

CMPE468 Speech Processing

For this coding assignment, it was asked to analyze two speech files (should.wav and s5.vaw) and plot the measurements speech waveform, short-time energy, short-time magnitude, and short-time zero-crossing.

For these calculations, I implemented a MATLAB code (SpeechProcessing.m) and ran the code to get the asked calculations as output. Here are the screenshots for MATLAB code and its outputs:

```

1  windowSize = 1024;
2  windowShift = 512;
3
4  should = 'should.wav';
5  s5 = 's5.wav';
6
7  analyze(should, windowSize, windowShift);
8  analyze(s5, windowSize, windowShift);
9
10
11 function analyze(filePath, windowSize, windowShift)
12     [speech, sampleRate] = audioread(filePath);
13
14     energy = rms(buffer(speech, windowSize, windowSize - windowShift, 'nodelay'));
15
16     magnitude = abs(spectrogram(speech, hann(windowSize), windowSize - windowShift));
17
18     zeroCross = sum(abs(diff(sign(buffer(speech, windowSize, windowSize - windowShift, 'nodelay'))))) / windowSize;
19
20     figure('Position', [100, 100, 800, 800]);
21
22     subplot(4, 1, 1);
23     plot(speech);
24     xlabel('Sample');
25     ylabel('Amplitude');
26     title('Entire Speech Waveform');
27
28     subplot(4, 1, 2);
29     time = (windowSize/2 : windowShift : windowSize/2 + windowShift * (size(energy, 2)-1)) / sampleRate;
30     plot(time, energy);
31     xlabel('Time (s)');
32     ylabel('Energy');
33     title('Short-time Energy (En)');
34

```

Figure 1: Screenshot of MATLAB Code (1)

```

34
35     subplot(4, 1, 3);
36     time = (windowSize/2 : windowShift : windowSize/2 + windowShift * (size(magnitude, 2)-1)) / sampleRate;
37     imagesc(time, [], 20*log10(magnitude));
38     set(gca, 'YDir', 'normal');
39     xlabel('Time (s)');
40     ylabel('Frequency');
41     title('Short-time Magnitude (Mn)');
42     colorbar();
43
44     subplot(4, 1, 4);
45     time = (windowSize/2 : windowShift : windowSize/2 + windowShift * (length(zeroCross)-1)) / sampleRate;
46     plot(time, zeroCross);
47     xlabel('Time (s)');
48     ylabel('Zero-crossing rate');
49     title('Short-time Zero-crossing (Zn)');
50
51     sgtile(sprintf('Speech Analysis: %s', filePath));
52 end
53

```

Figure 2: Screenshot of MATLAB Code (2)

