



國立陽明交通大學
NATIONAL YANG MING CHIAO TUNG UNIVERSITY

[IIAI30003] Digital Speech Processing

Homework 2

Po-Chuan, Chen
Student ID: 311511052
present90308.ee11@nycu.edu.tw

NATIONAL YANG MING CHIAO TUNG UNIVERSITY

October 6, 2023

Contents

1	Speech Analysis &Feature Extraction	2
2	Audio Data Augmentation	10
2.1	Nlpaug	11
2.2	Audiomentations	15

1 Speech Analysis & Feature Extraction

First, I record the audio file, and the program will output the waveform.

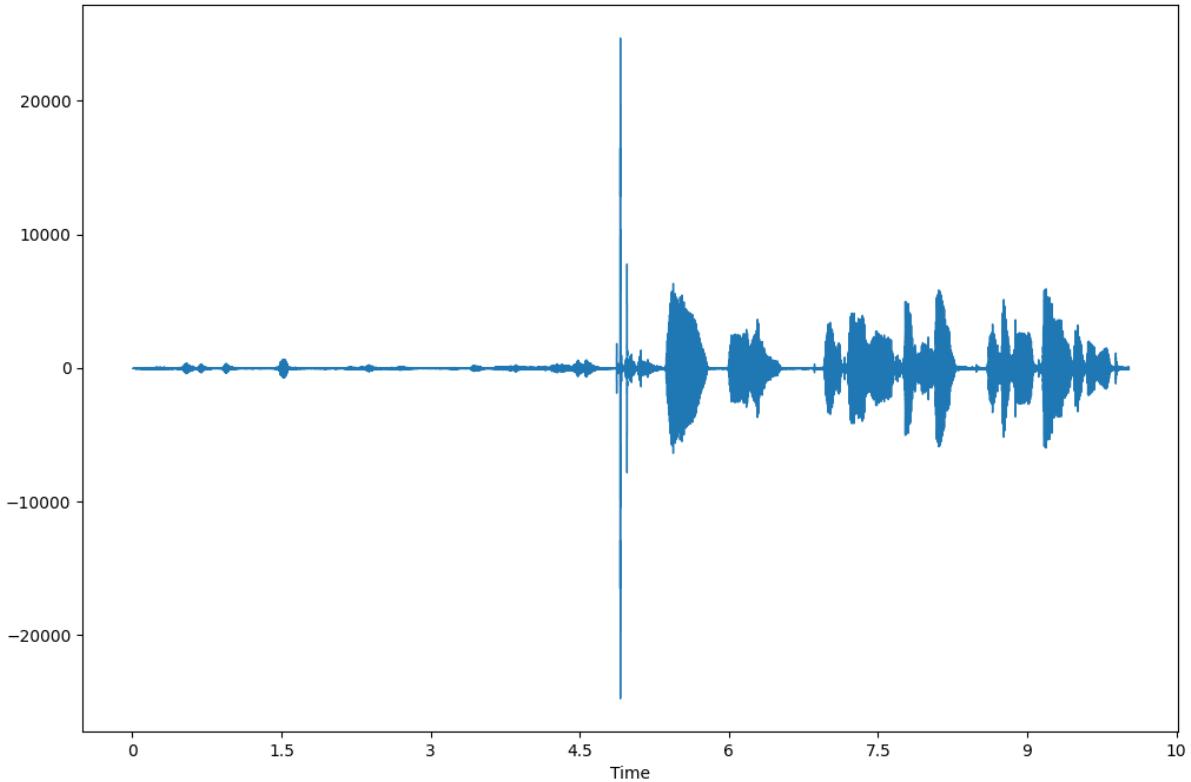


Figure 1: Waveform

With the `librosa.stft`, we can use Short-Time Fourier Transform convert a time-domain signal into its frequency-domain representation.

```
1 x = librosa.stft(audio, n_fft=2048, hop_length=480) # sampling rate  
    =48000, frame_size=32 ms, frame_shift=10 ms  
2 plt.figure(figsize=(12, 8))  
3 librosa.display.specshow(librosa.amplitude_to_db(np.abs(x)), sr=sr,  
    y_axis='linear', x_axis='time')
```

Listing 1: Code section for Spectrogram

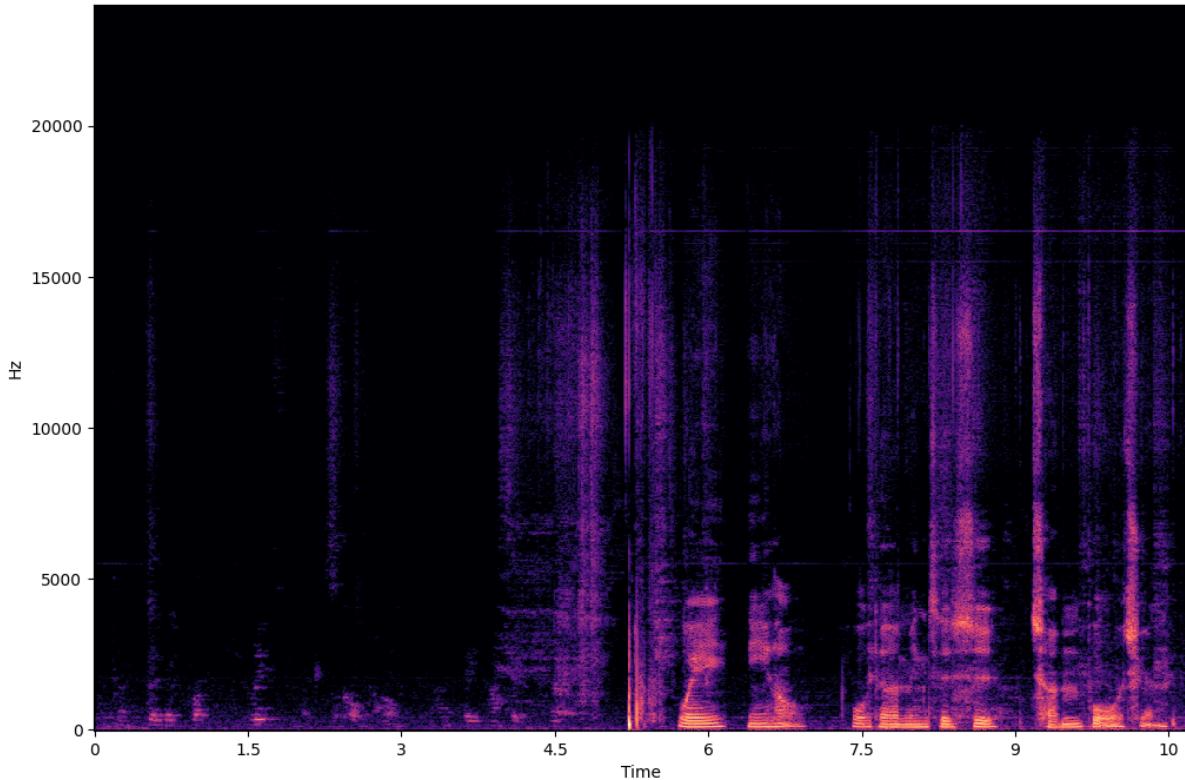


Figure 2: Spectrogram

The left subplot displays the entire pitch information estimated from the audio signal, while the right subplot shows a closer view of the first 100 rows of the pitch information.

```

1 pitches, magnitudes = librosa.piptrack(y=audio, sr=sr)
2 plt.figure(figsize=(28, 8))
3 plt.subplot(1,2,1)
4 plt.imshow(pitches[:, :], aspect="auto", interpolation="nearest",
   origin="lower")
5 plt.subplot(1,2,2)
6 plt.imshow(pitches[:100, :], aspect="auto", interpolation="nearest",
   origin="lower")

