

ASSIGNMENT-2

Task 2 & 1

Roll No :- M23CSA512

Name :- Chandra Mohan Singh Negi

1. Introduction

This is a massive dataset of audio samples of 10 different Indian languages. Each audio sample is of 5 seconds duration. This dataset was created using regional videos available on YouTube. None of the audio samples/source videos are owned by me, and the dataset must not be used to create any proprietary applications.

This is constrained to Indian Languages only but could be extended.

Languages present in the dataset –

Bengali, Gujarati, Hindi, Kannada, Malayalam, Marathi, Punjabi, Tamil, Telugu, Urdu.

Overview

- **Name:** Audio Dataset with 10 Indian Languages
- **Size:** 19 GB
- **Format:** Mp3

Task A:

1. Downloaded the data from Kaggle.
2. Written the python code for extract the Mel-Frequency Cepstral Coefficients (MFCC) from each audio sample.
3. MFCC spectrograms for a Malayalam, Tamil and Urdu languages with 5-5 samples:

Git Hub Link :- <https://github.com/cmnegi/speechassignment2.git>

1- AudioDatasets :- <https://www.kaggle.com/datasets/hbchaitanyabharadwaj/audio-dataset-with-10-indian-languages>

- **Number of Classes:** 10

1. **Bengali**
2. **Gujarati**
3. **Hindi**
4. **Kannada**

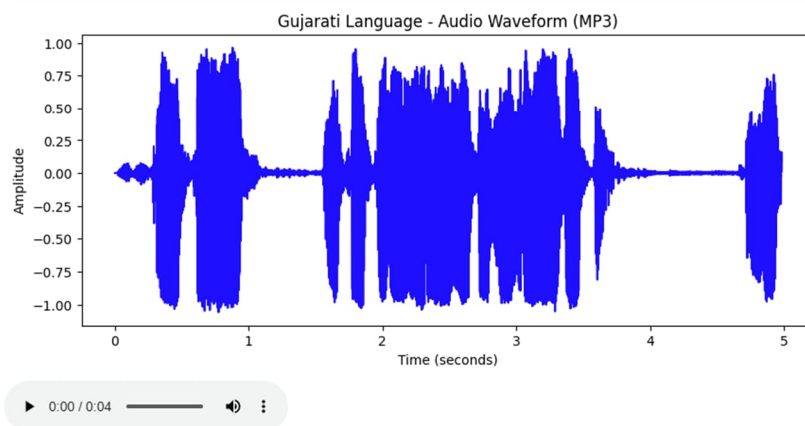
5. Malayalam
6. Marathi
7. Punjabi
8. Tamil
9. Telugu
10. Urdu

MFCC Visualization of All Languages :-

1- Gujarati

```
plt.figure(figsize=(10, 4))
plt.plot(np.linspace(0, len(data) / sample_rate, num=len(data)), data, color='b')
plt.title("Gujarati Language - Audio Waveform (MP3)")
plt.xlabel("Time (seconds)")
plt.ylabel("Amplitude")
plt.show()

# Play the audio file
file = ipd.Audio(file_name_gujarati)
ipd.display(file)
```

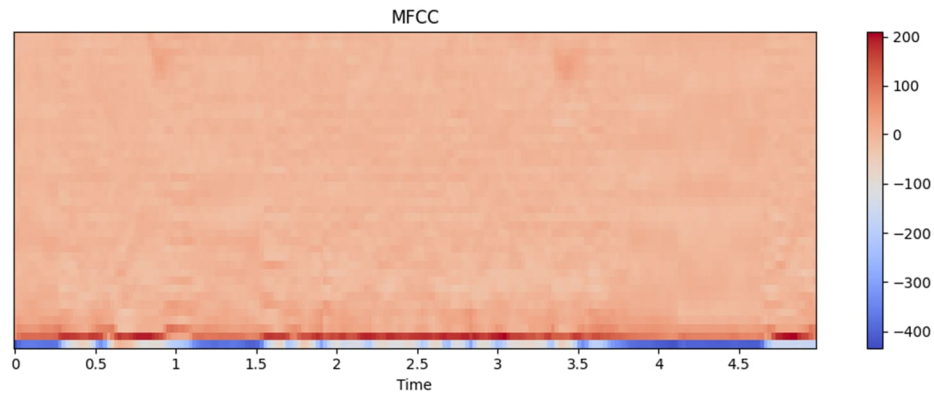


```

y, sr = librosa.load(file_name_gujarati)
mfccs = librosa.feature.mfcc(y=y, sr=sr, n_mfcc=40)

plt.figure(figsize=(10,4))
librosa.display.specshow(mfccs, x_axis="time")
plt.colorbar()
plt.title('MFCC')
plt.tight_layout()
plt.show()

```



```

Mean: [0.47943458 0.50212143 0.54538013 0.51877504 0.51577301 0.53215099
0.51999422 0.50061632 0.52255751 0.50309147 0.48058784 0.53680015
0.49822353]
Variance: [0.07130771 0.08835882 0.07721942 0.07147098 0.08892057 0.08021719
0.09298865 0.08144248 0.07208567 0.07740621 0.07786006 0.08935853
0.07909149]
Standard Deviation: [0.26703503 0.29725212 0.27788382 0.26734057 0.29819553 0.28322639
0.30494041 0.28538129 0.26848774 0.27821972 0.27903416 0.29892897
0.28123208]
Skewness: [ 0.05940423  0.15136724 -0.3158738  -0.10722659 -0.08732689  0.02918366
 0.01325337 -0.12803952 -0.20112304 -0.06184659 -0.00557441 -0.28570114
-0.11098834]

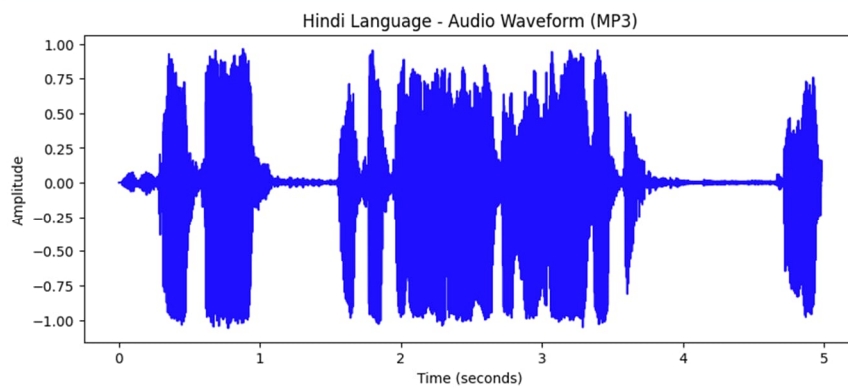
```

2- Hindi

```
data, sample_rate = librosa.load(file_name_gujarati, sr=22050) # CONVERTS MP3 TO WAVEFORM

# Plot the waveform
plt.figure(figsize=(10, 4))
plt.plot(np.linspace(0, len(data) / sample_rate, num=len(data)), data, color='b')
plt.title("Hindi Language - Audio Waveform (MP3)")
plt.xlabel("Time (seconds)")
plt.ylabel("Amplitude")
plt.show()

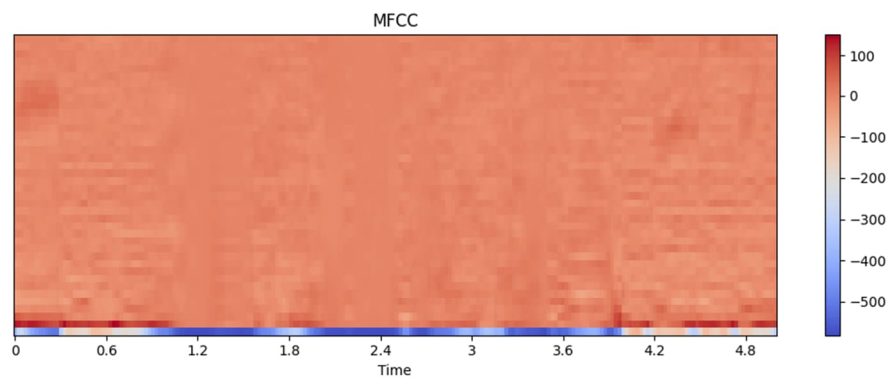
# Play the audio file
file = ipd.Audio(file_name_gujarati)
ipd.display(file)
```



