

Spectral Analysis

Shrey Mehta (200580)

Abstract—Given an audio file, use python to plot:

- the waveform $x(t)$
- spectrogram $X(f,t)$, as a function of window size (in s), hop length (in s), N, etc.
- spectrum $X(f)$, as a function of start time, N, etc.

We use Python language to take as input the .wav file of the audio signal and output the required plots using the following code:

- Import the required libraries

```
1 import pyaudio
2 import wave
3 import matplotlib.pyplot as plt
4 import numpy as np
5 from scipy.fftpack import fft
```

- Process the audio signal using pyaudio

```
1 audio_file = 'audio_file.wav'
2 p = pyaudio.PyAudio()
3 wf = wave.open(audio_file, 'rb')
4 sample_width = wf.getsampwidth()
5 sample_rate = wf.getframerate()
6 audio_data = wf.readframes(-1)
7 audio_data = np.frombuffer(audio_data, dtype=np.
    int16)
8 print(sample_rate)
9 print(len(audio_data))
10 wf.close()
11 p.terminate()
```

- Plotting the waveform $x(t)$

```
1 duration = len(audio_data) / sample_rate #
    Calculate audio duration
2 time = np.linspace(0, duration, num=len(
    audio_data)) # Create time axis
3 plt.figure(figsize=(10, 4))
4 plt.plot(time, audio_data)
5 plt.xlabel("Time (s)")
6 plt.ylabel("Amplitude")
7 plt.title("Audio Waveform")
8 plt.grid()
9 plt.show()
```

- Plotting the spectrogram $X(f,t)$ for the given waveform

```
1 # Define the window size and overlap
2 window_size = 1024 # Adjust this as needed
3 overlap = int(window_size * 0.5)
4
5 # Calculate the number of FFTs
6 num_ffts = int((len(audio_data) - overlap) / (
    window_size - overlap))
7
8 # Create an empty array to store the results
9 spectrogram = np.empty((num_ffts, window_size //
    2))
10
11 # Perform FFT on overlapping segments
12 for i in range(num_ffts):
13     start = i * (window_size - overlap)
14     end = start + window_size
15     segment = audio_data[start:end]
16     fft_result = np.fft.fft(segment)
```

```
17     frequencies = np.fft.fftfreq(len(fft_result)
    , 1.0 / sample_rate)
18     magnitude = 20 * np.log10(np.abs(fft_result
    [:len(fft_result) // 2]))
19     spectrogram[i, :] = magnitude
20
21 # Create a time axis
22 time_axis = np.arange(0, num_ffts) * (
    window_size - overlap) / sample_rate
23
24 # Create a frequency axis
25 frequency_axis = frequencies[:window_size // 2]
26
27 # Plot the spectrogram
28 plt.figure(figsize=(10, 4))
29 plt.pcolormesh(time_axis, frequency_axis,
    spectrogram.T, shading='auto', cmap='viridis')
30 plt.colorbar(label='Magnitude (dB)')
31 plt.xlabel('Time (s)')
32 plt.ylabel('Frequency (Hz)')
33 plt.title('Spectrogram (Time vs. Frequency)')
34 plt.show()
```

- Plotting the spectrum $X(f)$ for the given waveform (using decibel scale)

```
1 # Perform FFT on the entire audio signal
2 fft_result = np.fft.fft(audio_data)
3 frequencies = np.fft.fftfreq(len(fft_result),
    1.0 / sample_rate)
4
5 # Take the absolute value and scale to dB
6 magnitude = 20 * np.log10(np.abs(fft_result[:len
    (fft_result) // 2]))
7
8 # Plot the frequency spectrum
9 plt.figure(figsize=(10, 4))
10 plt.plot(frequencies[:len(frequencies) // 2],
    magnitude)
11 plt.xlabel('Frequency (Hz)')
12 plt.ylabel('Magnitude (dB)')
13 plt.title('Frequency Spectrum (FFT)')
14 plt.grid()
15 plt.show()
```