```

Listing 2: Code section for Pitch Tracking

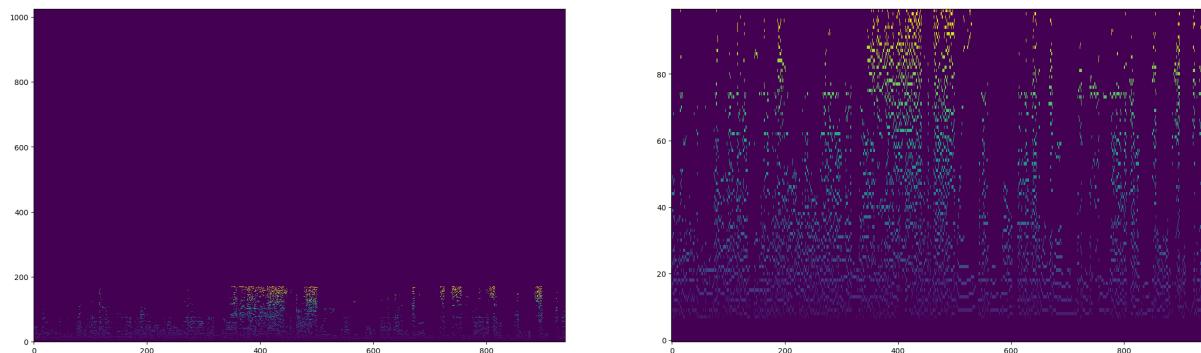


Figure 3: Pitch Tracking

Here, we use the Librosa library to perform some audio processing and visualization tasks. It generates and displays two spectrograms side by side: one for the original audio signal and one for the resampled audio signal.

```

1 plt.figure(figsize=(24, 8))
2
3 plt.subplot(1,2,1)
4 x = librosa.stft(audio, n_fft=2048, hop_length=480) # # sampling rate
      =16000, frame_size=32 ms, frame_shift=10 ms
5 librosa.display.specshow(librosa.amplitude_to_db(np.abs(x)), sr=sr,
      y_axis='linear', x_axis='time')
6
7 # resampling
8 target_sr = 16000
9 audio = librosa.resample(audio, orig_sr=sr, target_sr=target_sr)
10 sr = target_sr
11
12 plt.subplot(1,2,2)
13 x = librosa.stft(audio, n_fft=2048, hop_length=480) # # sampling rate
      =16000, frame_size=32 ms, frame_shift=10 ms
14 librosa.display.specshow(librosa.amplitude_to_db(np.abs(x)), sr=sr,
      y_axis='linear', x_axis='time')
```

Listing 3: Code section for Resampling

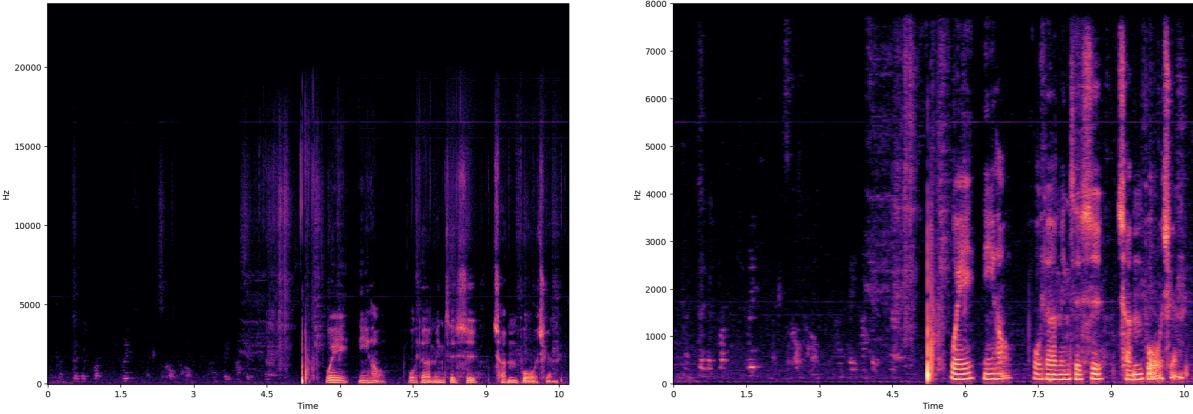


Figure 4: Resampling

The left subplot has a larger `n_fft` value, which provides higher frequency resolution but potentially lower time resolution, while the right subplot has a smaller `n_fft` value, providing higher time resolution but lower frequency resolution.

```

1 plt.figure(figsize=(24, 8))
2
3 plt.subplot(1,2,1)
4 x = librosa.stft(audio, n_fft=2048, hop_length=80) # # sampling rate
      =16000, frame_size=32 ms, frame_shift=10 ms
5 librosa.display.specshow(librosa.amplitude_to_db(np.abs(x)), sr=sr,
      y_axis='linear', x_axis='time')
6
```

```

7 plt.subplot(1,2,2)
8 x = librosa.stft(audio, n_fft=128, hop_length=80) # # sampling rate
# =16000, frame_size=32 ms, frame_shift=10 ms
9 librosa.display.specshow(librosa.amplitude_to_db(np.abs(x)), sr=sr,
y_axis='linear', x_axis='time')

```

Listing 4: Code section for NarrowBand and WideBand Spectrogram

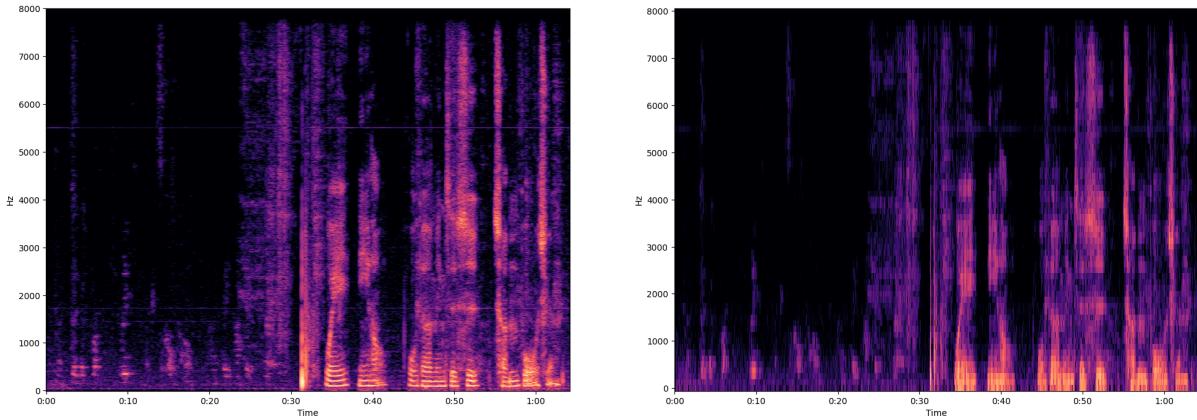


Figure 5: NarrowBand and WideBand Spectrogram

Using the Parselmouth library to perform some audio analysis and visualization.

```

1 snd = parselmouth.Sound(human_sound_file)
2 snd.resample(new_frequency=16000)
3 pitch = snd.to_pitch()
4 # If desired, pre-emphasize the sound fragment before calculating the
# spectrogram
5 pre_emphasized_snd = snd.copy()
6 pre_emphasized_snd.pre_emphasize()
7 spectrogram = pre_emphasized_snd.to_spectrogram(window_length=0.03,
maximum_frequency=8000)
8
9 plt.figure(figsize=(12, 8))
10 draw_spectrogram(spectrogram)
11 plt.twinx()
12 draw_pitch(pitch)
13 plt.xlim([snd.xmin, snd xmax])
14 plt.show() # or plt.savefig("spectrogram_0.03.pdf")

