, while also factoring in global  
106 variance parameters.

107 Models are constructed to convert CMU's arctic database of voices to two native Indian  
108 speakers' voices, both male. The test utterances of 30 speakers (23 standard American  
109 newsreaders, 5 from the CMU arctic database, and the 2 native Indian speakers) through  
110 these models, and discard the unintelligible outputs.

111 On empirical analysis, it is observed that a male voice passed through a male-male  
112 transformation is usually intelligible. Similarly, a female voice passed through a female-  
113 male transformation is usually intelligible.

114

#### 115 **5 Transformation Selection for De-Identification**

116 The best transformation is defined as the most de-identifiable yet fully comprehensible  
117 transformation of the input utterance.

118 By clustering the non-transformed original speaker's utterances using a variation of the k-  
119 means algorithm, a transformed utterance is chosen whose clusters are furthest away from it.  
120 We start with two clusters formed uniformly distributed about the mean of the input data,  
121 and successively split high population clusters, killing the low population clusters.

122 Initially, for a known speaker, it is empirically determined which transformations are  
123 intelligible, and only their clusters are used to find the most de-identifiable transformation.

124 The datasets of MFCCs and F0s for transformed voices were analyzed to obtain some trend  
125 or metric to gauge the intelligibility. After thorough experimentation, it was found that the  
126 variance of the MFCC coefficients for each frame had smooth transitions for intelligible  
127 voices. However, considerable transients were noticed for un-intelligible voices. Therefore,  
128 some preliminary conclusions can be drawn regarding a possible correlation between the  
129 variance of the MFCCs and the intelligibility of a given voice. We used a summed derivative  
130 along the frames axis to eliminate a subset of all those transformations that were  
131 unintelligible.

132

#### 133 **6 Results**

134 For 20 non-transformed utterances of 30 speakers, of which 23 are standard American  
135 newsreaders, the speaker ID system proved accurate 91.83% of the time. For de-identified  
136 transformed voices, the speaker ID system gave an accuracy of only 4.5%, and thus was  
137 sufficiently confused. Some preliminary results linking MFCC variance to intelligibility  
138 were also established.

139

140

141

142

Table 1: Accuracy of Speaker ID on non-transformed voices

Speakers	Accuracy of Speaker ID on non-transformed voices
23 standard American voices	90.43%
5 CMU Arctic voices	98%
2 native Indian speakers	92.5%
Total Accuracy	91.83%

143

144

Table 2: Accuracy of Speaker ID on transformed voices

Speakers	Accuracy of Speaker ID on transformed voices
5 CMU Arctic voices	4.72%
2 native Indian speakers	60%
Total Accuracy	12%

145

146

## 7 Future Work

147

148

149

150

151

152

153

154

155

## 8 Acknowledgements

156

157

158

159

160

161

162

163

164

165

- Prof. Alan W. Black for his ideas, suggestions and support.
- Esfandiar Zavarehei (Brunel University) for the MATLAB code for K-means with splitting clusters.
- Anil Alexander, Andrzej Drygajlo (Swiss Federal Institute of Technology, Laussane) for ‘Speaker Identification: A Demonstration using MATLAB’.
- Malcolm Slaney for his mfcc.m function from his Auditory Toolbox.

## References

- [1] Toda, T., Black, A., and Tokuda, K. (2007) "Voice Conversion Based on Maximum-Likelihood Estimation of Spectral Parameter Trajectory", *IEEE Transactions on Audio, Speech and Language Processing*, 15(8), pp 2222-2236.