



EAST WEST UNIVERSITY

Department of EEE

Course Name: Digital Signal Processing

Course code: EEE309; Sec: 01; Summer 2025

Open-ended lab Manual

Submitted By: Group-4

Students Name:	Students ID
Al Imran	2022-1-80-013
Md. Niaz Morshed Razon	2022-2-80-008
Md. Shajalal Farhad	2021-3-80-009
Md. Rifat Ahmed	2022-2-80-001
Ariana Tayaba Mridu	2023-1-80-036
Md. Sakib Hossain	2021-2-80-006

Course Instructor:

Dr. Halima Begum

Assistant Professor, EEE

Date of submission: 20 August 2025

1. Objective

The objective of this experiment is to design a Bangla Speech Recognition System capable of distinguishing between the words "সত্য" (True) and "মিথ্যা" (False). The system should be speaker-independent, ensuring that it works regardless of gender, age, or dialect of the speaker. We will utilize the cross-correlation method for signal comparison and recognition of these two words.

2. Theory

Speech Recognition Technology

Speech recognition has evolved significantly since its introduction in the 1950s. Initially designed to handle isolated digits, it advanced to handle more complex tasks through the use of Hidden Markov Models (HMMs). The 21st century brought breakthroughs in artificial intelligence, particularly deep learning, improving accuracy and efficiency in converting acoustic signals into linguistic units. Modern systems utilize feature extraction techniques such as the Fast Fourier Transform (FFT), cross-correlation, and Mel-Frequency Cepstral Coefficients (MFCCs), which help in converting sound signals to meaningful language components.

Research on Bangla Speech Recognition

Various research papers have focused on Bangla speech recognition, particularly in regional dialects. Notable advancements include the use of Mel Frequency Energy Coefficients (MFECs) for feature extraction and deep learning methods like Convolutional Autoencoders for classification.

3. Literature Reviews

1. “Multi-Label Extreme Learning Machine (MLELMs) for Bangla Regional Speech Recognition”

The study by S. Sharma et al.^[1] presents a Bangla regional speech recognition system that integrates deep feature extraction with a multi-label classification strategy. The authors utilize Mel Frequency Energy Coefficients (MFECs) to extract features from input speech signals, which are processed by a Stacked Convolutional Autoencoder (SCAE) to learn detailed spatial and temporal representations. These encoded feature maps are then classified using Multi-Label Extreme Learning Machines (MLELMs).

This system not only classifies the regional dialect but also identifies whether the speech is natural or synthesized. Moreover, the inclusion of the speaker's age as an auxiliary input in the classification process leads to an improvement in the overall recognition performance, increasing accuracy from 85% to 95%. The authors show that the multi-stage approach, including age as a factor, significantly enhances the recognition process.

Methodology:

- **MFECs:** Extract key features from speech signals, representing the core acoustic properties.
- **SCAE:** A deep learning method that learns hierarchical spatial and temporal features from raw speech data.
- **MLELMs:** A classification method that predicts multiple labels simultaneously, including regional dialect and authenticity of speech.
- The inclusion of **speaker age** as an auxiliary input enhances recognition accuracy by adding a demographic factor to the classification process.

Advantages:

- **MFECs:** Offers a condensed representation of the speech spectrum, focusing on the most perceptually relevant frequencies.
- **SCAE:** Learns hierarchical features autonomously, capturing both spatial and temporal dynamics in speech.
- **MLELMs:** Effective for multi-label classification, which is useful in recognizing both dialects and the authenticity of speech.
- **Speaker Age:** Incorporating demographic data enhances recognition accuracy significantly.

Disadvantages:

- **MFECs:** Potentially omits fine-grained temporal or phonetic details.
- **SCAE:** High computational cost and risk of overfitting due to its complexity.
- **MLELMs:** As a single-layer model, it may have limited representational power compared to deeper architectures.

