

## Introduction

With the rapid growth of computational power and complex algorithms, we propose a novel approach to detect arrhythmias in Phonocardiograms (PCGs). Typically, Electrocardiograms are used to diagnose arrhythmias, requiring medical-grade equipment to recognize cardiac illnesses accurately. However, PCGs provide ease of access to everyone who has a device capable of recording audio, allowing medical professionals to treat arrhythmias in the developmental stages. The new design comprises two subsystems; one is based on the relationship between Electrocardiograms (ECGs) and PCGs, and the other between PCGs and arrhythmias. The association between ECGs and PCGs is amended to translate from one space to another, where ECGs become dimensionally reduced, then reconstructed into a PCG signal. The second subsystem uses a Generative Adversarial Networks (GAN), in which both arbitrary PCG signals and generated signals are fed into a discriminator that detects if an arrhythmia is present or if the signal is false.

## Background

An estimated three million cases of arrhythmia occur in the United States yearly (Mayo Clinic), with 300,000 sudden deaths per year – an incidence rather higher than stroke, lung cancer, or breast cancer (American Heart Association). Traditionally, non-invasive arrhythmia analysis is based on multiple electrodes that reflect the electrical activity on ECGs. This method, despite being accurate, limits the use case to hospitals and clinics with specialized equipment; thus, limiting the portability of diagnosing, let alone classification of the type of pathology.

Phonocardiograms (PCGs) are sounds that are created by the mechanical movement of the heart. This physical movement produces four distinct sounds: S1, S2, S3, S4, and murmurs. S1 and S2 are sounds created by a healthy heart; whereas, S3, S4, and murmurs refer to diseases or anomalies.

The first heart sound, S1, marks the start of Systole. Systole occurs when the heart muscle contracts and pumps blood from the chambers into the arteries. The second heart sound, S2, marks the end of Systole and the start of Diastole. Diastole is a phase of the heartbeat when the heart muscle relaxes and allows the chambers to fill with blood.

Figure 1 illustrates the representation of heart sound recording positions. Although heart sound databases do exist, these datasets are still limited by the number of pathologies that are collected, often having to divide the dataset into two categories: normal and abnormal. Currently, only three major PCG datasets exist: PhysioNet Classification of Heart Sound Recording Challenge dataset, PASCAL Heart Sound Challenge dataset, and the Heart Sound and Murmur Library. These datasets are all anonymized and de-identified for the safety of their subjects, and thus includes no personal information such as name, income, age, etc.

The PhysioNet Classification of Heart Sound Recording Challenge dataset was produced as a part of the 2016 PhysioNet Computing in Cardiology Challenge. The heart sounds were collected from both clinical and non-clinical environments (in-home visits). The challenge focused on creating an accurate dataset of normal and abnormal heart sound recordings, especially in real-world (extremely noisy and low signal quality) scenarios. These recordings were sourced from nine independent databases and in total, contain 4,593 heart sound recordings from 1072 subjects, lasting from 5-120 seconds. Of which, 409 recordings that were collected from 121 patients contain one PCG lead and one simultaneously recorded ECG. Though, all recordings are resampled to 2,000 Hz using an anti-alias filter. Furthermore, the dataset is comprised of 3 classes: normal, abnormal, and unsure (this is due to poor recording quality), and have the following proportion respectively: 77.1%, 12.0%, 10.9%.

The PASCAL Classifying Heart Sounds Challenge dataset was released to the general public in 2011. The challenge consisted of two sub-challenges: heart sound segmentation, and heart sounds classification; these sub-challenges corresponded with dataset A, and dataset B respectively. Both datasets have recordings of varying lengths, between 1 second and 30 seconds. Dataset A was collected via the iSethoscope Pro iPhone app, and contained 176 heart sound recordings, 124 of which are divided into four classes: Normal (31 recordings), Murmur (34 recordings), Extra heart sound (19 recordings), and Artifact (40 recordings); the rest of the records are unlabeled for testing purposes. Dataset B was collected using a DigiScope (a digital stethoscope), and included 656 heart sounds. All except 370 were separated into three classes: Normal (320 recordings), Murmur (95 recordings), and Extra-systole (46 recordings). Both datasets A and B vary in sound recordings between lengths of 1 second and 30 seconds.

More than 300 million ECG recordings are analyzed yearly, and thus create an exceptional tool for arrhythmia classification. Coupled with the recent surge in research interest in 2015, many massive publicly available datasets have been published, notable by PhysioNet – the moniker of the Research Resource for Complex Physiologic Signals. Numerous datasets ECG exist, however, many are limited to few classes (Normal and Abnormal). At present, three public datasets exist that have more than 4 classes: AF Classification Challenge 2017, PTB Diagnostic ECG, and PTB-XL dataset. Additionally, iRhythm Technologies have developed a semi-public dataset, that is available upon request, that contains 12 classes.

The main challenge of ineffective heart sound detection stems from an analysis of noisy heartbeats, e.g., background noise. For clean datasets, e.g., the PhysioNet Challenge dataset, a varieties time and frequency of methods converged on localization accuracy of 96.9% (Fernando et al.) and 96.0% classification accuracy (Mostafa et al.). For large datasets with noisy signals, e.g., the PASCAL Challenge dataset, the performance of time and frequency methods remained inconsistent at a localization accuracy of 93.3% (Singh et al.) and 93.3% classification accuracy (Singh et al.).