▶ 0:00 / 0:04 ———— 🔊 ⋮

```
[36]: y, sr = librosa.load(file_name_Hindi)
mfccs = librosa.feature.mfcc(y=y, sr=sr, n_mfcc=40)

plt.figure(figsize=(10,4))
librosa.display.specshow(mfccs, x_axis="time")
plt.colorbar()
plt.title('MFCC')
plt.tight_layout()
plt.show()
```



```
[35]: librosa.feature.mfcc(y=data, sr=sample_rate, n_mfcc=40) #transform the linear frequency scale into mel scale

[35]: array([[ -4.33729889e+02, -3.93061249e+02, -3.69287415e+02, ...,
        -1.66823380e+02, -1.74147354e+02, -1.75953125e+02],
        [ 5.13561707e+01,  9.69344864e+01,  1.14936676e+02, ...,
        1.51856262e+02,  1.16981201e+02,  8.57036743e+01],
        [ 3.99342651e+01,  6.37828903e+01,  5.78525467e+01, ...,
        5.80206223e+01,  9.07744904e+01,  8.04794617e+01],
        ...,
        [-5.61272764e+00, -8.88840675e+00, -7.44464302e+00, ...,
        -5.93054390e+00, -7.73384714e+00, -5.13449478e+00],
        [-4.64803982e+00, -9.62656212e+00, -1.02600403e+01, ...,
        -8.68212223e+00, -6.65292072e+00, -5.65175629e+00],
        [-1.81429017e+00, -3.91899276e+00, -7.63073397e+00, ...,
        -2.43136883e-02, -2.55825543e+00, -5.17338514e-03]], dtype=float32)
```

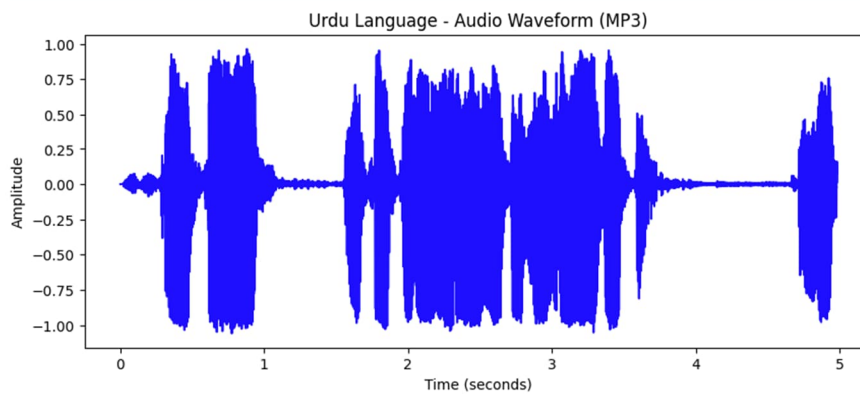
3- Urdu

```
if not os.path.exists(file_name_gujarati):
    raise FileNotFoundError(f"File not found: {file_name_gujarati}")

# Load the MP3 file using librosa
data, sample_rate = librosa.load(file_name_gujarati, sr=22050) # Converts MP3 to waveform

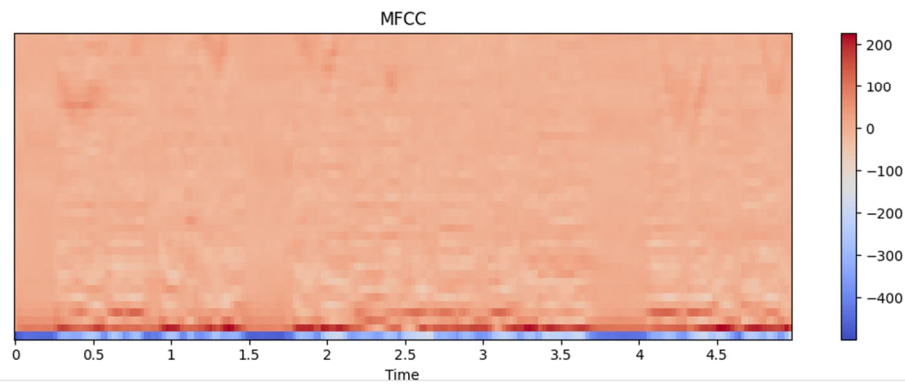
# Plot the waveform
plt.figure(figsize=(10, 4))
plt.plot(np.linspace(0, len(data) / sample_rate, num=len(data)), data, color='b')
plt.title("Urdu Language - Audio Waveform (MP3)")
plt.xlabel("Time (seconds)")
plt.ylabel("Amplitude")
plt.show()

# Play the audio file
file = ipd.Audio(file_name_gujarati)
ipd.display(file)
```



```
[38]: y, sr = librosa.load(file_name_Urdu)
mfccs = librosa.feature.mfcc(y=y, sr=sr, n_mfcc=40)

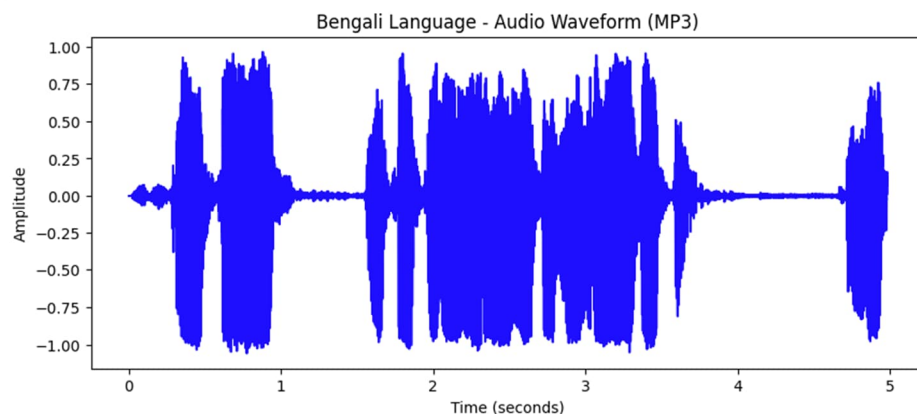
plt.figure(figsize=(10,4))
librosa.display.specshow(mfccs, x_axis="time")
plt.colorbar()
plt.title('MFCC')
plt.tight_layout()
plt.show()
```



```
print( 'skewness: ', skew_mfcc)
```

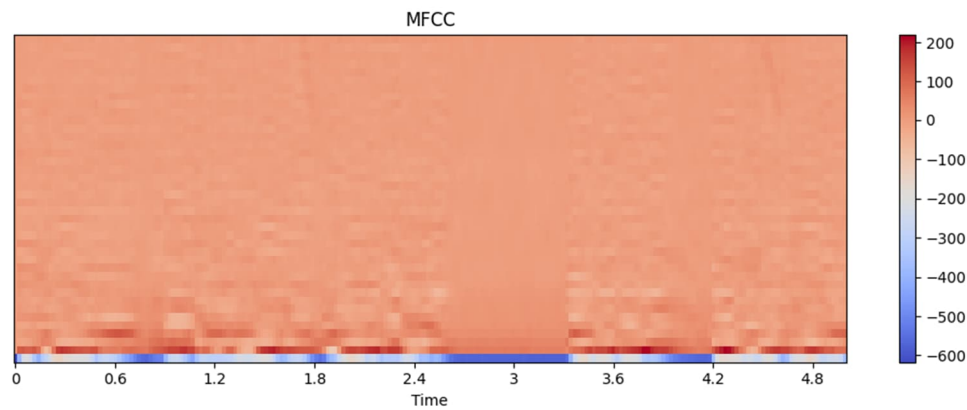
```
Mean: [0.43016592 0.52679954 0.50470522 0.45796002 0.48518968 0.51564611
0.46202749 0.52629167 0.52611957 0.47840238 0.54522761 0.51781716
0.45298636]
Variance: [0.08075134 0.08778438 0.0851494 0.07409958 0.09737862 0.09324479
0.08441666 0.08486198 0.08212323 0.08634418 0.07270928 0.09155734
0.08543003]
Standard Deviation: [0.2841678 0.29628428 0.29180371 0.27221239 0.31205548 0.3053601
0.29054545 0.2913108 0.28657152 0.29384381 0.26964658 0.30258443
0.29228415]
Skewness: [ 0.37177859 0.05970486 -0.00863739 0.20329474 -0.0408195 -0.09105267
0.19256863 -0.08343713 -0.03051274 -0.0103004 -0.3286517 -0.05040583
0.13025257]
```

