

**UC Davis
EEC 201
Final Report**

Speaker Recognition

Qian, Jianshu

Nguyen, Kevin

Human Hearing Test

| Training | Testing |
|----------|---------|
| 1 | 1 |
| 2 | 2 |
| 3 | 3 |
| 4 | 4 |
| 5 | 5 |
| 6 | 6 |
| 7 | 7 |
| 8 | 8 |
| 9 | 9 |
| 10 | 10 |
| 11 | 11 |

Randomly Rearranged Files

| Training | Testing |
|----------|---------|
| 1 | 8 |
| 2 | 4 |
| 3 | 7 |
| 4 | 9 |
| 5 | 1 |
| 6 | 10 |
| 7 | 2 |
| 8 | 5 |
| 9 | 3 |
| 10 | 11 |
| 11 | 6 |

The testing data was randomly rearranged in order to prevent bias then referred to the table directly above to determine if the human test was accurate. There is 100% accuracy from the human test.

Feature Extraction

The features of each of the training speakers must be extracted using a MFCC processor. The received speech signal must be divided into frames, windowed, then converted to the frequency domain using FFT. Mel-frequency wrapping involves processing the signal using mel-spaced filter banks that scale linearly at low frequencies and logarithmically at higher values in order to emphasize the frequencies that are important to human hearing. By doing so, cepstrum coefficients from the mel spectrum of each frame can be obtained to be later used in vector quantization.

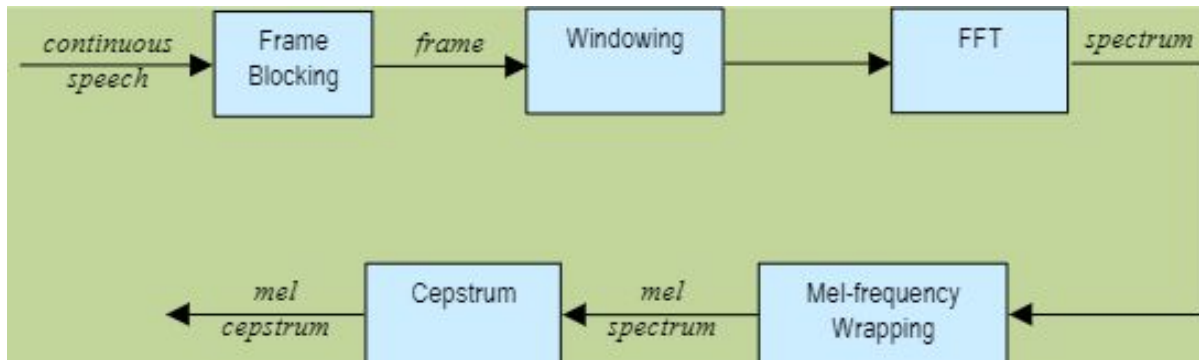


Figure 1: Block Diagram of MFCC extraction

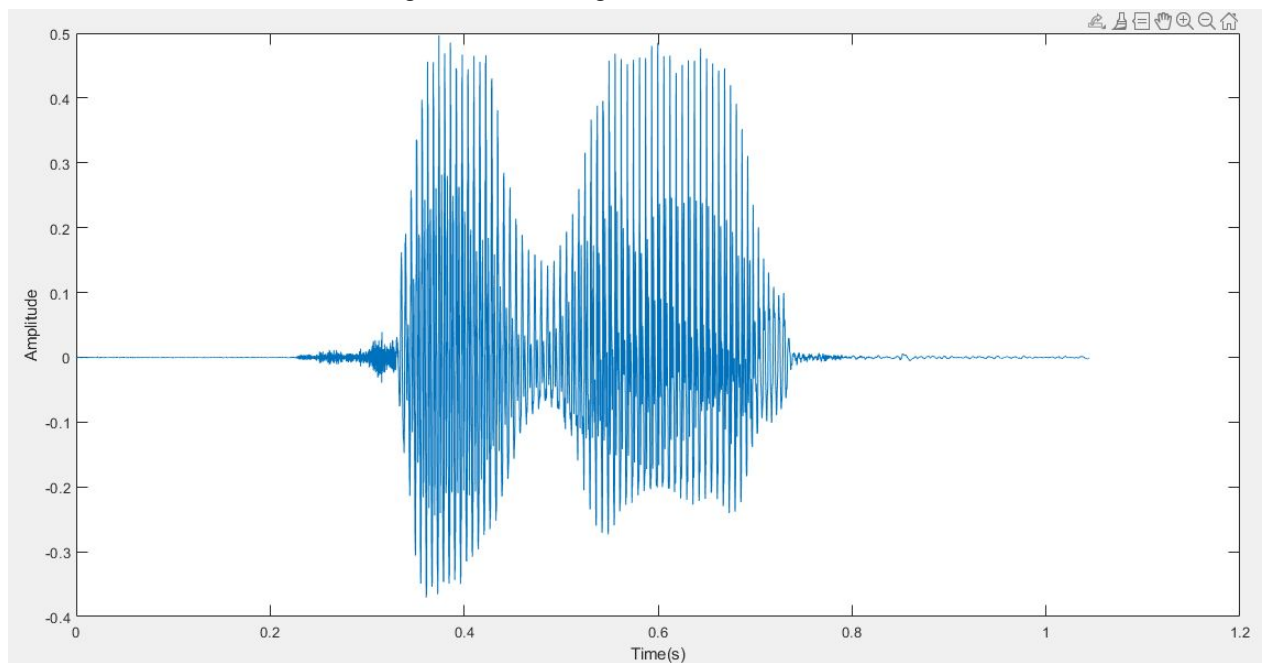


Figure 2: Speech Signal of Speaker 1 Unmodified

Frame Blocking

For testing each frame consists of $N = 256$ samples with each frame starting $M = 100$ samples after the other for an overlap of 156. The sampling frequency used was 12500 Hz, so each frames equates to 0.02048 seconds

Windowing

Each frame was processed through a hamming window of equal size to the number of samples in the frame to reduce spectral distortion. The hamming window used the form

$$w(n) = 0.54 - 0.46 \cos\left(\frac{2\pi n}{N-1}\right), \quad 0 \leq n \leq N-1$$