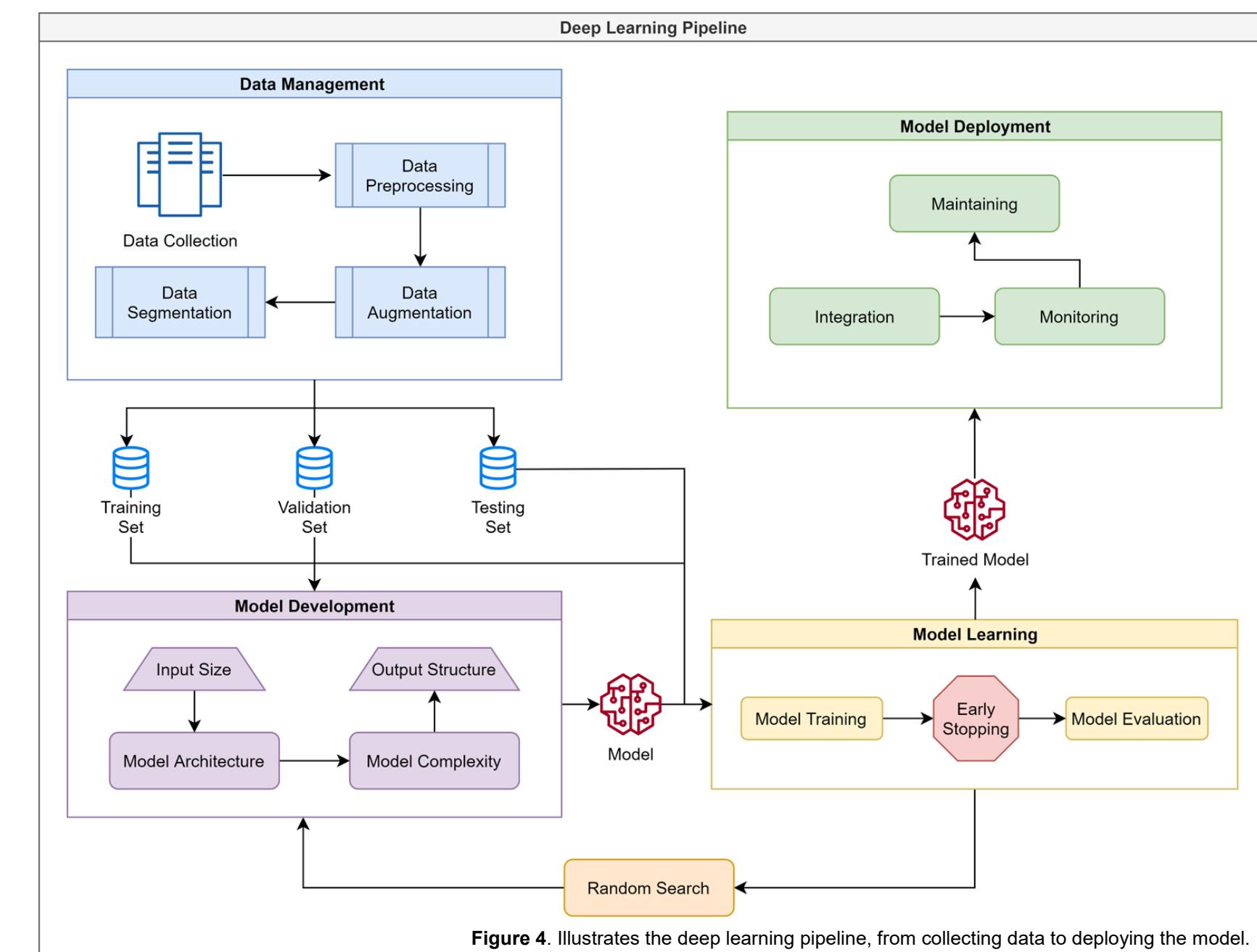
From the viewpoint of practical applications, the development of computationally efficient solutions is extremely important to the success of a model's deployment. Many studies have negated to comment on the practicality of their proposed methods. From our research, we have concluded only two studies have noted their time efficiency, (Fernando et al.) and (Messner et al.). The fastest model processed 1000 heart state classifications in 56.89 seconds (Fernando et al.), suggesting the model can process 18 bps. The average heart rate of a human heart is around 60–100 bps; thus, current models need severe optimization to achieve near to real-time analysis. These results are excluding the classification of heart arrhythmias.

Thus, the problem of computationally efficient and accurate classification of noisy heartbeats, especially with datasets with a variety of pathologies still remains a problem.

## Engineering Goals

- 1) Develop a System for End-to-End Heart Sound Arrhythmia Detection—Create a system that is able to record and analyze heart sounds for Cardiovascular modalities. The system should implement an adversarial model that is both time and space-efficient, and accurate
- 2) Increase the Number of Cardiovascular Pathologies—Develop a model to construct heart sounds from pre-existing data.
- 3) Real-World Testing — Test the end-to-end system in a real-world environment to ensure practicality and generality of the system.

## Methodology



# PCG-Net & VQGAN: Generative Adversarial Networks for PCG Arrhythmia Detection

## Data Management

### Data Collection

Although PCG signals are analyzed less often than ECG signals, these signals are rather analyzed in real-time by physicians and healthcare workers. Preliminary studies were done on PCG segmentation and classification primarily used private datasets. Hence, there existed no publicly available datasets until recently. Since then, many public datasets have been developed aiding researchers in their studies and creating open benchmarks for researchers to use in comparing similar findings. However, these datasets are still limited by the number of classes that are collected, when compared to ECG datasets.

Currently, only three major supervised PCG datasets exist: PhysioNet Classification of Heart Sound Recording Challenge dataset, PASCAL Heart Sound Challenge dataset, and the Heart Sound and Murmur Library. These datasets are all anonymized and de-identified for the safety of their subjects, and thus includes no personal information such as name, income, age, etc.

The PhysioNet Classification of Heart Sound Recording Challenge dataset was produced as a part of the 2016 PhysioNet Computing in Cardiology Challenge.

The heart sounds were collected from both clinical and non-clinical environments (in-home visits). The challenge focused on creating an accurate dataset of normal and abnormal heart sound recordings, especially in real-world (extremely noisy and low signal quality) scenarios. These recordings were sourced from nine independent databases and in total, contain 4,593 heart sound recordings from 1072 subjects, lasting from 5-120 seconds. Of which, 409 recordings that were collected from 121 patients contain one PCG lead and one simultaneously recorded ECG. Though, all recordings are resampled to 2,000 Hz using an anti-alias filter.

Furthermore, we use noise injection directly to preprocessed PCG recordings (Messner et al.). This process is identical to the process of synthetically adding noise to PCG recordings described in the preprocessing step. A variety of noises, like white noise, is added to the signal to increase the sample of recordings per class. This method is extremely beneficial for training on small datasets, like the PASCAL dataset.

We use to resample the heart sound recordings to different frequencies to simulate slower and faster beats per minute (bpm). The normal bpm for a human is between 60-100 bpm. Thus, measuring the sample distance between the first S1 (the start of systole) and the second S1, we calculate the bps and resample accordingly.

Figure 5 shows plots of the classes of the PhysioNet 2016 heart sound dataset (Normal and Abnormal).

The PASCAL Classifying Heart Sounds Challenge dataset was released to the general public in 2011. The challenge consisted of two sub-challenges: heart sound segmentation, and heart sounds classification; these sub-challenges corresponded with dataset A, and dataset B respectively. Both datasets have recordings of varying lengths, between 1 second and 30 seconds. Dataset A was collected via the iSethoscope Pro iPhone app, and contained 176 heart sound recordings, 124 of which are divided into four classes: Normal (31 recordings), Murmur (34 recordings), Extra heart sound (19 recordings), and Artifact (40 recordings); the rest of the records are unlabeled for testing purposes. Dataset B was collected using a DigiScope (a digital stethoscope), and included 656 heart sounds. All except 370 were separated into three classes: Normal (320 recordings), Murmur (95 recordings), and Extra-systole (46 recordings). Both datasets A and B vary in sound recordings between lengths of 1 second and 30 seconds.

Figure 6 shows plots of the classes present in the PASCAL dataset (Normal, Abnormal, and Artifact).

