**Project: Feeble Audio-Based Transcript Generation**
**Course: IE 643**
**Team: NLPeeps**
**Members: Aryash Srivastava, Dhruv Garg**
**Guide: Prof. Balamurugan**
**TAs: Mr. Rahul, Mr. Bheeshm**

---

## Introduction

This project aims to generate accurate transcriptions from low-quality audio with minimal loss of information, addressing challenges in recognizing speech affected by noise or missing segments.

## Abstract

This project addresses the challenge of transcribing feeble audio, focusing on handling low-quality and noise-filled audio inputs. By applying deep learning methods and pre-trained models, the aim was to enhance transcription accuracy, even when audio is distorted or contains missing segments. Key results show improved transcription accuracy with minimal semantic loss, despite audio noise.

## Introduction

Transcribing feeble or distorted audio is crucial for accessibility in communication and information retrieval. Poor audio quality can lead to information loss and misinterpretation, especially in automated systems.

## Workflow

1. Initial research and model selection.
2. Decision to utilise pre-trained models for efficiency and accuracy.
3. Creation of synthetic audio distortions (e.g., Gaussian and white noise) for model testing.
4. Experimentation with masked language models for transcription.
5. Final integration of functions into a unified model.

## Related Work

1. **Raskinov et al., 2013** - Discussed multi-label anomaly detection in speech, providing insight into dealing with noise.
2. **SpecAugment, 2024** - A technique that enhances robustness of audio models against distortions, referenced for spectrogram analysis.
3. **Spectrogram-Based Approaches** - Leveraged in the detection of muted audio sections, highlighting the value of spectral patterns in transcription models.
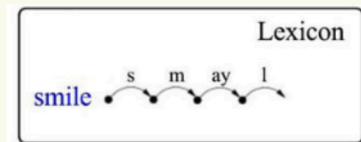
# Methods and Approaches

**Before Prep Presentation:**

- Conceptualised a custom speech recognition model based on Hidden Markov Models and Lexicon Models.
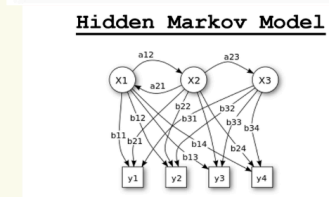- Initial experimentation with n-gram predictions for filling muted segments in audio.

**Lexicon Model:**
- breaks down word pronunciations into phonemes
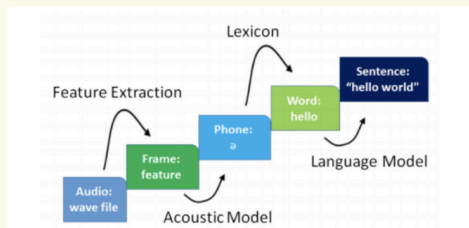- helping the system recognize and map to words