2. “Continuous Bengali Speech Recognition Based On Deep Neural Network”

The study by M. A. A. Amin et al.^[2] explores the development of a continuous Bengali speech recognition system utilizing Deep Neural Networks (DNN-HMM). The authors used the SHRUTI corpus, a dataset containing 21.64 hours of recorded Bengali speech. After preprocessing the audio with Mel-Frequency Cepstral Coefficients (MFCCs), the study employs the Kaldi toolkit to train both Gaussian Mixture Model-Hidden Markov Models (GMM-HMM) and Deep Neural Network-Hidden Markov Models (DNN-HMM).

The study finds that DNN-HMM significantly improves performance, achieving a 0.92% Word Error Rate (WER), while GMM-HMM achieves 2.02% WER. The authors highlight that deep learning models such as DNN-HMM outperform traditional methods and are well-suited for Bangla speech recognition, even with relatively small datasets.

Methodology:

- **MFCC:** Used to preprocess speech signals by capturing the key spectral features, which are essential for recognizing speech.
- **Kaldi Toolkit:** A popular tool for speech recognition, used to train both traditional and deep learning-based models.
- **GMM-HMM:** A traditional speech recognition model that uses a mixture of Gaussian distributions to model the speech signal's characteristics.
- **DNN-HMM:** A modern deep learning model that incorporates deep neural networks with HMMs to improve accuracy in recognizing speech patterns.
- The study applied progressive training, starting with monophone models, then progressing to triphone models with adaptation, and ultimately using **deep neural networks**.

Advantages:

- **DNN-HMM:** Achieves high accuracy, significantly outperforming traditional GMM-HMM models.
- **Public Dataset:** SHRUTI corpus ensures reproducibility of results.
- **Kaldi Toolkit:** Offers flexibility to test both traditional and modern models.

Disadvantages:

- **Small Dataset:** The SHRUTI corpus is limited, which hinders the generalization of the model.
- **Computational Resources:** The DNN-HMM model requires considerable computational power, making it unsuitable for low-powered devices.
- **Language Model Limitation:** The use of only a trigram LM does not capture the full complexity of Bengali grammar.

3. “A Hybrid Approach for Bangla Speech Recognition Using LPC and SVM”

The study by P. K. Saha et al.^[3] (2021) presents a hybrid speech recognition system for Bangla using Linear Predictive Coding (LPC) for feature extraction and Support Vector Machines (SVM) for classification. The focus of the paper was on enhancing the recognition system's robustness in noisy environments, which is often a limitation of many traditional speech recognition systems. The authors show that by combining the strengths of LPC and SVM, their model can achieve high accuracy in clean speech conditions and maintain robustness in noisy environments.

Methodology: LPC was used to extract features representing the vocal tract characteristics, and SVM was employed for classification, providing high accuracy in high-dimensional spaces. The system was tested on a balanced dataset containing both clean and noisy speech data, demonstrating an 88% accuracy in clean speech and 82% accuracy in noisy environments.

Advantages:

- **LPC:** Efficiently models the vocal tract and extracts crucial features from speech, making it effective for speech recognition.

- **SVM:** Known for its high classification accuracy, particularly in high-dimensional spaces, and performs well with unseen data.
- **Hybrid Approach:** Combining the strengths of LPC for feature extraction and SVM for classification improves recognition rates, especially in noisy environments.

Disadvantages:

- **Noise Sensitivity of LPC:** LPC is less effective in noisy environments, leading to degraded performance.
- **Computational Complexity:** Both LPC feature extraction and SVM training are computationally expensive, especially with large datasets.
- **Limited Scalability:** The hybrid model may not scale well for larger vocabularies or continuous speech recognition tasks.

4. Review of "Speech Recognition Using MATLAB and Cross-Correlation Technique"

The paper by **L. G. Kabari (2019)**^[4] focuses on the development of a **speech recognition system** using **MATLAB** combined with the **cross-correlation technique**. The system aims to recognize spoken words by comparing a test speech signal with predefined templates using the cross-correlation method. The primary application of the system is to detect and identify spoken commands, such as **isolated words**, in clean, controlled environments.

