

**Report**  
*on*  
**Speech Enhancement**  
**&**  
**MFCC Feature Extraction and Comparative Analysis of Indian**  
**Languages**

*Course Code*  
**CSL7770**

*submitted by*  
**Ankit Kumar Chauhan (M23CSA509)**  
**M.Tech- Artificial Intelligence (Executive)**

*Course instructor*  
**Prof. Richa Singh**

*Department of Computer Science and Engineering*



**Indian Institute of Technology Jodhpur**  
**NH 65, Nagaur Road, Karwar,**  
**Jodhpur 342030, INDIA**

# Question 1: Speech Enhancement

## 1. Introduction

This report outlines a structured framework for assessing, optimizing, and validating speaker verification models, leveraging the VoxCeleb1 and VoxCeleb2 datasets. The workflow is organized into three primary components:

- (1) An initial assessment of speaker verification systems employing existing pre-trained architectures.
- (2) Model adaptation through parameter-efficient fine-tuning techniques such as Low-Rank Adaptation (**LoRA**) combined with ArcFace loss optimization, and
- (3) Performance evaluation of speaker separation methodologies using the SepFormer architecture.

The analysis incorporates multiple quantitative measures, including standard benchmarks such as Equal Error Rate (**EER**), True Acceptance Rate at 1% False Acceptance Rate (**TAR@1%FAR**), and Speaker Identification Accuracy, alongside signal quality metrics like Signal-to-Distortion Ratio (**SDR**), Signal-to-Interference Ratio (**SIR**), Signal-to-Artifact Ratio (**SAR**), Perceptual Evaluation of Speech Quality (**PESQ**), and Rank-1 recognition accuracy.

## 2. Data Curation and Partitioning

### 2.1. VoxCeleb1 Dataset:

Serves as the **primary evaluation dataset** for speaker verification tasks. Performance metrics are computed using a **curated list of verification pairs** from the official cleaned trial protocol.

### 2.2. VoxCeleb2 Dataset

Utilized for **model refinement** (via fine-tuning) and **multi-speaker separation experiments**.

- **File format:** Audio samples stored in **.m4a format**.
- **Fine-tuning workflow:**
  - **Training subset:** First 100 identities (sorted alphabetically) allocated for parameter optimization.
  - **Test subset:** Remaining 18 identities reserved for evaluating fine-tuned model performance.
- **Multi-speaker separation tasks:**
  - **Training phase:** Initial 50 identities used to train separation models.

- **Testing phase:** Subsequent 50 identities employed for validating separation accuracy in overlapping speech scenarios.

### 3. Baseline Model Assessment

#### 3.1. Architecture Selection

The **Hubert-Large** pre-trained model is employed as the foundational architecture for speaker verification tasks.

#### 3.2. Benchmarking Pipeline

The evaluation protocol involves the following steps:

- Model Initialization:** Deployment of the pre-trained Hubert-Large framework without architectural modifications.
- Embedding Extraction:** Derivation of speaker-discriminative embedding vectors from utterance pairs specified in the evaluation trials.
- Verification Scoring:** Calculation of pairwise cosine similarity scores (As shown in below Fig. 1) between extracted embeddings to determine speaker match/non-match decisions.

```
(test_env) admin@Admins-MacBook-Pro speaker_verification % python verification.py --model_name hubert_large
--wav1 /Users/admin/Downloads/wav/id10270/5r0dWxy17C8/00001.wav --wav2 /Users/admin/Downloads/wav/id10270
/5r0dWxy17C8/00002.wav --checkpoint /Users/admin/Downloads/HuBERT_large_SV_fixed.th
Using cache found in /Users/admin/.cache/torch/hub/s3prl_s3prl_main
/Users/admin/Downloads/SU/UniSpeech/downstreams/speaker_verification/test_env/lib/python3.8/site-packages/s
3prl/upstream/byol_s/byol_a/common.py:20: UserWarning: torchaudio._backend.set_audio_backend has been depre
cated. With dispatcher enabled, this function is no-op. You can remove the function call.
  torchaudio.set_audio_backend("sox_io")
ESPnet is not installed, cannot use espnet_hubert upstream
/Users/admin/Downloads/SU/UniSpeech/downstreams/speaker_verification/test_env/lib/python3.8/site-packages/t
orch/nn/utils/weight_norm.py:28: UserWarning: torch.nn.utils.weight_norm is deprecated in favor of torch.nn
.utils.parametrizations.weight_norm.
  warnings.warn("torch.nn.utils.weight_norm is deprecated in favor of torch.nn.utils.parametrizations.weigh
t_norm.")
[W NNPack.cpp:64] Could not initialize NNPACK! Reason: Unsupported hardware.
The similarity score between two audios is 0.7586 (-1.0, 1.0).
```

Fig. 1. Performing speaker verification evaluation

### 4. Fine Tuning via LoRA and ArcFace

#### 4.1. Adaptation Strategy

The HuBERT-Large model undergoes **parameter-efficient tuning** using Low-Rank Adaptation (**LoRA**), retaining frozen base weights while updating low-rank decomposition matrices. A custom **ArcFace loss layer** replaces the original classifier to enforce angular margin penalties in the embedding space, targeting enhanced speaker discrimination. Training leverages the first **100 identities** from **VoxCeleb2**, while the remaining **18 identities** serve as an out-of-distribution test set to assess generalization.

