# Research Paper - A Hybrid Approach for Voice Recognition Using MFCC, VQ and K-Means Algorithm

Nisha Saini
M.Tech Scholar, Department of Computer Science & Engineering
Deenbandhu Chhotu Ram University of Science & Technology (DCRUST), Sonepat, India

Dr. Dinesh Singh
Asst. Professor, Department of Computer Science & Engineering
Deenbandhu Chhotu Ram University of Science & Technology (DCRUST), Sonepat, India

*Abstract:* The speech or voice of any human being is his/her unique personal characteristics. No two have almost identical voice there are some features which must be present in one voice and missing from another voice. In order to identify this uniqueness a robust and efficient technique is required so that we can accurately identify the genuine voice from the bunch of fake voices. The process of recognizing a voice of a given speech from the group of given speakers is called speaker identification. This paper proposes a new improved and efficient technique for voice recognition system by applying Mel Frequency Cepstral Coefficients (MFCC), Vector Quantization, K-means and Euclidean distance technique.

*Keywords:* Voice Recognition, MFCC, Vector Quantization, Euclidean Distance

## I. INTRODUCTION

Biometrics refers to "automatic recognition of individuals based on their physiological and behavioral characteristics." Today there is need for simple and secure access control mechanism for authentication of a person. Biometrics technique for authentication is the solution for this. With biometric security we need not carry different types of ID cards or remembering various passwords or keys, our body parts are sufficient for authentication purpose. Even more secure & robustness can be achieved if we combine biometric security with other security modes such as key, password or any identification card.

There are various application areas where we can use biometric security such as banking application where signature recognition identify our authentication or our voice in case of phone banking, driving license and passport where our image (picture) is used for identification.

There are many advantages of using biometric security over normal security for authentication [1]. These advantages are listed below:

1. Reliability: The biometric cannot be lost and stolen, therefore it provides more reliability.

2. Uniqueness: Another important feature of biometric is that it is unique for everyone. For example the fingerprint of one person is different from another person.

3. Nothing to remember or carry: With biometric we need not carry different types of ID cards or remembering various passwords or keys, our body parts are sufficient for authentication purpose.

The biometric security works as follows:

First capture the biometric image, extract features & create database of these features. Now when any person comes for biometric verification then his/her biometric image is collected again features are extracted and then these features are matched with the features already stored in database and then matching is performed. One of important biometric system is human voice.

The main aim of this paper is voice or speaker identification. In this method we compare voice of any unknown person with database of known voices, It then returns the message whether voice matches or not. In this paper, we use hybrid approach of "Mel Frequency Cepstral Coefficients (MFCC) technique, Vector Quantization, K-means algorithm and Euclidean distance."

## II. OVERVIEW OF WORK

Voice recognition identification [2] is the process of recognizing the speaker automatically. The speech waves of individual voice form the basis of identification of speaker. We can use voice identification in multiple application areas such as telephone banking, shopping through telephone, access to database information and voice mail. One of the powerful applications of voice recognition is for security purpose where a person can enter his/her voice for authentication.

The voice recognition system [3] is broadly divided into following two categories- "Speaker identification & Speaker verification". The process of recognizing a voice of a given speech from the group of given speakers is called speaker identification. The speaker whose maximum voice characteristics are matches with the stored voice is identified & the speaker whose voice characteristics are not matched is eligible for new entry in the database.

The main building blocks [4] of any speaker or voice recognition systems are "*feature extraction* and *feature matching*." As the name suggests the feature extraction block of voice recognition system is to extract the unique characteristics of any voice signal. These characteristics represent the identity of any person. On the other hand feature matching block is use to match the unknown voice signal with the database of given voice and identify the claim person.