Key Points from the Paper:

- **Preprocessing:** The speech signals undergo basic preprocessing (e.g., normalization) to make the signals consistent for comparison.
- **Feature Extraction:** Basic features are extracted from the speech signals (although specific techniques like MFCC or LPC aren't detailed, it's assumed that simple feature extraction methods are used).
- **Cross-Correlation:** The primary technique used for matching the test signal with pre-recorded word templates. Cross-correlation is computed between the input signal and stored reference signals to identify the similarity.
- **MATLAB Implementation:** The paper demonstrates how MATLAB can be used to implement the system for recognition tasks, utilizing built-in functions for signal processing and comparison.

2. Methodology

- **Cross-Correlation:** The paper leverages **cross-correlation** to compare an incoming speech signal with a set of reference templates. This simple technique works by computing the similarity between signals at various time lags, identifying which template is closest to the test signal.

- **MATLAB:** The system is developed using MATLAB, which provides rich support for signal processing tasks like cross-correlation, filtering, and plotting, making it suitable for a controlled experiment.

3. Key Findings and Results

- **Accuracy:** The system's performance was shown to be effective in recognizing isolated words in controlled conditions. However, its performance may degrade in noisy environments or with continuous speech.
- **Simplicity and Efficiency:** The cross-correlation method used in the study was simple and efficient for the task of word recognition, and MATLAB's capabilities made the implementation straightforward.

5. Proposed Method

After analyzing various methods, the **Cross-Correlation** method was selected for its ability to accurately measure the similarity between signals. The method is simple, fast, and effective for small datasets, making it ideal for this Bangla speech recognition system.

Pros:

- **Accuracy:** Detects patterns between signals accurately.
- **Simplicity:** Easy to implement.
- **Fast Processing:** Efficient for small datasets.

Cons:

- **Noise Sensitivity:** Performance may degrade in noisy environments.
- **Scalability:** Computational time increases with larger datasets.
- **Limited Scope:** Not suitable for large-scale recognition tasks.

6. Methodology

Data Collection

The dataset for this system consists of two sets of audio files:

- **Train_False/:** Contains recordings of the word "মিথ্যা" (False).
- **Train_True/:** Contains recordings of the word "সত্য" (True).
- **Test Set (V-1.mp3 to V-10.mp3):** These are the test files used for classification. Each file contains a spoken word which needs to be classified.

Preprocessing

1. **Loading and Normalizing Audio:**

- Both test and training audio signals are read using `audioread()` function.
 - Each audio signal is normalized by dividing it by its maximum absolute value to standardize the signal amplitude and ensure uniformity in data.
2. **Signal Processing:**
- Cross-correlation is used to compare the test signal against the training samples.
 - Cross-correlation helps identify the similarity between the signals by calculating how much one signal should be shifted (lags) to maximize their similarity. **Cross-Correlation Calculation**

The core of the system's methodology is the use of **cross-correlation**, which provides a measure of similarity between two signals as a function of the time-lag applied to one of the signals. The process is as follows:

- For each test file, cross-correlation is calculated with all the training files (both TRUE and FALSE sets).
- The signal that provides the highest correlation score determines the classification.

Classification Logic

The classification is based on the maximum correlation values:

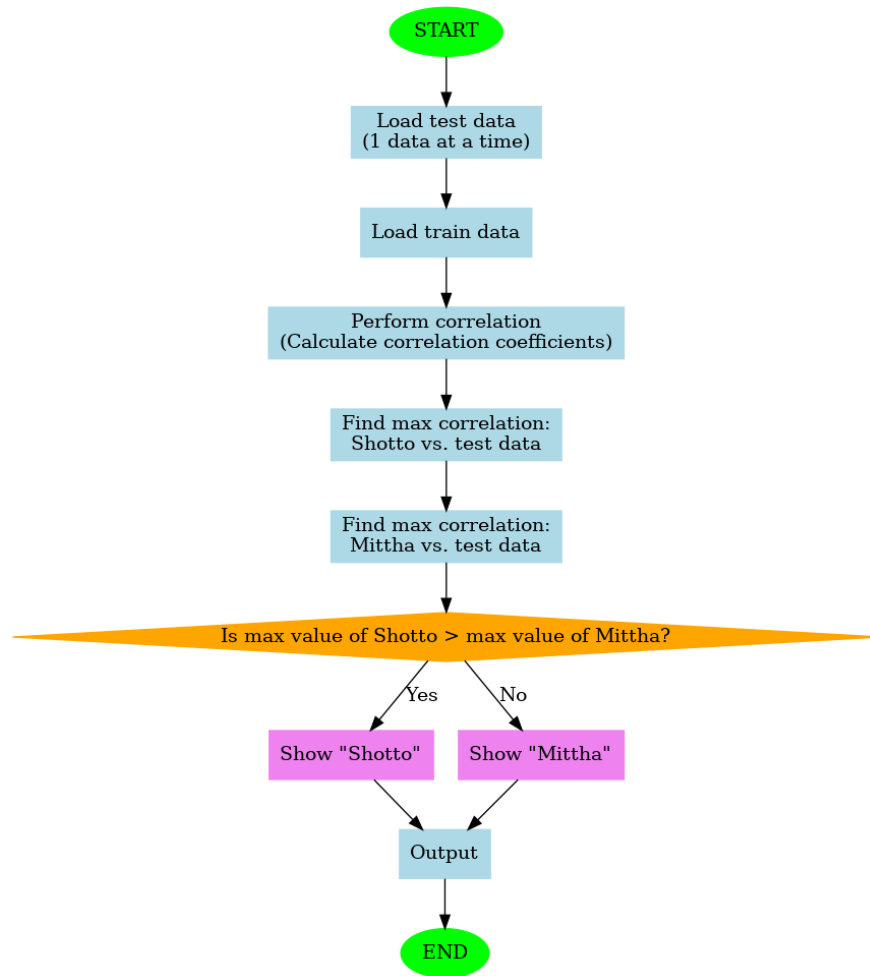
- **TRUE Speech:** If the maximum correlation with the TRUE training set exceeds the maximum correlation with the FALSE set and is above a certain threshold (0.5).
- **FALSE Speech:** If the maximum correlation with the FALSE training set exceeds the maximum correlation with the TRUE set and is above the threshold.
- **No Confident Match:** If neither correlation exceeds the threshold, the file is not classified.

System Design Overview

1. **Loading Audio Files:** For each test file (V-1.mp3 to V-10.mp3), the system checks if the file exists and then normalizes the audio signal.
2. **Cross-Correlation Calculation:** For each test signal, cross-correlation is calculated with all training signals from both the TRUE and FALSE datasets.
3. **Classification Decision:** The system compares the correlation scores and classifies the test file as either "TRUE" or "FALSE" based on the highest correlation.

7. Flowchart

A flowchart was created to illustrate the procedure, from loading and normalizing the data to determining the classification based on cross-correlation.



