

An Insight into Audio Signal Processing for Musical Information Retrieval

Aditya Tatwawadi - Wren's



Westminster School
Computational Engineering Project
John Locke
August 2021

Foreword: My main goal of this John Locke has been to make an application with Python. Please reach out if you would like a live demonstration. Much of this report is supplementary as an overview of Audio Signal Processing. I also have attempted this with very little to no background in music.

[Consider accessing GitHub repository here \[1\]](#)

Contents

1	Introduction	2
2	The Parson's Code	3
2.1	Theory	3
2.2	Practical Implementation of Parson's Code	3
2.3	Brief Overview of Application	4
3	Scientific Methods	5
3.1	Physics behind ASP	5
3.2	Mathematics behind ASP	6
3.3	Practical Implementation of Fourier Transform	7
3.4	Database	7
4	Appendix	8
5	Postface	15
6	Bibliography	16

List of Python Code:

1	main.py	8
2	compute.py	9
3	webscrape.py	12
4	Fourier Transform Demonstration	13

An Insight into Audio Signal Processing

1 Introduction

Audio Signal Processing (ASP) has addressed many problems previously faced across various industries. Being able to identify songs was one such notoriously difficult issue. Companies such as *Shazam* have found ingenious methods to solve this and create a user-friendly product. Today, *Shazam* is valued at over \$400 million, and has been acquired by *Apple* [2]. Modern day applications of ASP include data-compression, speech processing, noise-cancellation and more. This is an industry to look out for.

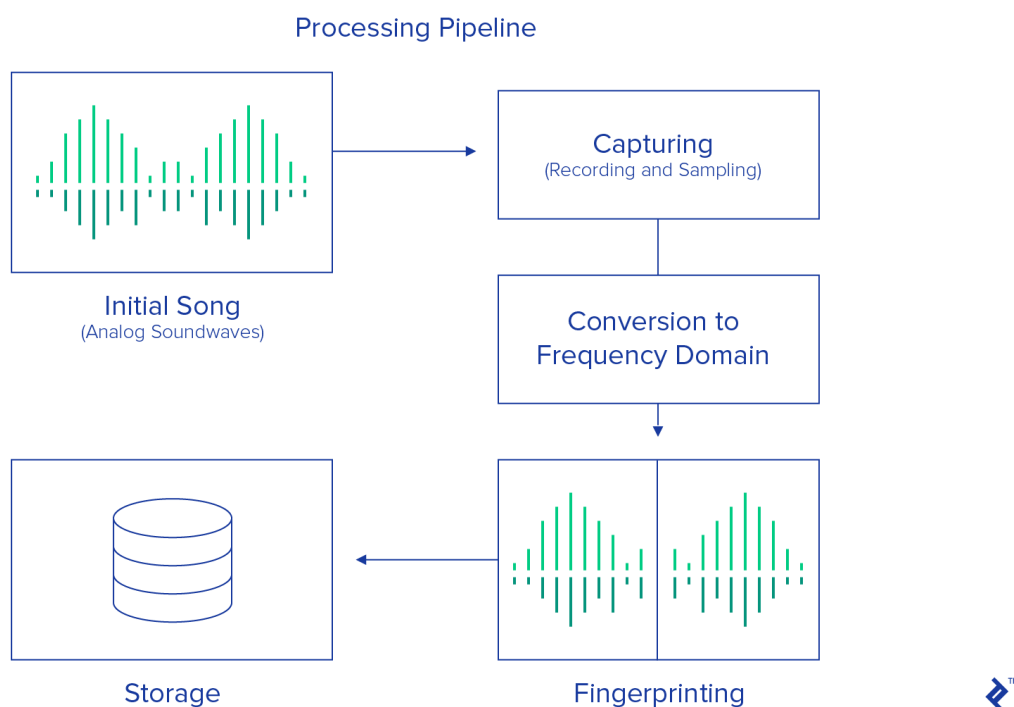


Figure 1: Shazam application timeline [3]

The theory and applications behind ASP have historically been championed at *Bell Labs* since the start of the 20th century after the introduction of telephones, radios and phonographs. Previously, analogue signals - which relied on continuous signals being "analogous" to the sound waves in air - needed to be physically altered by electric circuits. This made the process cumbersome and cost inefficient. Thanks to the advent of computational power, digital signals can now be used to express audioforms in binary numbers. This makes modern day signal processing more powerful, reliable and cheaper to implement [4].