The PTB-XL is the largest publicly available dataset for ECGs and contains 21,837 clinical 12-lead ECG recordings from 18,885 patients of 10 second length. These recordings are separated into 5 super-classes: Normal, Myocardial Infarction, Hypertrophy, ST/T-Change, and Conduction Disturbance. These super-classes are further split into 71 sub-classes that range from AV Block to Posterior Myocardial Infarction. The raw signal data were downsampled to 100 Hz and annotated by up to two cardiologists, who assigned potentially multiple ECG statements to each record.

iRhythm Technologies developed a large, 12 classes ECG dataset using raw single-lead ECG inputs. The 12 classes include Atrial fibrillation and flutter, AVB, Bigeminy, EAF, IFR, Junctional rhythm, Noise, Sinus rhythm, SVT, Trigeminy, Ventricular tachycardia, and Wenckebach. The dataset consists of 91,232 ECG recordings from 53,549 patients. This training dataset is available upon request under license from iRhythm Technologies, Inc. The publicly available test dataset contains 328 records collected from 328 unique patients, split between 6 classes. Both datasets A and B vary in sound recordings between lengths of 1 second and 30 seconds.

More than 300 million ECG recordings are analyzed yearly, and thus create an exceptional tool for arrhythmia classification. Coupled with the recent surge in research interest in 2015, many massive publicly available datasets have been published, notable by PhysioNet – the moniker of the Research Resource for Complex Physiologic Signals. Numerous datasets ECG exist, however, many are limited to few classes (Normal and Abnormal). At present, three public datasets exist that have more than 4 classes: AF Classification Challenge 2017, PTB Diagnostic ECG, and PTB-XL dataset.

Additionally, iRhythm Technologies have developed a semi-public dataset, that is available upon request, that contains 12 classes.

The main challenge of ineffective heart sound detection stems from an analysis of noisy heartbeats, e.g., background noise. For clean datasets, e.g., the PhysioNet Challenge dataset, a varieties time and frequency of methods converged on localization accuracy of 96.9% (Fernando et al.) and 96.0% classification accuracy (Mostafa et al.).

For large datasets with noisy signals, e.g., the PASCAL Challenge dataset, the performance of time and frequency methods remained inconsistent at a localization accuracy of 93.3% (Singh et al.) and 93.3% classification accuracy (Singh et al.).

From the viewpoint of practical applications, the development of computationally efficient solutions is extremely important to the success of a model's deployment. Many studies have negated to comment on the practicality of their proposed methods. From our research, we have concluded only two studies have noted their time efficiency, (Fernando et al.) and (Messner et al.). The fastest model processed 1000 heart state classifications in 56.89 seconds (Fernando et al.), suggesting the model can process 18 bps. The average heart rate of a human heart is around 60–100 bps; thus, current models need severe optimization to achieve near to real-time analysis. These results are excluding the classification of heart arrhythmias.

Thus, the problem of computationally efficient and accurate classification of noisy heartbeats, especially with datasets with a variety of pathologies still remains a problem.

## Model Deployment

Model Deployment is one of the last stages of any machine learning project and involves releasing the model to the public.

### Integration

Integration consists of implementing the model in a system, whether it happens on the client-side or the backend. The most popular backend model integration tools involve Flask, Azure, and FastAPI. These tools create APIs that encapsulate the model prediction, given a GET request with the desired input.

### Monitoring & Maintaining

Following model integration and deployment, we move onto the next phase, monitoring and maintaining the system. As more and more data passes through the model, it increases the opportunity for the model to learn from a more generalized dataset. Though such data would be unsupervised, we could use unsupervised techniques to categorize the data. Based on the improvement of the model, the model is retrained and deployed. In essence, looping the whole process from data management to model learning.

## PCG-Net: Dense Generator

The input consists of a 2-dimensional tensor (batch size x signal) with a length of 128. This input is randomly generated.

The Dense layers upscale the length of the input vector until the length reaches that of a real input signal (2500). This allows for the output of the generator to be directly fed into the discriminator.

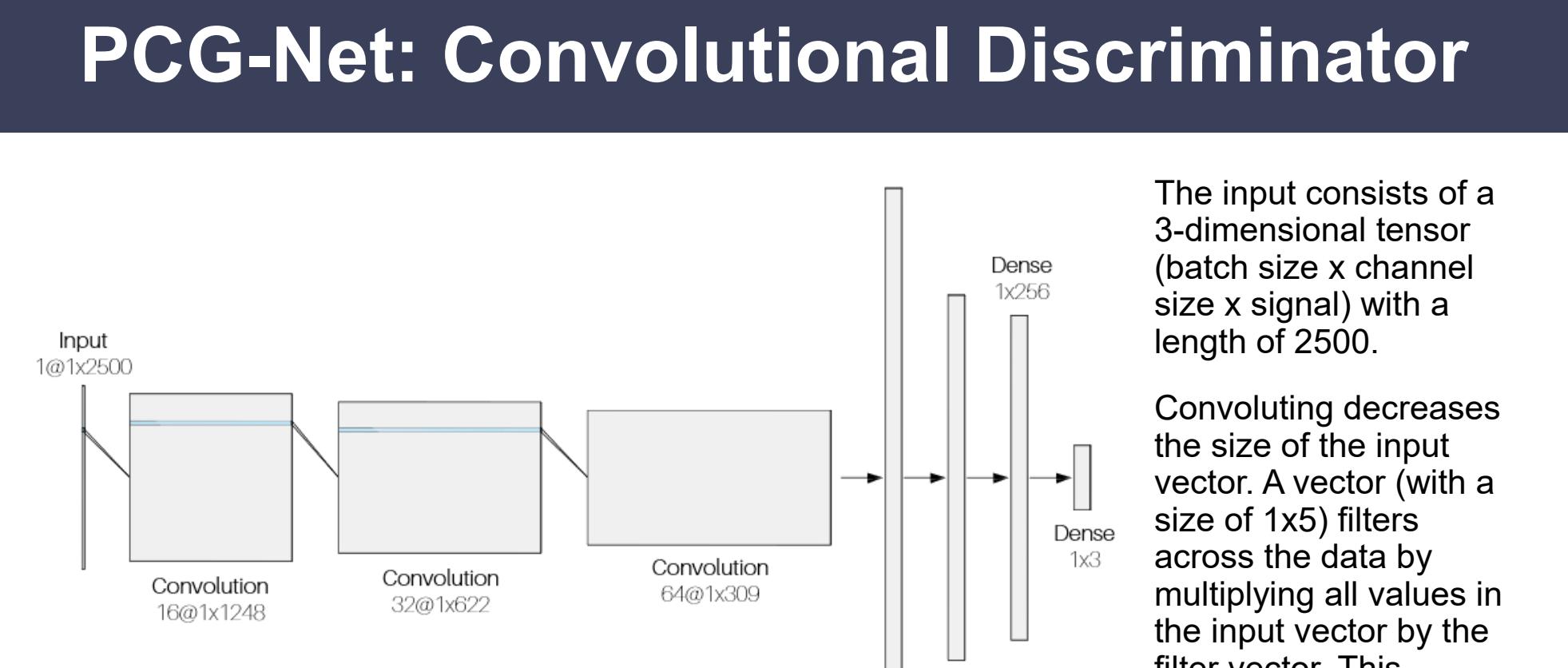
The Rectified Linear activation function alters the range of the incoming data by setting all numbers below 0 to 0 and leaving all positive numbers intact.