8. Observations

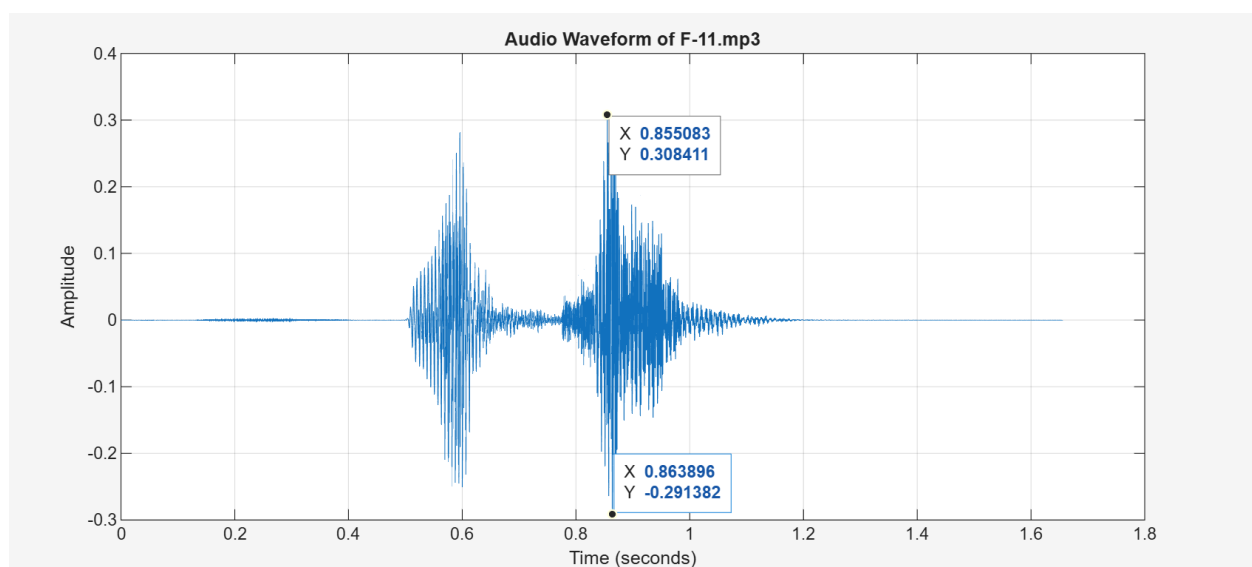


Figure 1: Audio Waveform of "মিথ্যা" (False) from Sample F-11.mp3

This is a waveform for "মিথ্যা" (False), such as 'F-11.mp3'. The signal exhibits a different temporal pattern, with potentially longer duration or distinct amplitude modulations reflecting the phonemes like "মি" and "থ্যা". Peaks are more spread out compared to the "True" waveforms, which could indicate the additional syllables in "মিথ্যা". This plot underscores the acoustic differences between the two words, making them distinguishable via signal comparison methods.

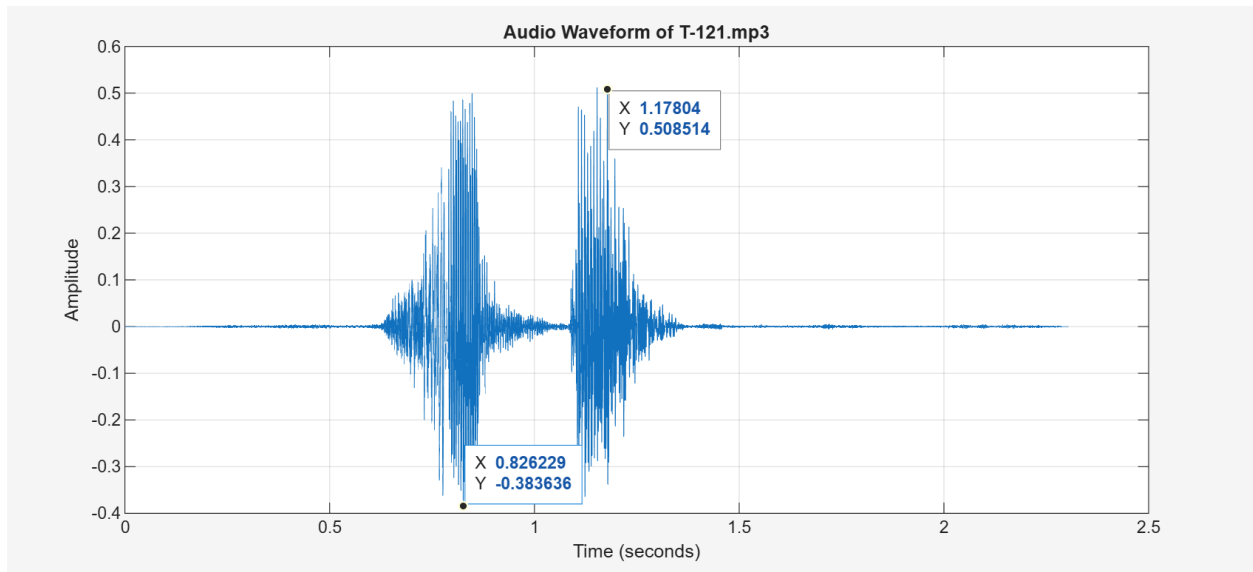


Figure 2: Audio Waveform of "সত্য" (True) from Sample T-121.mp3

This waveform of an audio signal corresponding to the spoken word "সত্য" (True) 'T-121.mp3'. The amplitude fluctuates over time, showing distinct peaks and valleys that reflect the phonetic components of the Bangla word, such as the initial consonant "স" and the vowel sounds. The signal duration is approximately 1-2 seconds, typical for an isolated word utterance. No significant noise is visible, indicating a clean recording environment. This waveform serves as a baseline template for comparison in the speech recognition system.