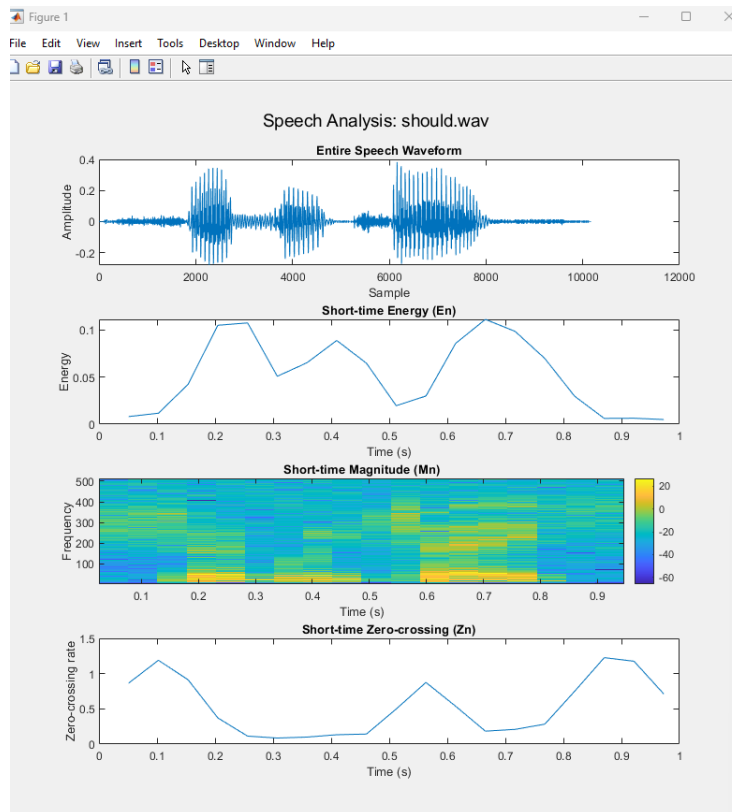


Figure 3: Calculations for `should.wav`

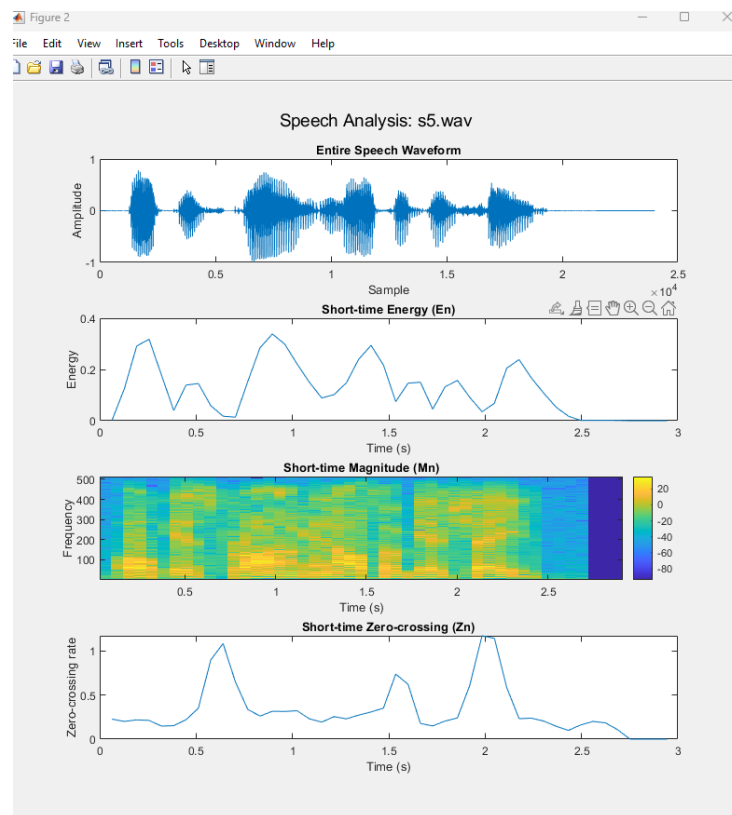


Figure 4: Calculations for `s5.wav`

The reasons behind chosen values of window sizes, window shifts, and window for the analysis can be explained as following:

1. Window Sizes: The window size defines how long each analytic window will be. The function uses a 1024-sample window size. The window size should be carefully selected to achieve a balance between temporal and frequency resolution. A wider window size provides better frequency resolution but worse time resolution, whereas a smaller window size provides better time resolution but lower frequency resolution. The 1024-sample size is a good compromise between time and frequency resolution for many speech analysis tasks.
2. Window Shifts: The window shift determines how many samples are shifted in successive analysis windows. The code uses a hop length of 512 samples. Choosing the proper hop length is essential for striking a balance between the amount of information acquired in the analysis and the processing performance. A shorter hop length offers finer temporal detail, but at the expense of additional windows and processing. While enhancing the temporal detail, a longer hop length reduces the cost of computing. Because it strikes a balance between maintaining computing efficiency and gathering enough temporal information, the 512 sample hop length was chosen.
3. Window: Each analysis window receives a window function to lessen spectral leakage and artifacts that may result from the window's abrupt start and end locations. A Hann window is used in the code (`hann(window_size)` used in calculation of magnitude, line 16 of the code). Due to its effective sidelobe attenuation and high frequency resolution, the Hann window is a common choice in speech analysis. It offers a reasonable middle ground between sidelobe suppression and main lobe breadth.