Figure 7: Representation of PCG-Net dense generator structure.

Generated Noisy Data

The graph to the right shows the input values (seeds) and corresponding generated noise. The variance in amplitudes and periods in the signal mimics that of real heartbeats changing. At a closer examination of the results, the generator appears to resample the seed values to that of the output (2500). This is expected as the generator consists of dense layers, which linearly increase the signal length and complexity.

Figure 8: Plot of the noise (artifact class) generated from the generator.



The Rectified Linear activation function alters the range of the incoming data by setting all numbers below 0 to 0 and leaving all positive numbers intact.

The liner function flattens the incoming result (batch size x 64 x 309) into a 2D tensor (batch size x 1976).

The Dense layers downscale the length of the input vector until the length reaches that of the number of classes. This allows for the output of the discriminator to be directly interpreted by the cost function.

The Linear Output layer transforms the Linear layer output into a 1x3. Each column in the tensor represents a class's likelihood of being the correct class in the dataset. Thus, the column with the largest values is the model prediction for the input.

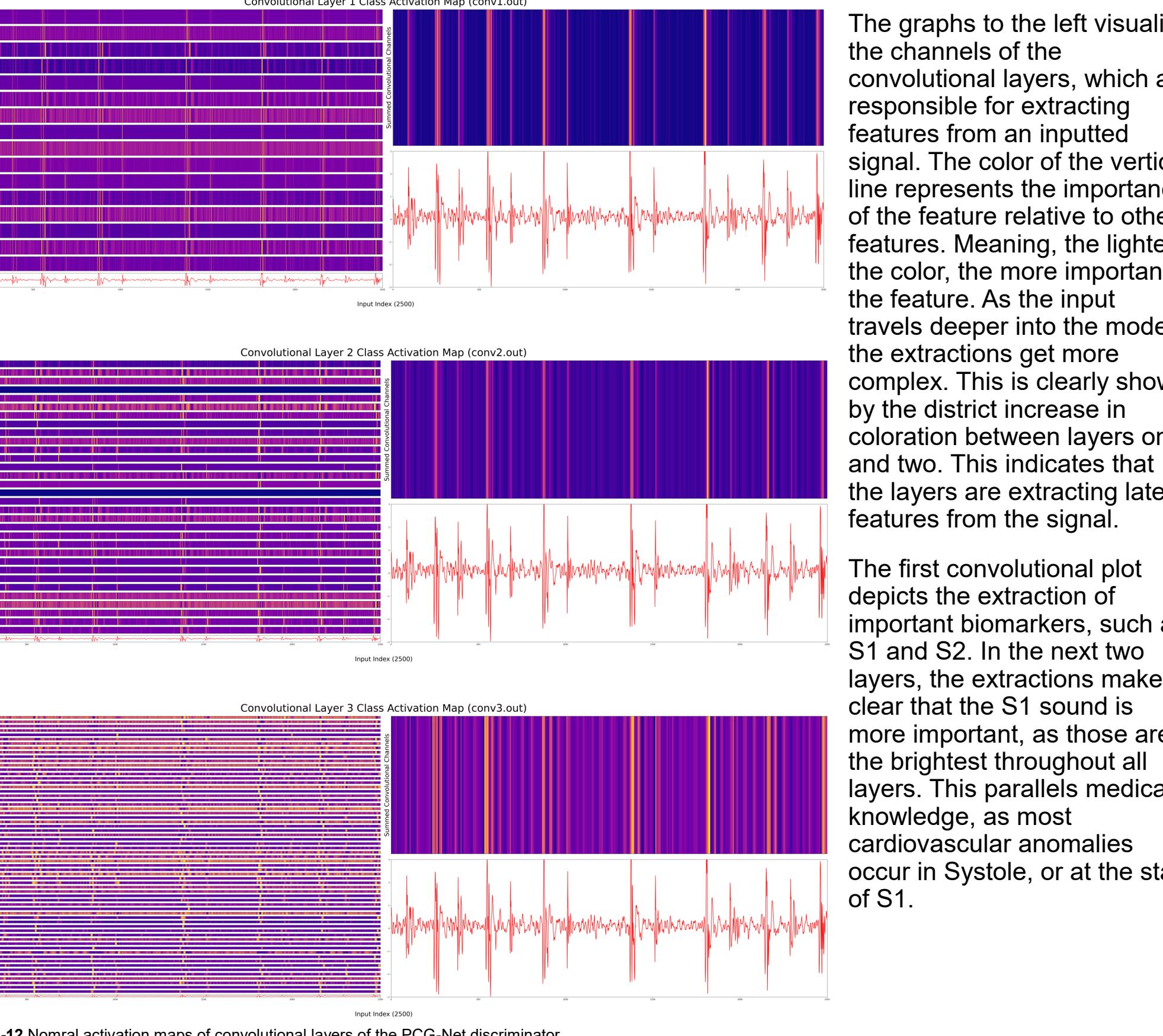
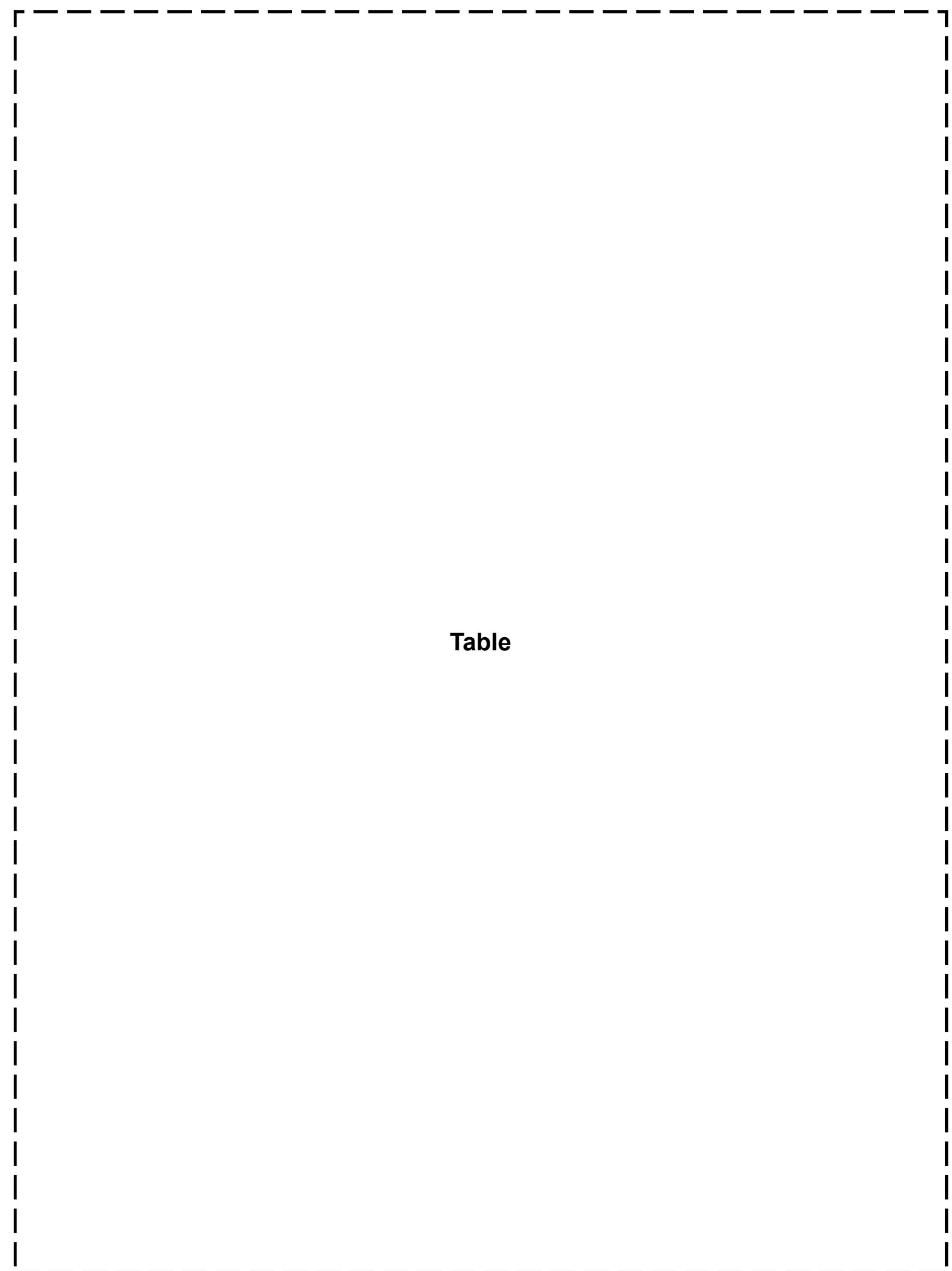