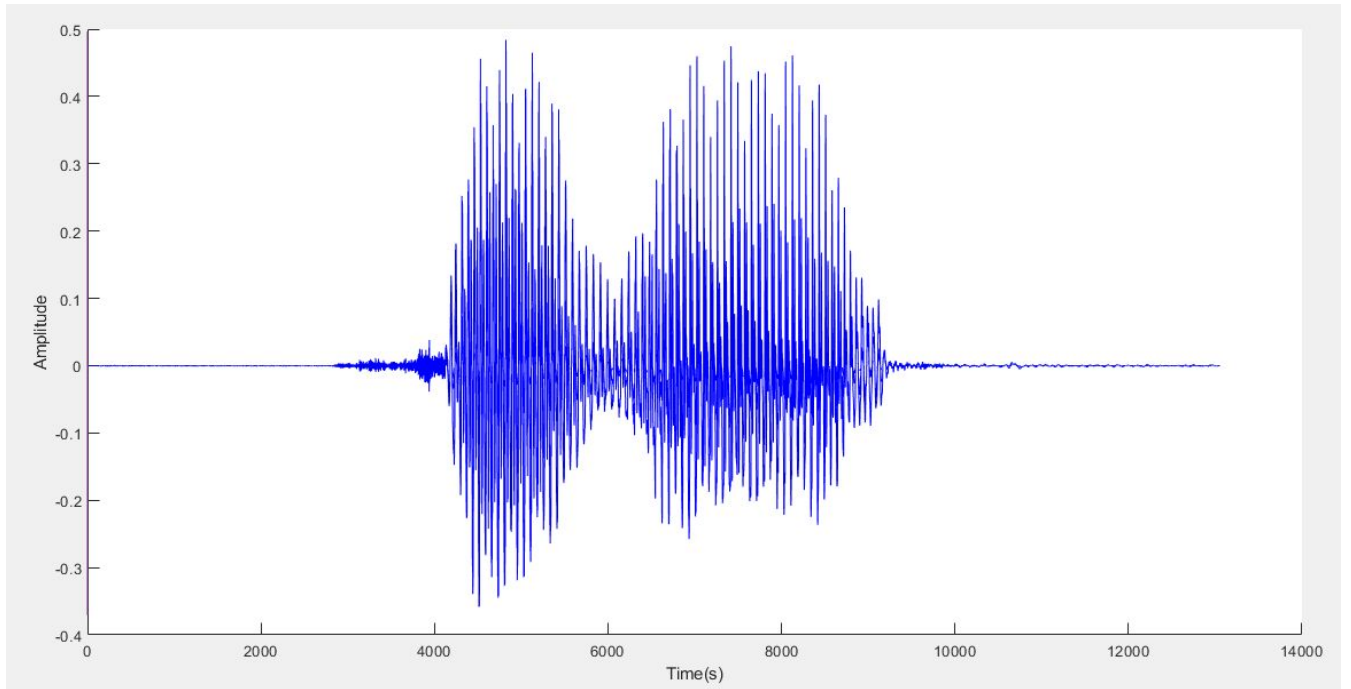


Figure 3: Speech Signal of Speaker 1 Framed and Windowed with $N = 256$

FFT

The FFT of the signal is taken to get the periodogram.