### A. Voice Feature Extraction [5, 6]
The voice feature extraction block of voice recognition system is to extract the unique characteristics of any voice signal. These characteristics represent the identity of any person. It converts voice into speech wave as shown in figure 1. It uses the concept of digital signal processing (DSP) tool for converting voice signal into speech wave.
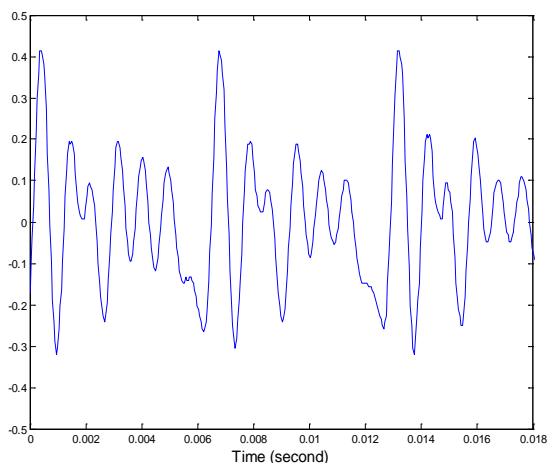
Figure 1: Example of speech wave

The voice signal can be converted into speech wave for extracting unique features by using one of the following methods:

1. Linear Prediction Coding (LPC)
2. Mel-Frequency Cepstrum Coefficients (MFCC)

In this work we use MFCC technique for feature extraction because it is well studied and popular method for voice feature extraction

MFCC [7] described by authors as "MFCC's are based on the known variation of the human ear's critical bandwidths with frequency, filters spaced linearly at low frequencies and logarithmically at high frequencies have been used to capture the phonetically important characteristics of speech." The MFCC is generraly expressed as *mel-frequency* scale.

## B. Voice Feature Matching

Voice feature matching [8] block is use to match the unknown voice signal with the database of given voice and identify the claim person. Voice feature matching works as - When any person comes for voice verification then his/her voice signal is collected, then features are extracted and then these features are matched with the features already stored in database and then matching is performed.

The voice signal can be matched with voices already stored in database by using one of the following methods:

1. Dynamic Time Warping (DTW)
2. Hidden Markov Modeling (HMM)
3. Vector Quantization (VQ).

In this work, we use the VQ approach [9] because it gives high accuracy & its implementation is easy.

## III. PROPOSED WORK

The main aim of this paper is voice or speaker identification. In this method we compare voice of any unknown person with database of known voices, It then returns the message whether voice matches or not. In this paper, we use hybrid approach of "Mel Frequency Cepstral Coefficients (MFCC) technique, Vector Quantization, K-means algorithm and Euclidean distance."

MFCC [10] is used to extract unique features from voice entered by the user. For estimation of probability distributions of the computed feature vectors, we use the concept Vector Quantization (VQ). VQ represents the

centroid of extracted features. The k-means algorithm is used for dividing the features into n clusters; each cluster represents the unique characteristics of any feature of voice [11].

The steps performed in K-means algorithm are shown in figure 2 below:

1. Clusters the data into k groups where k is predefined.
2. Selects k points at random as cluster centers.
3. Assigns objects to their closest cluster center according to the Euclidean distance function.
4. Calculates the centroid or mean of all objects in each cluster.
5. Repeats steps 2, 3 and 4 until the same points are assigned to each cluster in consecutive rounds.

Figure 2: Steps in K-means Algorithm

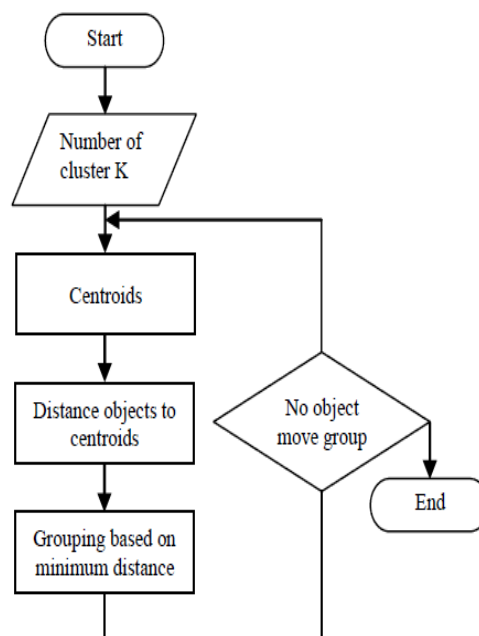The flowchart of above steps is shown in Figure 3 below:



Figure 3: Flow chart of the K-means algorithm

Finally to identify the unknown speaker, we use Euclidean distance. The Euclidean distance measure the alteration distance of given two vector sets. It chooses one of the speakers with the smallest alteration distance for identification as the unknown person.