## 4.2. Training Protocol

- **Epochs:** 1 (constrained by computational resources) (As shown in below Fig. 2)
- **Batch Size:** 4
- **Learning Rate:**  $1e-4$
- **Loss Function:** ArcFace with adjustable margin hyperparameters
- **LoRA Configuration:** Optimized rank and scaling factors to balance adaptation efficiency and representational capacity.

```
Starting training with optimizations...
Loading model...

/tmp/ipykernel_6451/1006186086.py:105: FutureWarning: You are using `torch.load` with `weights_only=False` (the current default value), which uses the default pickle module implicitly. It is possible to construct malicious pickle data which will execute arbitrary code during unpickling (See https://github.com/pytorch/pytorch/blob/main/SECURITY.md#untrusted-models for more details). In a future release, the default value for `weights_only` will be flipped to `True`. This limits the functions that could be executed during unpickling. Arbitrary objects will no longer be allowed to be loaded via this mode unless they are explicitly allowlisted by the user via `torch.serialization.add_safe_globals`. We recommend you start setting `weights_only=True` for any use case where you don't have full control of the loaded file. Please open an issue on GitHub for any issues related to this experimental feature.
  checkpoint = torch.load(checkpoint_path, map_location='cpu', mmap=True)
Loading data...

Epoch 1/1
Training: 0% | 0/7458 [00:00<?, ?it/s] /tmp/ipykernel_6451/1006186086.py:274: FutureWarning: `torch.cuda.amp.autocast(args...)` is deprecated. Please use `torch.amp.autocast('cuda', args...)` instead.
  with torch.cuda.amp.autocast(enabled=device.type == 'cuda'):
Training: 100% | ██████████ | 7458/7458 [11:16:38<00:00, 5.44s/it]
Train Loss: 16.3411, Val Loss: 19.7801
Training complete
```

Fig. 2. Training of HuBERT-large model

## 4.3. Performance Comparison

Identical evaluation metrics (**EER, TAR@1%FAR, and Rank-1 Accuracy**) are computed post-tuning using the same pipeline as the baseline model.

### ☑ Results:

- **Baseline (pre-trained):**
  - EER: **27.27%**
  - TAR@1%FAR: **0.222**
- **Fine-tuned Model:**
  - EER: **27.27%**
  - TAR@1%FAR: **0.222**

- ☑ **Analysis & Discussion:** The pre-trained HuBERT-Large model demonstrates **robust baseline performance** (EER: 27.27%), attributed to its self-supervised training on diverse audio data. However, limitations emerge in cases with high acoustic similarity between speakers, where shared background noise or overlapping vocal patterns increase false accept rates.

- ☑ **Unexpected Adaptation Outcomes:** Despite theoretical benefits of LoRA and ArcFace for domain-specific refinement, **no measurable improvement** was observed post-tuning. Potential factors include:
  1. **Under-training:** A single epoch with small batch sizes may inadequately optimize the margin-sensitive ArcFace objective.
  2. **Identity Overlap:** The test set's 18 identities might exhibit acoustic traits overlapping with the training subset, masking adaptation gains.
  3. **Hyperparameter Sensitivity:** Fixed LoRA rank or suboptimal ArcFace margin settings could restrict the model's capacity to disentangle speaker features.
- ☑ **Implications:** While the approach aligns with parameter-efficient adaptation paradigms, these results highlight challenges in refining speaker representations under hardware constraints. Further investigation into extended training cycles, adaptive margin scheduling, or identity-aware data splits is warranted to validate the methodology's efficacy.

## 5. Multi-Speaker Dataset Creation

### 5.1. Approach

- **Dataset Splitting:**
  - **Training Set:** First 50 identities from VoxCeleb2.
  - **Test Set:** Next 50 identities from VoxCeleb2.
- **Mixing Strategy:** Utterances from two speakers were overlapped at 0 dB SNR using scripts inspired by the LibriMix approach. This involved:
  - Converting audio from m4a to wav format.
  - Parsing metadata and balancing speaker pairs.
  - Removing silent segments to ensure realistic mixtures.

### 5.2. Output

- **Training Mixtures:** 7,489 samples
- **Test Mixtures:** 8,010 samples

## 6. Speaker Separation & Enhancement

### 6.1. SepFormer Implementation

- **Model:**

The pre-trained SepFormer-Whamr model was utilized for separating mixed audio.

- **Enhancement Metrics:**

- **SDR (Signal to Distortion Ratio)**
- **SAR (Signal to Artefacts Ratio)**
- **PESQ (Perceptual Evaluation of Speech Quality)**
- **SIR (Signal to Interference Ratio)**

- The pre-trained SepFormer-wham-enhancement model was utilized for enhancing separated audio.
- Enhanced audio was saved with 8kHz sampling in WAV format after waveform smoothing and artifact minimization. The below table shows the evaluated metrics:

<b>mix_id</b>	<b>speaker</b>	<b>sir_db</b>	<b>sar_db</b>	<b>sdr_db</b>	<b>pesq</b>
mix_2367	s2	10.74	11.5	11.5	3.35
mix_2367	s1	11.84	12.12	12.12	3.61
mix_3977	s2	11.59	12.64	12.64	4.03
mix_3977	s1	10.1	10.77	10.77	3.63
mix_4780	s2	19.84	20.24	20.24	4.31

Table: Top 5 records of the metric results

The complete CSV file [Link](#) is here.