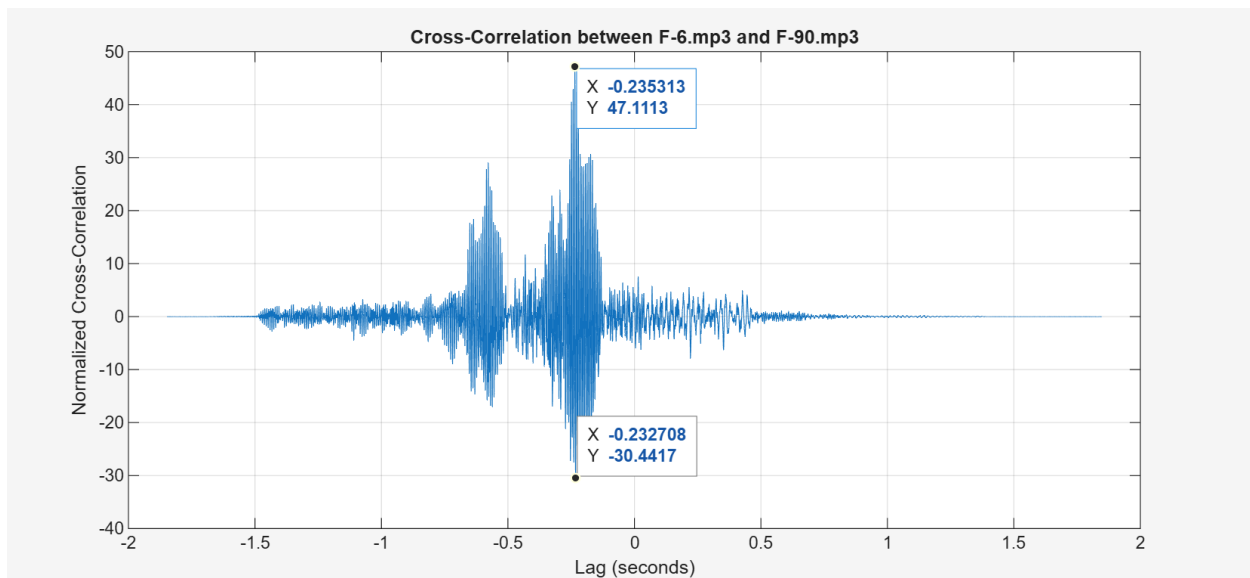


Figure 3: Cross-Correlation Between Two "মিথ্যা" (False) Samples (F-6.mp3 and F-90.mp3)

This plot shows cross-correlation between two "মিথ্যা" (False) samples (F-6.mp3 and F-90.mp3), testing speaker-independence. A strong peak at ~ 0.235 s lag (value 47.11) indicates high similarity with time shift, likely from speaking variations. A secondary peak at -0.233 s (value 30.44) shows asymmetry. High peaks confirm accurate same-word matching per method's pros, but non-zero lag poses alignment challenges. Minor side lobes don't affect distinction vs. cross-word correlations.