In this report, the focus will be placed on music information retrieval, as well as the insights into the mathematics and physics to achieve this. Additionally, I have also built my own application to give song predictions using the *Python 3* programming language. This can be found referring to section 2.2 on page 3.

2 The Parson's Code

2.1 Theory

Developed by Denys Parsons in his 1975 book *The Directory of Tunes and Musical Themes*, the Parson's code is a surprisingly simple yet accurate way of identifying melodies; this makes it well-suited towards music information retrieval [5].

One of the key benefits of using the Parson's code is that no technical music knowledge is needed to implement it. It works through melodic motion: that is by identifying movements of pitch varying up and down. The first tone of reference is denoted with an asterix: *. If the following note is higher, a U (up) is used. If the following note is lower a D (down) is used. If the following note is repeated, a "R" (repeat) is used.

e.g. The Parson's code to display *Fur Elise* would be:

*DUDUDUDDDUUU.



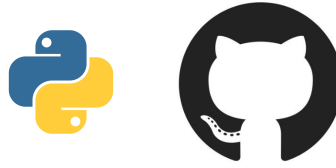
Figure 2: Parson's contour for "Für Elise" using my application

The Parson's code is only reliable in music retrieval if there is an identifiable, but not simultaneously varied pitch. Due to this, extraction is only feasible for classical/ folk/ distinct melodies, and not so much for hip-hop/ jazz songs or anything containing mixed timbres. Other limitations include difficulty in storing data, due to the Parson's code being relatively unknown and not well documented.

2.2 Practical Implementation of Parson's Code

On page 8 is the code of my project written to suggest to a user the top 3 most likely song predictions of a given MIDI audio file. See Listings 1,2,3 from Appendix. To access a cleaner version of my code, consider forking my repository on *Github* [1]. An explanation can be found on page 4.

2.3 Brief Overview of Application



There are 3 main files which run the code - `main.py`, `compute.py` and `webscrape.py`.

- The `main.py` file imports the code from the other two files, and produces an image from *Musipedia*: which has the information about the top 3 most likely song names and their composer.
- The `compute.py` file is an algorithm to extract the Parson's code from a MIDI file, produce a visual contour, and store it as a string to be accessed later.
- The `webscrape.py` file hits the internet and uses the *Musipedia* database to extract the songs corresponding to the identified Parson's Code.