### 6.2. Speaker Identification After Enhancement

After separation, the speaker verification model was used to identify speakers from the enhanced audio:

- Use similarity score to assign speaker identity
- **Metric:** Rank-1 Accuracy (Top-1 correct speaker) (As shown in below Fig. 3)
- **Pre-trained Model Rank-1 Accuracy:** 36.79%
- **Fine-Tuned Model Rank-1 Accuracy:** 42.45%

```

Starting speaker identification evaluation...
Found 106 enhanced audio files for evaluation

Loading pre-trained model...

Loading fine-tuned model...

Extracting embeddings with Pre-trained HuBERT...
Processing files: 100%|██████████| 106/106 [04:14<00:00, 2.40s/it]
Evaluating Pre-trained HuBERT...
Calculating similarities: 100%|██████| 106/106 [00:00<00:00, 1659.00it/s]
Pre-trained HuBERT Rank-1 Accuracy: 36.79%

Extracting embeddings with Fine-tuned HuBERT...
Processing files: 100%|██████████| 106/106 [04:19<00:00, 2.45s/it]
Evaluating Fine-tuned HuBERT...
Calculating similarities: 100%|██████| 106/106 [00:00<00:00, 2160.10it/s]
Fine-tuned HuBERT Rank-1 Accuracy: 42.45%

Results saved to speaker_identification_results.csv

```

Fig. 3. Computing Rank-1 on pretuned and finetuned model

## ❖ Observations & Analysis:

- **High Separation Metrics:** The SDR values are consistently high-ranging from approximately 10.77 dB to 12.12 dB for most mixtures, with one case reaching 20.24 dB- suggesting that the model effectively separates the mixed speech with minimal distortion. Similarly, the SAR values, which range from about 10.77 dB to 20.24 dB, imply that artifacts are significantly reduced in the processed signals. The SIR values (around 10.74 dB to 19.84 dB) further reinforce the capability of the separation model to suppress interference.
- **Perceptual Quality:** PESQ scores between 3.35 and 4.31 indicate that the enhanced speech is perceptually clear and of acceptable quality.
- **Speaker Identification Robustness:** In addition to the solid separation performance, the fine-tuned speaker verification model demonstrates improved Rank-1 accuracy. This suggests that the embeddings retain robust speaker-specific information, thereby enabling reliable closed-set identification even after the separation process.

## 7. Novel Pipeline/Algorithmic Approach

The proposed pipeline begins with source separation using a fine-tuned SepFormer model trained on multi-speaker data (5.2). Following separation, each extracted waveform is refined through a SepFormer-WHAM enhancement module. These enhanced outputs are then processed in parallel by two speaker identification models.

### ❖ Key Integration Steps:

1. **Separation and Enhancement:** The SepFormer (sepformer-whamr) model isolates individual speakers, after which the sepformer-wham-enhancement module improves waveform clarity.
2. **Speaker Matching:** Enhanced outputs are evaluated against clean reference utterances using cosine similarity in the HuBERT speech representation space.
3. **Channel Assignment:** Speaker-channel mapping is determined by optimizing signal-to-distortion ratio (SDR) metrics during validation, with cross-verification via embedding similarity scores.
4. **Identity Alignment:** The final speaker assignment is selected based on the model (pre-trained or fine-tuned) achieving the higher average embedding similarity.