```

Listing 5: Code section for Pitch Contour Extraction

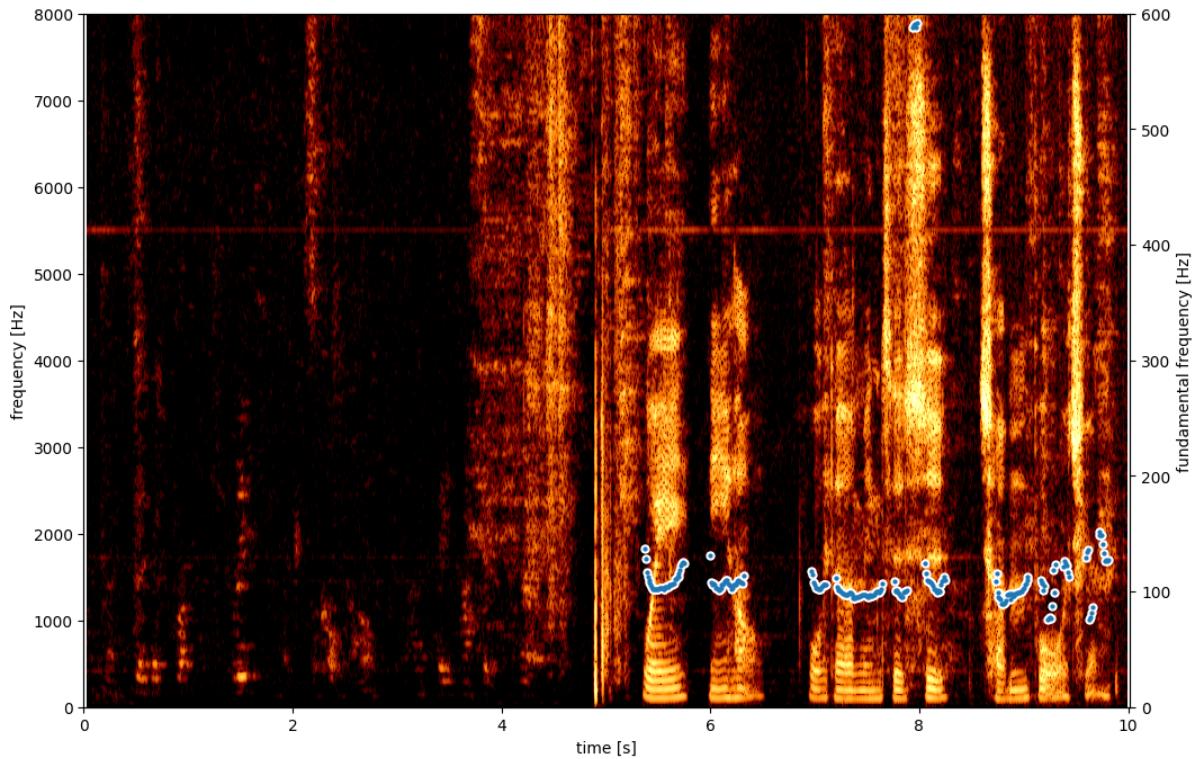


Figure 6: Pitch Contour Extraction

Here, the code generates and displays two spectrograms side by side: a linear power spectrogram on the left and a mel-scaled power spectrogram on the right.

```

1 plt.figure(figsize=(24, 8))
2
3 plt.subplot(1,2,1)
4 x = librosa.stft(audio, n_fft=2048, hop_length=80) # # sampling rate
        =16000, frame_size=32 ms, frame_shift=10 ms
5 librosa.display.specshow(librosa.amplitude_to_db(np.abs(x)), sr=sr,
    y_axis='linear', x_axis='time')
6 # Put a descriptive title on the plot
7 plt.title('linear power spectrogram')
8 # draw a color bar
9 plt.colorbar(format='%+02.0f dB')
10
11
12 plt.subplot(1,2,2)
13 # Let's make and display a mel-scaled power (energy-squared)
        spectrogram
14 S = librosa.feature.melspectrogram(y=audio, sr=sr, n_mels=128)
15 # Convert to log scale (dB). We'll use the peak power (max) as
        reference.
16 log_S = librosa.power_to_db(S, ref=np.max)
17 # Display the spectrogram on a mel scale
18 # sample rate and hop length parameters are used to render the time
        axis
19 librosa.display.specshow(log_S, sr=sr, x_axis='time', y_axis='mel')
20 # Put a descriptive title on the plot
21 plt.title('mel power spectrogram')
```

```

22 # draw a color bar
23 plt.colorbar(format='%+02.0f dB')
24
25 # Make the figure layout compact
26 plt.tight_layout()

```

Listing 6: Code section for Mel-Scaled Spectrogram

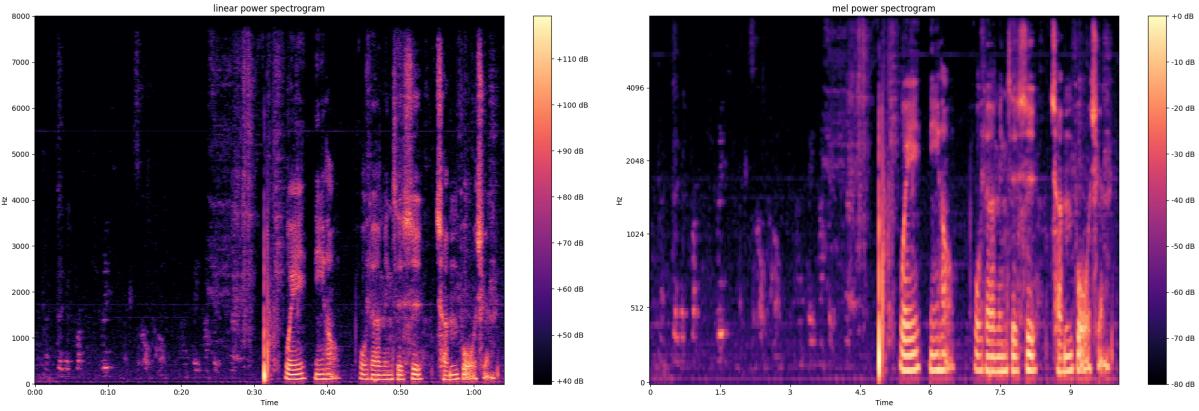


Figure 7: Mel-Scaled Spectrogram

Now it focusing on extracting Mel-frequency cepstral coefficients (MFCCs) and their derivatives. The resulting coefficients are standardized, and visualizations of these coefficients are displayed in separate subplots.

```

1 # Next, we'll extract the top 13 Mel-frequency cepstral coefficients (
2 # MFCCs)
3 mfcc      = librosa.feature.mfcc(S=log_S, n_mfcc=13)
4 print(mfcc.shape)
5 print(mfcc)
6
6 # zero mean, unit variance
7 mfcc      = preprocessing.scale(mfcc, axis=1)
8
9 # Let's pad on the first and second deltas while we're at it
10 delta_mfcc = librosa.feature.delta(mfcc)
11 delta2_mfcc = librosa.feature.delta(mfcc, order=2)
12
13 # How do they look? We'll show each in its own subplot
14 plt.figure(figsize=(12.5, 12))
15
16 plt.subplot(3,1,1)
17 librosa.display.specshow(mfcc)
18 plt.ylabel('MFCC')
19 plt.colorbar()
20
21 plt.subplot(3,1,2)
22 librosa.display.specshow(delta_mfcc)
23 plt.ylabel('MFCC-$\Delta$')
24 plt.colorbar()

```

```

25
26 plt.subplot(3,1,3)
27 librosa.display.specshow(delta2_mfcc, sr=sr, x_axis='time')
28 plt.ylabel('MFCC-$\Delta^2$')
29 plt.colorbar()
30
31 plt.tight_layout()
32
33 # For future use, we'll stack these together into one matrix
34 M = np.vstack([mfcc, delta_mfcc, delta2_mfcc])