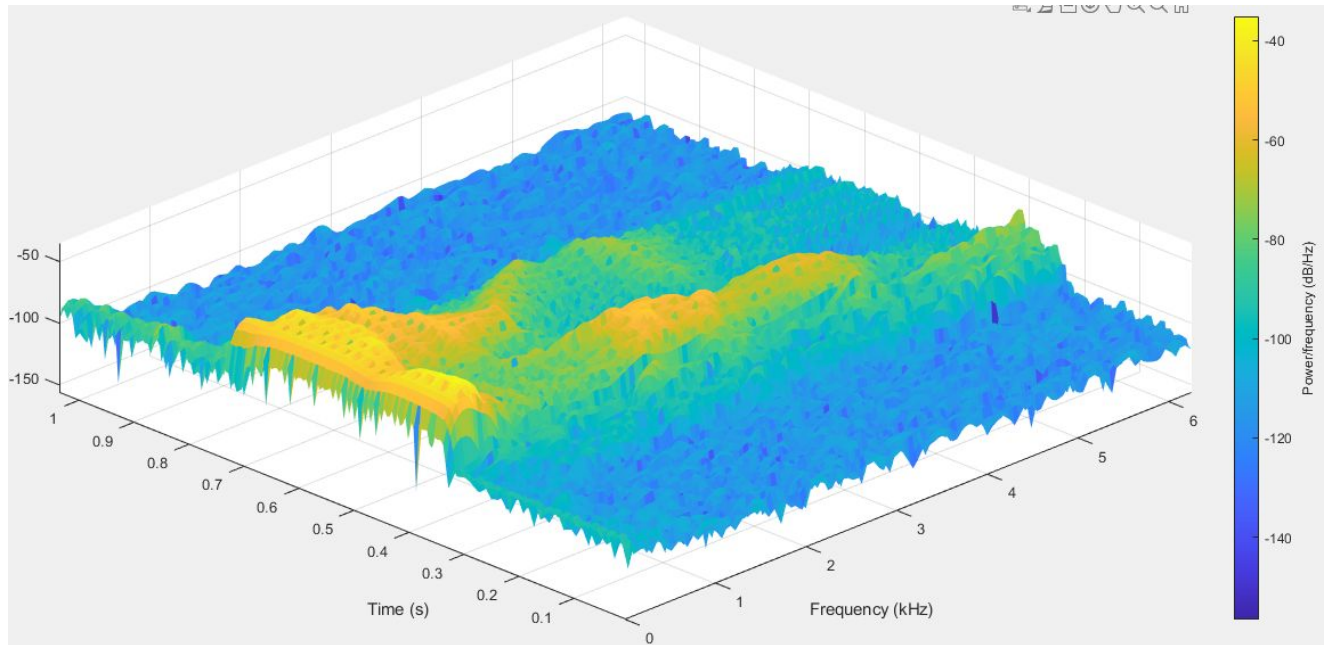


Figure 4: Periodogram of Speaker 1, $N = 128$, $M = 42$

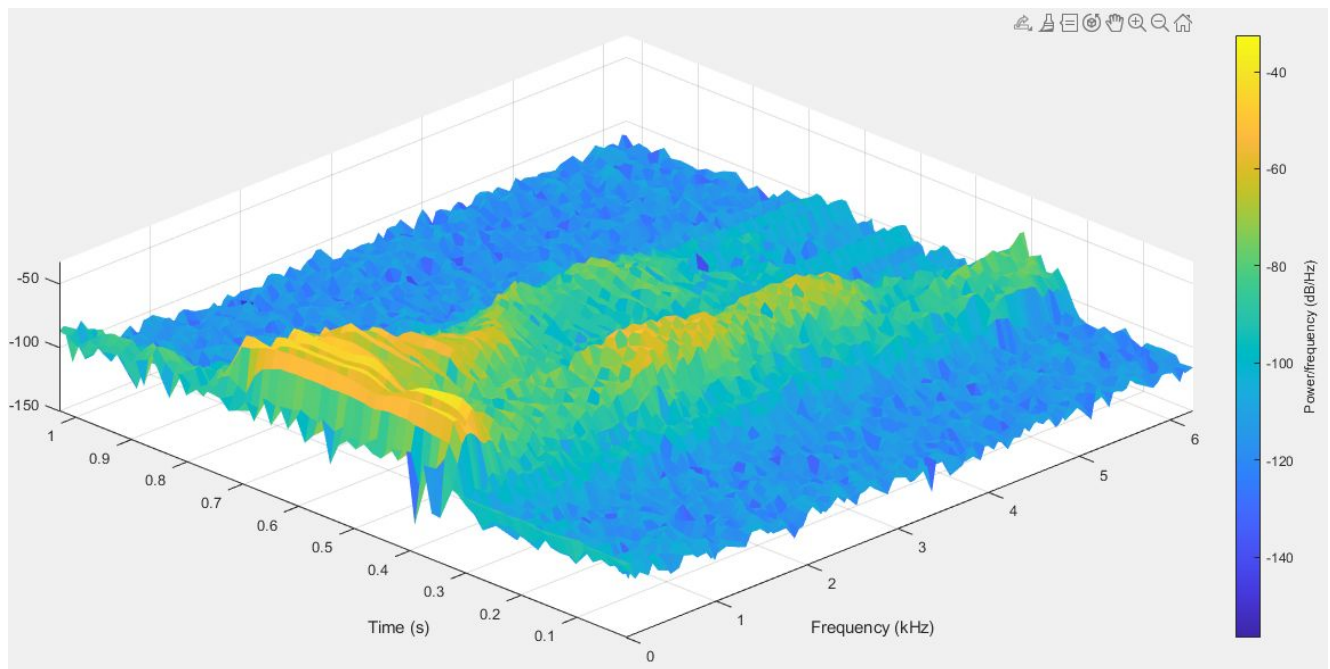


Figure 5: Periodogram of Speaker 1, $N = 256$, $M = 85$

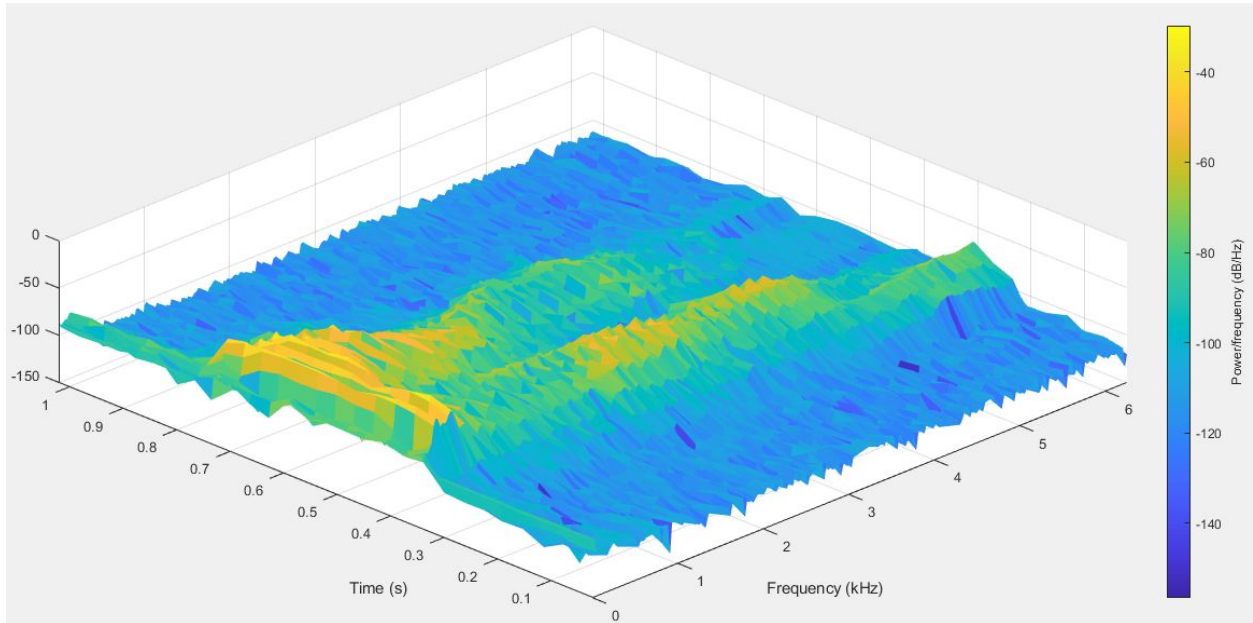


Figure 6: Periodogram of Speaker 1, $N = 512$, $M = 170$

The frequencies that contain most of the power/frequency are located below 1 kHz which reveals the dominant frequencies in human speech, thus they are emphasized by the mel-spaced filter banks.

Mel Frequency Wrapping

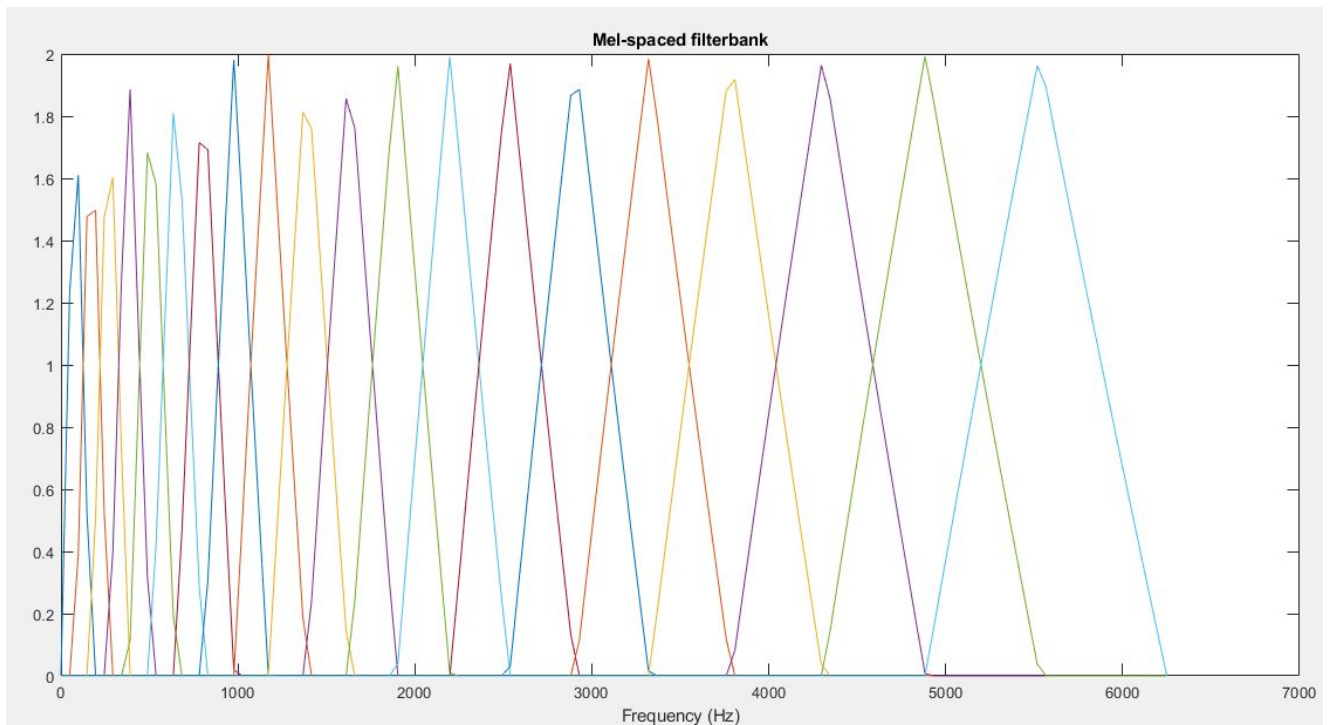


Figure 7: Mel-spaced Filter Bank Responses

20 filter banks were used resulting in 20 mel spectrum coefficients for each frame.

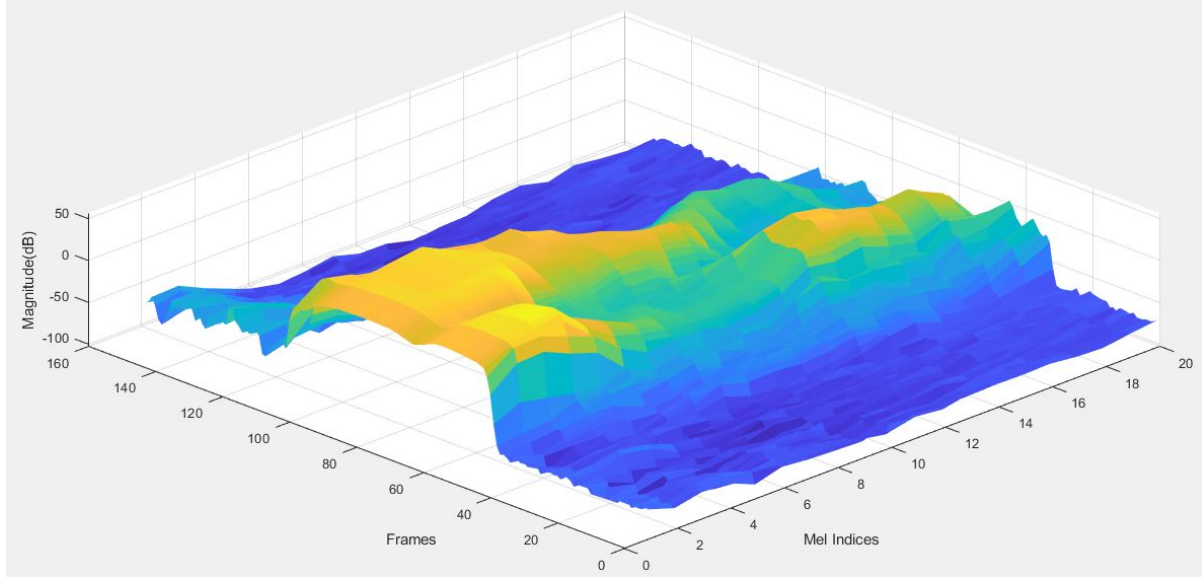


Figure 8: Periodogram of Speaker 1 after Mel Frequency Wrapping, N = 256, M = 85

Each filter bank produces one coefficient for a total of 20 per frame which results in a smoother periodogram.