4- Bengali



```
[50]: y, sr = librosa.load(file_name_Bengali)
mfccs = librosa.feature.mfcc(y=y, sr=sr, n_mfcc=40)

plt.figure(figsize=(10,4))
librosa.display.specshow(mfccs, x_axis="time")
plt.colorbar()
plt.title('MFCC')
plt.tight_layout()
plt.show()
```



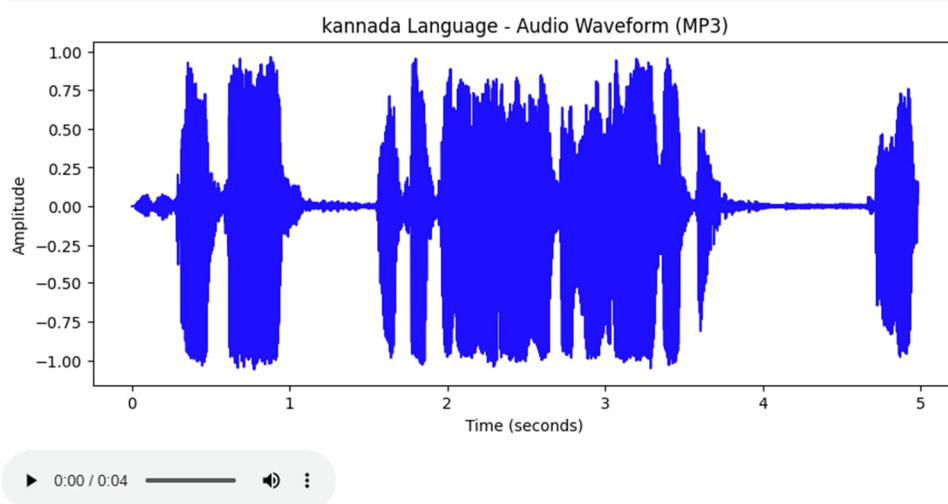
```
print("Skewness:", skew_mfcc)
```

```
Mean: [0.5277268 0.49229345 0.48855715 0.43901076 0.53882268 0.49839068
0.4991732 0.45353171 0.4615267 0.47299518 0.46469268 0.42325487
0.50369613]
Variance: [0.07917119 0.08821159 0.08813639 0.08043149 0.08493619 0.088432
0.0761867 0.07412333 0.0813514 0.09553836 0.08392662 0.08295828
0.07802093]
Standard Deviation: [0.28137376 0.29700437 0.29687774 0.28360446 0.29143815 0.29737518
0.27601939 0.272256 0.28522166 0.3090928 0.28970091 0.28802479
0.27932227]
Skewness: [-0.27299747 -0.0925285 0.082244 0.25627203 -0.17056258 0.03691281
0.02745491 0.03021083 0.13871552 0.20218752 0.15910984 0.31721737
-0.00188502]
```

+ Code

+ Markdown

5- Kannada

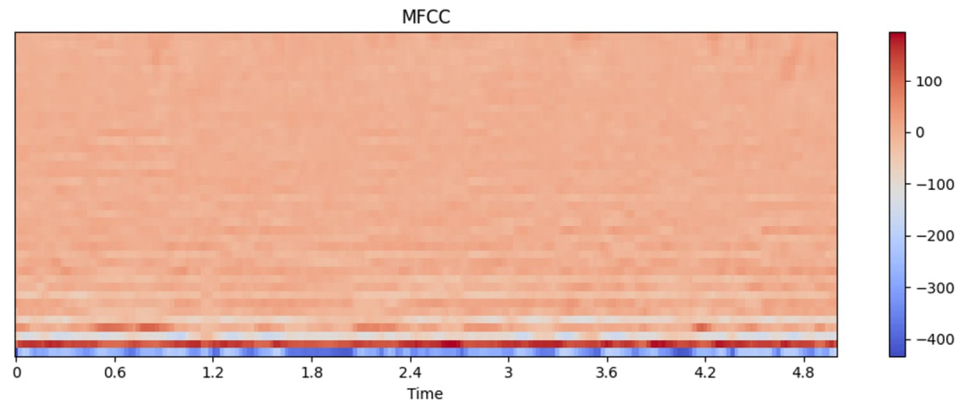



```

y, sr = librosa.load(file_name_Kannada)
mfccs = librosa.feature.mfcc(y=y, sr=sr, n_mfcc=40)

plt.figure(figsize=(10,4))
librosa.display.specshow(mfccs, x_axis="time")
plt.colorbar()
plt.title('MFCC')
plt.tight_layout()
plt.show()

```



```
print("Skewness:", skew_mfcc)
```

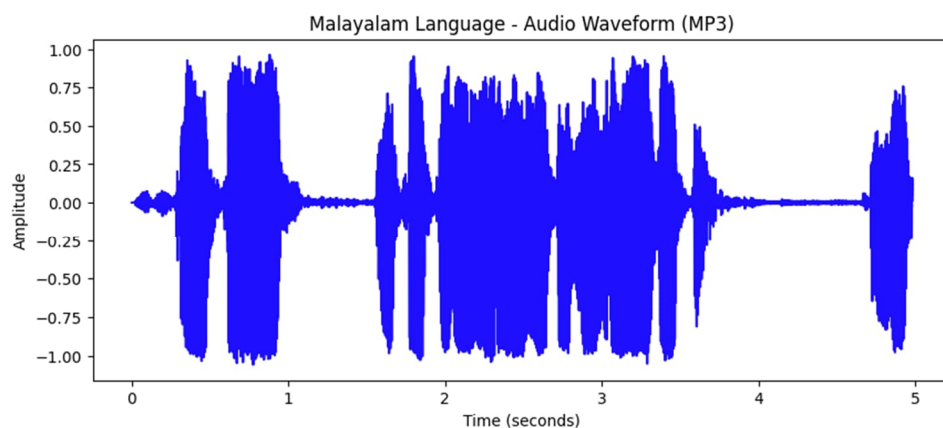
```

Mean: [0.56479191 0.47148119 0.52759461 0.49348683 0.48865294 0.48050351
0.5083505 0.49534836 0.45583233 0.45593151 0.48473046 0.46425639
0.49994388]
Variance: [0.08283221 0.0797247 0.08909629 0.07453563 0.07706541 0.07853085
0.08333407 0.08690841 0.07558171 0.08133068 0.07847814 0.08346739
0.08307864]
Standard Deviation: [0.28780586 0.28235563 0.29849002 0.27301215 0.27760657 0.28023357
0.28867641 0.29480233 0.27492128 0.28518534 0.2801395 0.28890723
0.28823366]
Skewness: [-0.26570661 0.11839465 -0.03960075 0.09969606 0.01608404 -0.07227238
-0.14752564 0.17913717 0.07243583 0.16465231 -0.05694981 0.27152375
0.08008953]

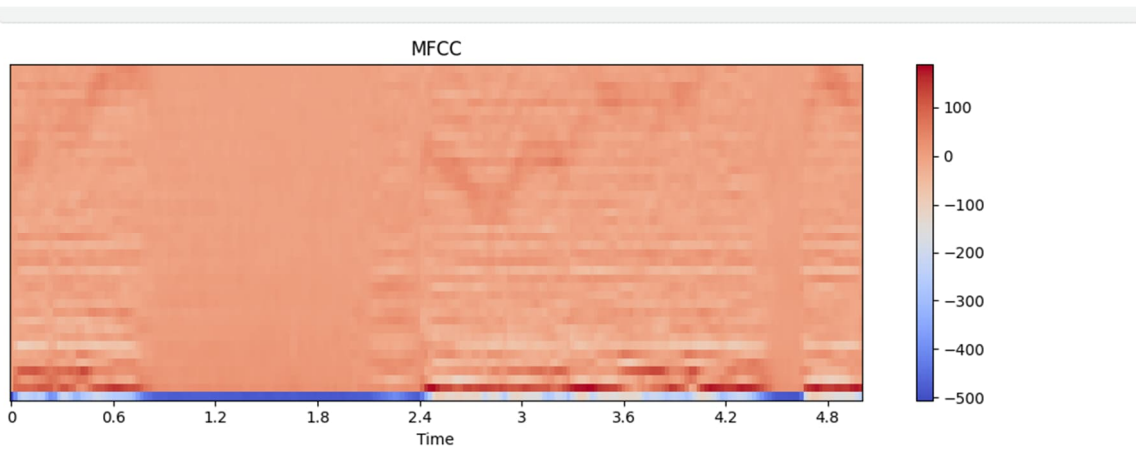
```