```

Listing 7: Code section for MFCCs

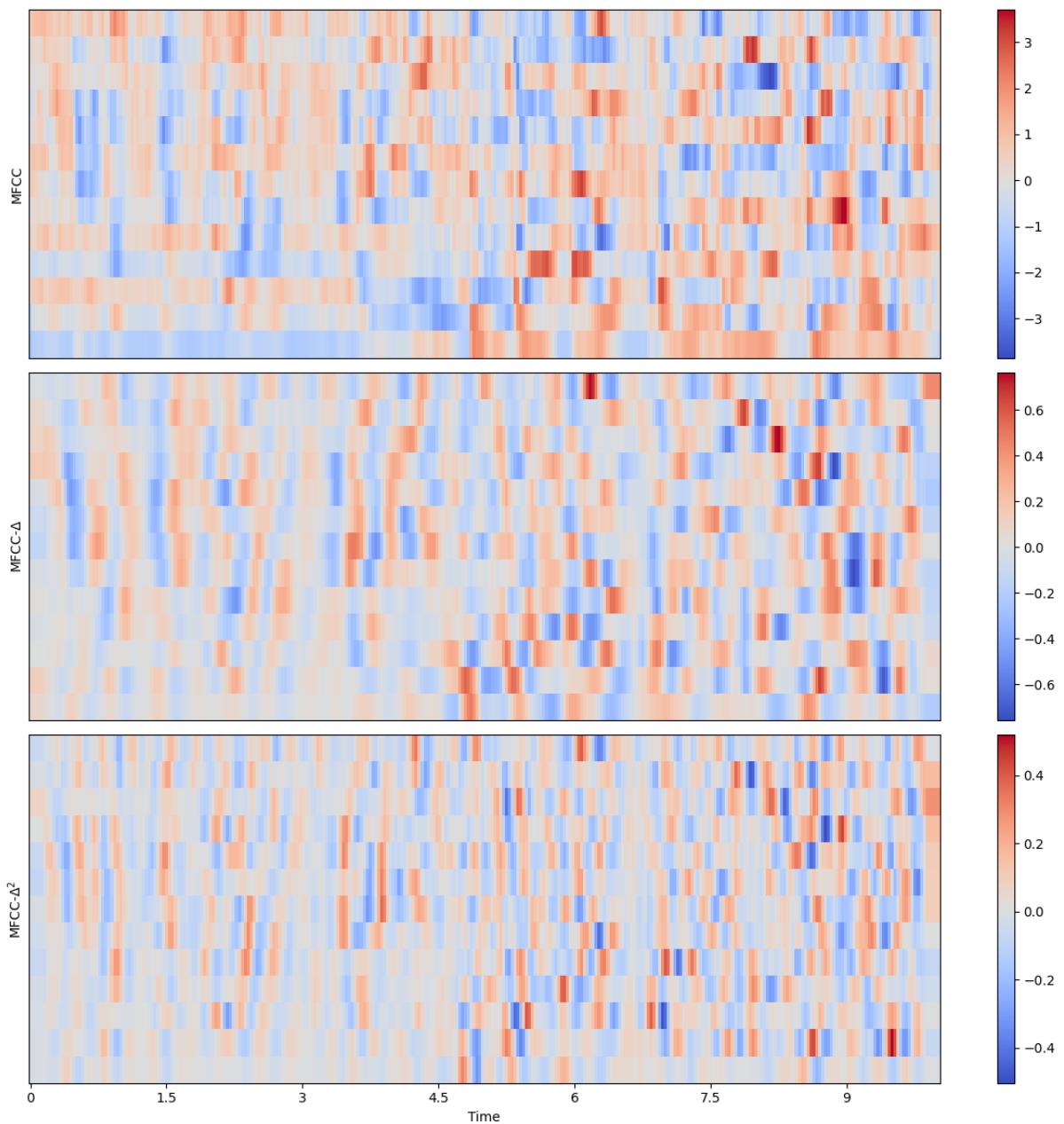


Figure 8: MFCCs

The last step, it uses the Parselmouth library to perform audio analysis and visualization for two different audio files. Two of the plots are a spectrogram and pitch contour for a different audio file.

```

1 plt.figure(figsize=(28, 8))
2
3 plt.subplot(1,2,1)
4 snd = parselmouth.Sound(human_sound_file)
5 pitch = snd.to_pitch()
6 # If desired, pre-emphasize the sound fragment before calculating the
    spectrogram
7 pre_emphasized_snd = snd.copy()
8 pre_emphasized_snd.pre_emphasize()
9 spectrogram = pre_emphasized_snd.to_spectrogram(window_length=0.03,
    maximum_frequency=8000)
10
11 draw_spectrogram(spectrogram)
12 plt.twinx()
13 draw_pitch(pitch)
14 plt.xlim([snd.xmin, snd xmax])
15
16 plt.subplot(1,2,2)
17 snd = parselmouth.Sound('output.wav')
18 pitch = snd.to_pitch()
19 # If desired, pre-emphasize the sound fragment before calculating the
    spectrogram
20 pre_emphasized_snd = snd.copy()
21 pre_emphasized_snd.pre_emphasize()
22 spectrogram = pre_emphasized_snd.to_spectrogram(window_length=0.03,
    maximum_frequency=8000)
23
24 draw_spectrogram(spectrogram)
25 plt.twinx()
26 draw_pitch(pitch)
27 plt.xlim([snd.xmin, snd xmax])

```

Listing 8: Code section for Rubber Band

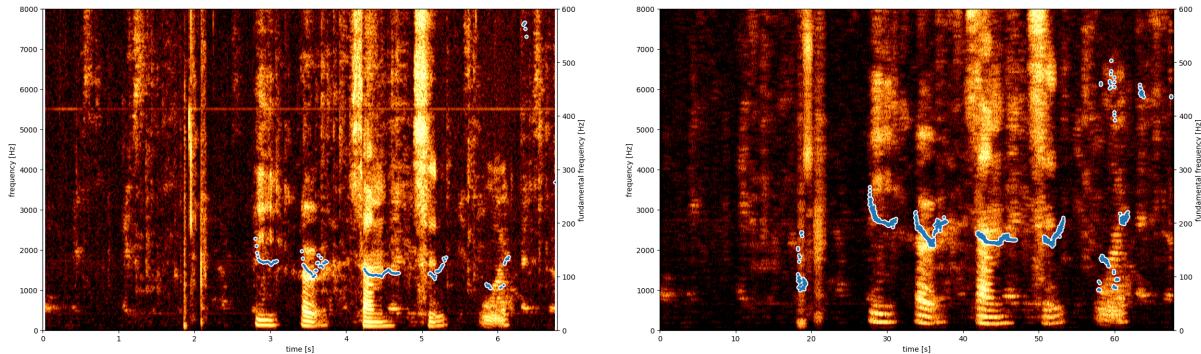


Figure 9: Rubber Band

2 Audio Data Augmentation

The source of the audio is Jay Chou's song, In the Name of the Father.