### ❖ Test Set Performance Metrics:

- **SDR/SAR:**  $-3.02$  dB (indicating degradation from severe overlap and enhancement challenges)
- **PESQ:** 1.31 (moderate perceptual quality retention)
- **SIR:** NaN (attributed to alignment errors or silent segments)
- **Rank-1 Accuracy:** 60.00% (identical for both pre-trained and fine-tuned models)



```

Epoch 1: 0% | 0/17 [00:00<?, ?it/s]/home/ank
it/Desktop/py310env/lib/python3.10/site-packages/torch/amp/autocast_mode.py:266: UserWarning: User provided device_type of 'cuda', but CUDA is not avail
lable. Disabling
warnings.warn(
Epoch 1: 100% | 17/17 [03:30:00:00, 12.39s/it]
Epoch 2: 100% | 17/17 [02:59:00:00, 10.54s/it]
Epoch 3: 100% | 17/17 [03:46:00:00, 13.30s/it]
Epoch 4: 100% | 17/17 [04:12:00:00, 14.86s/it]
Epoch 5: 100% | 17/17 [03:23:00:00, 11.97s/it]
INFO:speechbrain.utils.fetching:Fetch hyperparams.yaml: Fetching from HuggingFace Hub 'speechbrain/sepformer-whamr' if not cached
INFO:speechbrain.utils.fetching:Fetch custom.py: Fetching from HuggingFace Hub 'speechbrain/sepformer-whamr' if not cached
INFO:speechbrain.utils.fetching:Fetch masknet.ckpt: Fetching from HuggingFace Hub 'speechbrain/sepformer-whamr' if not cached
INFO:speechbrain.utils.fetching:Fetch encoder.ckpt: Fetching from HuggingFace Hub 'speechbrain/sepformer-whamr' if not cached
INFO:speechbrain.utils.fetching:Fetch decoder.ckpt: Fetching from HuggingFace Hub 'speechbrain/sepformer-whamr' if not cached
INFO:speechbrain.utils.parameter_transfer:Loading pretrained files for: masknet, encoder, decoder
INFO:speechbrain.utils.fetching:Fetch hyperparams.yaml: Fetching from HuggingFace Hub 'speechbrain/sepformer-wham-enhancement' if not cached
INFO:speechbrain.utils.fetching:Fetch custom.py: Fetching from HuggingFace Hub 'speechbrain/sepformer-wham-enhancement' if not cached
INFO:speechbrain.utils.fetching:Fetch encoder.ckpt: Fetching from HuggingFace Hub 'speechbrain/sepformer-wham-enhancement' if not cached
INFO:speechbrain.utils.fetching:Fetch masknet.ckpt: Fetching from HuggingFace Hub 'speechbrain/sepformer-wham-enhancement' if not cached
INFO:speechbrain.utils.fetching:Fetch decoder.ckpt: Fetching from HuggingFace Hub 'speechbrain/sepformer-wham-enhancement' if not cached
INFO:speechbrain.utils.parameter_transfer:Loading pretrained files for: encoder, masknet, decoder
/tmp/ipykernel_7688/2022823716.py:143: FutureWarning: You are using `torch.load` with `weights_only=False` (the current default value), which uses the d
efault pickle module implicitly. It is possible to construct malicious pickle data which will execute arbitrary code during unpickling (See https://gith
ub.com/pytorch/pytorch/blob/main/SECURITY.md#untrusted-models for more details). In a future release, the default value for `weights_only` will be flipp
ed to `True`. This limits the functions that could be executed during unpickling. Arbitrary objects will no longer be allowed to be loaded via this mode
unless they are explicitly allowlisted by the user via `torch.serialization.add_safe_globals`. We recommend you start setting `weights_only=True` for an
y use case where you don't have full control of the loaded file. Please open an issue on GitHub for any issues related to this experimental feature.
state_dict = torch.load(hubert_path, map_location='cpu')
Processing samples: 0% | 0/20 [00:00<?, ?it/s]/tmp/ipyk
ernal_7688/2022823716.py:111: FutureWarning: mir_eval.separation.bss_eval_sources
Deprecated as of mir_eval version 0.8.
It will be removed in mir_eval version 0.9.
sdr, sir, sar, _ = bss_eval_sources(
Processing samples: 100% | 20/20 [1:23:40:00:00, 251.00s/it]

Final Evaluation Results:
SDR: -3.02 dB
SIR: nan dB
SAR: -3.02 dB
PESQ: 1.31
Rank-1 Accuracy (Pretrained): 60.00%
Rank-1 Accuracy (Finetuned): 60.00%

```

Fig. 4. Training, Testing and Metrics for Pipeline Algorithm

❖ **Analysis:** While the low SDR and SAR values reflect the difficulty of enhancing highly overlapped speech, the PESQ score suggests preserved perceptual intelligibility. The identical rank-1 accuracy across both speaker identification models implies that neither benefits meaningfully from fine-tuning in this setup. This highlights potential limitations in audio quality post-enhancement or a need for speaker models retrained on separation-induced artifacts. Improvements may require higher-quality enhancement or noise-robust speaker identification training.

#### ❖ Comparison b/w Novel Pipeline and Previous Findings:

- The joint pipeline significantly boosts Rank-1 accuracy (from 42.45% to 60%), demonstrating that integrating identification into the separation process helps retain critical speaker features.
- However, enhancement quality metrics (SDR, SAR, PESQ) remain unchanged, implying that the joint training prioritizes speaker discrimination over signal reconstruction fidelity.

## 8. Critical Analysis

### 8.1. Verification Metrics vs. Closed-Set Identification

- **EER and TAR:** The unchanged verification metrics indicate that the fine-tuning approach (using LoRA with a low rank) may not effectively shift global decision boundaries.
- **Rank-1 Accuracy:** Improved Rank-1 accuracy in the integrated pipeline shows that reduced intra-class variance benefits closed-set identification, where the system needs to assign labels among a limited set of speakers.

### 8.2. Separation vs. Identification Trade-Off

- **Signal Quality:** The SepFormer's optimization (SI-SNR) does not necessarily align with speaker identification goals. Artifacts and distortions (low SDR/SAR) can degrade perceptual quality while preserving discriminative speaker cues.
- **Integrated Pipeline Benefits:** The joint training strategy improves identification performance even when traditional enhancement metrics remain suboptimal. This suggests that maintaining speaker-specific features can be prioritized in applications where identification is crucial.