Cepstrum

In order to turn the mel spectrum into the mel cepstrum, the DCT of the natural log of the mel spectrum was taken for the purpose of converting it to the time domain as shown below. S_k are the coefficients of the mel spectrum.

$$\tilde{c}_n = \sum_{k=1}^K (\log S_k) \cos \left[n \left(k - \frac{1}{2} \right) \frac{\pi}{K} \right], \quad n = 0, 1, \dots, K-1$$

The result is 20 MFCC's for each frame.

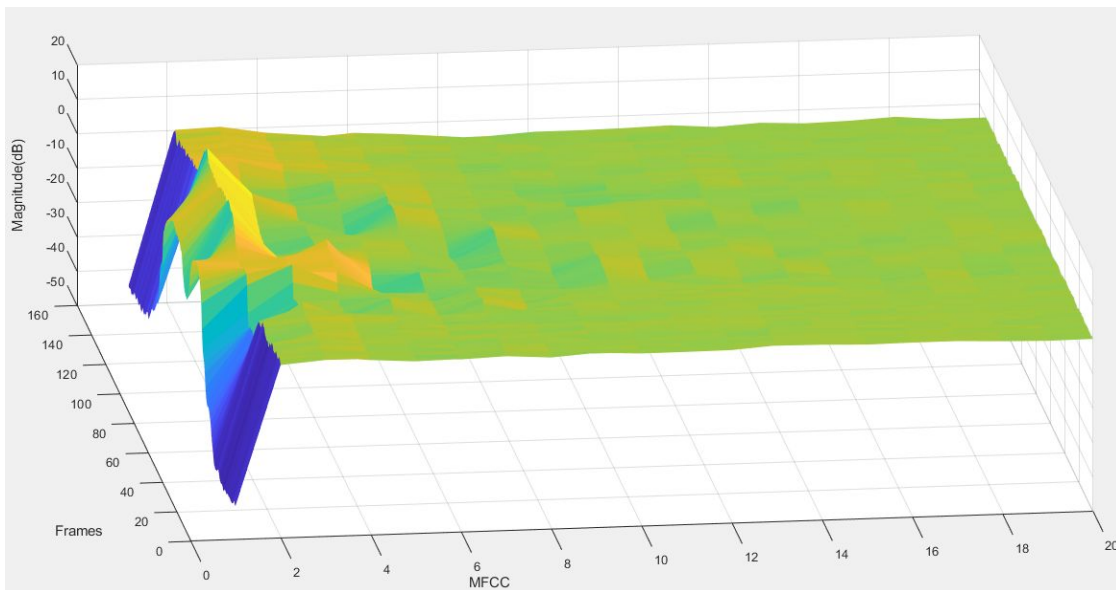


Figure 9: Periodogram of Speaker 1 MFCC, N = 256, M = 85

Vector Quantization

In order to match the extracted features to the test speaker, the training vectors from the MFCC must be used to create a VQ codebook for each speaker. The LBG algorithm was used to find the centroids and determine if the distortion between the centroid and samples is sufficiently small. The splitting parameter, $\epsilon = 0.01$, determines the position where each portion of the centroid divides, then the distance between the centroid and the closer codeword is used to update the centroid. This continues until the distortion is smaller than the splitting parameter and a codebook of size M is formed. Because each centroid splits into two for each iteration, the codebook will always be a size of the power of 2. A codebook of size 32 was used to increase the robustness of the system as one of size 16 had some errors in detecting the speaker. This will be discussed in the notch filter.