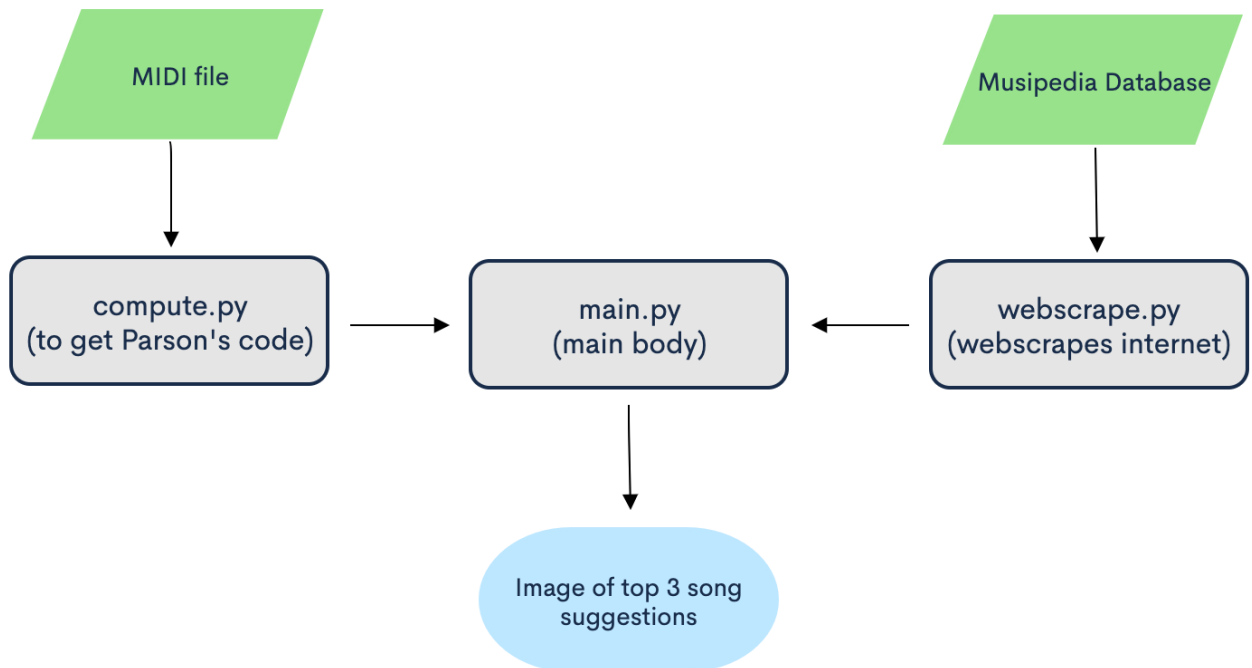


Figure 3: Flowchart of events

It is important to note that this is a somewhat abstraction of what is actually happening - the true computation lies in the *Mido* python library. The `MidiFile` object is able to read, write and play back any MIDI files. Although humans cannot hear .mid files like they would a .wav file, midi files are more accurate in transferring information between different places.

The next section will consider some of the mathematics and physics behind how exactly an audio file is decomposed, and an insight into how more sophisticated methods are used to generalize audio extraction to any type of music.

3 Scientific Methods

3.1 Physics behind ASP

What is sound? Sound can be described as a vibration of pressure that propagates as longitudinal waves. After hitting the eardrum, the vibration is then transmitted to hair cells in the cochlea, which produce electrical signals to the brain.

This process of converting to an electrical signal from the pressure of a sound wave (continuous signals) is imitated by recording devices. In a microphone, the Operational Amplifier (Op-Amp) translates this into an analog voltage signal. For the continuous signal to be of use, it must be translated into a discrete signal that can be stored digitally. This is done by capturing a digital value that represents the amplitude of the signal. The quantization of the input means that an analog-to-digital converter performs many conversions on very small pieces of the signal is a process known as sampling [6].

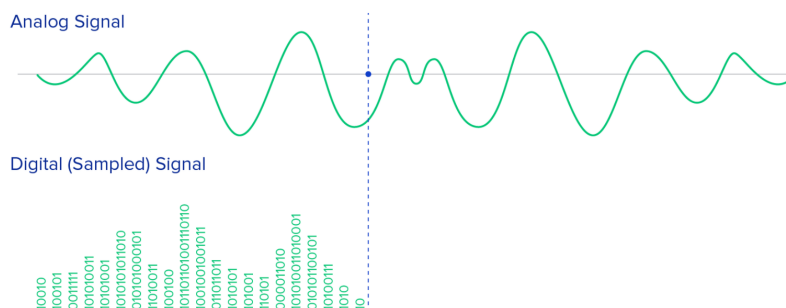


Figure 4: Change of signals [3]

Every sound in the world can be represented as the sum of multiples of pure tones in different amplitudes. As tones are expressed in sinusoidal forms, every sound in the world can therefore be expressed in terms of sinusoidal waves. Music is played by various instruments and singers, producing a combination of sinewaves at varying frequencies, creating an overall more obscure sine wave. This can be visualized using a Spectrogram: a 3D graph. Time on the horizontal X axis. Frequency on the vertical Y axis. A colour representing the amplitude of a defined frequency (in decibels, dB) as the 3rd dimension. Spectrograms can be computed using Fourier Transforms [7].