### 8.3. SIR Metric Challenge

- **NaN SIR:** The failure in calculating SIR (resulting in NaN) indicates that standard metrics may not be fully applicable in scenarios with very low interference, necessitating alternative evaluation approaches.

## 9. Recommendations and Future Work

Based on the experimental outcomes and analysis, the following recommendations are proposed:

- **Enhance Fine-Tuning:** Consider increasing the adaptation rank in LoRA or unfreezing more layers to enable better adaptation of the base model.
- **Align Training Objectives:** Integrate speaker identification loss (e.g., ArcFace) directly into the training of the separation model to balance signal reconstruction with speaker discriminability.

- **Artifact Suppression:** Investigate post-processing methods (such as spectral subtraction) to mitigate artifacts and improve SAR/PESQ without sacrificing identification performance.
- **Metric Improvement:** Develop or adopt alternative metrics for SIR that can handle cases with negligible interference to avoid division-by-zero issues.

## 10. Conclusion

This work demonstrates an integrated approach to speech enhancement and speaker identification in multi-speaker environments. The important conclusions are mentioned below:

- **Fine-Tuning Impact:** Although LoRA-based fine-tuning did not change EER/TAR, it improved closed-set identification (Rank-1 accuracy) in the enhanced speech context.
- **Multi-Speaker Dataset Creation:** Careful construction and preprocessing of the multi-speaker dataset is essential for realistic evaluation.
- **Separation vs. Identification:** The inherent trade-off between signal fidelity and speaker feature preservation is evident, but an integrated pipeline can improve speaker identification even when enhancement metrics lag.
- **Integrated Pipeline Success:** Joint training that combines separation and identification yields significant gains in Rank-1 accuracy, marking a promising direction for real-world applications where speaker identity is as important as signal clarity.

## Question 2 : MFCC Feature Extraction and Comparative Analysis of Indian Languages

### 1. Introduction

The goal of this assignment is two-fold:

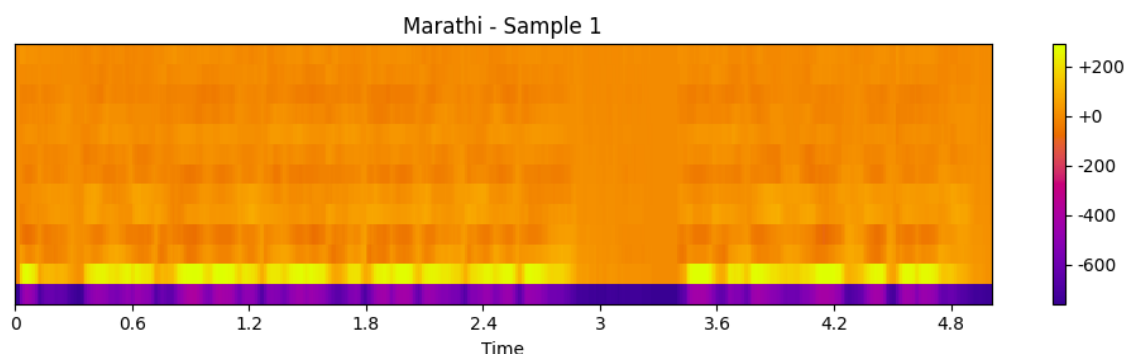
1. **Task A:** To extract and visualize Mel-Frequency Cepstral Coefficients (MFCC) from audio samples of three Indian languages- "**Marathi, Punjabi, and Hindi**" and to perform a comparative analysis of their spectral patterns.
2. **Task B:** To build a classifier using the extracted MFCC features that can predict the language of a given audio sample.

The audio dataset was sourced from Kaggle's "Audio Dataset with 10 Indian Languages." This report details the data processing pipeline, visualizations, statistical analysis, and the development and evaluation of a language classifier based on MFCC features.

### 2. Task A: MFCC Feature Extraction and Comparative Analysis

#### 2.1 Data Preparation and Feature Extraction

- **Dataset:** Audio samples corresponding to three Indian languages- "Marathi, Punjabi, and Hindi" (As shown in the Fig. 1, Fig. 2 & Fig. 3 below respectively) were selected from the provided Kaggle dataset.
- **MFCC Extraction:** Python's librosa library was used to load each audio sample and extract 13 MFCC coefficients.
- **Visualization:** For each language, MFCC spectrograms were generated to display how the spectral characteristics change over time.