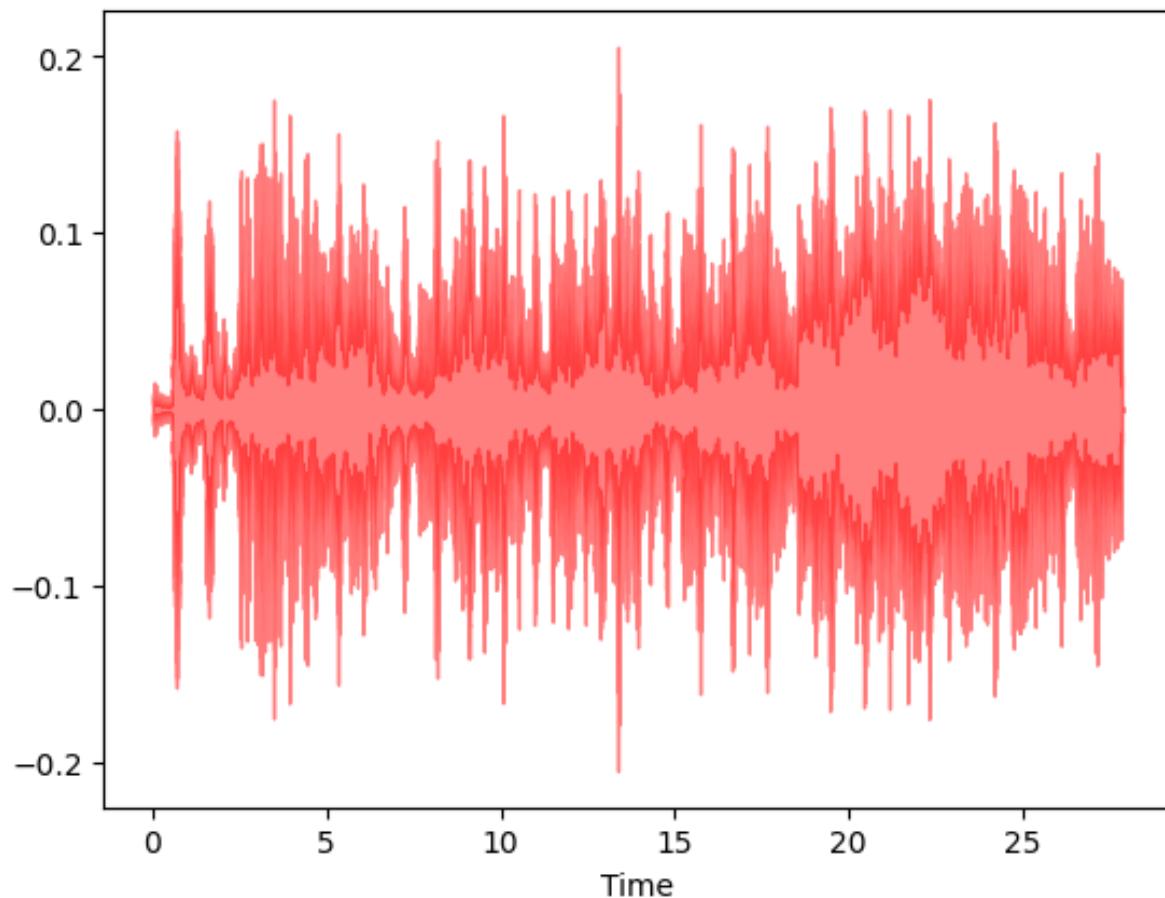


Figure 10: Waveform

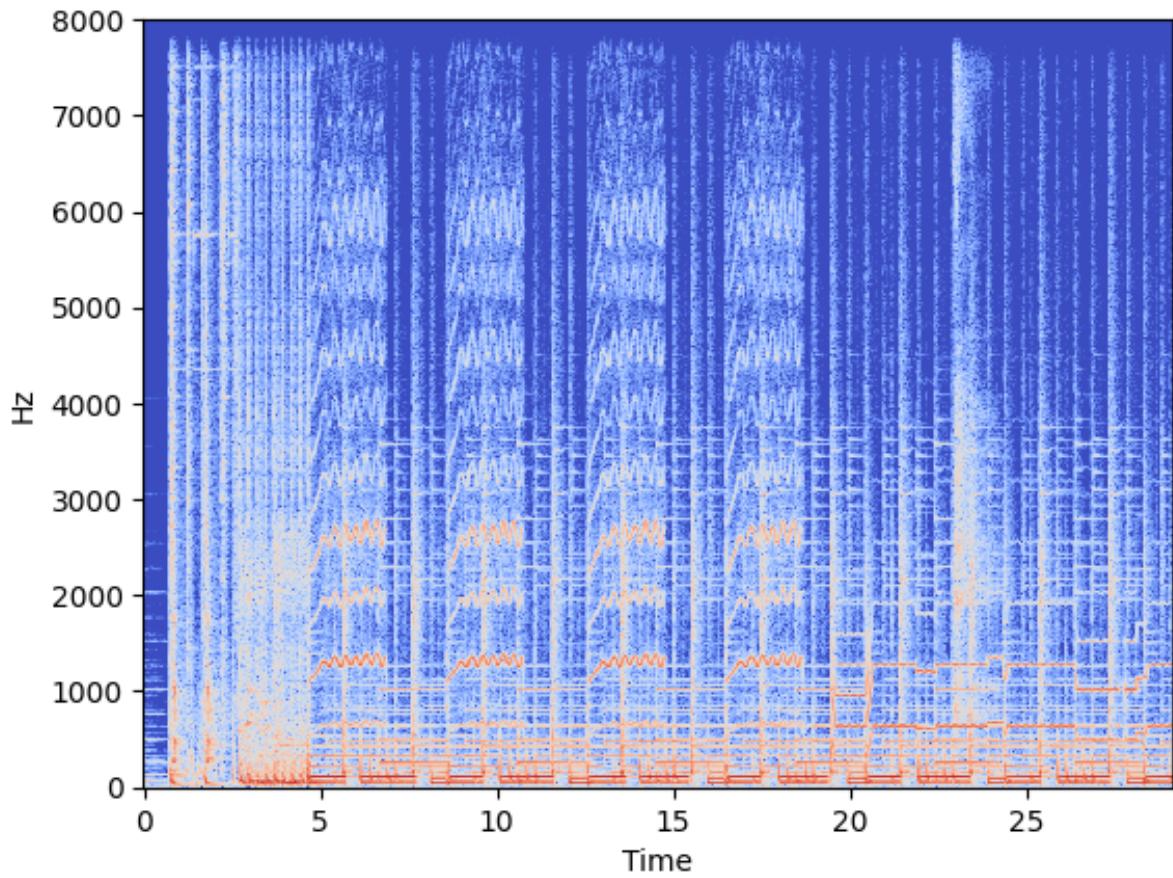


Figure 11: Spectrogram

2.1 Nlpaug

PitchAug

```

1 aug = naa.PitchAug(sampling_rate=sr, factor=(2,3))
2 augmented_data = np.array(aug.augment(data))
3
4 librosa_display.waveshow(data, sr=sr, color='b', alpha=0.5)
5 librosa_display.waveshow(augmented_data, sr=sr, color='g', alpha=0.25)
6 write_wav("jay-PitchAug.wav", sr, augmented_data)
7 plt.tight_layout()
8 plt.show()

```

Listing 9: Code section for PitchAug

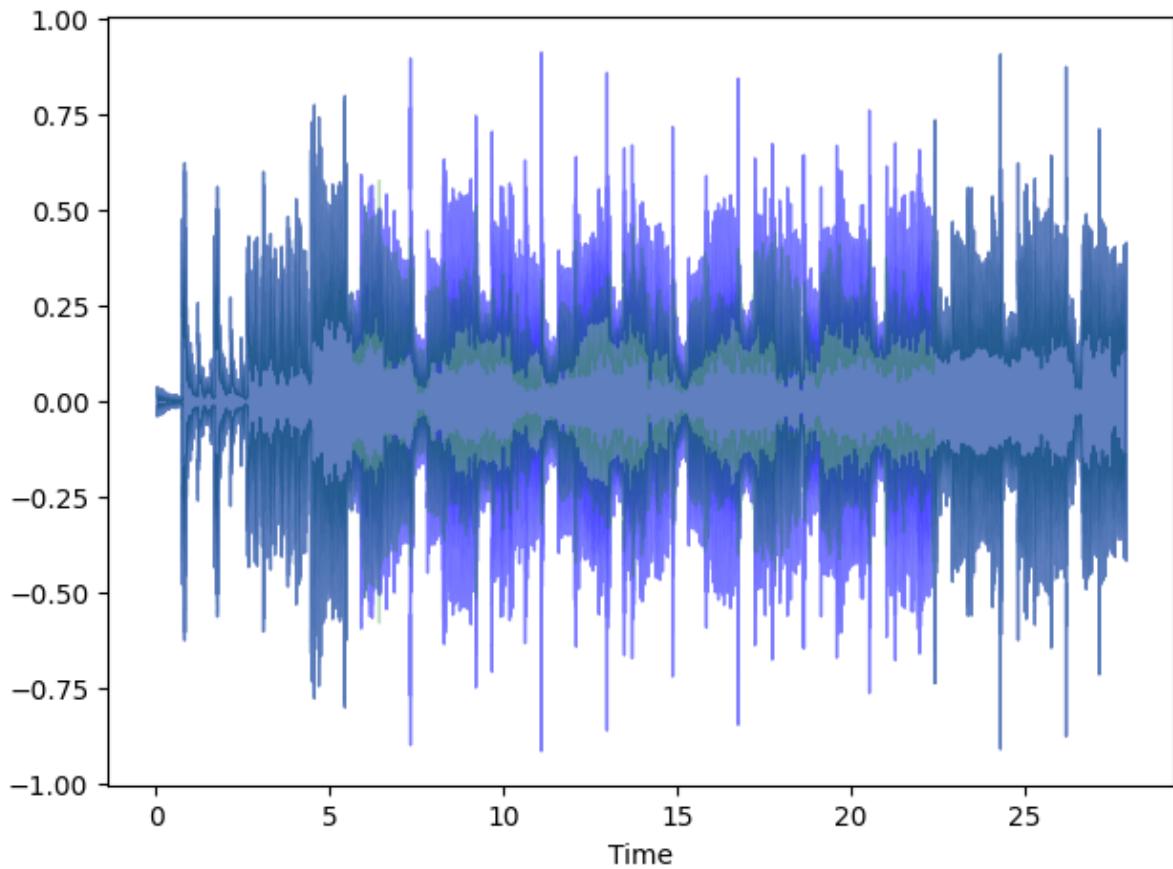


Figure 12: PitchAug

ShiftAug

```

1 aug = naa.ShiftAug(sampling_rate=sr)
2 augmented_data = np.array(aug.augment(data))
3
4 librosa_display.waveshow(data, sr=sr, color='b', alpha=0.5)
5 librosa_display.waveshow(augmented_data, sr=sr, color='g', alpha=0.25)
6 write_wav("jay-ShiftAug.wav", sr, augmented_data)
7 plt.tight_layout()
8 plt.show()