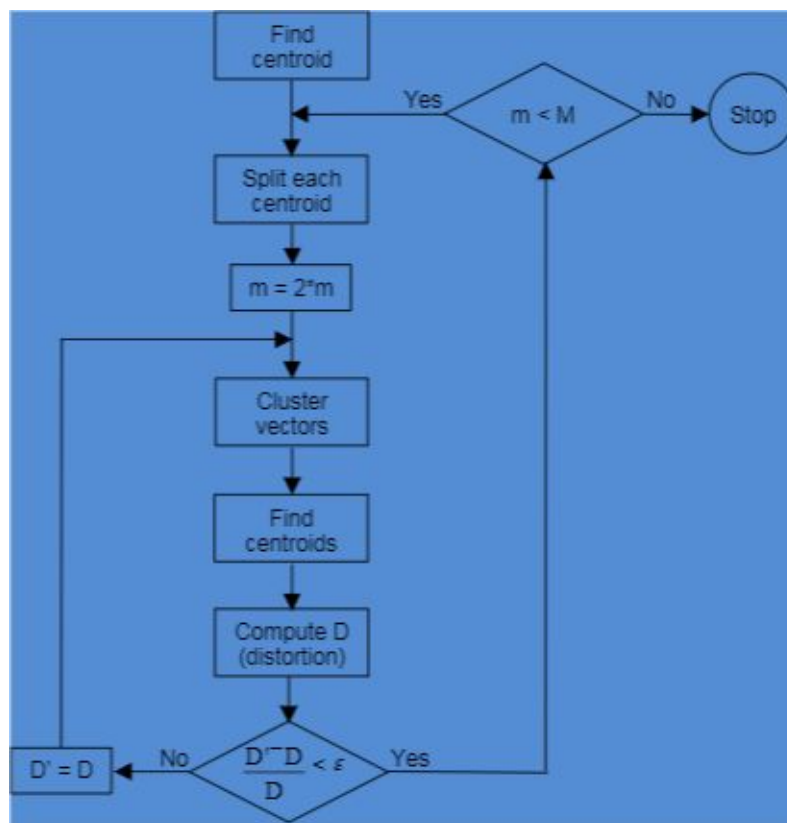


Figure 10: Flow Chart of LBG Algorithm

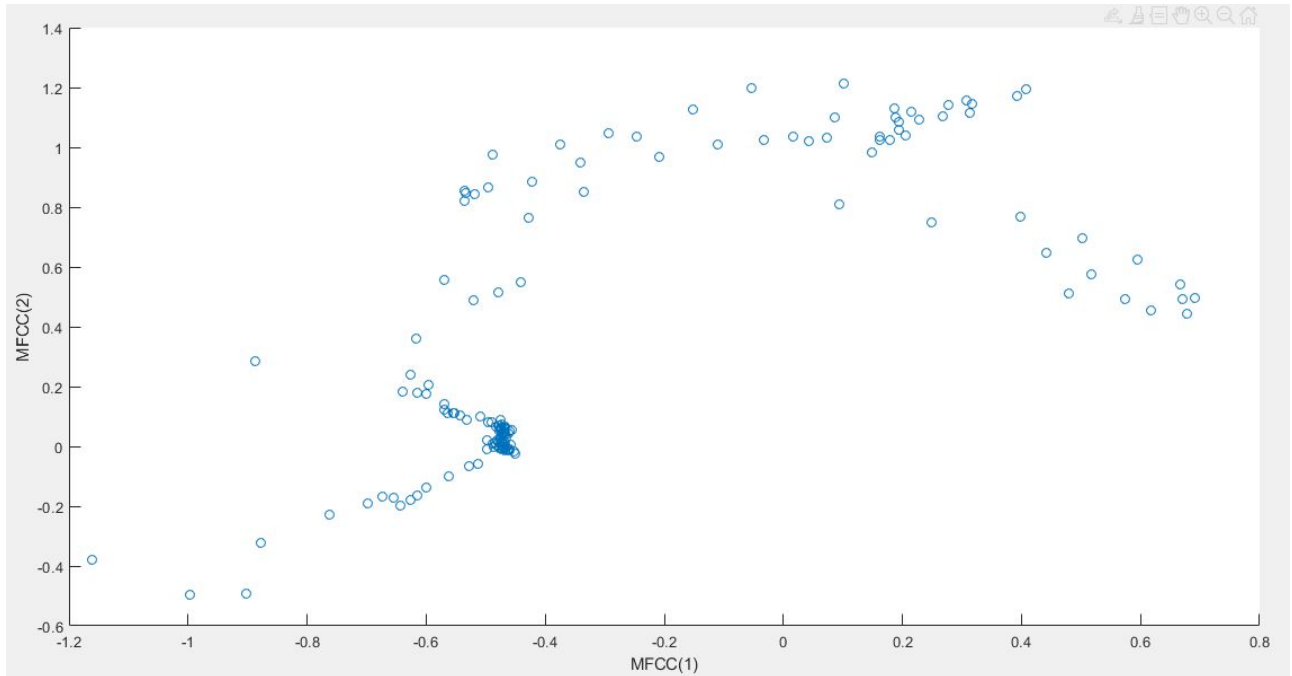


Figure 11: Clustering of MFCC's 1 and 2 of Speaker 1 before Vector Quantization

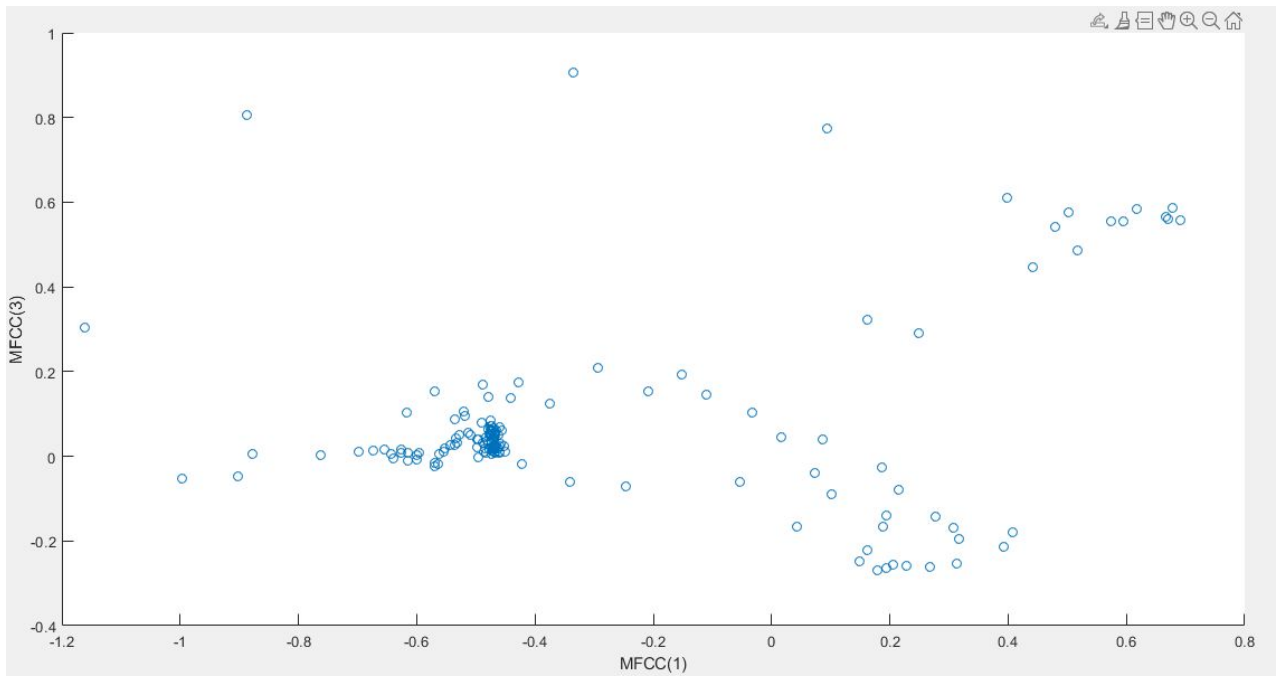


Figure 12: Clustering of MFCC's 1 and 3 of Speaker 1 before Vector Quantization

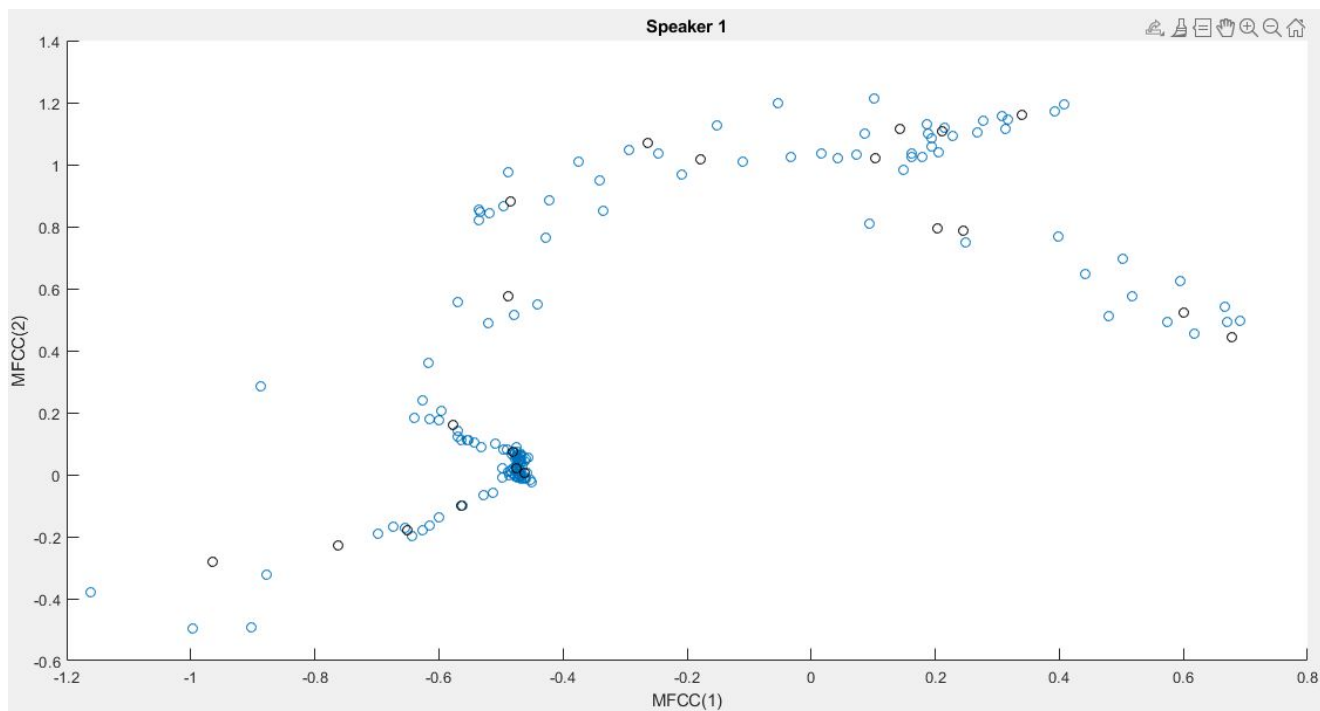


Figure 13: Clustering of MFCC's 1 and 2 of Speaker 1, Centroids(Black), Samples(Blue)