**Hidden Markov Model (HMM):**
- used to generate an encoded context vector
- predicting the next phoneme

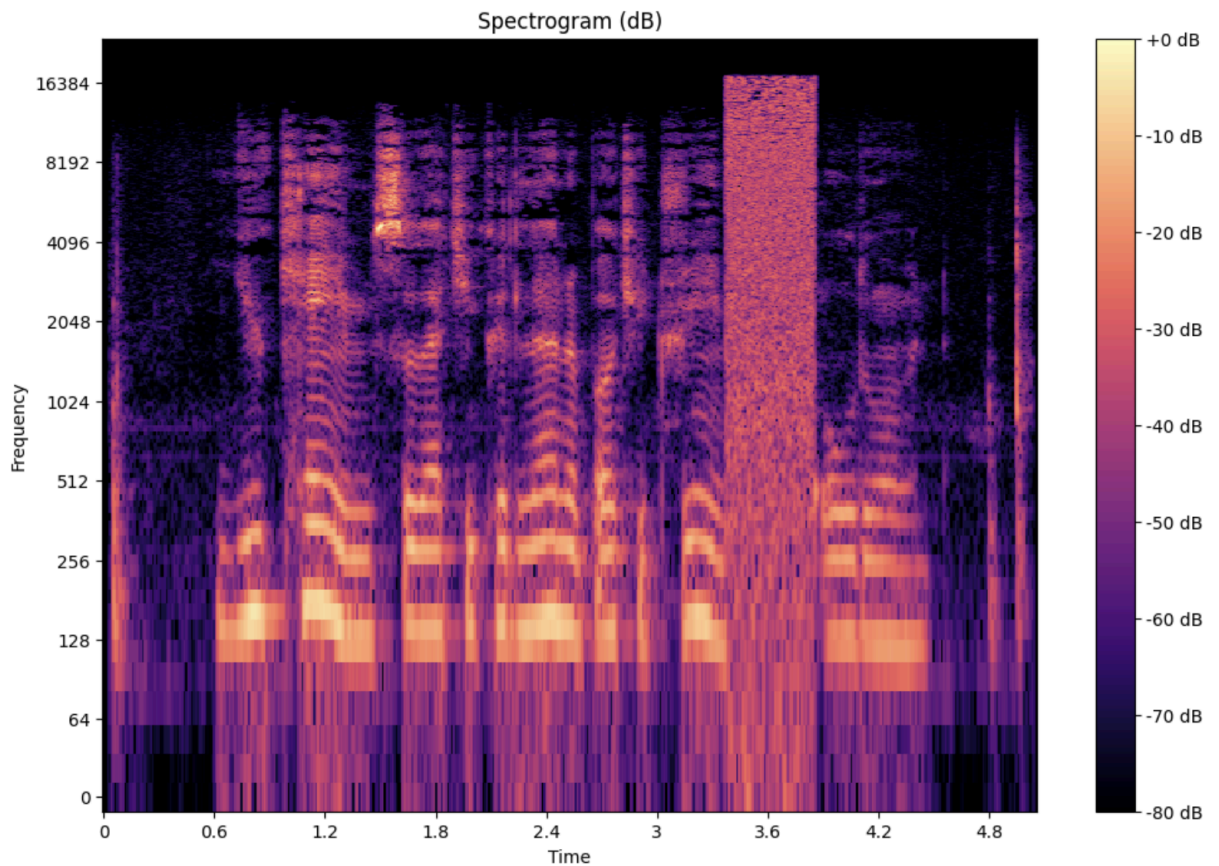**Transformer – Encoder Self-Attention:**
- self-attention mechanisms to focus on key parts of the audio features
- ensuring better context understanding for transcription

**After Prep Presentation:**

- Shifted to pre-trained models to address scalability.
- Implemented audio distortion simulation (Gaussian and white noise) using NumPy.
- Utilised a fill-mask model for predicting missing audio segments based on spectrogram analysis.

```
audio_path = '/kaggle/working/clip_audio/sample-000002.mp3'
plot_spectrogram(audio_path)
play_audio(audio_path)
```



- 
- Integrated all modules into a comprehensive "Final_Predictor" model for end-to-end transcription.
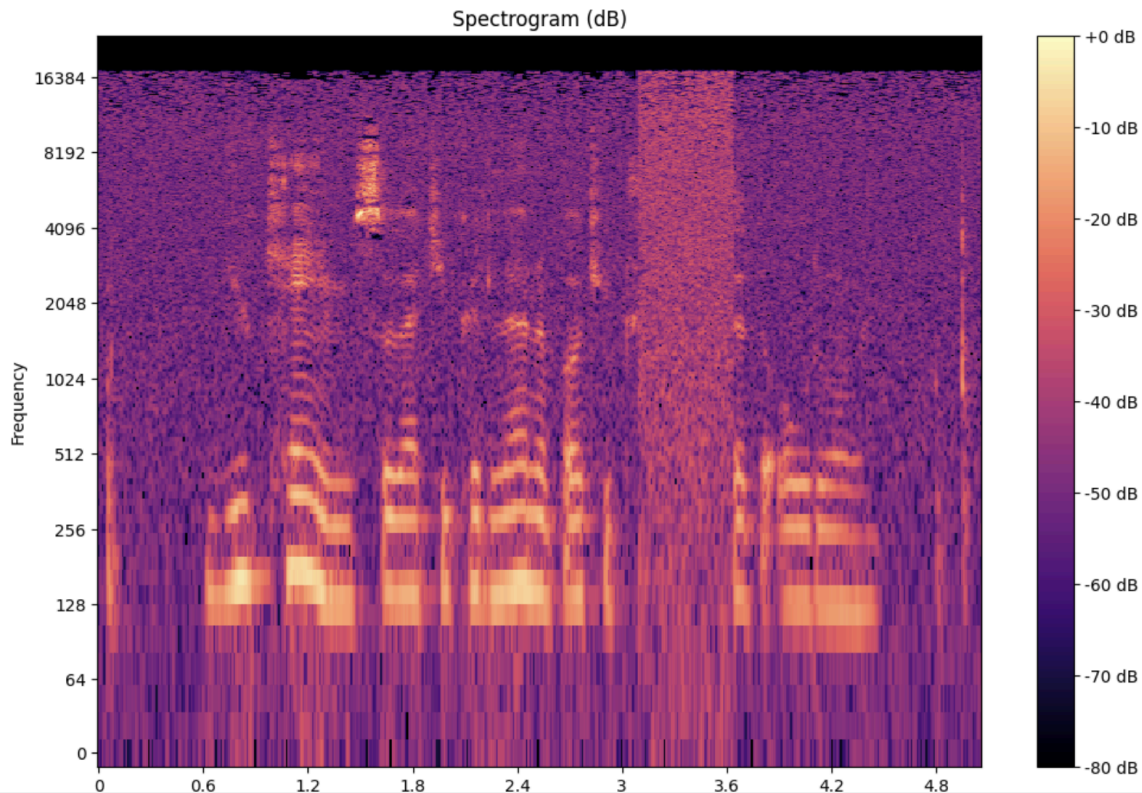
## Data