Figure 10-12 Nominal activation maps of convolutional layers of the PCG-Net discriminator.

## PCG-Net: Raw Results



Table

Table 1 The table shows the accuracy, sensitivity, and specificity of 150 trials of training the GAN model on the dataset. The metrics displayed are the results of the testing set.

**Accuracy:** percent of correctly identified normal and abnormal heart sounds.

**Sensitivity:** percent of correctly identified abnormal heart sound recording from a sample of abnormal heart signals.

**Specificity:** percent of correctly identified normal heart sounds recordings from a sample of normal heart signals.

## PCG-Net: Testing Sample Distributions

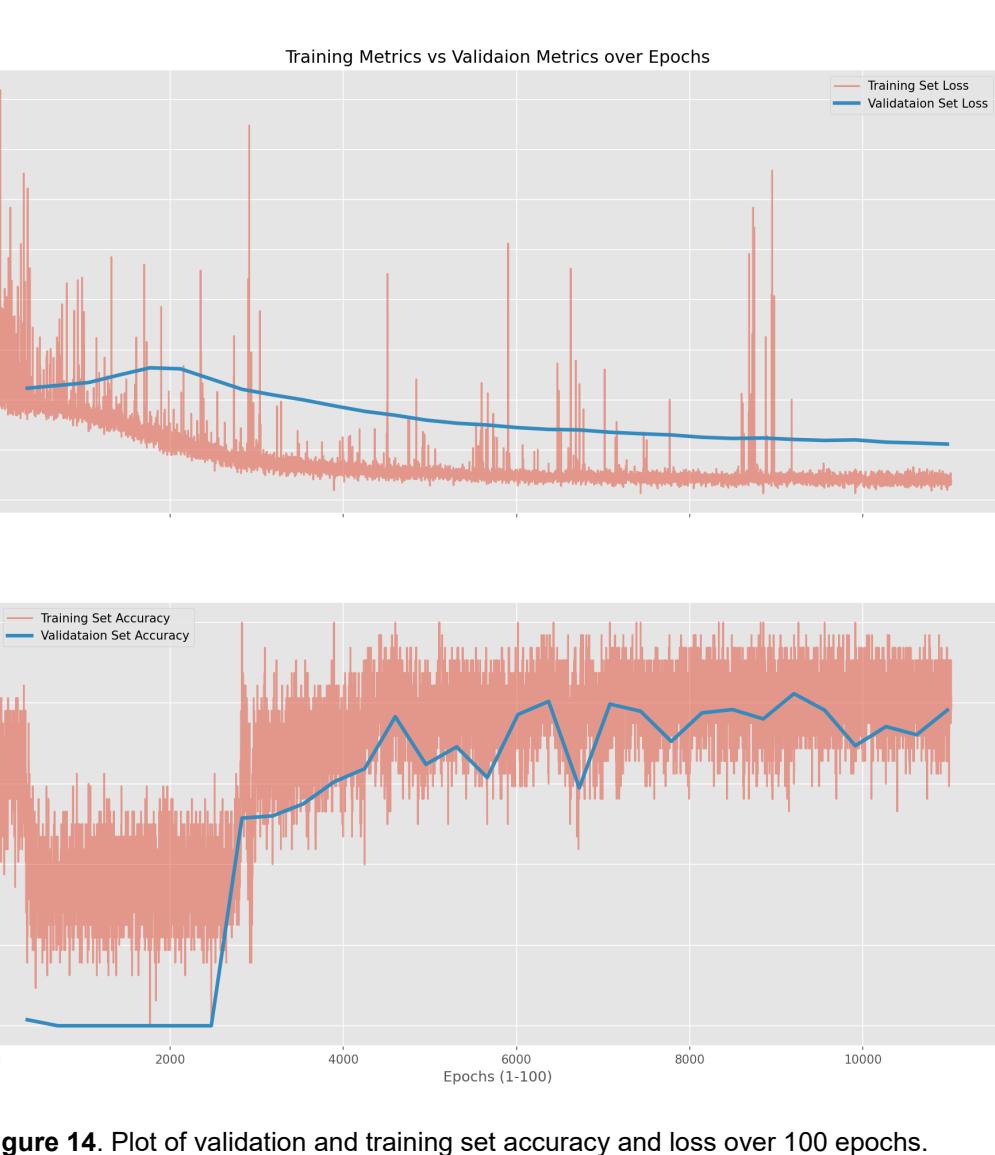
The graphs on the right illustrate that the GAN's accuracy, specificity, and sensitivity are skewed to the left. This suggests that the model is extremely successful at differentiating between abnormal and normal heart sounds.