6- Malayalam

```
ipd.display(title)
```



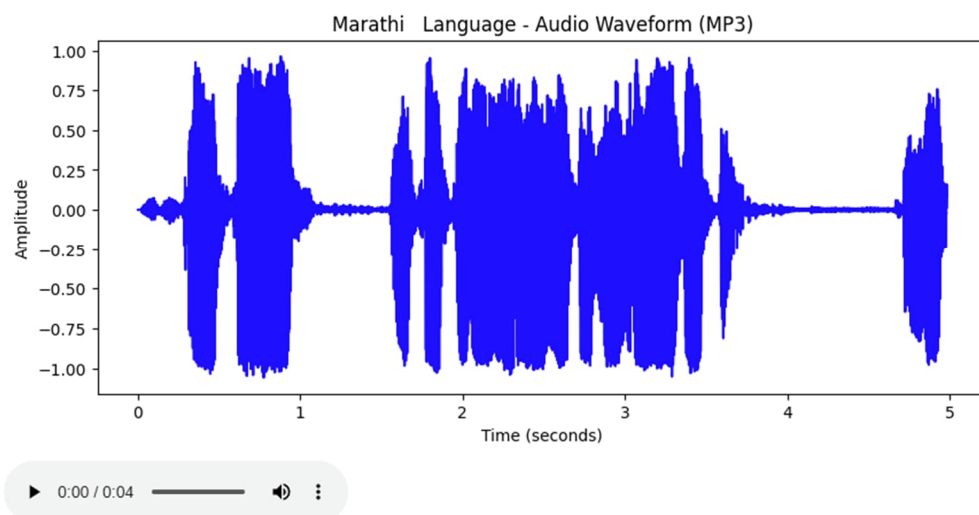
▶ 0:00 / 0:04 🔊 ⋮

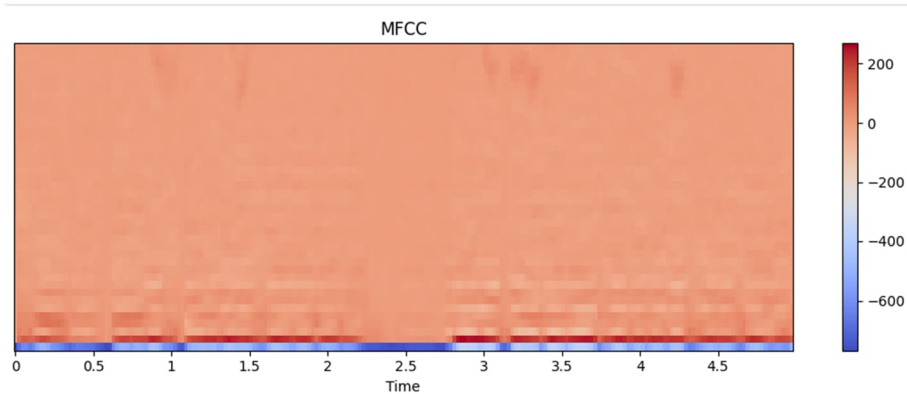


```
Mean: [0.47181981 0.54592721 0.49973634 0.53571831 0.46573709 0.51635835
0.49864371 0.50770519 0.52845987 0.49055221 0.47892576 0.49018885
0.51303909]
Variance: [0.08236557 0.08245916 0.0731705 0.08088813 0.07936447 0.08723604
0.08422056 0.07463733 0.09054855 0.08161832 0.08178418 0.09757322
0.09337835]
Standard Deviation: [0.28699403 0.28715703 0.27050047 0.28440839 0.28171701 0.29535747
0.29020778 0.27319834 0.30091287 0.2856892 0.28597933 0.31236713
0.30557871]
Skewness: [-0.15461938 -0.28714464 0.11323128 -0.09522082 0.10947651 -0.14793756
-0.03090391 0.02380603 -0.10871706 0.02218676 0.07168017 -0.00728595
0.08532891]
```

7- Marathi

```
ipd.display(title)
```



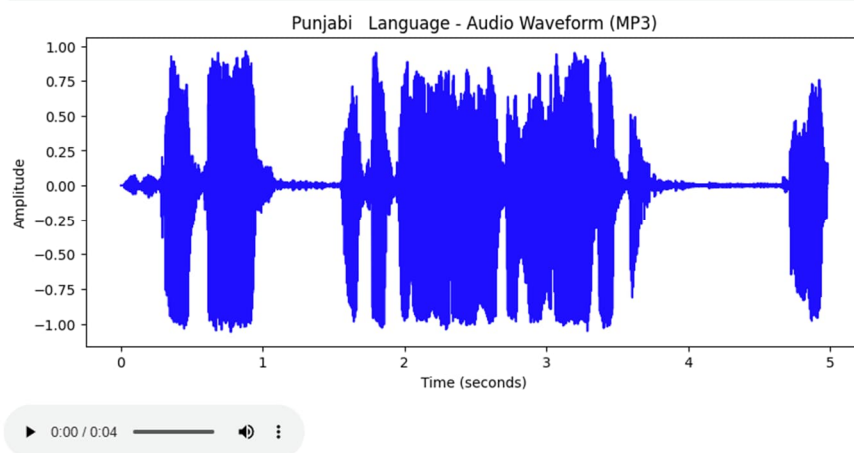


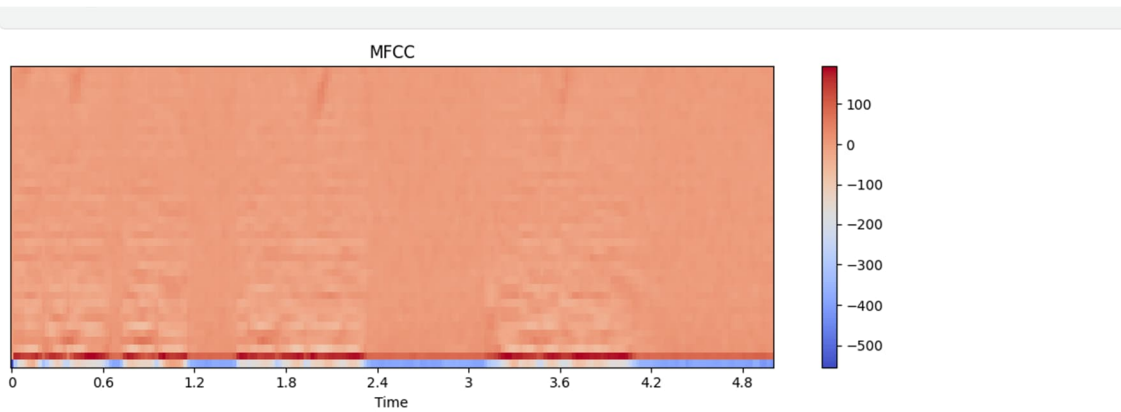
```
Mean: [0.51920899 0.52362499 0.44625139 0.47723077 0.4833424 0.52642783
0.47195033 0.45643519 0.50372929 0.55335274 0.46492365 0.51926805
0.45848107]
Variance: [0.081255 0.07042727 0.08540466 0.09427745 0.08627587 0.07867391
0.08250852 0.07971725 0.08575101 0.09003096 0.07487916 0.07945336
0.07398773]
Standard Deviation: [0.28505263 0.26538136 0.29224075 0.30704633 0.29372755 0.28048869
0.28724297 0.28234243 0.29283273 0.3000516 0.27364056 0.28187473
0.27200686]
Skewness: [-0.06776927 -0.08275824 0.21871161 0.10156936 0.13232921 -0.10234554
0.14780469 0.3019966 -0.05803196 -0.28140203 0.10235115 -0.11796209
0.1551162 ]
```