```

Listing 10: Code section for ShiftAug

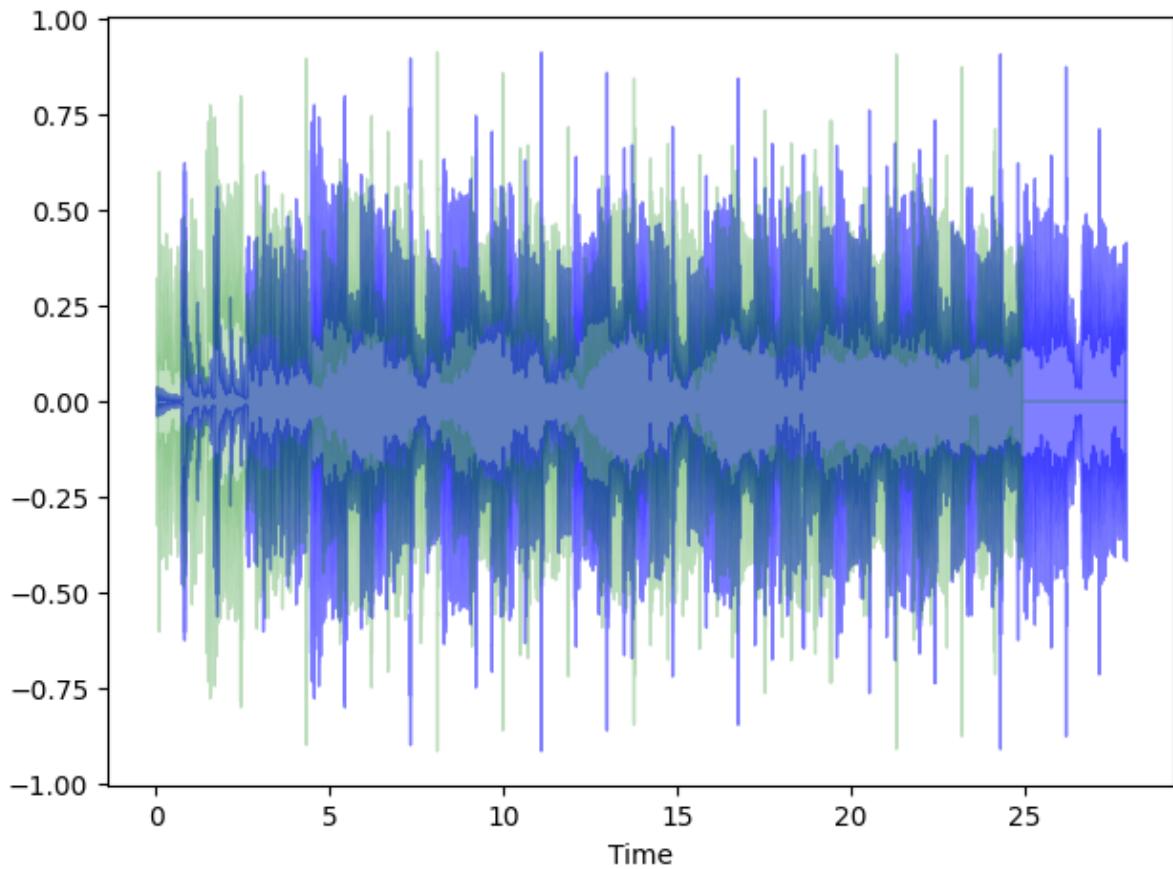


Figure 13: ShiftAug

SpeedAug

```

1 aug = naa.SpeedAug()
2 augmented_data = np.array(aug.augment(data))
3
4 librosa_display.waveshow(data, sr=sr, color='b', alpha=0.5)
5 librosa_display.waveshow(augmented_data, sr=sr, color='g', alpha=0.5)
6 write_wav("jay-SpeedAug.wav", sr, augmented_data)
7 plt.tight_layout()
8 plt.show()

```

Listing 11: Code section for SpeedAug

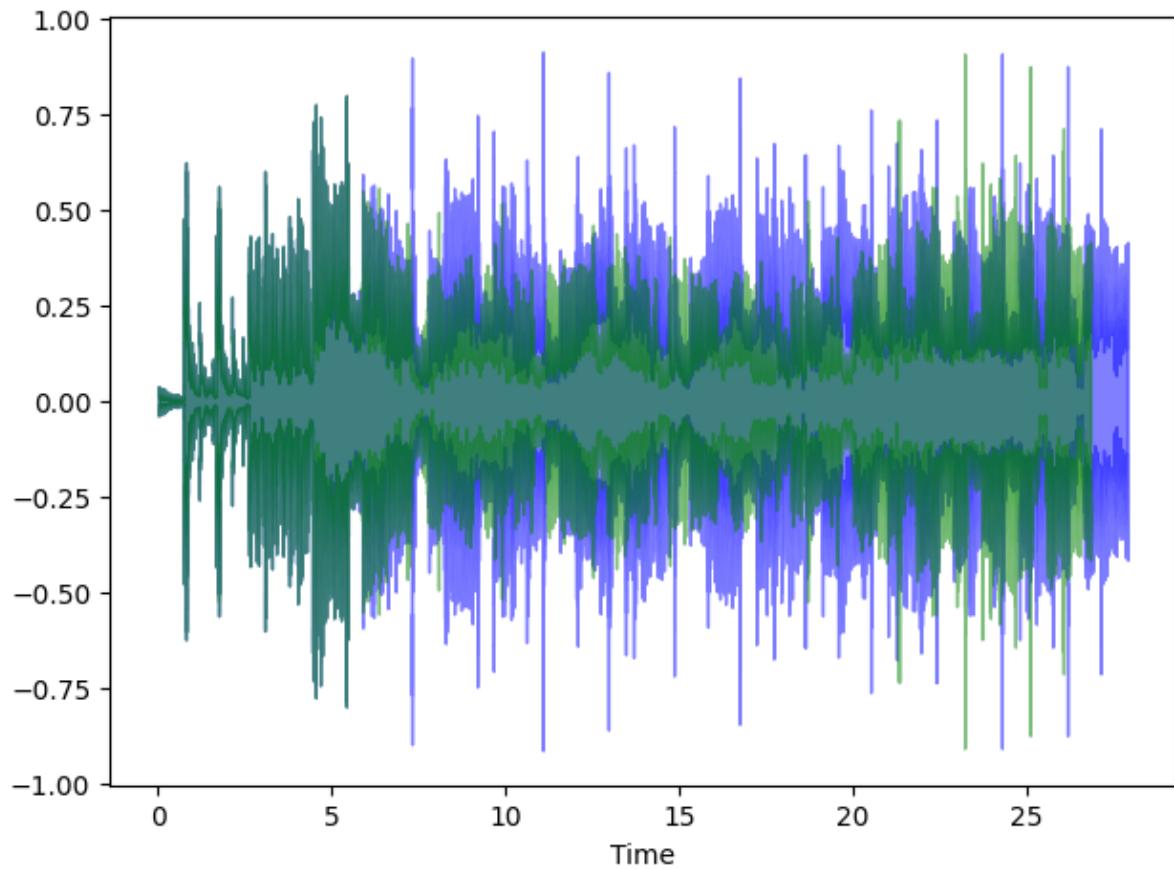


Figure 14: Rubber Band

VtlpAug

```

1 aug = naa.VtlpAug(sampling_rate=sr)
2 augmented_data = np.array(aug.augment(data))
3 write_wav("jay-VtlpAug.wav", sr, augmented_data)
4 AudioVisualizer.freq_power('VTLP Augmenter', data, sr, augmented_data)