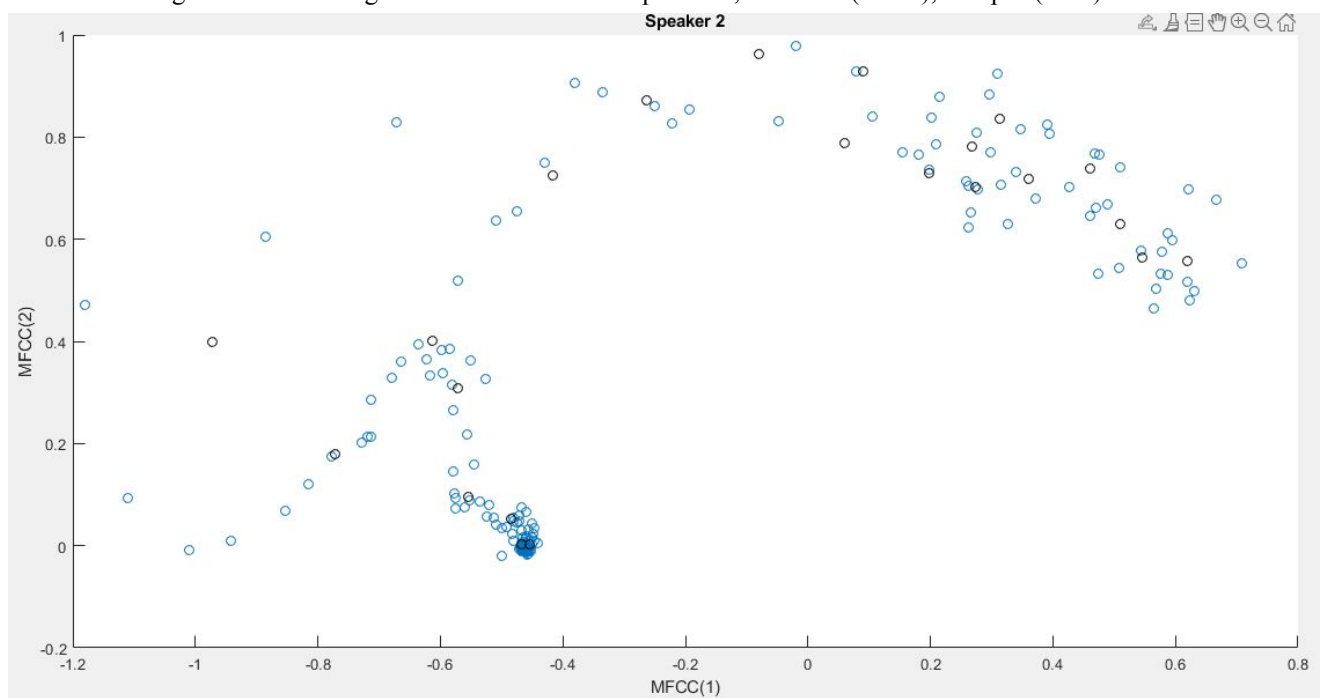


Figure 14: Clustering of MFCC's 1 and 2 of Speaker 2, Centroids(Black), Samples(Blue)

Testing

After the codebooks of each training speaker is initialized, the codebook of the testing speaker is taken. Whoever's training data has the lowest distortion with the testing data, the speaker is matched. Using a GUI, the results can be quickly matched.

| | | | |
|--|---|--|---|
| Input File Speaker ID: <input type="text" value="1"/> | Matched Speaker Speaker ID: <input type="text" value="1"/> | Input File Speaker ID: <input type="text" value="2"/> | Matched Speaker Speaker ID: <input type="text" value="2"/> |
| Input File Speaker ID: <input type="text" value="3"/> | Matched Speaker Speaker ID: <input type="text" value="3"/> | Input File Speaker ID: <input type="text" value="4"/> | Matched Speaker Speaker ID: <input type="text" value="4"/> |
| Input File Speaker ID: <input type="text" value="5"/> | Matched Speaker Speaker ID: <input type="text" value="5"/> | Input File Speaker ID: <input type="text" value="6"/> | Matched Speaker Speaker ID: <input type="text" value="6"/> |
| Input File Speaker ID: <input type="text" value="7"/> | Matched Speaker Speaker ID: <input type="text" value="7"/> | Input File Speaker ID: <input type="text" value="8"/> | Matched Speaker Speaker ID: <input type="text" value="8"/> |
| Input File Speaker ID: <input type="text" value="9"/> | Matched Speaker Speaker ID: <input type="text" value="9"/> | Input File Speaker ID: <input type="text" value="10"/> | Matched Speaker Speaker ID: <input type="text" value="10"/> |
| Input File Speaker ID: <input type="text" value="11"/> | Matched Speaker Speaker ID: <input type="text" value="11"/> | Input File Speaker ID: <input type="text" value="13"/> | Matched Speaker Speaker ID: <input type="text" value="12"/> |

Figure 14: Matching of Testing and Training Speakers

The corresponding testing speaker has the same ID as the training speaker which matches with what is shown above, so all the data is 100% accurate with a codebook size of 32. Training speaker 12 is a new speaker(teammate) which matches correctly to the added testing speaker 13.

Notch Filter

In order to test the robustness of the code, a notch filter of a specific frequency was used.