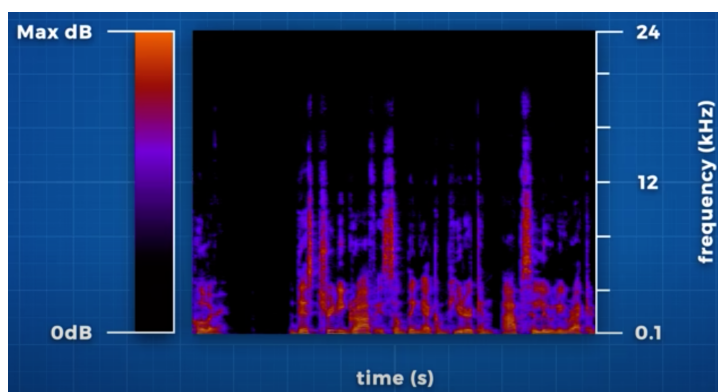


Figure 5: Spectrogram [8]

3.2 Mathematics behind ASP

Sound is made up of a lot of complex wavefunctions. One analogy which could be used to represent the essence of ASP is like blending up a smoothie, and then trying to reverse-engineering the process to get back its original ingredients.

Though of course impossible in real life, in the 19th century Jean-Baptiste Joseph Fourier made a trailblazing advancement that the time domain can be represented as the sum of multiple sinusoidal signals (as aforementioned on page 5). This is known as the Fourier Series.

The Fourier Series is used to represent a periodic function by a discrete sum of complex exponentials, while the Fourier transform is then used to represent a general, non-periodic function by an integral of complex exponentials. The Fourier Transform can be viewed as the limit of the Fourier series of a function as the period approaches infinity, so the limits of integration are $(-\infty, \infty)$. Instead of trying to crudely fit a signal to model a wave as $A \sin(2\pi(ft + \phi))$, any complex wavefunction can be represented as the sum of sines and cosines [9]. It is important to note that this is one of the most elegant, yet conceptually difficult mathematical topics to understand, and so cannot be given full justice in just one section of a report. We are concerned with Discrete FT in ASP.

The Fourier Transform (FT) applies to continuous time giving a continuous spectrum:

$$\hat{\mathcal{F}}(\xi) = \int_{-\infty}^{+\infty} f(t) e^{-2\pi i \xi t} dt = \int_{-\infty}^{+\infty} f(x) e^{-i\omega t} dt$$

where $\omega \in (-\infty, +\infty)$

The Discrete Fourier Transform (DFT) applies to finite time giving a discrete spectrum:

$$X(k) = \sum_{n=0}^{N-1} x(n) e^{-i(2\pi kn/N)}$$

where $n = 0, 1, \dots, N - 1$

N is the number of samples that composed the signal. $X(K)$ represents the kth bit of frequencies. $x(n)$ is the nth sample of the audio sample.

The use of the complex number $i = \sqrt{-1}$ is purely graphical: the algebra tends to be nicer with geometrically related problems to do with winding or rotation. This can be visualized on the Argand Diagram; instead of x,y axis on the Cartesian Plane, there are real and imaginary axis. As the Fourier Transform pertains to circular paths, Euler's Formula is a suitable way to generate one:

$$e^{i\theta} = \cos(\theta) + i \sin(\theta)$$

In industry, applications like *Shazam* need to optimise processing time to compute a spectrogram, and produce a match to the recorded audio, all within the constraints of the phone's processor. For this reason, the Fast Fourier Transform (FFT) - which serves the same function as the DFT, but only much faster, - is used. The simplest and most efficient version of the FFT (also used in the numpy package in the implementation of my code in page 7) is the Cooley-Tukey algorithm. The FFT has a time-complexity of $O(n\log(n))$ which is far superior to DFT's $O(n^2)$ [10]:

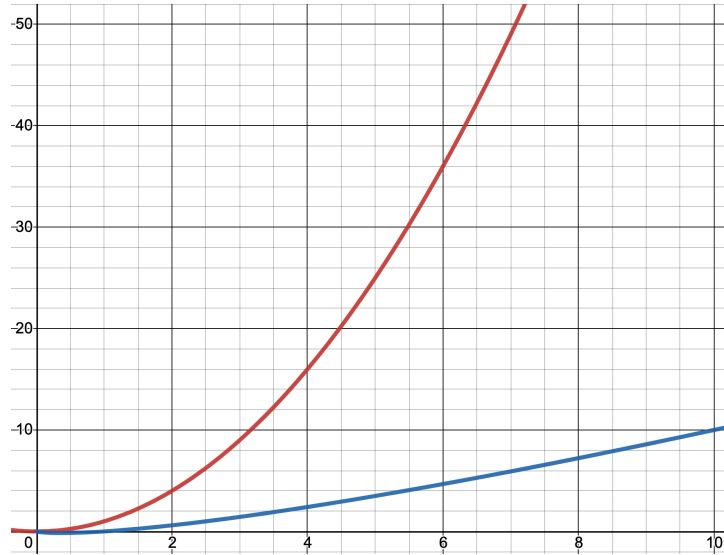


Figure 6: DFT - Red - $O(n^2)$ — FFT - Blue - $O(n\log(n))$

Such is the importance of the FFT, that the Institute of Electrical & Electronics Engineers (IEEE) deemed it as "the most ubiquitous algorithm in use to analyze digital data", and has been heralded as the top 10 most influential algorithms of the century by countless publishers and universities [11].

3.3 Practical Implementation of Fourier Transform

Python Code which performs an FFT can be located in the Appendix under Listing 4 on page 13. This takes a time domain signal, applies a FFT and converts it into a frequency domain signal. Like my previous application, this can be found on my GitHub repository.

3.4 Database

As a sidenote, it is intriguing to consider how *Shazam* stores this data. For each individual song, after the FFT is applied and a spectrogram is computed, a *fingerprint* - i.e. a starmap is used to reduce the data from 3D to 2D. Hash functions and tables are then used in searching algorithms. This goes beyond the scope of this report, but more detail can be located in this article by the founder of Shazam - Avery Wang [12].

4 Appendix

A compilation of all code. Everything hosted on GitHub [1].

Listing 1: main.py

```
1  """
2  Aditya Tatwawadi - Summer 2021 - Westminster School - John
   Locke
3
4  Application to take a song in a .mid file, convert it into
   Parson's Code
5  and display a photo of possible song suggestions
6
7  n.b.
8  - Generally only successful for Classical music with a single
   melody line in isolation (Though does work for most songs
   on musipedia).
9  - .mid file should be very accurate (ideally taken from
   musipedia).
10
11  Some test songs & their corresponding Parson's codes:
12
13  Fur Elise                / DUDUDUDDDUUU          / Giscard Rasquin
   & Ludwig van Beethoven
14  Eine kleine Nachtmusik / DUDUDUUDDUDUDDUD / Wolfgang Amadeus
   Mozart
15  Jana Gana Mana          / UURRRRR              / Rabindranath
   Tagore
16  """
17
18  #Import general dependency
19  import time
20
21  #Import Pillow to show screenshot (will need to be done with
   HomeBrew)
22  from PIL import Image
23
24
25  def show_suggestions():
26      music_suggestions = Image.open("musipedia_recommendations
   .png")
27      music_suggestions.show()
28
29  if __name__ == "__main__":
30
31      #Import file written to compute Parson's Code
32      import compute
33
```

```

34     #Delay needed to read txt file
35     time.sleep(1)
36
37     #Import file written to scrape Musipedia
38     import webscrape
39     time.sleep(1)
40     show_suggestions()

```

Listing 2: compute.py

```

1  #Import dependency to read a .mid file. RECOMMENDED - mido
   documentation
2  import mido
3
4  #Enter path of file
5  path_name_of_mid_file = "//Users/adityatatwawadi/Desktop/
   Projects/Python/JohnLocke2021/Midi Files/FurElise.mid"
6
7  #Compute stuff
8  class ComputeParsons:
9
10     def __init__(self):
11         pass
12
13     def midi_to_parsons(midifile, limit=13, offset=0):
14
15         """
16         Takes a midi file path, outputs a string
17         """
18
19         #Info to be added to empty variable string parsons
20         global parsons
21         count = 0
22         parsons = ""
23
24         #Reads midifile - a mid file is a series of 8-bit
           bytes
25         for message in mido.MidiFile(midifile):
26
27             if "note_on" in str(message) and offset == 0:
28
29                 if parsons == "":
30                     # initialise list
31                     prev = message.note
32                     parsons += "*"
33
34                 #comparison to determine u d or r (As per
           Parson's Code)