To achieve the desire result, first of all one has to divide the project in two parts: Training and Testing. In Training session the feature vectors are extracted from the training set and then the classifier is trained using these features. At testing session the classifier is tested with some unknown data. The following figure 4 illustrates the Speaker Recognition System concept.
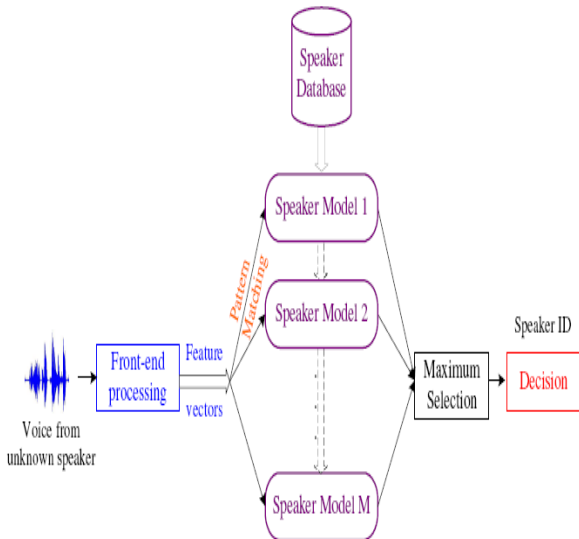
Figure 4: Block Diagram of the algorithm for Speaker recognition

In order to make a speaker identification system, main part is to find right feature.

## IV. FRAMEWORK FOR VOICE RECOGNITION

The basic framework for voice recognition system consists of two phases – training phase and testing phase.

In the training phase, voice samples of different person are trained & stored in a database for creating reference model of voices. In the testing phase, actual decision for voice recognition is performed with the input voice sample. Decision is depend upon feature matching values of reference model of database and input voice.

The flowchart for training phase is shown in figure 5 & the flowchart for testing phase is shown in figure 6 below.
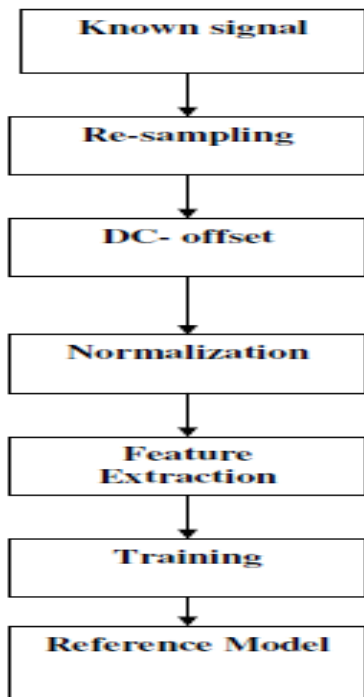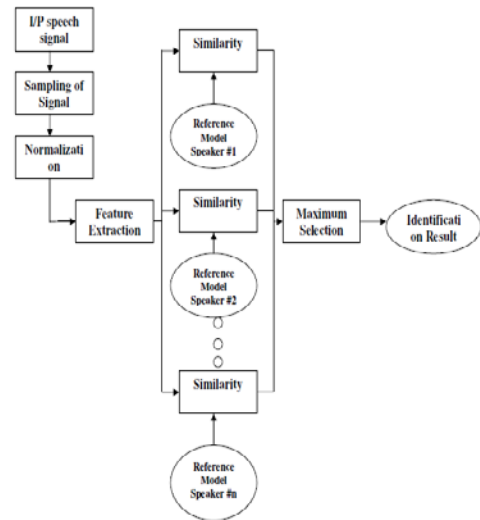