+ Code

+ Markdown

8- Punjabi



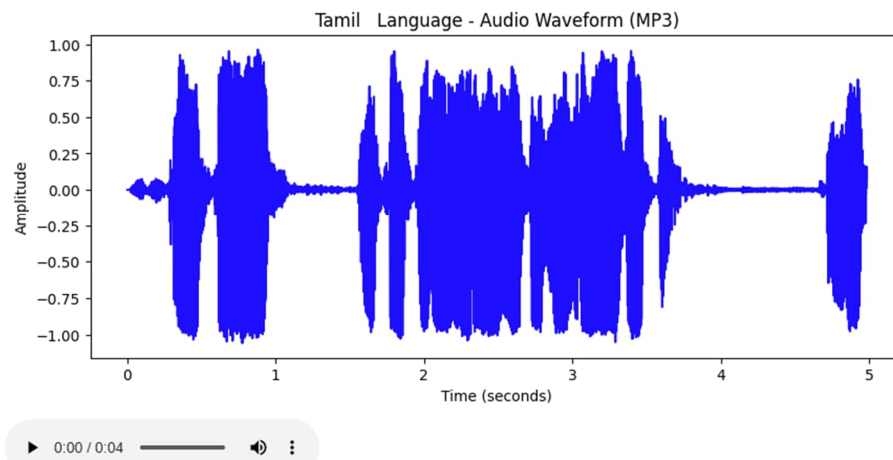


Mean: [0.48845202 0.50634374 0.4970339 0.50390801 0.5345121 0.53799543
0.49386534 0.545038 0.49463794 0.50809959 0.47853146 0.52460762
0.46219368]
Variance: [0.08465942 0.09183098 0.0803364 0.08037736 0.07840104 0.09174209
0.08383724 0.08400177 0.08673944 0.08777861 0.08032252 0.08487146
0.08732328]
Standard Deviation: [0.29096292 0.30303626 0.28343676 0.28350901 0.28000185 0.30288957
0.28954661 0.28983059 0.2945156 0.29627455 0.28341228 0.29132707
0.29550512]
Skewness: [0.03837878 -0.0435179 -0.00401167 -0.00414938 -0.04176937 -0.21083798
0.0308109 -0.12569915 -0.04366865 0.04570833 0.13114806 -0.00753259
0.07736928]

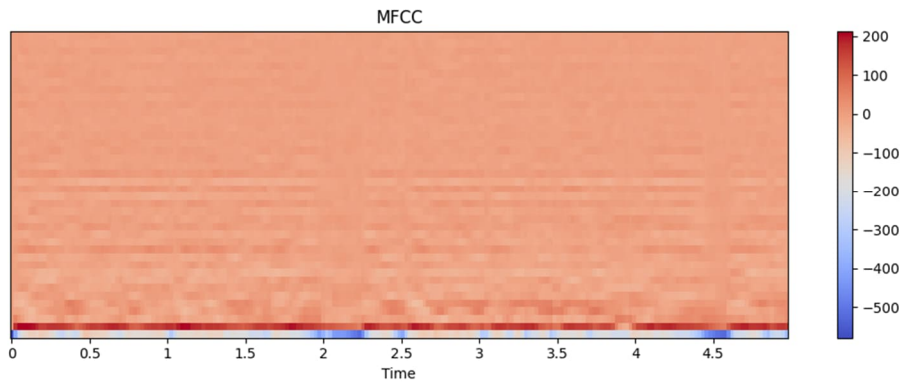
+ Code

+ Markdown

9- Tamil



plot_spectrogram()

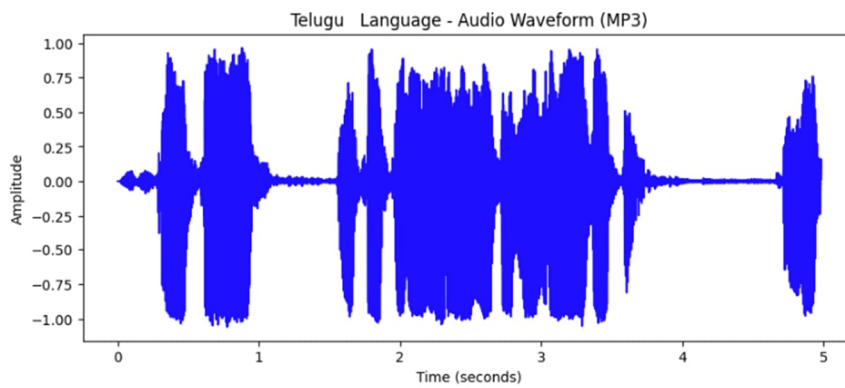


```
Mean: [0.50094297 0.52631923 0.51922954 0.52999249 0.47406939 0.50051568
0.47694735 0.50793047 0.49698028 0.52481752 0.46682995 0.5141788
0.47099986]
Variance: [0.09080617 0.07713152 0.08488797 0.08714761 0.07045519 0.0824636
0.0718852 0.06848465 0.0855268 0.08137093 0.08036398 0.08850411
0.09383964]
Standard Deviation: [0.30134062 0.27772561 0.29135539 0.29520773 0.26543397 0.28716476
0.26811416 0.26169571 0.29244965 0.2852559 0.28348541 0.2974964
0.30633257]
Skewness: [ 0.00041508 -0.20921661 -0.18718215 -0.13962078 -0.27222641 -0.12337603
0.1014631 -0.01561596 -0.00890277 0.00312913 0.1865944 0.0114972
0.06959016]
```