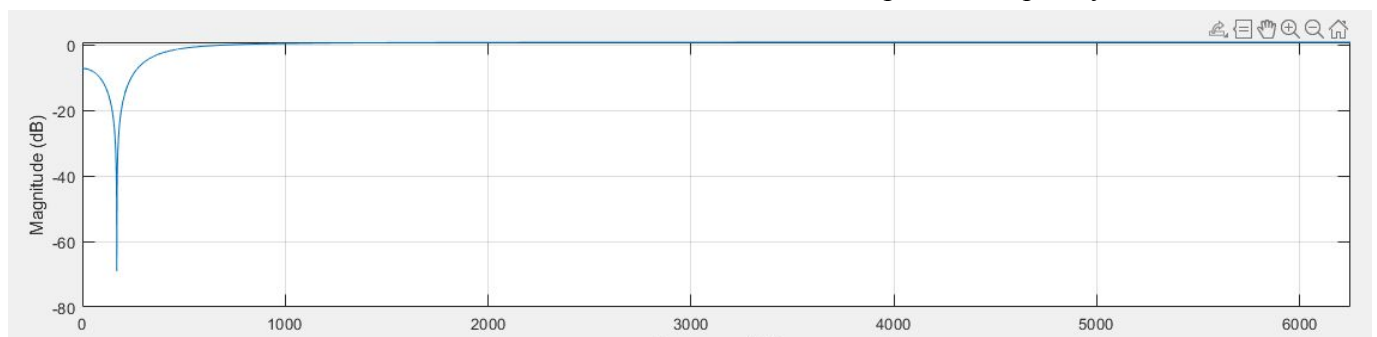


Figure 16: Notch Filter with Null Frequency at 157 Hz

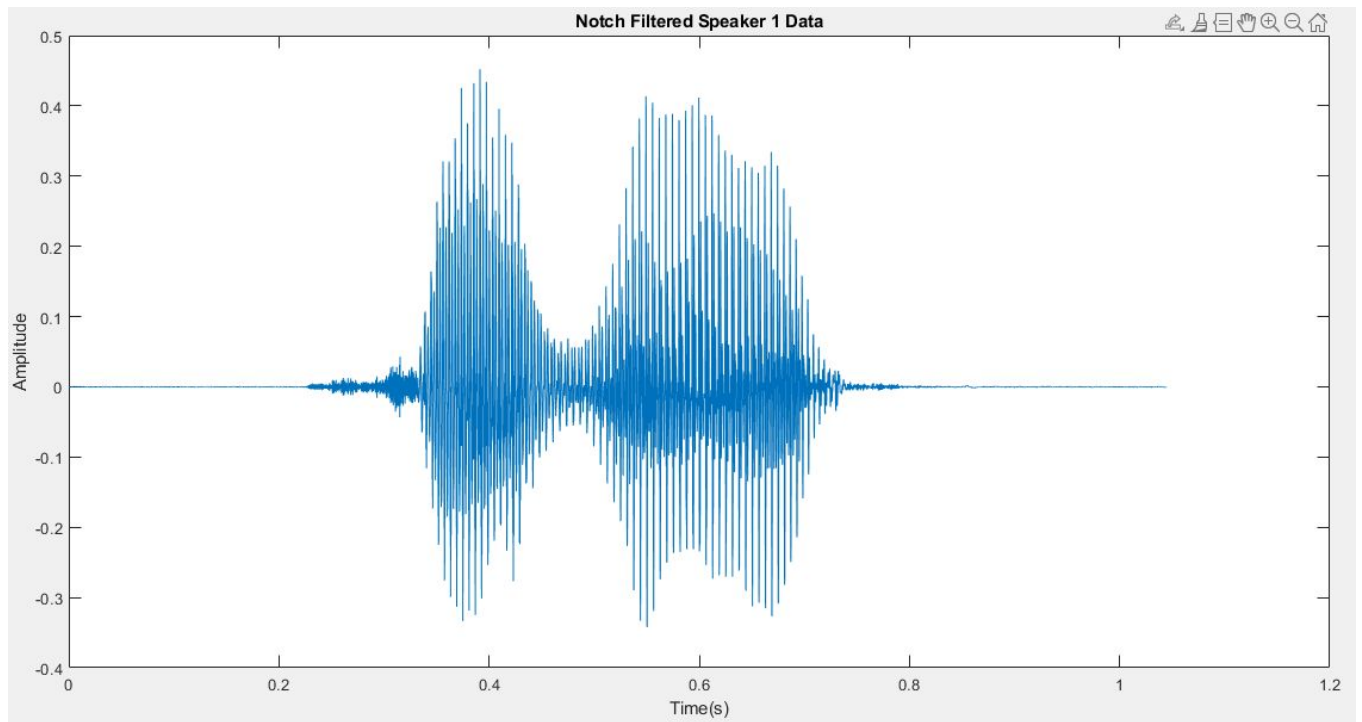


Figure 17: Notch Filtered Data of Speaker 1

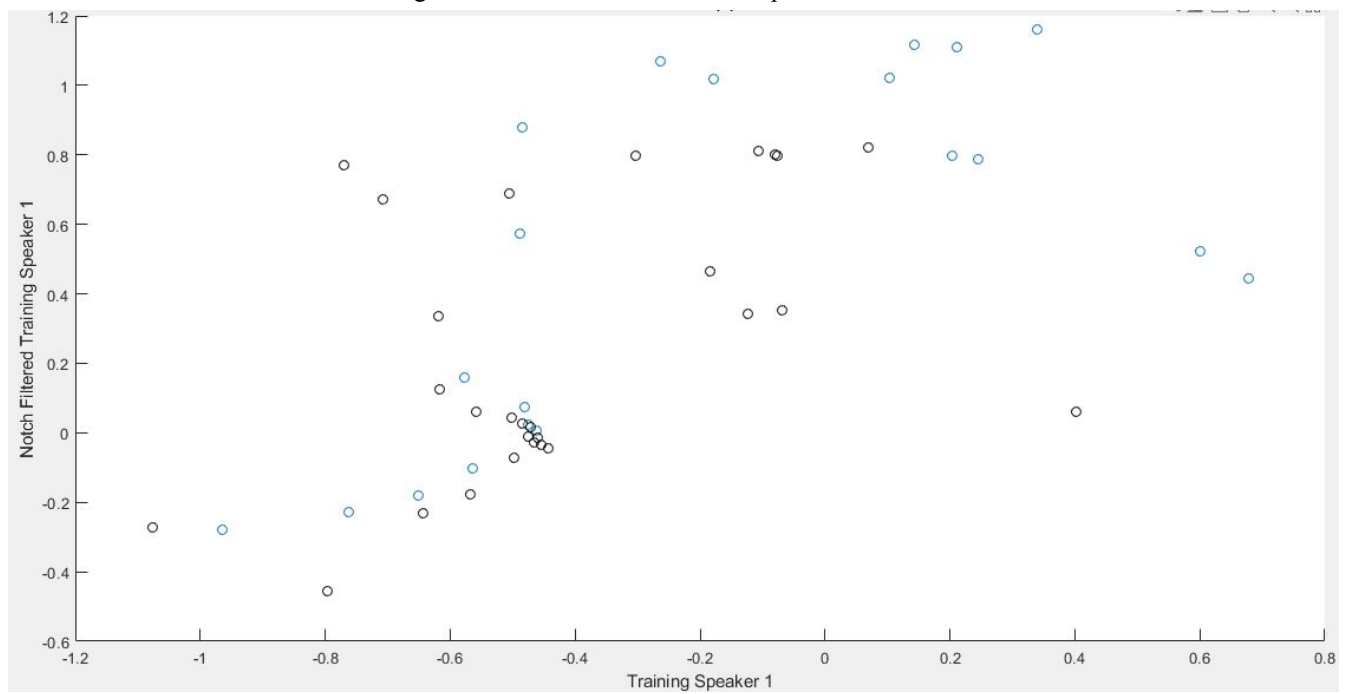


Figure 18: MFCC(1) and MFCC(2) Centroids of Training Speaker 1 Unmodified(Blue) and Notch Filtered(Black)

| Input File | Matched Speaker |
|---|--|
| Speaker ID: <input type="text" value="12"/> | Speaker ID: <input type="text" value="1"/> |

157 Hz is a more dominant frequency where more of the power in the signal is located and the system was able to correctly detect the corresponding speaker.

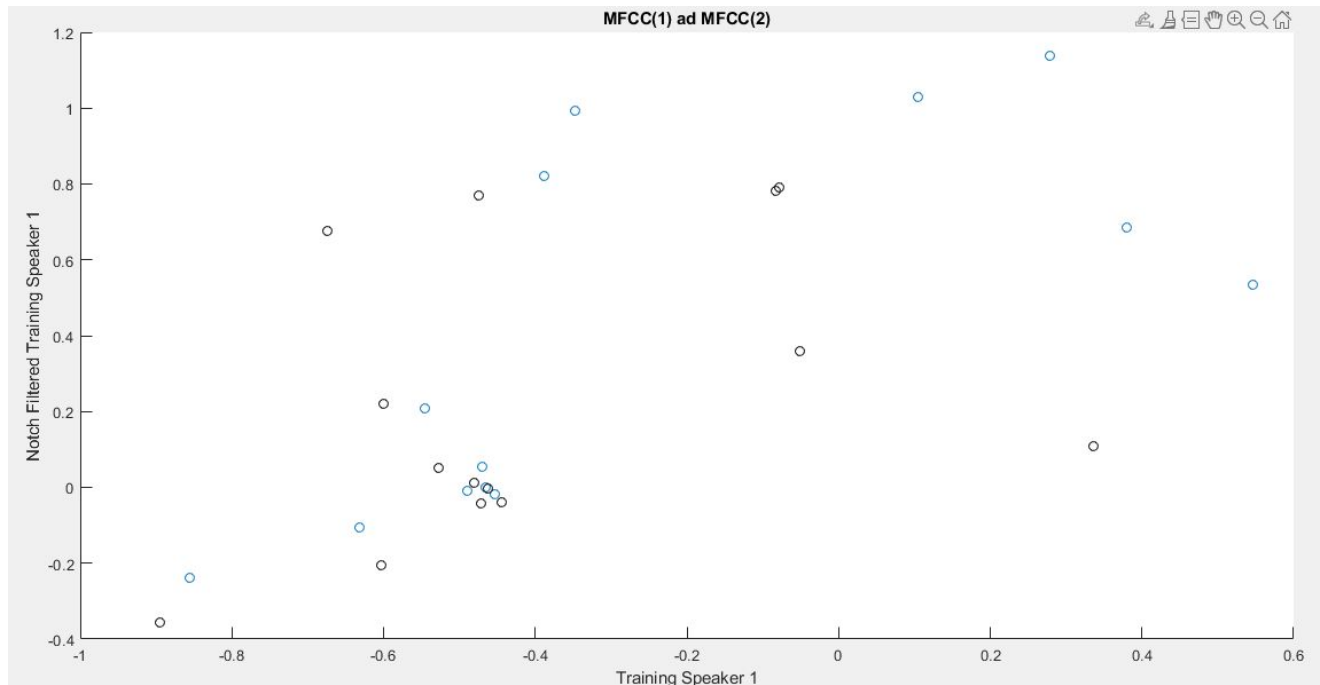


Figure 19: MFCC(1) and MFCC(2) Centroids of Training Speaker 1 Unmodified(Blue) and Notch Filtered(Black) with Codebook Size 16

| Input File | Matched Speaker |
|---|--|
| Speaker ID: <input type="text" value="12"/> | Speaker ID: <input type="text" value="2"/> |

When using a codebook of size 16, the output matched training speaker 2 to the notch filtered training speaker 1, so the distortion between the centroids and the samples was too large compared to the unfiltered version. The system could not properly match the correct speaker compared to a codebook of size 32.