```

```

35         elif message.note > prev:
36             parsons += "u"
37             prev = message.note
38         elif message.note < prev:
39             parsons += "d"
40             prev = message.note
41         elif message.note == prev:
42             parsons += "r"
43             prev = message.note
44
45         #increment count
46         count += 1
47         if count >= limit:
48             break
49
50         elif "note_on" in str(message):
51             offset -= 1
52
53     return(f"The Parson's code for the file you have
54           inputted with the path {path_name_of_mid_file} is
55           {parsons}")
56
57 def contour(code):
58     """
59     Creates a contor image using the Parson's code
60     Output: UI display
61     """
62
63     #Not needed - only for individual file testing
64     if code[0] != "*":
65         raise ValueError("Parsons Code must start with
66                           '*'")
67
68     #Initialize an empty dictionary
69     contour_dict = {}
70     pitch = 0
71     index = 0
72
73     max_pitch = 0
74     min_pitch = 0
75
76     contour_dict[(pitch, index)] = "*"
77
78     for character in code:
79         if character == "r":
80             index += 1
81             contour_dict[(pitch, index)] = "-"

```

```

80
81         index += 1
82         contour_dict[(pitch, index)] = "*"
83     elif character == "u":
84         index += 1
85         pitch -= 1
86         contour_dict[(pitch, index)] = "/"
87
88         index += 1
89         pitch -= 1
90         contour_dict[(pitch, index)] = "*"
91
92         if pitch < max_pitch:
93             max_pitch = pitch
94     elif character == "d":
95         index += 1
96         pitch += 1
97         contour_dict[(pitch, index)] = "\\\"
98
99         index += 1
100        pitch += 1
101        contour_dict[(pitch, index)] = "*"
102
103        if pitch > min_pitch:
104            min_pitch = pitch
105
106    for pitch in range(max_pitch, min_pitch+1):
107        line = [" " for _ in range(index + 1)]
108        for pos in range(index + 1):
109            if (pitch, pos) in contour_dict:
110                line[pos] = contour_dict[(pitch, pos)]
111
112        print("".join(line))
113
114    def save_parsons():
115
116        """
117        Saves Parson's code to a txt file so webscrape.py can
118        access it
119        """
120
121        #Writes code to a txt file
122        with open("parson.txt", 'w') as f:
123            f.write(parsons)
124
125    print(ComputeParsons.midi_to_parsons(path_name_of_mid_file))
126    ComputeParsons.contour(parsons)

```

```
127 ComputeParsons.save_parsons()
```

Listing 3: webscrape.py

```
1  #Import all dependencies to webscrape
2  from selenium import webdriver
3  from selenium.webdriver.common.by import By
4
5  #Scrape stuff
6  class WebScrape():
7
8      def __init__(self):
9          pass
10
11     def read_parsons():
12
13         """
14         Reads in Parson's code written by compute.py to use
15         as variable
16         """
17
18         #So that variable can be accessed elsewhere in the
19         file
20         global parsons_code
21         with open("parson.txt", "r") as parsons_file:
22             parsons_code = parsons_file.read()
23
24     def scrape_musipedia():
25
26         """
27         - Uses Selenium & ChromeDriver to extract information
28           about top 3 most likely songs
29         - Requires scraping html
30         - Relies on Musipedia not being updated
31         """
32
33         #Path details
34         path_of_driver = "/Users/adityatatwawadi/Downloads/"
35         chromedriver"
36         driver = webdriver.Chrome(executable_path =
37             path_of_driver)
38
39         #Access musipedia's website
40         website_url = "https://www.musipedia.org/"
41         melodic_contour.html"
42         driver.get(website_url)
43         #print(driver.title)
```

```

39      #Enters in the Parson's code to search engine
40      driver.find_element(By.NAME, 'tx_mpsearch_pi1[pc]').
41          send_keys(parsons_code)
42
43      #Submits & accesses database
44      click_button = driver.find_element(By.NAME, '
45          tx_mpsearch_pi1[submit_button]')
46      click_button.click()
47
48      #Scrolls down to top 3 identified songs
49      driver.execute_script("window.scrollTo(0,400)")
50
51      #Reframes & takes a screenshot of top3
52      recommendtaions
53      S = lambda X: driver.execute_script('return document.
54          body.parentNode.scrollTo'+X)
55      driver.set_window_size(S('Width'), S('Height'))
56      driver.find_element_by_tag_name('body').screenshot("
57          musipedia_recommendations.png")
58
59      #Exits automated chrome
60      driver.quit()
61
62      WebScrape.read_parsons()
63      WebScrape.scrape_musipedia()