+ Code

+ Markdown

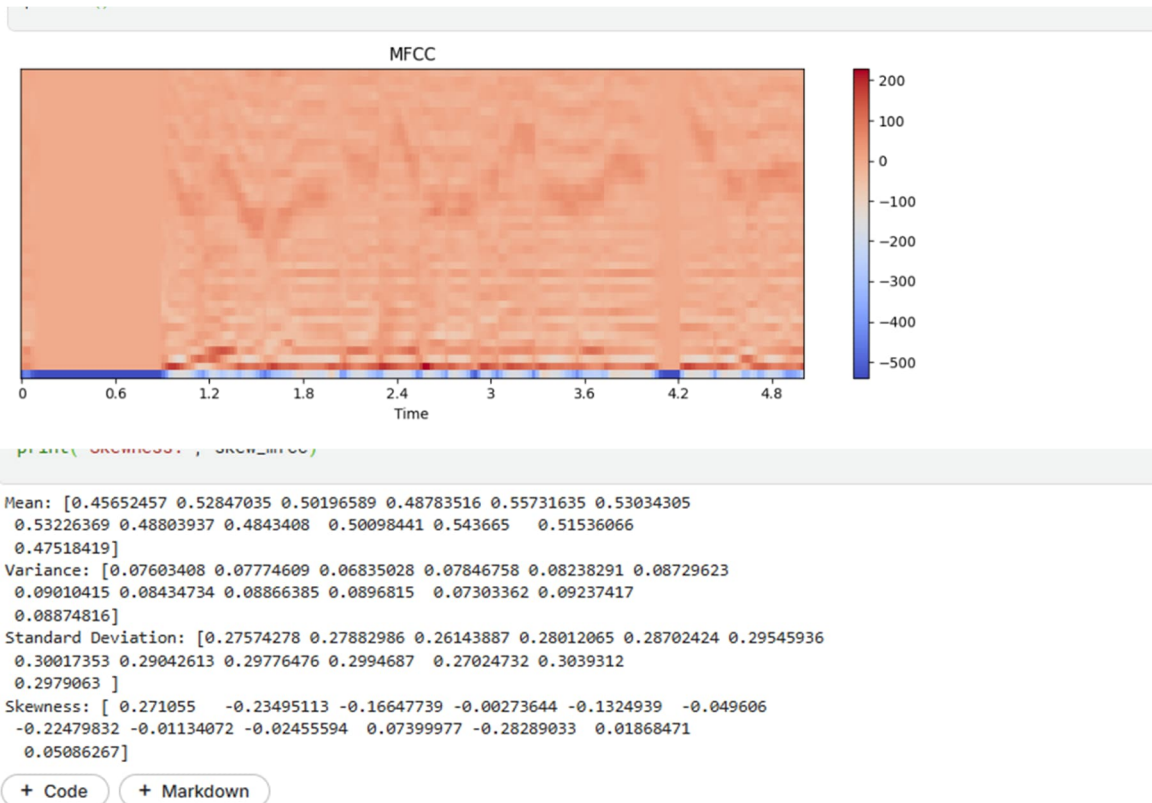
10- Telugu



▶ 0:00 / 0:04 🔊 ⋮

+ Code

+ Markdown

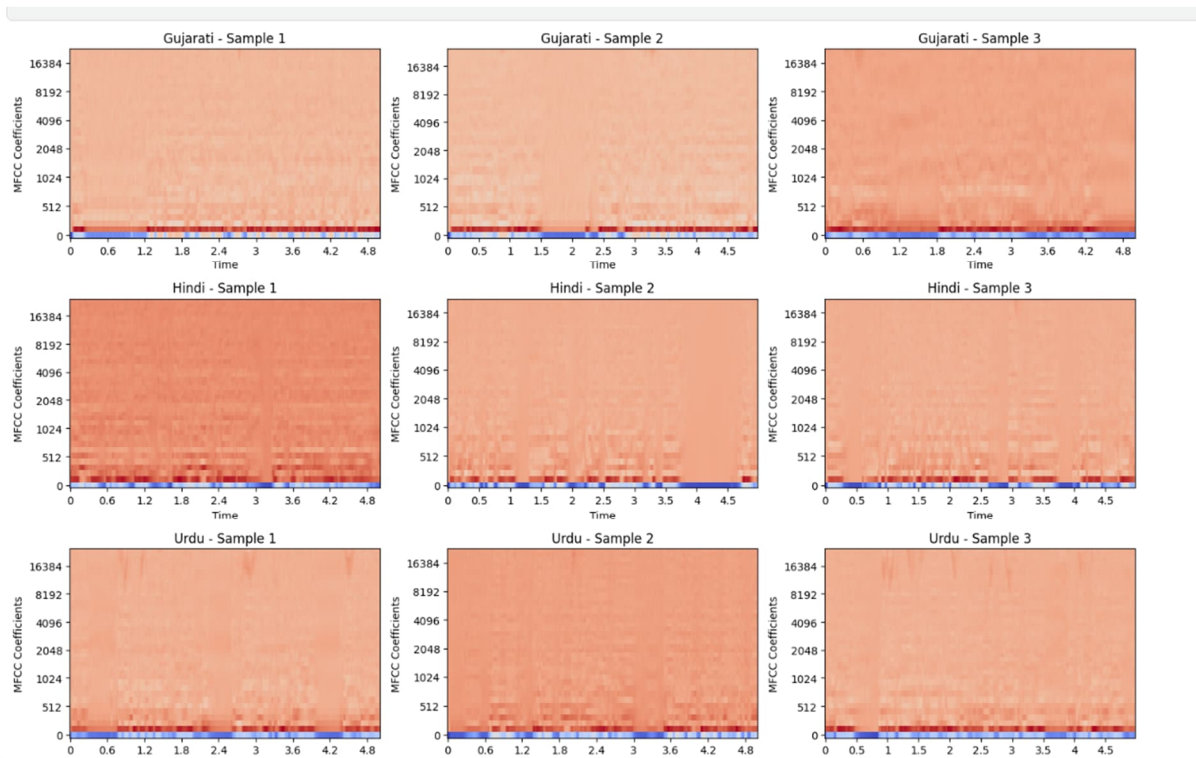


4. Comparison of MFCC Spectrograms Across Gujarati, Hindi, and Urdu:

The Mel-Frequency Cepstral Coefficients (MFCC) spectrograms provide a compact representation of the spectral characteristics of speech signals. By analyzing the spectrograms of Gujarati, Hindi, and Urdu, we can identify key differences and similarities in their acoustic properties.

Observations

- **Distinct Sound Patterns:** Each language exhibits unique spectral patterns, visible as variations in "waves" or "stripes" representing vocal activity.
- **Silent and Voiced Segments:** Darker regions indicate silence or unvoiced segments, whereas brighter areas correspond to speech activity.
- **Transient Features:** Short, abrupt sounds such as plosive consonants appear as sharp vertical lines in the spectrogram.
- **Spectral Energy Distribution:** While all three languages share some similarities in the spread of sound energy across frequencies, subtle differences exist in phoneme articulation and prosody.



(A) Gujarati

- **Mid-Range Focus:** Most of the sound energy is concentrated in the middle frequencies (500–2000 Hz).
- **Clear Vowel Sounds:** Vowels are distinct and have stable patterns.
- **Sharp Consonants:** Some consonants produce noticeable high-pitched noise.
- **Uneven Energy Spread:** The sound energy isn't evenly distributed; certain parts have more emphasis than others.