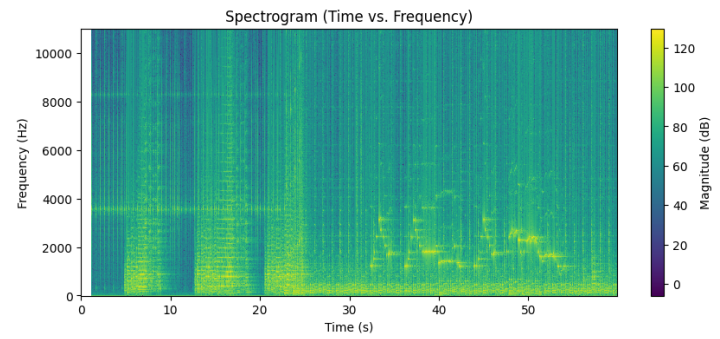
- Plotting the spectrogram for the given waveform (using mel scale)

```
1 # Create a function to compute Mel filterbanks
2 def mel_filterbank(num_filters, fft_size,
    sample_rate):
3     mel_min = 0
4     mel_max = 2595 * np.log10(1 + (sample_rate /
    2) / 700)
5     mel_points = np.linspace(mel_min, mel_max,
    num_filters + 2)
6     hz_points = 700 * (10**((mel_points / 2595) -
    1))
7     bin_points = np.floor((fft_size + 1) *
    hz_points / sample_rate).astype(int)
8     filters = np.zeros((num_filters, fft_size //
    2 + 1))
9     for i in range(1, num_filters + 1):
10         filters[i - 1, bin_points[i - 1]:
            bin_points[i]] = (np.arange(bin_points[i -
            1], bin_points[i]) - bin_points[i - 1]) / (
                bin_points[i] - bin_points[i - 1])
```

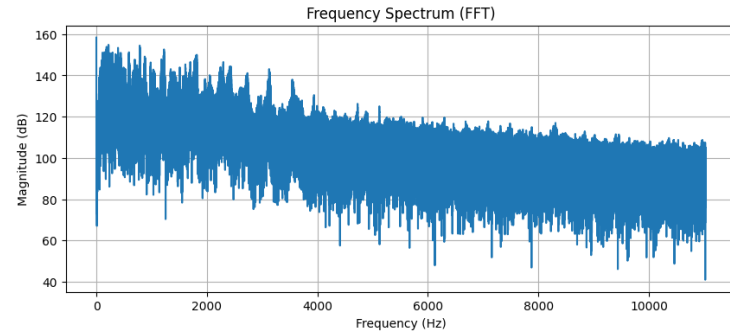
```

11     filters[i - 1, bin_points[i]:bin_points[
12         i + 1]] = 1 - (np.arange(bin_points[i],
13         bin_points[i + 1]) - bin_points[i]) / (
14         bin_points[i + 1] - bin_points[i])
15     return filters
16
17 # Define the parameters for the spectrogram
18 window_size = 512 # Adjust this as needed
19 overlap = window_size // 2 # Adjust this as
20     needed
21 hop_size = window_size - overlap
22
23 # Calculate the FFT frequencies
24 fft_frequencies = np.fft.fftfreq(window_size,
25     1.0 / sample_rate)
26
27 # Create an empty list to store the spectrogram
28     frames
29 spectrogram_frames = []
30
31 # Create the Mel filterbank
32 num_filters = 40 # Adjust this as needed
33 mel_filters = mel_filterbank(num_filters,
34     window_size, sample_rate)
35
36 # Iterate through the audio data with
37     overlapping windows
38 for i in range(0, len(audio_data) - window_size,
39     hop_size):
40     audio_segment = audio_data[i:i + window_size
41     ]
42     fft_result = fft(audio_segment)
43     mel_spectrum = np.dot(mel_filters, np.abs(
44     fft_result[:window_size // 2 + 1]))
45     spectrogram_frames.append(mel_spectrum)
46
47 # Convert the list of frames to a numpy array
48 spectrogram = np.array(spectrogram_frames).T
49
50 # Plot the spectrogram
51 plt.figure(figsize=(10, 4))
52 plt.imshow(spectrogram, aspect='auto', cmap='
53     viridis', origin='lower')
54 plt.colorbar(label='Magnitude')
55 plt.xlabel('Time Frame')
56 plt.ylabel('Mel Filterbank Index')
57 plt.title('Mel Spectrogram')
58 plt.show()

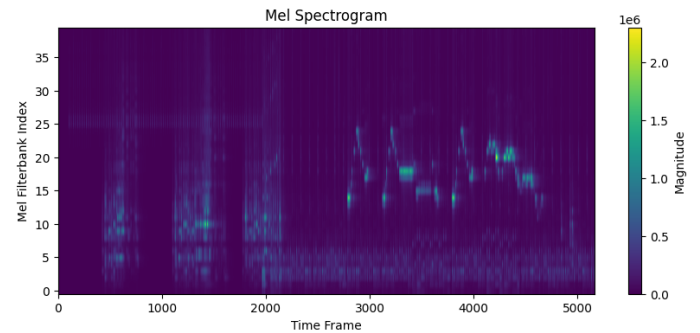
```