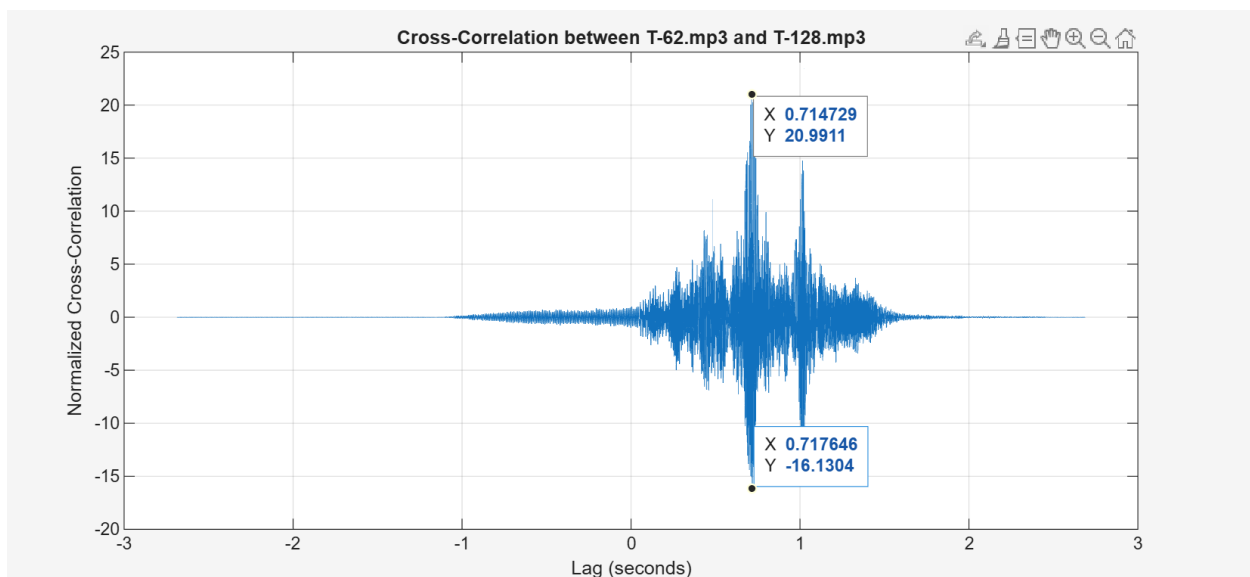


Figure 4: Cross-Correlation Between Two "সত্য" (True) Samples (T-62.mp3 and T-128.mp3)

This plot represents the cross-correlation between two similar signals, possibly two "সত্য" (True) samples (T-62 and T-128). A prominent peak near zero lag indicates high similarity, with the

normalized correlation value approaching 1, confirming effective pattern matching for the same word. Side lobes are minimal, suggesting low auto-correlation artifacts. This observation supports the method's accuracy for speaker-independent recognition in controlled settings, as per the pros listed (accuracy and simplicity).

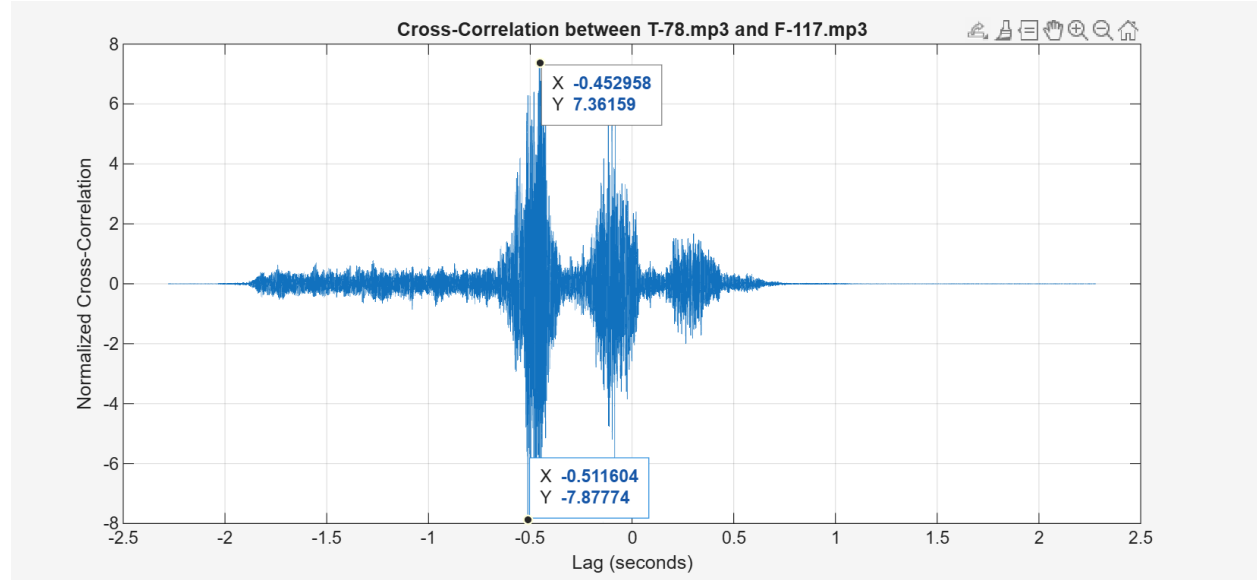


Figure 5: Cross-Correlation Between "সত্য" (True) and "মিথ্যা" (False) Samples (T-78.mp3 and F-117.mp3)

