# Speech Processing Assignment - 2

## Dataset Description :

- This LJ Speech dataset has been taken from the Kaggle.
- It contains 13,100 short audio clips of a single speaker
- These clips vary in length from 1 to 10 sec and have a total length of around 24 hours.

## Objective :

- ‣ Load and preprocess a speech signal (convert to mono and resample to 16kHz).
  ‣ Use a pre-trained Wav2Vec2 model to recognize phonemes from the speech signal.
  ‣ Extract specific phoneme segments based on their time intervals.
  ‣ Visualize the phoneme waveforms and label them according to the recognized phonemes.

## Code :

```python
import torchaudio
import torchaudio.transforms as T
import librosa.display
import torch
from transformers import Wav2Vec2Processor, Wav2Vec2ForCTC
import matplotlib.pyplot as plt
import numpy as np
import nltk
from nltk.corpus import cmudict

nltk.download('cmudict')
cmu_dict = cmudict.dict()

audio_path = "/content/LJ001-0030.wav"
waveform, sample_rate = torchaudio.load(audio_path)

resampler = T.Resample(orig_freq=sample_rate, new_freq=16000)
waveform = resampler(waveform.mean(dim=0, keepdim=True))

processor = Wav2Vec2Processor.from_pretrained("facebook/wav2vec2-large-960h")
model = Wav2Vec2ForCTC.from_pretrained("facebook/wav2vec2-large-960h")
```

```python
24  input_values = processor(waveform.squeeze().numpy(), return_tensors="pt",
    sampling_rate=16000).input_values
25  with torch.no_grad():
26      logits = model(input_values).logits
27
28
29  predicted_ids = torch.argmax(logits, dim=-1)
30  transcription = processor.batch_decode(predicted_ids)[0].upper()
31  print("Recognized Text:", transcription)
32
33  words = transcription.split()
34  phonemes = []
35  for word in words:
36      if word in cmu_dict:
37          phonemes.extend(cmu_dict[word][0])
38      else:
39          phonemes.append(word)
40
41  print("Recognized Phonemes:", " ".join(phonemes))
42
43  total_duration = waveform.shape[1] / 16000
44  time_intervals = np.linspace(0, total_duration, num=len(phonemes))
45
46  plt.figure(figsize=(12, 4))
47  librosa.display.waveshow(waveform.squeeze().numpy(), sr=16000)
48  plt.xlabel("Time (s)")
49  plt.ylabel("Amplitude")
50  plt.title("Phoneme Visualization")
51
52
53  for i, phoneme in enumerate(phonemes):
54      if i < len(time_intervals):
55          plt.text(time_intervals[i], 0, phoneme, fontsize=10, ha='center',
           color='red', rotation=45)
56
57  plt.show()
```
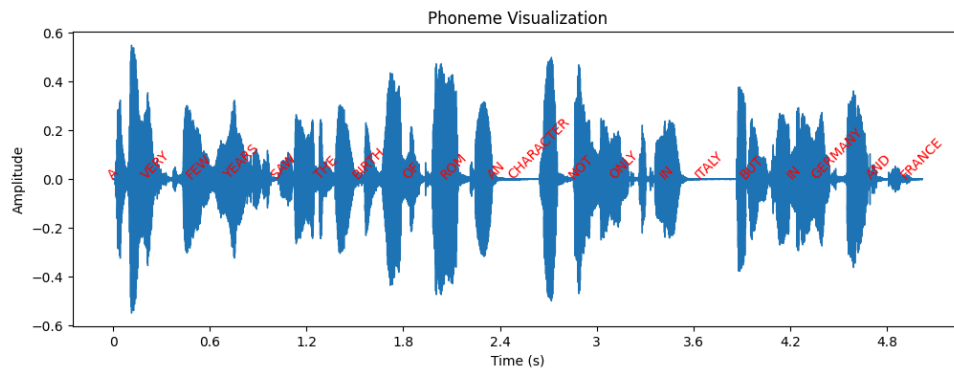
Figure 1: Final Plot with recognized text .