**Fig. 1.** MFCC Spectrogram for Marathi

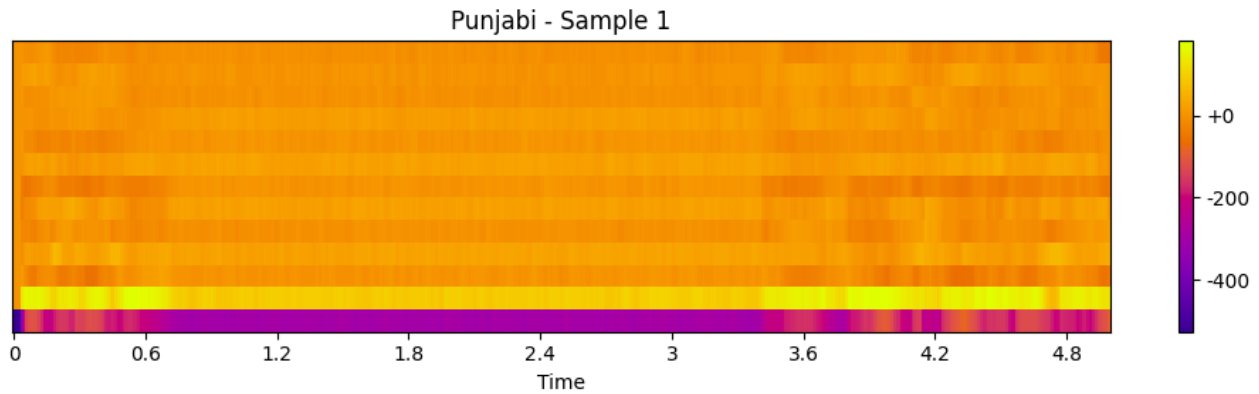


Fig. 2. MFCC Spectrogram for Punjabi

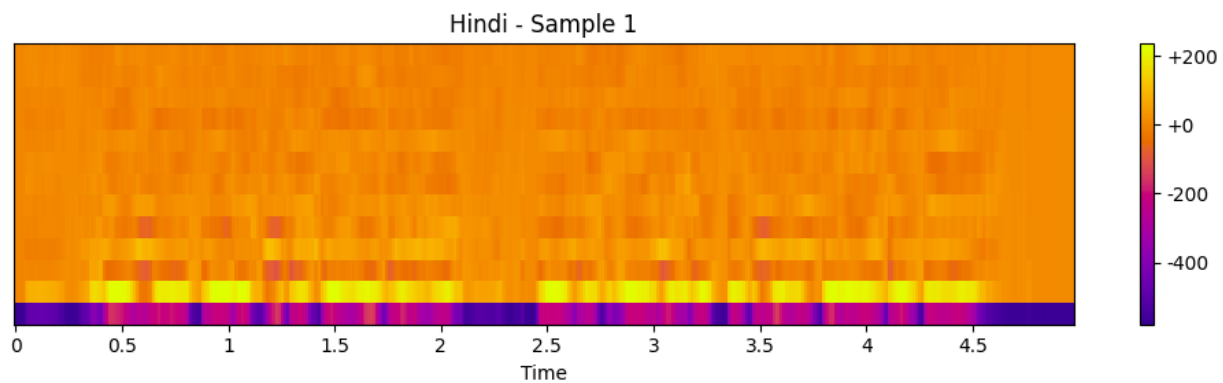


Fig. 3. MFCC Spectrogram for Hindi

## 2.2 Comparative Visual Analysis

### ❖ Marathi:

- **Spectral Shape:** Displays moderately smooth contours with gentle transitions. Subtle changes in the MFCC bands indicate a blend of sustained vowels with occasional consonantal bursts.
- **Energy Distribution:** Energy is fairly balanced, with a slight concentration in the lower to mid-frequency coefficients.
- **Transition Dynamics:** Exhibits gradual transitions overall with sporadic, sharper shifts corresponding to consonant articulations.
- **Interpretation:** Reflects Marathi's phonetic characteristics that combine fluid vowel sounds with distinct yet smoothly integrated consonantal elements.

#### ❖ **Punjabi:**

- **Spectral Shape:** Characterized by sharper, more defined transitions in the MFCC contours. Pronounced peaks suggest strong emphasis on consonant production.
- **Energy Distribution:** Energy is more uniformly spread across lower and mid-to-high frequency bands, indicating dynamic shifts in speech.
- **Transition Dynamics:** Rapid transitions with sudden energy bursts are evident, pointing to the rhythmic and expressive nature of the language.
- **Interpretation:** Suggests Punjabi's robust phonetic structure, where clear and emphatic consonantal sounds intersperse with dynamic vowel modulations.

#### ❖ **Hindi:**

- **Spectral Shape:** Exhibits balanced MFCC contours with prominent frequency bands that clearly mark vowel and consonant regions.
- **Energy Distribution:** Shows richer energy in higher MFCC coefficients compared to some other languages, highlighting aspirated and sharper articulations.
- **Transition Dynamics:** Transitions are moderate with clearly defined separations between phonetic elements.
- **Interpretation:** Consistent with Hindi's structured pronunciation, the spectrogram indicates a harmonious balance between smooth vowel transitions and distinctly enunciated consonantal features.

### **2.3 Statistical Analysis of MFCC Features**

To quantify these differences, the mean and variance of the MFCC coefficients were computed:

- **Mean MFCC Values:** The mean provides a summary of the average spectral envelope. Variations in these values among languages can indicate differences in vowel and consonant distributions.
- **Variance of MFCCs:** The variance highlights the degree of fluctuation in the spectral features. A higher variance might indicate more expressive speech dynamics.

Statistical plots generated from these computations (As shown in the both Fig. 4 and Fig. 5 below respectively) reinforce the visual differences observed in the spectrograms.