- Spectrum



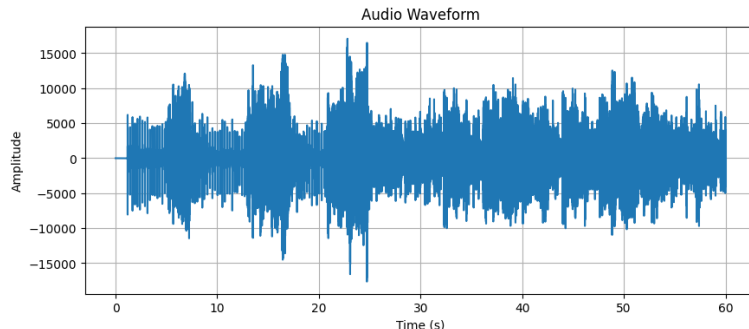
- Spectrogram (using mel scale)



Now, we will analyse these plots for various waveforms and draw conclusions from them.

Example: We take a sample audio file as an example and plot the corresponding four plots for the given audio signal. The plots obtained are as shown:

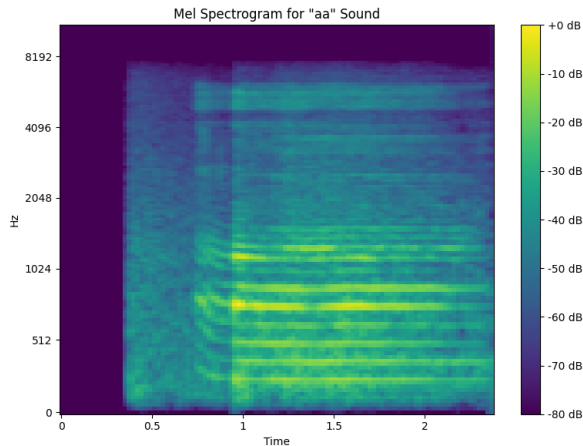
- Waveform



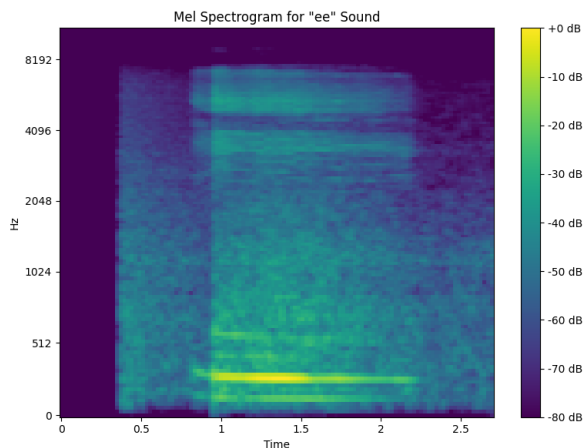
- Spectrogram (using dB scale)

Analysis of vowels: We will analyse the sounds such as aa, ee, o, u, etc., and study the characteristic features of their spectra/spectrograms. On plotting the spectrograms of these sounds, we obtain the following plots:

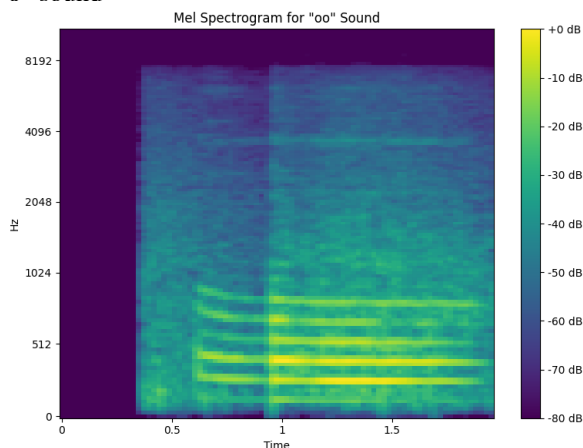
• **'aa' sound**



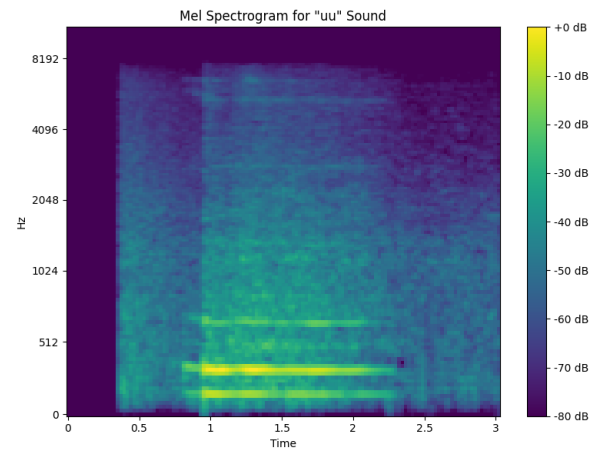
• **'ee' sound**



• **'o' sound**



• **'u' sound**



Observations:

- **'aa'** sound is a low, open vowel sound. The spectrogram shows strong fundamental frequency (the lowest frequency component) and several harmonics (integer multiples of the fundamental). The formants (resonant frequencies) are relatively evenly spaced and typically low in frequency.
- **'ee'** sound is a high, front vowel sound. The spectrogram shows strong fundamental frequency and harmonics, similar to "aa" but the formants are higher in frequency. The first formant is relatively high, and the second formant is even higher, making it a front vowel.
- **'o'** sound is a mid-to-high, back vowel sound. The spectrogram shows strong fundamental frequency and harmonics, but the formants are shifted towards the lower-mid frequencies. The first formant is relatively low, and the second formant is moderate in frequency, placing it in the back vowel category.
- **'u'** sound is a high, back vowel sound. The spectrogram shows strong fundamental frequency and harmonics like the others, but the formants are positioned at relatively high frequencies. The first formant is low, and the second formant is high, making it a high, back vowel.

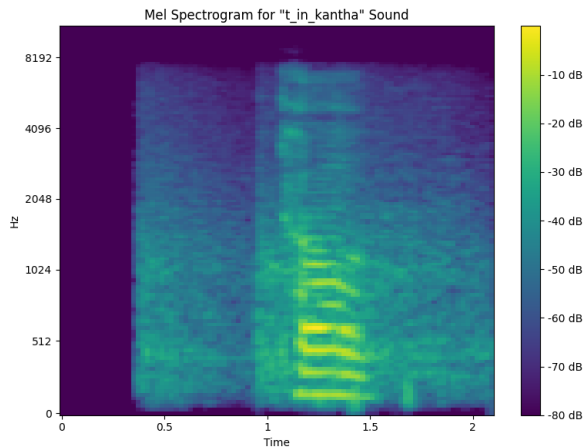
In spectrograms of these vowel sounds:

- "aa" will show a strong, dark band in the low-frequency range.
- "ee" will have formants concentrated in the high-frequency range.
- "o" will show formants in the mid-frequency range.
- "u" will exhibit formants in the high-frequency range but shifted towards the back.

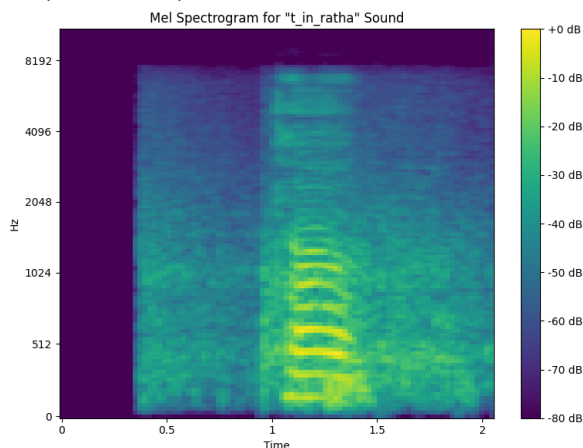
These spectral and spectrographic characteristics are essential in phonetics and speech analysis as they help differentiate between different vowel sounds and contribute to our understanding of speech production and perception.