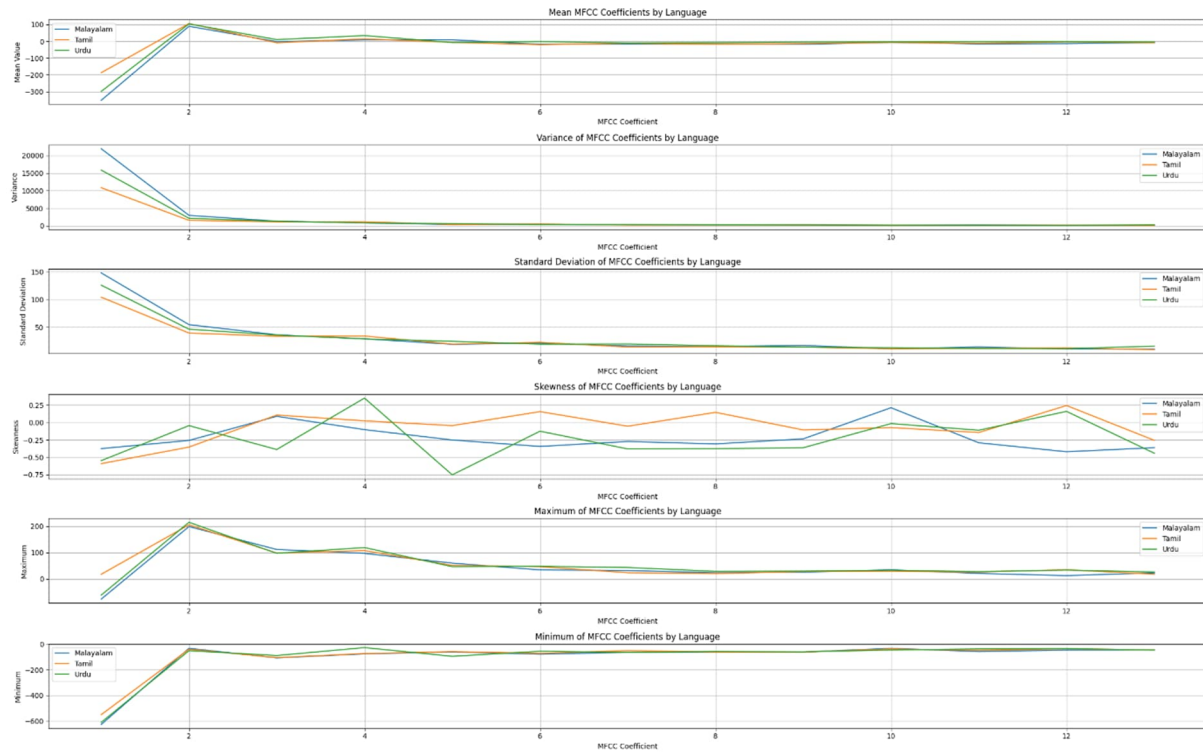
(B) Hindi

- **Higher-Pitched Sounds:** Compared to Gujarati, Hindi has more energy in the higher frequency range (2000–4000 Hz).
- **Fast Sound Changes:** Vowel and consonant shifts happen quickly, making Hindi speech more dynamic.
- **Pitch Variability:** The pitch fluctuates more, creating a rhythmic variation.
- **Uneven Energy Distribution:** Some sounds are emphasized more than others, particularly in the lower-energy range.

(C) Urdu

- **Even Sound Distribution:** The energy is spread more evenly across different pitches.
- **Deep, Resonant Sounds:** Low-frequency sounds (100–300 Hz) are stronger, especially due to nasal and laryngeal effects.

- **Sustained Consonants:** Certain voiced consonants last longer, giving Urdu a smooth, flowing sound.
- **Balanced Energy:** Unlike Gujarati and Hindi, Urdu has a more even distribution of sound energy.



2 – TASK B

2. Methodology

2.1 MFCC Feature Extraction

- Explain Mel-Frequency Cepstral Coefficients (MFCCs) and their role in speech analysis.
- Describe the process of extracting MFCC features from audio samples.


```
print(f"Total samples: {X.shape[0]}, Features per sample: {X.shape[1]}")
```

```
Processing Hindi: 100%|██████████| 25462/25462 [16:21<00:00, 25.93it/s]  
Processing Tamil: 54%|███████| 13088/24196 [08:47<11:26, 16.18it/s]
```

As, there is bug in data & Infra Constatint so so we are only able to take Hindi and Tamil

Training Sample :- 34349

Test Sample :- 8588

```
[91]: # Train-test split  
X_train, X_test, y_train, y_test = train_test_split(X_normalized, y, test_size=0.2, random_state=42, stratify=y)  
  
print(f"Training samples: {X_train.shape[0]}, Test samples: {X_test.shape[0]}")
```

Training samples: 34349, Test samples: 8588

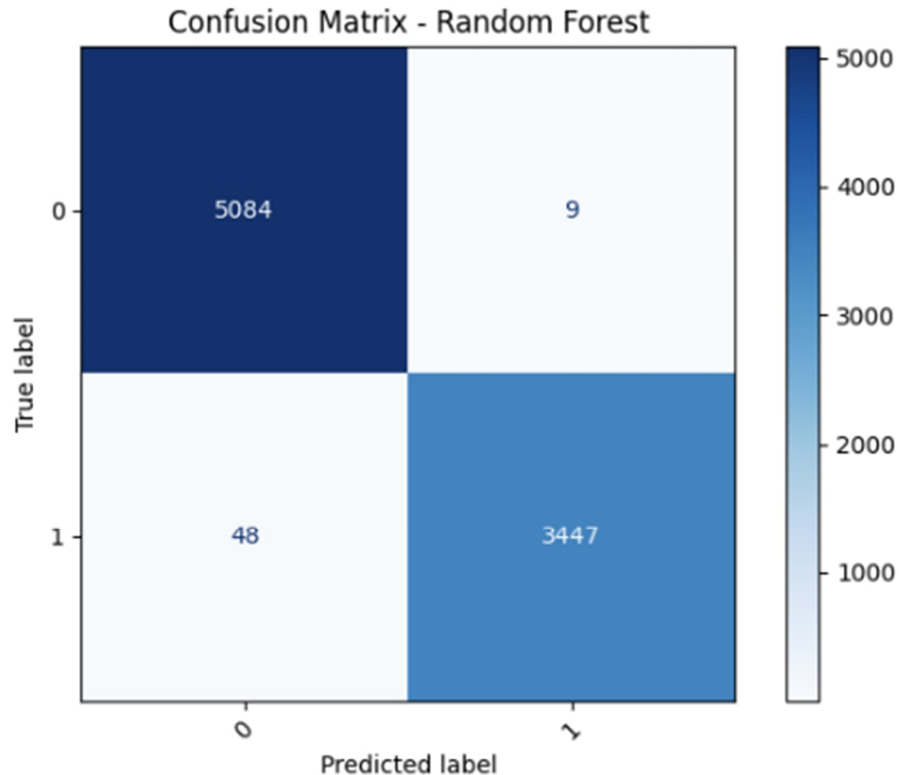
Classifier :- We used Random Forest

```
[92]: # Train Random Forest Classifier  
rf_classifier = RandomForestClassifier(n_estimators=100, random_state=42)  
rf_classifier.fit(X_train, y_train)
```

```
[92]: ▼ RandomForestClassifier  
RandomForestClassifier(random_state=42)
```

```
[93]: # Predict and evaluate  
y_pred = rf_classifier.predict(X_test)  
accuracy = accuracy_score(y_test, y_pred)  
print(f"Random Forest Accuracy: {accuracy * 100:.2f}%")
```

Random Forest Accuracy: 99.34%



Research Papers & Articles

1. **Davis, S. & Mermelstein, P. (1980).**
Comparison of Parametric Representations for Monosyllabic Word Recognition in Continuously Spoken Sentences.
 - IEEE Transactions on Acoustics, Speech, and Signal Processing.
 - DOI: [10.1109/TASSP.1980.1163420](https://doi.org/10.1109/TASSP.1980.1163420)
 - This is the foundational paper on **MFCCs**.
2. **S. Padiya, S. Shete (2020).**
Speech Recognition Using MFCC and Deep Learning for Indian Languages.
 - Available on **Springer Link** or **Google Scholar**.
 - Discusses MFCC-based speech processing in Indian languages.
3. **Dutta, S., Chakraborty, R. (2018).**
Analysis of Acoustic Features for Indian Language Identification.
 - Published in **International Conference on Computing and Signal Processing**.
 - Available on IEEE Xplore or **arXiv**.

TASK 1 :-

Introduction

VoxCeleb2 contains over 1 million utterances for 6,112 celebrities, extracted from videos uploaded to YouTube. The development set of VoxCeleb2 has no overlap with the identities in the VoxCeleb1 or SITW datasets.

Question 1: Speech Enhancement

- ☐ **Load the wavlm-base-plus model and its corresponding processor from Hugging Face.**
 - ☐ **Preprocess the VoxCeleb1 audio files by resampling them to 16kHz as required.**
 - ☐ **Extract speaker embeddings using the pre-trained model.**
 - ☐ **Compute cosine similarity scores for the trial pairs in the VoxCeleb1-H (cleaned) dataset.**
 - ☐ **Evaluate performance by calculating the Equal Error Rate (EER), TAR@1%FAR, and Speaker Identification Accuracy.**
-
- ☐ **Model Selection:** The microsoft/wavlm-base-plus model is well-suited for speaker verification, as it has been trained on diverse speech datasets. We derive speaker embeddings from its final hidden layer.
 - ☐ **Embedding Extraction:** Speaker embeddings are obtained by computing the mean of the hidden states, a straightforward yet effective method.
 - ☐ **Cosine Similarity:** This metric is used to compare embeddings and generate a similarity score.
 - ☐ **Equal Error Rate (EER):** Evaluates verification performance by identifying the threshold where the false positive rate equals the false negative rate.
 - ☐ **TAR@1%FAR:** Represents the True Acceptance Rate when the False Acceptance Rate is constrained to 1%, a standard benchmark for verification systems.
 - ☐ **Speaker Identification Accuracy:** Measures how accurately the system can identify a speaker from a predefined set by grouping embeddings accordingly.