The specificity distribution has a mean of 0.9030, with a standard deviation of 0.0547. The best model reached a specificity of 0.9672. This suggests that ~90.3% of normal heart sounds were correctly identified.

The sensitivity distribution has a mean of 0.9952, with a standard deviation of 0.0197. The best model reached a sensitivity of 1.0. This suggests that 99.5% of abnormal heart sounds were correctly identified. In detecting pathologies in medicine, we often attempt to maximize sensitivity, the rate at which a subject with a disease is correctly identified amongst other ill subjects. This is because we want to ensure that all potential subjects with a disease are sent for further examination. Essentially, weighting sensitivity higher than specificity, the rate of correctly identified normal or healthy patients from a sample of healthy patients. Thus, from the results, the proposed model is robust, in that it can detect abnormalities nearly ~100% of the time.

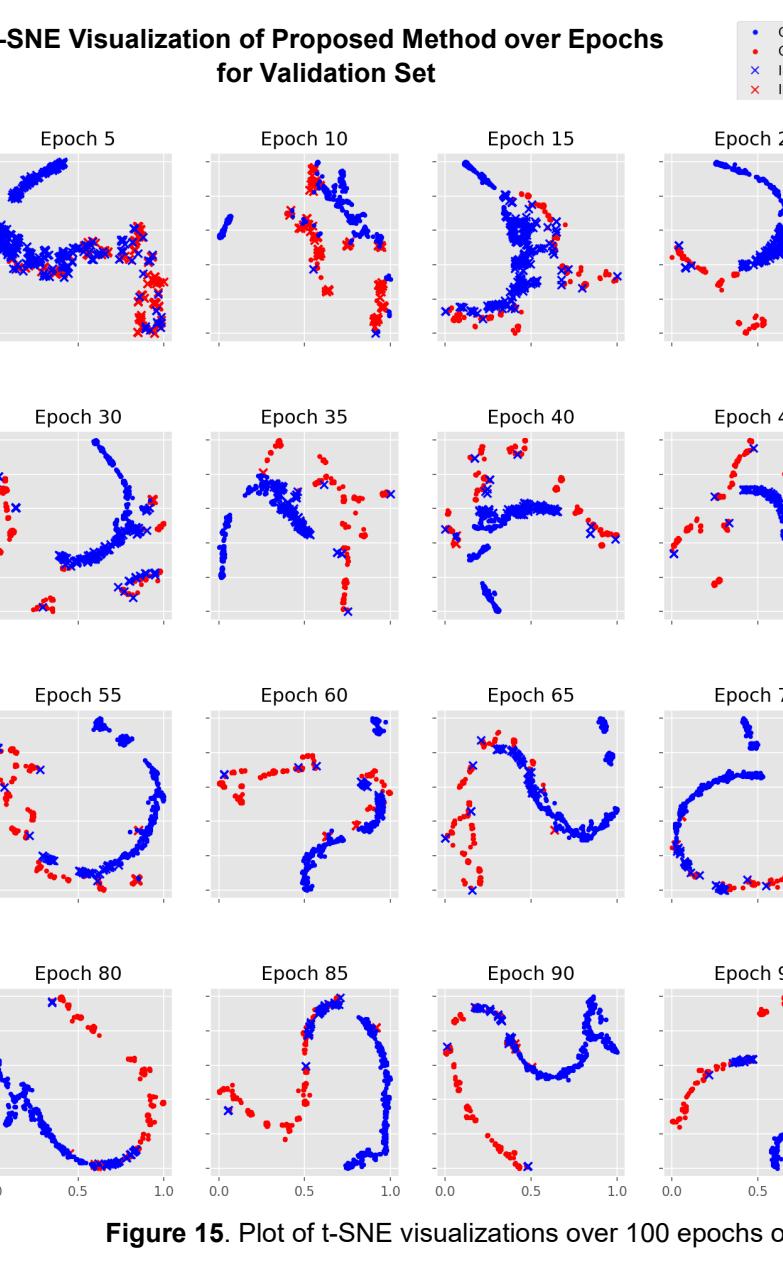
The accuracy distribution has a mean of 0.9498, with a standard deviation of 0.0293. The best model reached an accuracy of 0.9875. This suggests that ~95.0% of both abnormal and normal sounds were classified correctly.

## PCG-Net: Generalization Statistics



Generalization is important in creating accurate predictions, as it establishes that the model is learning meaningful features that are not just applicable to the training data, but signals overall. The graphs to the left plot the loss and accuracy of the training and validation set over each epoch. The large fluctuations in the training set metrics are caused by logging the metrics after each step in the training set (on every change in the model parameters). Thus, the optimizer is bound to decrease gradients in the wrong direction, thus correcting for such variations cause those fluctuations. Both lines on the loss plot resemble an exponential curve, which suggests the model continues to learn as training progresses. The average deviation between the validation and training set for each epoch is 0.4817; though this deviation is high, it is due to the lack of meaningful surface-level features that would lead to accurate detection. Meaning, the model's deep feature extraction layers are responsible for the gap. Furthermore, the accuracy plots follow a logarithmic curve. In other words, as the training accuracy increases, the validation correspondingly increases, though at a slower rate. Both graphs illustrate the model has reached convergence by the end of the training phase. This is confirmed by the static change in metrics in both datasets.

## PCG-Net: Dataset Feature Visualization



Dataset visualization is critical in understanding the dataset's complexity and model's effectiveness. Here, we use t-distributed stochastic neighbor embedding (t-SNE), a statistical method for visualizing multi-dimensional data in less computationally expensive dimensions. The method is presented with the raw prediction values for each input of the validation set and maps the corresponding predictions into a 2-dimensional space. Tracked over epochs, the visualization allows us to view the progression of the model's competence while training. The visualization highlights clear clustering within the dataset, which suggests the model is stable. Though, from epochs 70 and onwards, it is evident that there is overlapping between abnormal and normal signals. Assuming these signals as ground truth, this implies that additional feature engineering is required to adequately classify heart sounds.