```

Listing 4: Fourier Transform Demonstration

```

1  """
2  Aditya Tatwawadi - Westminster School
3  Computes a Fast Fourier Transform using Numpy - Turns Time
4  domain signal into Frequency domain signal
5  Plots the frequency spectrum using matplotlib
6  """
7
8  #Import dependencies
9  import numpy as np
10 import matplotlib.pyplot as plt
11
12 def compute():
13     #Construct a time signal from sine wave:
14
15     #Sampling Frequency (Hz):
16     frequency_sample = 2000
17     #Sample time interval:
18     tstep = 1 / frequency_sample
19     #Signal Frequency:

```

```

19     signal_frequency = 100
20
21     #Number of samples:
22     N = int(10 * frequency_sample / signal_frequency)
23
24     #Time steps:
25     t = np.linspace(0, (N-1)*tstep, N)
26     #Frequency interval:
27     fstep = frequency_sample / N
28     #Frequency steps:
29     f = np.linspace(0, (N-1)*fstep, N)
30
31     #Create a chaotic wave - new one added is 2 times the
       original input frequency and has a magnitude of 4
32     y = 1 * np.sin(2 * np.pi * signal_frequency * t) + 3 * np
       .sin(2 * np.pi * 2 * signal_frequency * t)
33
34     #Perform FFT:
35     X = np.fft.fft(y)
36     X_mag = np.abs(X) / N
37
38     f_plot = f[0:int(N/2+1)]
39     X_mag_plot = 2 * X_mag[0:int(N/2+1)]
40     X_mag_plot[0] = X_mag_plot[0] / 2
41
42     #Plot
43     fig, [ax1, ax2] = plt.subplots(nrows=2, ncols=1)
44     ax1.plot(t, y, "-.")
45     ax2.plot(f_plot, X_mag_plot, "-.")
46     ax1.set_xlabel("Time (s)")
47     ax2.set_xlabel("Frequency (Hz)")
48     ax1.grid()
49     ax2.grid()
50
51     ax1.set_xlim(0, t[-1])
52     ax2.set_xlim(0, f_plot[-1])
53     plt.tight_layout()
54     plt.show()
55 compute()