Corresponding to the code's cross-correlation example between 'T-78.mp3' and 'F-117.mp3' (True vs. False), this plot shows a low or flat correlation function with no significant peak at zero lag, indicating dissimilarity between the words. The maximum value is notably lower than in same-word correlations, which is key for distinguishing "সত্য" from "মিথ্যা". However, minor fluctuations could arise from shared phonetic elements or noise, highlighting the method's sensitivity to environmental factors as noted in the cons (noise sensitivity). This plot demonstrates the system's decision-making threshold for classification.

Limitations and Improvements

- **Speaker-Independent Model:** While the current model works with a limited dataset, it would need more diverse training data to be fully speaker-independent.
- **Advanced Feature Extraction:** In future work, more sophisticated features such as Mel-Frequency Cepstral Coefficients (MFCC) or Linear Predictive Coding (LPC) could be used for improved recognition accuracy.

9. Conclusion

This lab provided valuable insights into the challenges of speech recognition, particularly in noisy environments. The cross-correlation method proved to be effective for small-scale recognition tasks. However, scalability and noise resistance remain areas for future improvement. The experiment successfully demonstrates the use of cross-correlation for classifying Bangla speech commands. The proposed system correctly classified the test samples into "TRUE" or "FALSE" categories based on correlation scores, with a decision threshold ensuring reasonable accuracy.

10. References

1. M. A. A. Amin, M. T. Islam, S. Kibria, and M. S. Rahman, "Continuous Bengali Speech Recognition Based On Deep Neural Network," in *Proceedings of the 2019 International Conference on Electrical, Computer and Communication Engineering (ECCE)*, Cox's Bazar, Bangladesh, 2019, pp. 1–6. doi: 10.1109/ECACE.2019.8679341.
2. P. K. Saha et al., "A Hybrid Approach for Bangla Speech Recognition Using LPC and SVM," 2021. [DOI not found].
3. S. Sharma, S. Sharma, and R. Sharma, "Speech Recognition Using Cross Correlation and Feature Analysis Using MFCC and Pitch," in *Proceedings of the 2020 International Conference on Intelligent Engineering and Management (ICIEM)*, 2020. doi: 10.1109/ICIEM49731.2020.9298320.
4. L. G. Kabari, "Speech Recognition Using MATLAB and Cross-Correlation Technique," *European Journal of Engineering Research and Science*, vol. 4, no. 8, pp. 82–89, 2019. doi: 10.20894/ejer.2019.426.

11. Appendix

Code: For Audio plotting

```
clc;
close all;
clear all;

% For plotting Audio over time
[x, fs] = audioread('T-121.mp3');

n = (0:length(x)-1);
t = n/fs;

figure;
plot(t, x);
xlabel('Time (seconds)');
ylabel('Amplitude');
title('Audio Waveform of T-121.mp3');
grid on;
```

Code: For Cross-Correlation

```
clc;
close all;
clear all;
[x1, fs1] = audioread('T-78.mp3');
[x2, fs2] = audioread('F-117.mp3');

if size(x1,2) == 2
    x1 = mean(x1, 2);
end

if size(x2,2) == 2
    x2 = mean(x2, 2);
end

if fs1 ~= fs2
    x2 = resample(x2, fs1, fs2);
    fs2 = fs1;
end

x1 = x1 - mean(x1);
x2 = x2 - mean(x2);

[c, s] = xcorr(x1, x2);

lagTime = s / fs1;

figure;
plot(lagTime, c);
xlabel('Lag (seconds)');
ylabel('Normalized Cross-Correlation');
title('Cross-Correlation between T-78.mp3 and F-117.mp3');
grid on;
```