Dataset: **Mozilla Common Voice**

- **Size**: ~200,000 audio samples (3–7 seconds each).
- **Attributes**: Audio clips with corresponding transcriptions.
- **Preprocessing**: Clipping and noise addition using `librosa`. Gaussian noise added randomly; certain audio segments replaced with white noise to simulate real-world distortions.

# Experiments

```
plot_spectrogram('/kaggle/working/noise_and_clipped_audio/sample-000002.mp3')
amplitudes=list(store_amplitude('/kaggle/working/noise_and_clipped_audio/sample-000002.mp3')
```



Spectrogram (dB)

- **Model Integration**: Integrated spectrogram, noise addition, and masked prediction functions into one cohesive model.

```
Final_Predictor('/kaggle/input/common-voice/cv-valid-train/cv-valid-train/sample-000019.mp3')
```

The Original Audio

▶ 0:00 / 0:03 🔊 ⋮

Then I got a hold of some dough and went goofy.
The Distorted File is saved in /kaggle/working/noise_and_clipped_audio/sample-000019.mp3
The Distorted Audio:

▶ 0:00 / 0:03 🔊 ⋮

▶ 0:00 / 0:03 🔊 ⋮

Then I got some joe when went goofy.

▶ 0:01 / 0:01 🔊 ⋮

▶ 0:00 / 0:03 🔊 ⋮

```
text1: Then I go
text2:  got some joe when went goofy
Masked: Then I go <mask> got some joe when went goofy
Then I go and got some joe when went goofy
```

- **Training and Testing**: Used Mozilla Common Voice and real-time audio recordings for model validation.
- **Hardware**: The project was tested on a standard GPU configuration for efficient processing.
- **GitHub Repository**: https://github.com/AryashSrivastava/IE_643-Feeble-Audio-Transcript-Generation

## Results

**Quantitative Results**:

- No Data. Because it was difficult to apply any performance metric due to being different and convey the same meaning.

**Qualitative Results**:

- **Example**:
  - Original: "Then I got a hold of some dough and went goofy."
  - Predicted: "Then got a hold of some joe and went goofy."

## Plan for Novelty Assessment

Actually, there are still a lot of improvements which can be made. After doing all those we can look upon to deploy this model on some private website.

## Conclusion

Hence, by Experimental results we can see that. It is not 100% accurate but it predicts the text almost correctly in most of the cases. Performance metric can not be used here, because even similar words will cover same semantic meaning but will not be classified same.

## References

1. **Steven Raskinov, Bob Dietterich, and Rita Barnard.** "SVMs for Multi-label Anomaly Detection," *Journal of Machine Learning*, 2013. https://arxiv.org/1303.1231223
2. **Chris Bisred.** https://blog.com/samplewebsite, Accessed on: 7th March 2024.