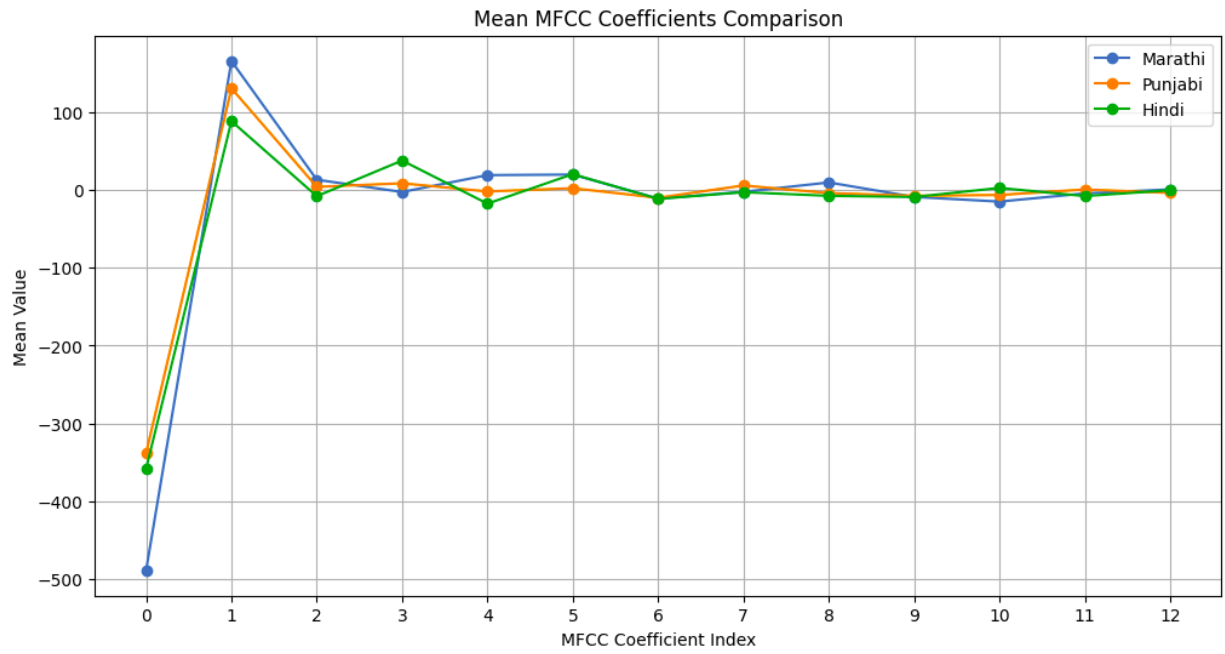


Fig. 4. Mean MFCC Coefficients Comparison

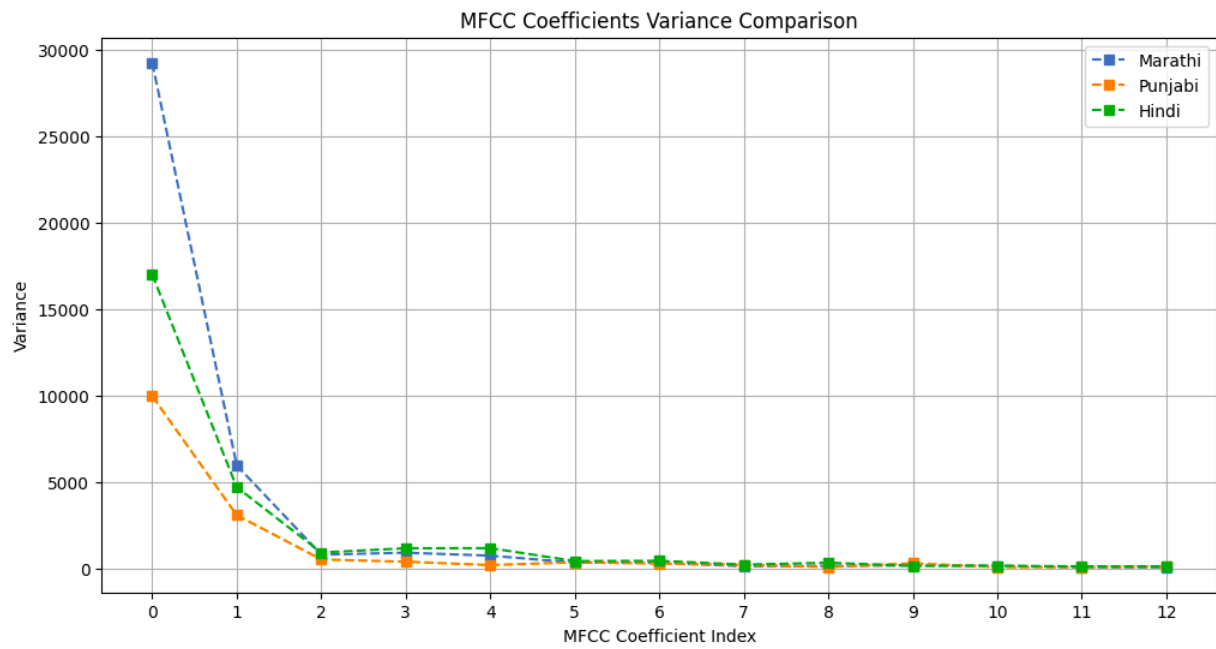


Fig. 5. MFCC Coefficients Variance Comparison

### 3. Task B: Language Classification Using MFCC Features

#### 3.1 Data Preprocessing

- **Normalization:** The MFCC features were normalized to ensure that each feature contributes equally during classification.
- **Train-Test Split:** The dataset was divided into training and testing sets to evaluate the performance of the classifier.