Figure 5: Testing Phase



Figure 6: Testing Phase

First of all a set of speech signal is collected. Then they are re-sampled at 22050 samples per second with 16 bit per sample. DC offset of signal is corrected to 0. Then signal is normalized to 60%. The Mel-frequency cepstral coefficient of each signal is calculated and stored in database. This is the training phase. After that the unknown signal is taken for recognition. Then they are re-sampled at 22050 samples per second with 16 bit per sample. DC offset of signal is corrected to 0. Then signal is normalized to 60%. Then Mel frequency for each signal is calculated. Then Euclidean distance of the unknown signal is calculated. Minimum distance found is the result sound.

MATLAB [12] has been used to implement the above procedure. First training sample belonging to each class is passed to the function speaker_training.m, which reads each samples and calculates MFCC. The values of these MFCC's are then written into a database file called cepstrum.mat along with the name of training signal called name.mat. After that the unknown signal is passed through speaker_test.m which reads the signals and calculates MFCC. After that the cepstral coefficients of training signal and test signal are passed through distmeasure.m, which calculates the Euclidean distance between test signal and each of the training signals.

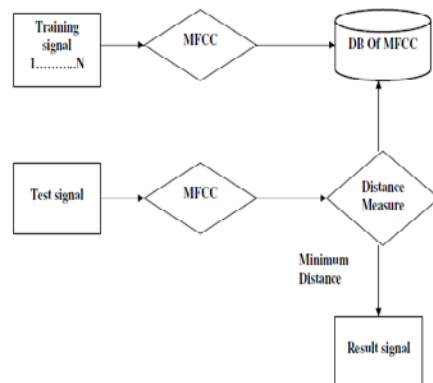Figure 7 shows the general framework for voice recognition system.



Figure 7: A framework for voice recognition system

## V. ANALYSIS OF RESULT

In this work we compare voice of any unknown person with database of known voices, It then returns the message whether voice matches or not. In this paper, we use hybrid approach of "Mel Frequency Cepstral Coefficients (MFCC) technique, Vector Quantization, K-means algorithm and Euclidean distance."

MFCC is used to extract unique features from voice entered by the user. For estimation of probability distributions of the computed feature vectors, we use the concept Vector Quantization (VQ). VQ represents the centroid of extracted features. The k-means algorithm is used for dividing the features into n clusters; each cluster represents the unique characteristics of any feature of voice.

Finally to identify the unknown speaker, we use Euclidean distance. The Euclidean distance measure the alteration distance of given two vector sets. "It chooses one of the speakers with the smallest alteration distance for identification as the unknown person".

The figure 8 & 9 below shows the results obtained using the concept of feature extraction and feature mapping.
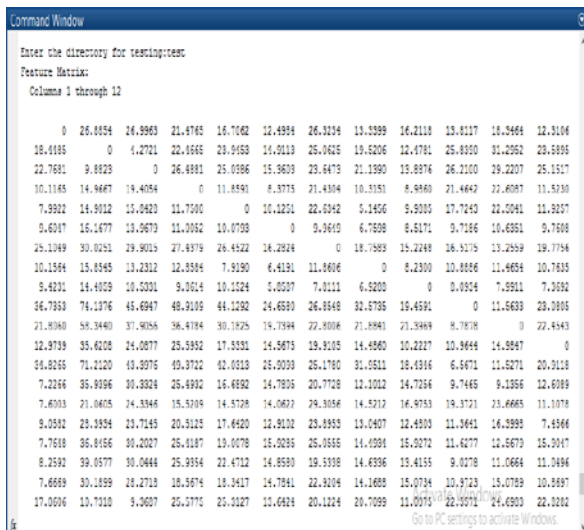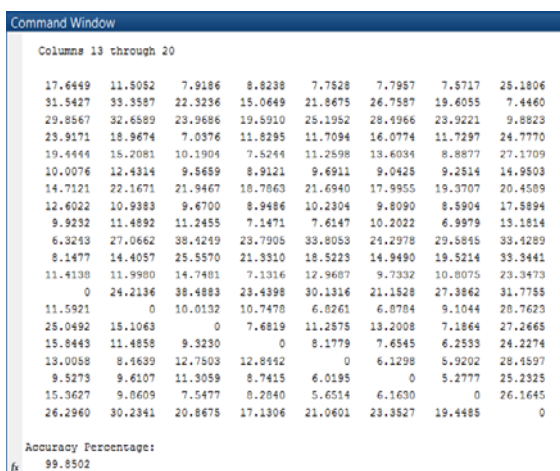