```

Listing 12: Code section for VtlpAug

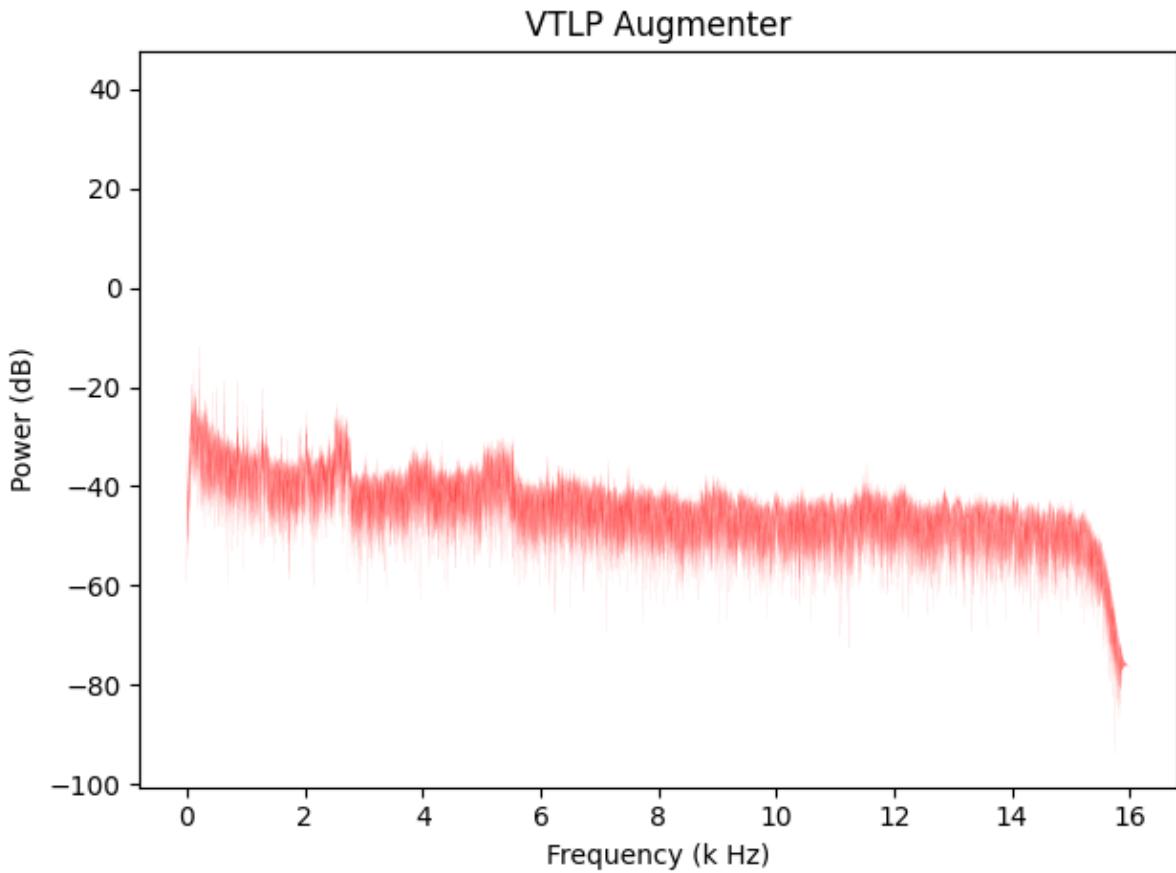


Figure 15: VtlpAug

2.2 Audiomentations

TimeMask

```

1 from audiomentations import TimeMask
2
3 transform = TimeMask(
4     min_band_part=0.1,
5     max_band_part=0.15,
6     fade=True,
7     p=1.0,
8 )
9
10 my_waveform_ndarray = np.array(data)
11
12 augmented_sound = transform(my_waveform_ndarray, sample_rate=16000)
13
14 librosa_display.waveform(augmented_sound, sr=sr, color='b', alpha=0.25)
15 write_wav("jay-TimeMask.wav", sr, augmented_sound)
16 plt.tight_layout()
17 plt.show()
```

Listing 13: Code section for TimeMask

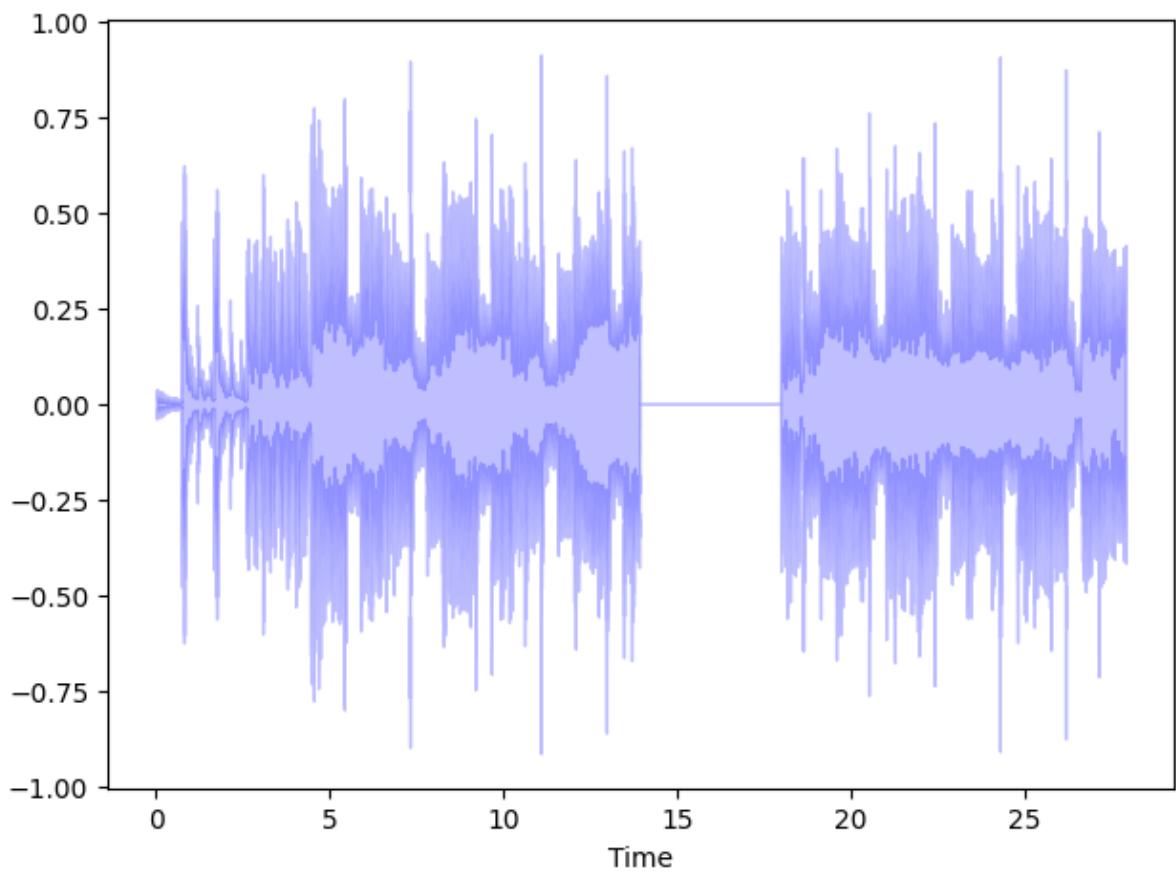


Figure 16: TimeMask

FrequencyMask

Here we use SpecFrequencyMask.

```

1 from audiomentations import SpecFrequencyMask
2 import numpy as np
3
4 transform = SpecFrequencyMask(p=0.5)
5
6 my_waveform_ndarray = np.array(data)
7
8 # Augment/transform/perturb the spectrogram
9 augmented_spectrogram = transform(my_waveform_ndarray)
10
11 librosa_display.waveshow(augmented_spectrogram, sr=sr, color='b', alpha
12 =0.25)
13 write_wav("jay-SpecFrequencyMask.wav", sr, augmented_spectrogram)
14 plt.tight_layout()
15 plt.show()
```

Listing 14: Code section for FrequencyMask

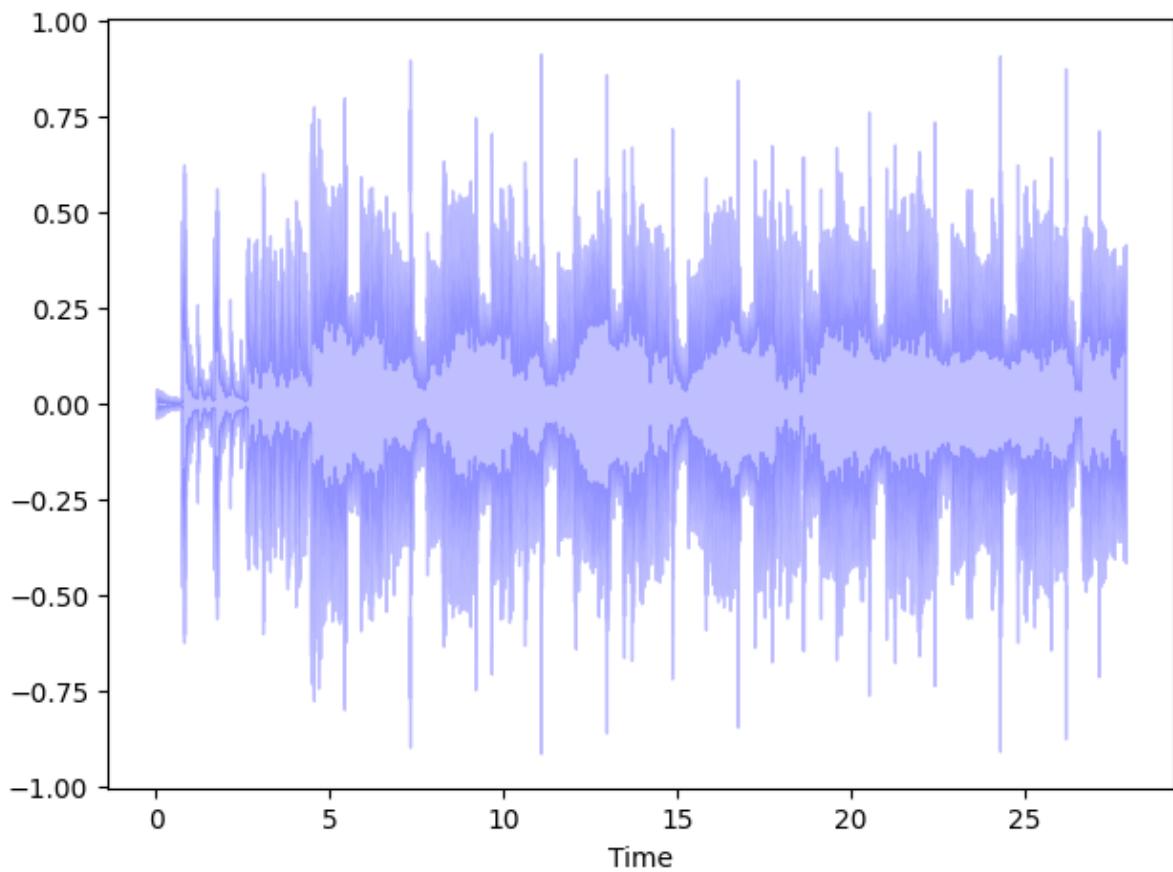


Figure 17: FrequencyMask

ApplyImpulseResponse

```

1 from audiomentations import ApplyImpulseResponse
2
3 my_waveform_ndarray = np.array(data)
4
5 transform = ApplyImpulseResponse(ir_path=path_2, p=1.0)
6
7 augmented_sound = transform(my_waveform_ndarray, sample_rate=48000)
8
9 librosa_display.waveshow(augmented_sound, sr=sr, color='b', alpha=0.25)
10 write_wav("jay-ApplyImpulseResponse.wav", sr, augmented_sound)
11 plt.tight_layout()
12 plt.show()
```