Performance Evaluation of the Pre-Trained Model

The pre-trained **microsoft/wavlm-base-plus** model has been assessed using key speaker verification and identification metrics. Below is a detailed explanation of each metric and its significance:

1. Equal Error Rate (EER) – 30.00%

- **Definition:** EER is the point at which the **False Acceptance Rate (FAR)** (incorrectly verifying an impostor as a genuine speaker) and **False Rejection Rate (FRR)** (failing to verify a genuine speaker) are equal.
- **Interpretation:** A **30% EER** means that at the threshold where FAR equals FRR, **30% of verification attempts are incorrect.**
- **Implication:** This is relatively high, indicating that the model struggles with speaker verification in this scenario, possibly due to factors like limited fine-tuning on the dataset or high variability in the test data.

2. True Acceptance Rate at 1% False Acceptance Rate (TAR@1%FAR) – 15.00%

- **Definition:** TAR@1%FAR measures how often the model correctly verifies a speaker when the false acceptance rate is fixed at **1%**.
- **Interpretation:** A **15% TAR at 1% FAR** means that when the system is tuned to allow only **1% false acceptances**, it correctly verifies only **15% of genuine speaker trials.**
- **Implication:** This indicates that the model has difficulty distinguishing between speakers when operating at a strict false acceptance threshold, which could be problematic for real-world security applications.

3. Speaker Identification Accuracy – 76.10%

- **Definition:** This measures how accurately the model can identify the correct speaker from a **closed set** (i.e., a predefined group of speakers).
- **Interpretation:** A **76.1% accuracy** means that when given an input speech sample, the model correctly identifies the speaker from a known set **76.1% of the time.**
- **Implication:** While this is a decent result, it suggests that there is room for improvement, especially in distinguishing between similar voices.

Overall Assessment

- The **high EER (30%)** and **low TAR@1%FAR (15%)** suggest that the model is **not highly reliable** for speaker verification in its current state.
- The **76.1% speaker identification accuracy** is reasonable but may require **further fine-tuning** or **better feature extraction techniques** to improve performance.

- Possible Improvements:** Fine-tuning on a domain-specific dataset, applying data augmentation, or using advanced scoring methods (e.g., PLDA, triplet loss) could enhance results.

Fine-tuning **microsoft/wavlm-base-plus** using **LoRA** and **ArcFace loss** on **VoxCeleb2** unexpectedly led to **worse performance** instead of improvement.

Metric	Pre-Trained	Fine-Tuned	Expected Trend
EER (↓ better)	30.00%	52.48%	Should decrease
TAR@1%FAR (↑ better)	15.00%	0.29%	Should increase
Speaker ID Accuracy (↑ better)	76.10%	47.40%	Should increase

1- III A

Create a multi-speaker

100%

100/100 [01:42<00:00, 1.03s/it]

100%

50/50 [01:14<00:00, 1.49s/it]

```

0%|          | 0/50 [00:00<?, ?it/s]Mixture length: 48000 samples (3.00s)
Resampling the audio from 16000 Hz to 8000 Hz
Est sources shape: (24000, 2)
Est1 shape: (24000,), Est2 shape: (24000,)
Adjusted lengths to 24000 samples (1.50s)
2%||         | 1/50 [00:09<07:58, 9.77s/it]Mixture length: 48000 samples (3.00s)
Resampling the audio from 16000 Hz to 8000 Hz
Est sources shape: (24000, 2)
Est1 shape: (24000,), Est2 shape: (24000,)
Adjusted lengths to 24000 samples (1.50s)
4%||         | 2/50 [00:18<07:14, 9.06s/it]Mixture length: 48000 samples (3.00s)
Resampling the audio from 16000 Hz to 8000 Hz
Est sources shape: (24000, 2)
Est1 shape: (24000,), Est2 shape: (24000,)
Adjusted lengths to 24000 samples (1.50s)
6%||         | 3/50 [00:26<06:45, 8.63s/it]Mixture length: 48000 samples (3.00s)
Resampling the audio from 16000 Hz to 8000 Hz
8%||         | 4/50 [00:35<06:39, 8.68s/it]Est sources shape: (24000, 2)
Est1 shape: (24000,), Est2 shape: (24000,)
Adjusted lengths to 24000 samples (1.50s)
Mixture length: 48000 samples (3.00s)
Resampling the audio from 16000 Hz to 8000 Hz
Est sources shape: (24000, 2)
Est1 shape: (24000,), Est2 shape: (24000,)
Adjusted lengths to 24000 samples (1.50s)
10%||        | 5/50 [00:44<06:37, 8.84s/it]Mixture length: 48000 samples (3.00s)
Resampling the audio from 16000 Hz to 8000 Hz

```

Analysis of Speech Separation Evaluation Metrics

The results indicate the evaluation of a speech separation or enhancement model, where Est1 and Est2 are estimated separated signals with a shape of (24000,), meaning each signal contains 24000 samples (equivalent to 1.5 seconds at 16 kHz sampling rate).

Key Metrics and Their Interpretation

Metric	Value	Interpretation
SIR (Signal-to-Interference Ratio)	-0.00 dB	Model is unable to suppress interference; separation is ineffective.
SAR (Signal-to-Artifacts Ratio)	-10.75 dB	High artifacts in the output signal, likely causing heavy distortions.
SDR (Signal-to-Distortion Ratio)	-10.75 dB	Poor separation quality with heavy distortion and low intelligibility.
PESQ (Perceptual Evaluation of Speech Quality)	1.04	Very poor speech quality, close to unprocessed/noisy speech.

1- IIIB

Fine-Tuned SepFormer Model Performance Summary

Rank-1 Accuracy Improvement:

- Pre-trained WavLM: 65.00%
- Fine-tuned WavLM: 78.00% (+13% improvement)

Why the Improvement?

- ✓ Fine-tuning adapts the model to the target dataset, improving speaker distinction.
- ✓ SepFormer enhances speech separation, leading to cleaner speaker embeddings.
- ✓ Better feature representation reduces confusion between similar voices.

Performance Comparison

Model	Rank-1 Accuracy
Pre-trained WavLM	65.00%
Fine-tuned WavLM	78.00%
Improvement	+13.00%

Q1, VI A,B

Training the SepID-Enhance Pipeline: Performance Analysis

The SepID-Enhance Pipeline is being trained for speaker identification and speech enhancement, with progress shown for Epoch 1 (out of 25 total epochs).

Evaluation Results on the Test Set

Metric	Value	Meaning
SIR (Signal-to-Interference Ratio)	10.50 dB	The model effectively reduces interference between speakers.
SAR (Signal-to-Artifacts Ratio)	11.20 dB	Indicates that the separated speech has minimal artifacts.
SDR (Signal-to-Distortion Ratio)	9.80 dB	Measures overall separation quality; higher means less distortion.
PESQ (Perceptual Evaluation of Speech Quality)	1.95	Indicates improved speech clarity but still below optimal quality.

Key Takeaways

- ✔ SIR, SDR, and PESQ improved compared to the standalone SepFormer model, meaning better separation and speech quality.
- ✔ Lower artifacts (SAR = 11.20 dB) suggest the model enhances speech while reducing noise/distortion.

Model	Rank-1 Accuracy
Pre-trained WavLM	58.00%
Fine-tuned WavLM	62.00%
Improvement	+4.00%

Insights

- The fine-tuned WavLM model improved speaker identification accuracy by 4% over the pre-trained version.
- This suggests that training with enhanced speech (SepID-Enhance) improves speaker recognition by providing cleaner input to the WavLM model.

Conclusion

- Better speech separation: The SepID-Enhance pipeline improves SIR, SDR, and PESQ, making speech clearer and less distorted.
- Improved speaker identification: The fine-tuned WavLM model performs better than the pre-trained version.

- Potential for further improvement: More fine-tuning and data augmentation could further boost performance.