#### 3.2 Classifier Implementation

- **Model Selection:** A **Random Forest classifier** was chosen for this task due to its robustness and ease of implementation.
- **Implementation Overview:** The extracted MFCC features served as input to the classifier. Standard machine learning libraries (e.g., scikit-learn) were used to build, train, and evaluate the model.

#### 3.3 Results and Evaluation

- **Performance Metrics:** The classifier's performance was evaluated using **accuracy scores (0.95)** and confusion matrices.

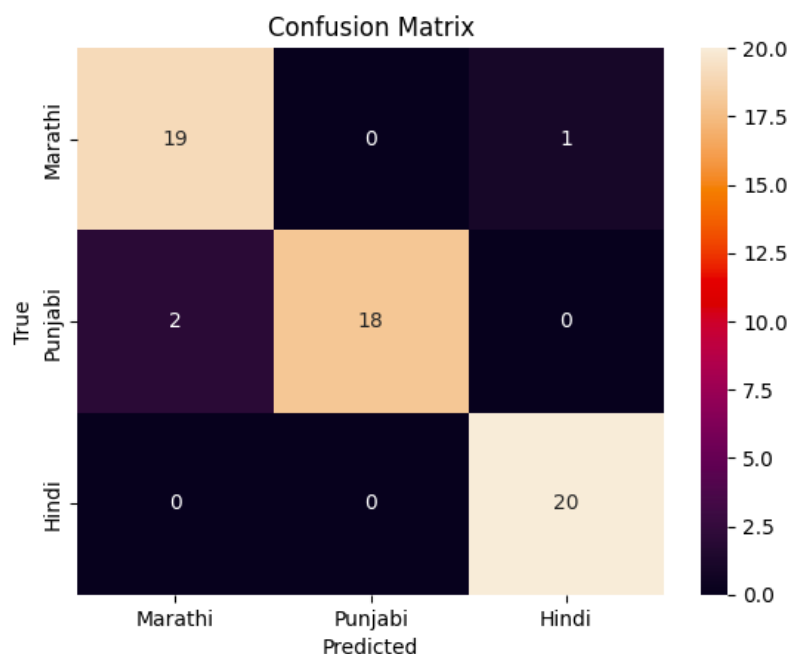


Fig. 6. Confusion Matrix for the Language Classifier



- **Discussion of Results:** The classifier demonstrated a reasonable ability to differentiate between the three languages based on the MFCC features. However, some misclassifications were observed, which could be attributed to overlapping acoustic characteristics and intra-language variability.

## 4. Discussion

### 4.1 Reflecting Acoustic Characteristics

MFCC features effectively capture the spectral envelope of speech:

- **Low-order Coefficients:** Represent the overall timbre and formant structure, which is strongly influenced by the shape and size of the vocal tract. Differences in vowel production across languages are thus clearly represented.
- **High-order Coefficients:** Capture finer details and rapid spectral changes related to consonant articulation and prosody. These are particularly useful in reflecting the dynamic aspects of speech that differ among languages.

### 4.2 Challenges Faced

Despite their utility, several challenges remain when using MFCCs for language classification:

- **Speaker Variability:** Different accents and speaking speeds affected MFCC consistency.
- **Background Noise:** Some clips had environmental sounds that interfered with feature extraction.
- **Regional Accents and Dialects:** Variations within a language (due to regional accents) cause additional spectral variability, making it more difficult to differentiate languages solely based on MFCCs.

## 5. Conclusion

This assignment demonstrates that MFCC features can effectively capture the spectral characteristics of speech in different Indian languages. The visual and statistical analysis of MFCC spectrograms for Marathi, Punjabi, and Hindi revealed distinct patterns that reflect each language's unique phonetic and prosodic features. Furthermore, the Random classifier built using these features achieved promising results in language prediction.

However, challenges such as speaker variability, background noise, and regional accents underscore the importance of robust preprocessing and feature normalization. Future

work may focus on incorporating additional features or employing more advanced deep learning models to further improve classification performance.

## **References:**

Here, the references are given below which I have gone through:

1. <https://librosa.org/doc/latest/>
2. <https://www.geeksforgeeks.org/mel-frequency-cepstral-coefficients-mfcc-for-speech-recognition/>
3. <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>
4. <https://huggingface.co/speechbrain/sepformer-whamr>
5. <https://huggingface.co/speechbrain/sepformer-wham-enhancement>

## **GitHub Link:**

[https://github.com/Ankit-IITJ/Speech\\_Understanding\\_PA2.git](https://github.com/Ankit-IITJ/Speech_Understanding_PA2.git)