Analysis of 't' sounds: We will analyse the sounds such as t (as in tiwari), th (as in ratha), t (as in topi), th (as in neelkantha), and study the characteristic features of their spectra/spectrograms. On plotting the spectrograms of these sounds, we obtain the following plots:

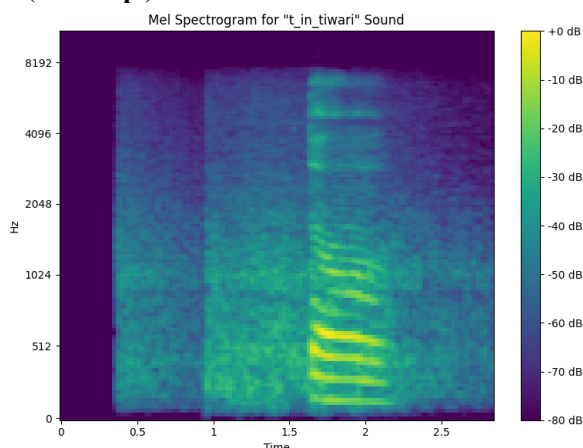
• **'t (as in tiwari)' sound**



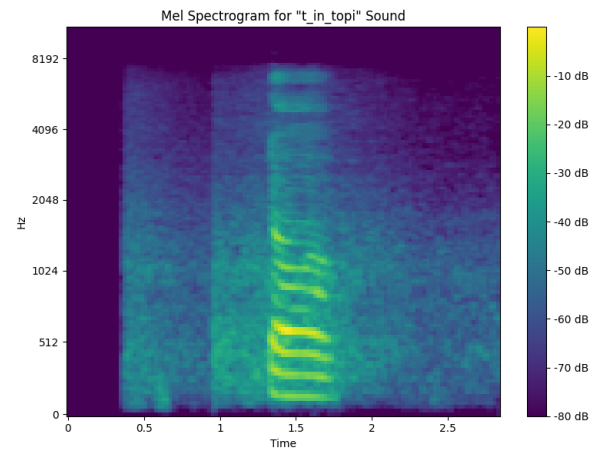
• **'th (as in ratha)' sound**



• **'t (as in topi)' sound**



• **'th (as in neelkantha)' sound**



Observations:

- **'t (as in tiwari)'**: The spectrogram shows burst of energy at a specific frequency corresponding to the release of the closure (the moment when the tongue is released).
- **'th (as in ratha)'**: The spectrogram shows a similar burst of energy as in the above "t," but with a longer duration.
- **'t (as in topi)'**: The spectrogram shows a burst of energy at a specific frequency corresponding to the release of the closure.
- **'th (as in neelkantha)'**: The spectrogram shows a similar burst of energy as in the above "t," but with a longer duration.

As seen from the plots on Page 2, we can also analyse the differences between the spectrograms plotted using the decibel scale and the mel scale:

- **Using the dB (Decibel) Scale:**
 - The dB scale is a logarithmic scale that measures the intensity or amplitude of a signal relative to a reference level. It represents power or energy ratios in a more human-perceptible way.
 - When the dB scale is used for spectra or spectrograms, the y-axis (intensity) is represented in decibels instead of linear amplitude.
 - The dB scale is particularly useful because it compresses the range of intensity values. In speech processing, this helps highlight the signal's relevant features, making it easier to identify and analyze formants, harmonics, and other spectral components.
- **Using the Mel Scale:**
 - Designed to mimic the way humans perceive differences in pitch. It's based on the fact that our perception of pitch is not linear but rather logarithmic.
 - In the Mel scale in spectrograms, the frequency axis is transformed so that it's no longer linear but follows the Mel scale. This means that the spacing between frequency bins is wider at lower frequencies and narrower at higher frequencies.
 - By using the Mel scale, we can represent the spectrogram in a way that's more relevant to how our ears perceive speech sounds.
 - It helps in capturing important acoustic features for speech recognition and synthesis.