Database

Additional voices were received from the Speech Accent Archives to further test the robustness of the code.

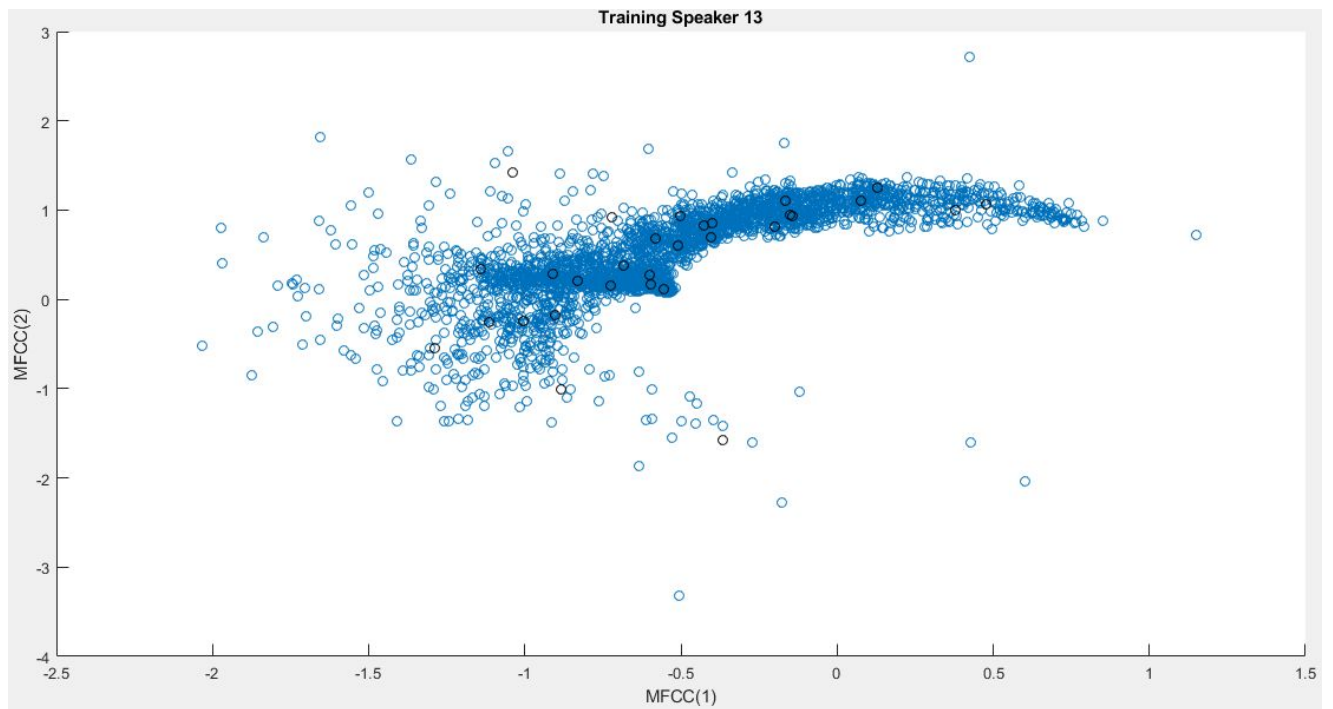


Figure 20: Clustering of MFCC's 1 and 2 of Training Speaker 13, Centroids(Black), Samples(Blue)

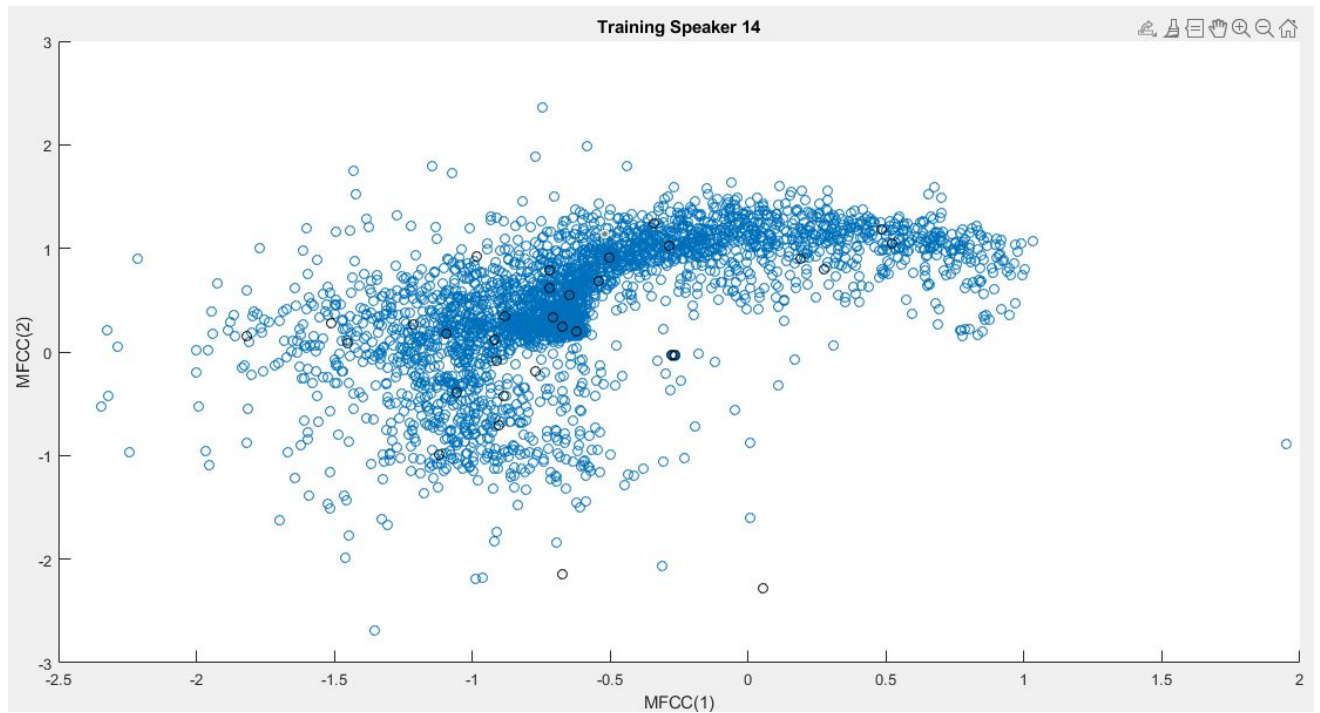


Figure 21: Clustering of MFCC's 1 and 2 of Training Speaker 14, Centroids(Black), Samples(Blue)

| Input File | Matched Speaker |
|----------------|-----------------|
| Speaker ID: 14 | Speaker ID: 13 |

| Input File | Matched Speaker |
|----------------|-----------------|
| Speaker ID: 15 | Speaker ID: 14 |

From the database the number of samples was much higher as the wav file lasted approximately 20 seconds. Using feature extraction and vector quantization, the codebook for each additional training speaker was formed, although the clusters are not very apparent and the appropriate testing speaker matched correctly to the training data.

Conclusions

Using the LBG algorithm for vector quantization accurately determines the centroids of each codebook with a given signal depending on the size of the codebook. The testing speaker mostly matched to the correct training speaker with the exception of using a codebook of size 16 and a notch filter with a null where a prominent frequency in the spectrum is located. This issue is fixed by increasing the codebook size to 32.

References

- [1] Professor Ding, "Speaker Recognition System", 2021
- [2] Y. Linde, A. Buzo & R. Gray, "An algorithm for vector quantizer design", IEEE Transactions on Commu 28, pp.84-95, 1980.
- [3] Y. Tiwari, "MFCC and its applications in speaker recognition", 10 Feb 2010