



Project Title: Heartbeat Sound Segmentation and Classification

Group Number: 19G11 Supervisor Name: Ellen De Mello Koch

Student Name A: Elias Sepuru Student Name B: Boikanyo Radiokana

Student Number A: 1427726 Student number B: 1386807

Ethics: ☐ Request for waiver (does not involve human participants or sensitive data)

☐ Copy of ethics application attached (Non-medical) – School Committee

Supervisor Signature ☒ Copy of ethics application attached (Medical) – University Committee

Project Outline: *(give a brief outline such that ethics reviewers understand what will be done, 100 words maximum)*

This project aims to create the first level screening of detecting signs of heart diseases in patients. This will aid medical practitioners in their field and possibly home use by patients. A method to locate lub (S1) and dub (S2) sounds in audio data of patients' heartbeats will be implemented. After location the heartbeat audio data will be segmented based on S1 and S2. Followed by segmentation, a method to classify a heartbeat into normal and diseased categories will be implemented.

## Project Specification:

### 1. Project Objectives:

#### 1.1 Primary Objectives:

- To successfully locate S1 and S2 and segment the audio heartbeat data into the categories S1 and S2.
- To train a machine learning model that is going to classify a heartbeat sound into either diseased or normal

**Success criteria:** Current existing solutions have an accuracy of less than or equal to 79%, this project aims to obtain an accuracy of 79% or higher.

#### 1.2 Secondary Objectives:

- If time allows, the project aims to create a user-interface, for home use or use by medical practitioners.

### 2. Project Work Breakdown Structure:

#### 2.1 Heartbeat Segmentation

For location and segmentation of the heartbeat sounds into S1 and S2, the following high-level methodology represented by illustration 1 will be used.

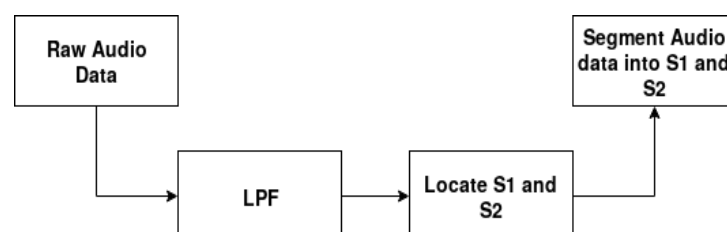


Figure 1: High level methodology to be followed for heartbeat audio segmentation

- The audio data is first filtered using a Low Pass Filter (LPF) to remove high frequency noise components. The prospective filters that might be used are the Discrete Waveform Transform Filter (DWT), Daubechies filter or any other wavelet filter.
- To locate and segment the heartbeat audio data into S1 and S2, the fact that the time from S1-S2 is shorter than the time from S2-S1 will be exploited.

## 2.2 Heartbeat Classification

- For heartbeat classification, methods for finding feature importance of the segmented heartbeat audio signal are going to be used.
- After finding the feature importance, numerous machine learning models are going to be trained and tested on the features. The model with the highest accuracy is to be selected.

### Milestones:

1. Literature review [0.5 Weeks]
2. Data Cleansing [2.5 weeks]
3. Model Training & Testing [2 weeks]
4. Documentation [1 week]

(If time permits the user interface will be built on the last week of Model Training and Testing)

### Preliminary Budget & Resources:

- MATLAB and Jupyter Notebook for signal processing and building models.
- KNIME and Anaconda building models.
- Two compatible computers for data cleansing, training and testing of the models.
- Internet for research.
- Digital Stethoscope.
- Funds for printing of technical report and the poster for open day.

### Risks / Mitigation:

- The two laptops that will be used for testing and training the data might not be able to process the big data which can result in the two machines crashing. This risk will be mitigated by using the machines provided by Professor Otto for running big data simulations.
- Failure to remove all the noise present in the data might compromise the success of the project and reduce the desired accuracy. This will be mitigated by using alternative secondary open source data from MIT. The data from MIT has consistent length and reduced/no noise.
- Mislocation of S1 and S2 sounds in the audio data will lead to misinterpretation, incorrect analysis and unacceptable prediction results. This will be mitigated by using multiple methods of locating S1 and S2 to check for any errors and correspondence of the results.