## PCG-Net: Time Complexity

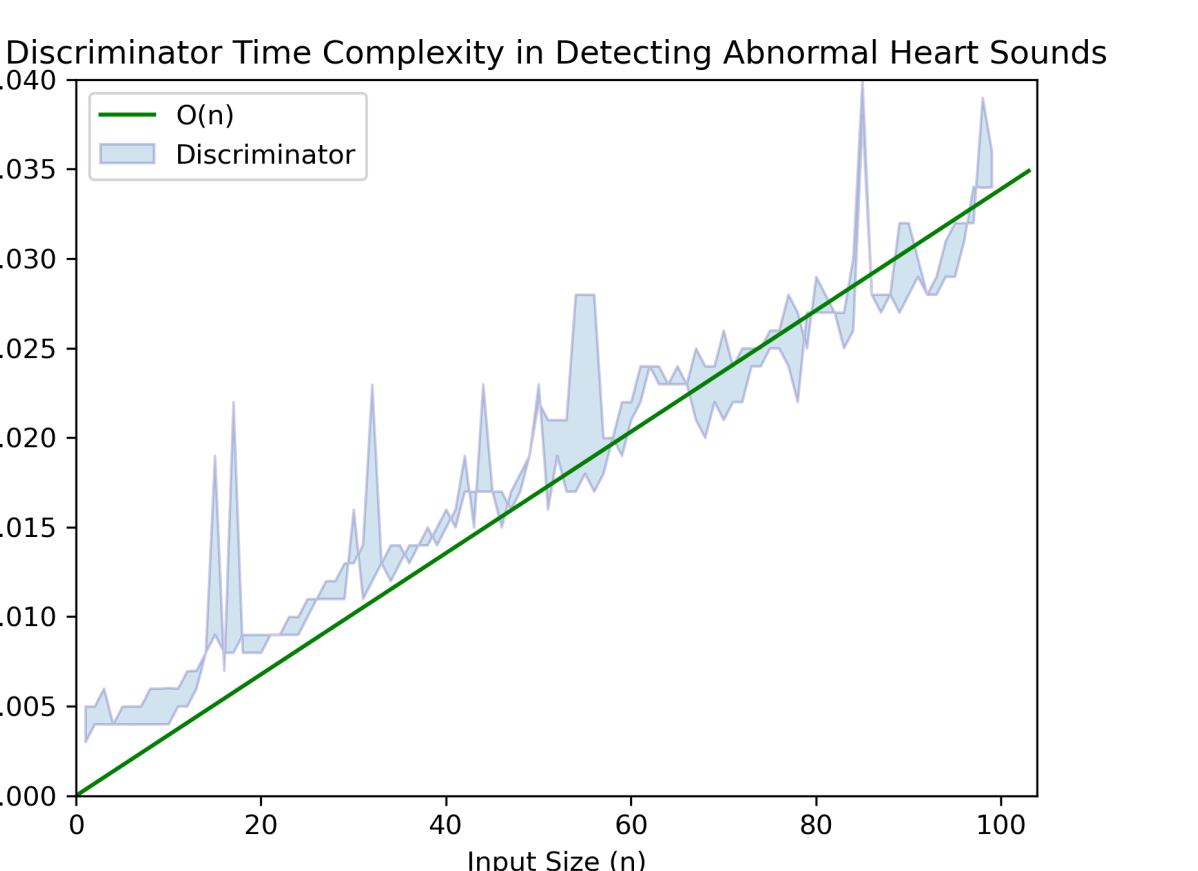


Figure 17. Time complexity of discriminator in classifying heart sounds.

Model complexity is used to gauge and evaluate the efficacy of a model against an increase in data ( $n$ ). We mainly focus on time complexity as it is most relevant to the problem at hand (space complexity is  $O(1)$ ). Depending on model deployment and integration, the complexity can vary. For example, GPUs have parallel processing capabilities, which allow them to process multiple signals at once, efficiently decreasing the model complexity to  $O(1)$ . For this reason, we use the worst-case scenario (a CPU), for analysis of the proposed method's time complexity. The graph to the left implies the model's time complexity is directly and linearly correlated to the input size, suggesting the complexity is  $O(n)$ . Thus, the model on average, can predict 2800 heart sounds in the worst case scenario. This will prove to be greatly helpful in real-time detection.

## PCG-Net: Confusion Matrix

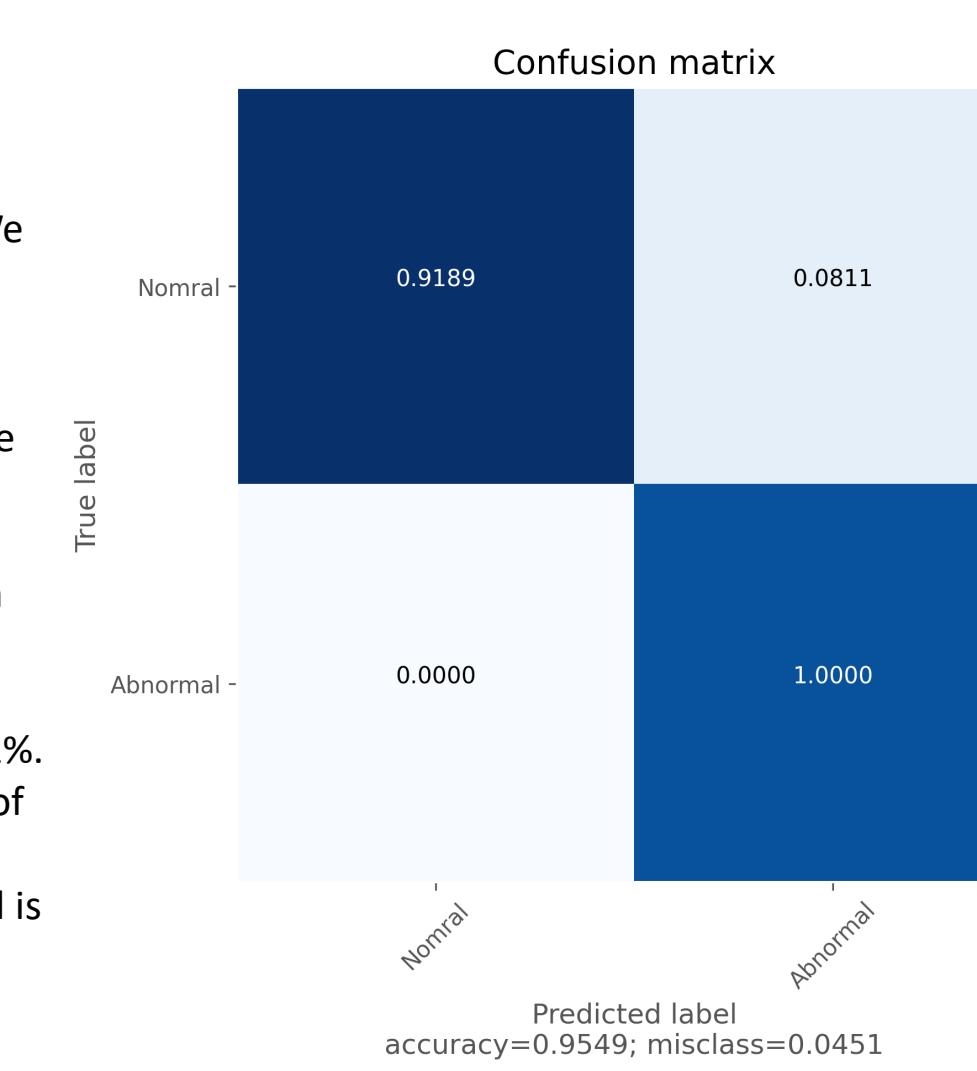


Figure 16. Matrix of accuracy between labels in the dataset.