```

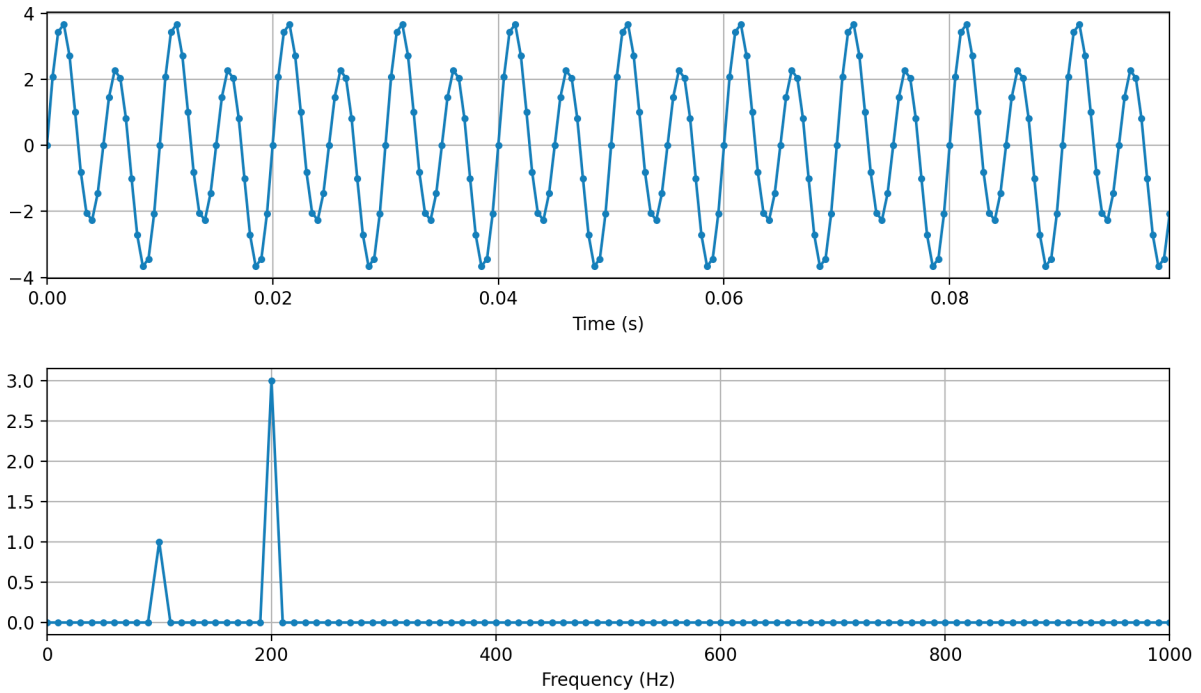


Figure 7: Matplotlib graph of FT taking Time to Frequency domain [1]

5 Postface

This project was incredibly useful in exposing me to the mathematics behind many first year EEE/ EIE degrees. Along the way I gained a lot of practical technical information in python, such as Object-Oriented Programming, file-handling and proficiency with multiple libraries. This let me perform a wide range of objectives from data analysis, web-scraping to using logic to come to a solution. With more time, I would have looked to either Tkinter or web-development to be able to host my application on a more User-Friendly level (and hope to come back to this later). I also learnt how to host software on GitHub, intricacies on VSCode IDE, and also LaTeX. Audio Signal Processing is a much broader topic than what hasa been covered, and I look forward to returning to it.

6 Bibliography

- [1] Aditya Tatwawadi. *John Locke - Audio Signal Processing GitHub Repo*, 2021. <https://github.com/aditya-tatwawadi/Audio-Signal-Processing>.
- [2] BBC News. *Apple to buy Shazam 2017*. <https://www.bbc.co.uk/news/business-42299207>. Accessed August 2021.
- [3] Jovan Jovanovic. *How Shazam Works*. <https://www.toptal.com/algorithms/shazam-it-music-processing-fingerprinting-and-recognition>. Accessed August 2021.
- [4] Udo Zölzer. “Digital Audio Signal Processing”. In: (1997). ISBN 0-471-97226-6.
- [5] S Brown Denys Parsons. *The Directory of Tunes and Musical Themes*. ISBN 0-904747-00-X. 1975.
- [6] Jonathan Allday Steven Adams. *General Classical Physics Theories. Ideas adapted from Advanced Physics*. ISBN 0-199-14680-2. 2000.
- [7] Allen B. Downey. *ThinkDSP*. 2014.
- [8] Real Engineering Youtube. *How Shazam Works*. <https://www.youtube.com/watch?v=kMNSAhsyiDg>.
- [9] 3Blue1Brown. *But what is the Fourier Transform? A visual introduction*. <https://www.youtube.com/watch?v=spUNpyF58BY>.
- [10] John W. Tukey James W. Cooley. “An Algorithm for the Machine Calculation of Complex Fourier Series.” In: (1965).
- [11] F. Sullivan J. Dongarra. “Computing in Science and Engineering.” In: (2000).
- [12] Avery Li-Chun Wang. “An Industrial-Strength Audio Search Algorithm .” In: (2003).