Figure 8: Result of feature comparisons



Figure 9: Result of feature comparisons (cont…)

As we can clearly infer all the speakers 1 to 11 (in our example) were correctly identified. The final matrix

contains the Euclidean distance from training to testing feature vectors. The value 0 in each row represents exact matching of voice. Hence accuracy of our algorithm achieves 99.85%.

## VI.  CONCLUSION

The process of recognizing a voice of a given speech from the group of given speakers is called speaker identification. This paper proposes a new improved and efficient technique for voice recognition system by applying Mel Frequency Cepstral Coefficients (MFCC), Vector Quantization, K-means and Euclidean distance technique. MFCC is used to extract unique features from voice entered by the user. For estimation of probability distributions of the computed feature vectors, we use the concept Vector Quantization (VQ). VQ represents the centroid of extracted features. The k-means algorithm is used for dividing the features into n clusters; each cluster represents the unique characteristics of any feature of voice. Finally to identify the unknown speaker, we use Euclidean distance.

### REFERENCES

[1] Parwinder Pal Singh, Er. Bhupinder Singh, "Speech Recognition as Emerging Revolutionary Technology", International Journal of Advanced Research in Computer Science and Software Engineering, Volume 2, Issue 10, October 2012.

[2] Campbell, J.P., "Speaker recognition: a tutorial", Proceedings of the IEEE Volume 85, Issue 9, Sept. 1997 Page(s):1437 – 1462.

[3] Surbhi Mathur, Choudhary S. K. and Vyas J. M., "Speaker Recognition System and its Forensic Implications", Open Access Scientific Reports, Volume 2, Issue 4, 2013.

[4] Douglas A. Reynolds, "An Overview of Automatic Speaker Recognition Technology", ©2002 IEEE

[5] Bhupinder Singh, Rupinder Kaur, Nidhi Devgun, Ramandeep Kaur, "The process of Feature Extraction in Automatic Speech Recognition System for Computer Machine Interaction with Humans: A Review",IJARCSSE, Volume 2, Issue 2, February 2012.

[6] Genevieve I. Sapijaszko, Wasfy B. Mikhael, "An Overview of Recent Window Based Feature Extraction Algorithms for Speaker Recognition", IEEE, pp 880-883, 2012.

[7] Ufuk Ülüg, Tolga Esat Özkurt, Tayfun Akgül, "Bispectrum Mel-frequency Cepstrum Coefficients for Robust Speaker Identification", 2006 IEEE

[8] Maider Zamalloa, Germacn Bordel, Luis Javier Rodriguez, Mikel Penagarikano, "Feature Selection Based on Genetic Algorithms for Speaker Recognition", 2006, IEEE.

[9] HUO ChunBao SHAO Yan, "The improved VQ algorithm for speaker recognition",© 2009 IEEE

[10] Jorge Martinez, Hector Perez, Enrique Escamilla, Masahisa Mabo Suzuki, "Speaker recognition using Mel Frequency Cepstral Coefficients (MFCC) and Vector Quantization (VQ) Techniques, IEEE, pp 248-251, 2012.

[11] Izuan Hafez Ninggal & Abdul Manan Ahmad, "The Fundamental Of Feature Extraction In Speaker Recognition: A Review", Proceedings of the Postgraduate Annual Research Seminar 2006.

[12] MATLAB Primer, MathWorks Inc, 2014.