Listing 15: Code section for ApplyImpulseResponse

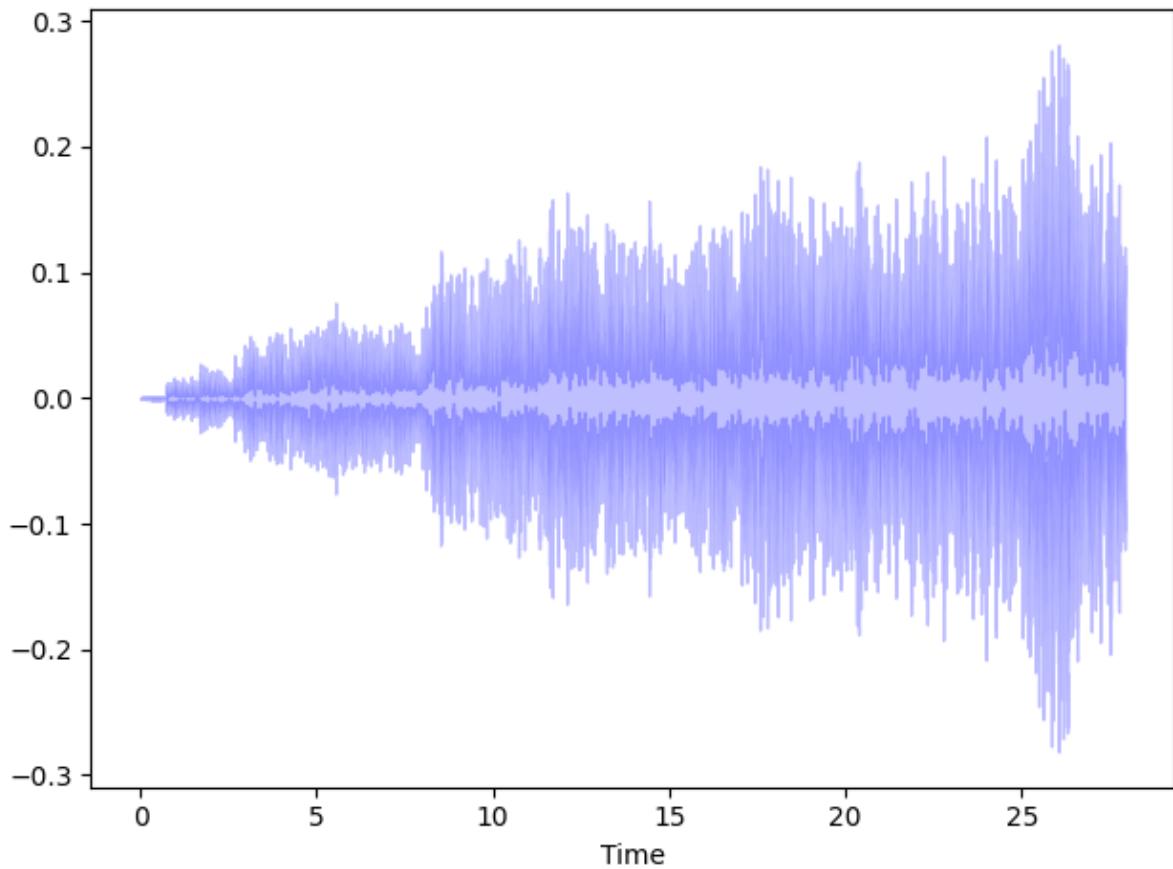


Figure 18: ApplyImpulseResponse

AddBackgroundNoise

```

1 from audiomentations import AddBackgroundNoise, PolarityInversion
2
3 my_waveform_ndarray = np.array(data)
4
5 transform = AddBackgroundNoise(
6     sounds_path=path_2,
7     min_snr_in_db=3.0,
8     max_snr_in_db=30.0,
9     noise_rms="relative",
10    noise_transform=PolarityInversion(),
11    p=0.4
12 )
13
14 augmented_sound = transform(my_waveform_ndarray, sample_rate=16000)
15
16 librosa_display.waveshow(augmented_sound, sr=sr, color='b', alpha=0.25)
17 write_wav("jay-AddBackgroundNoise.wav", sr, augmented_sound)
18 plt.tight_layout()
19 plt.show()

```

Listing 16: Code section for AddBackgroundNoise

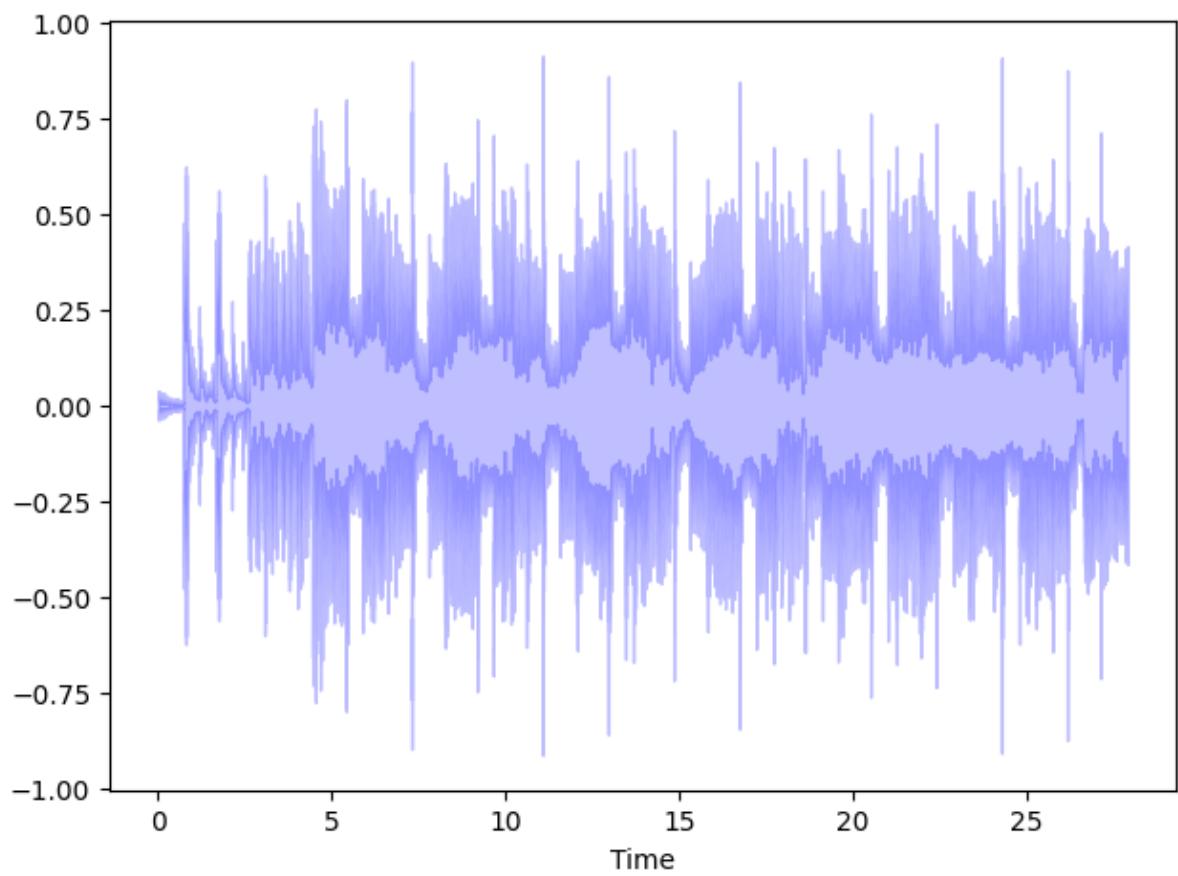


Figure 19: AddBackgroundNoise