## VQGAN: Raw Results

Table

Table 2. The table shows the loss of 150 trials of training the VQGAN model on the ECG and PCG datasets. The metrics displayed are the results of the testing set.

**Loss:** A vector that represents the correctness and confidence of the model's prediction. A low loss suggests the model was correct and confident in its answer.

## Real-World Test

Testing the model's viability is crucial for ensuring the model's success in the real world. Ideally, recording heart sounds are recorded with digital stethoscopes.

These tools use transducer technology to convert sound into an electrical signal. Over the past decade, this technology has grown immensely (by cause of speech recognition). Modern phones have the potential to record the sounds at a high resolution, given the microphone is located at the correct position relative to the heart. Such a device will prove to be extremely beneficial in providing diagnosis without the need for specialized equipment. The graph to the right shows a heart sound recording from a phone microphone. The plot shows import biomarkers, like S1 and S2, remain visible. This ensures that the recording doesn't contain excessive amounts of noise that may hinder the performance of the detection system. Hence, feeding it into the proposed method resulted in a normal classification.



## Conclusion

We proposed a Generative Adversarial Network (GAN), composed of a Dense Generator and a Convolutional Neural Network (CNN) Discriminator to detect abnormal heart sounds in a recording. The model achieved an accuracy of 94.98%, a specificity of 90.30%, and a sensitivity of 99.52% on the testing set. The previous state-of-the-art achieved an accuracy of 86.02%, a specificity of 77.81%, and a sensitivity of 94.24%. This data, along with results from the t-test revealed that the proposed alternative hypothesis was correct and that the null hypothesis should be rejected. This is because the proposed method reached better performance than the previous state-of-the-art methods. Additionally, the model attained a staggering ~2500 classification per second in the worst-case scenario. This is because of the nature of the CNN architecture; unlike other methods, the CNN's reduce the data dimensionality as it forward propagates through the model. Furthermore, the proposed method showed real-world deployment capabilities for autonomous heart sound abnormality detection with recordings collected from a phone microphone. This test shows extremely promising results for future applications and integrations.

We also set out to introduce new pathologies for increased arrhythmia labels in classification. We proposed using a VQGAN for constructing PCG signals from existing ECG datasets that contain a surplus amount of arrhythmia-specific data. The results were promising, in that the VQGAN discriminator was able to construct the general shape of the PCG spectrogram, but missed import details in the fluctuation of important biomarkers (S1 and S2). This caused the PCG waveform representation extracted from the PCG spectrogram to miss rapid oscillations present in the biomarkers.

The object of this study was to create a fast and accurate end-to-end heart sound arrhythmia detection system, capable of detecting abnormalities in real-time without specialized equipment. While also increasing the number of cardiovascular pathologies classified. With the data shown, our proposed method accomplishes exemplary statistics in abnormalities detection and shows promising results in increased arrhythmia construction. Hopefully, this study will shed light on PCG construction techniques and give birth to applications with autonomous abnormality detection.

## VQGAN: PCG Construction Visualization

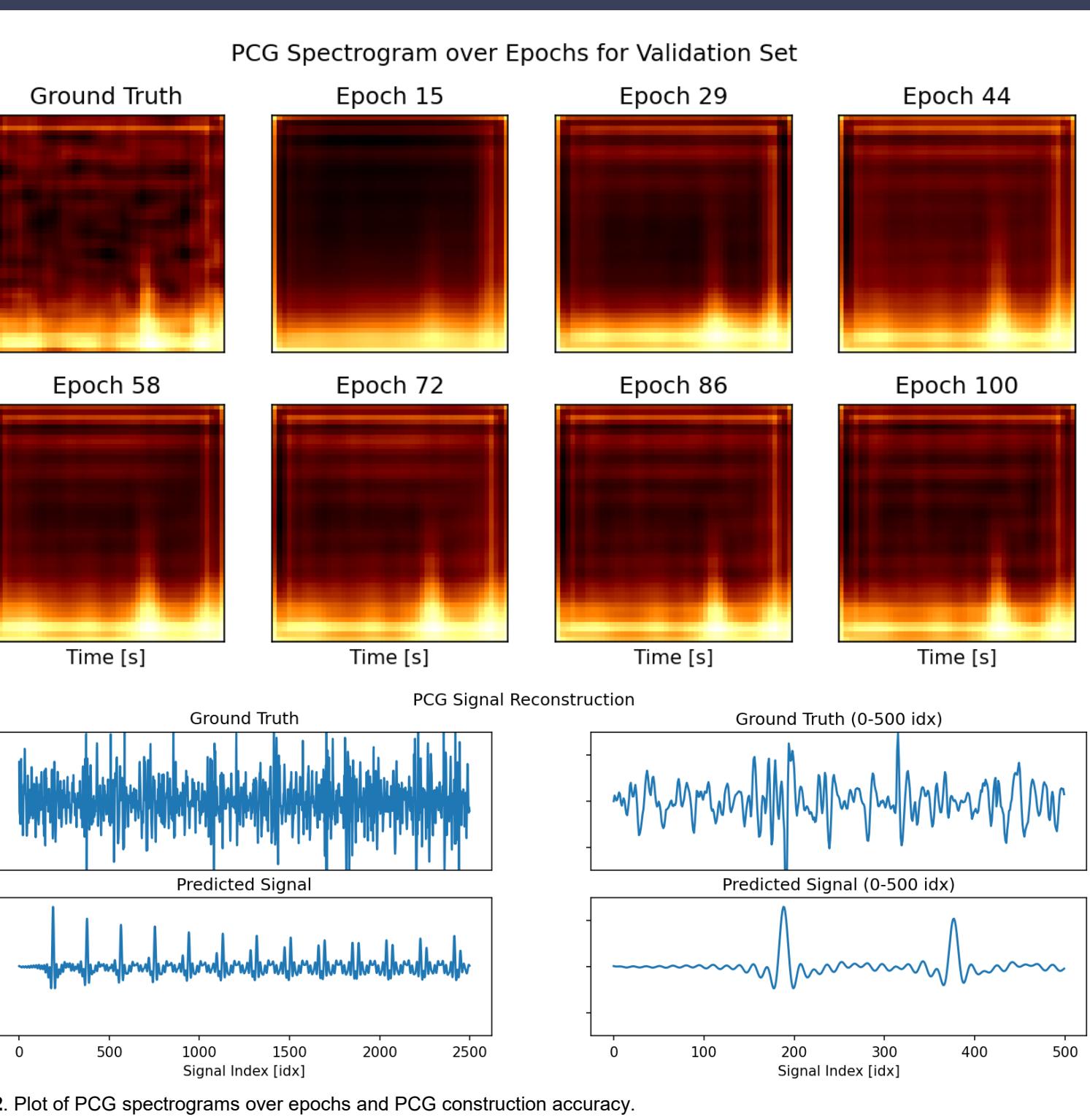
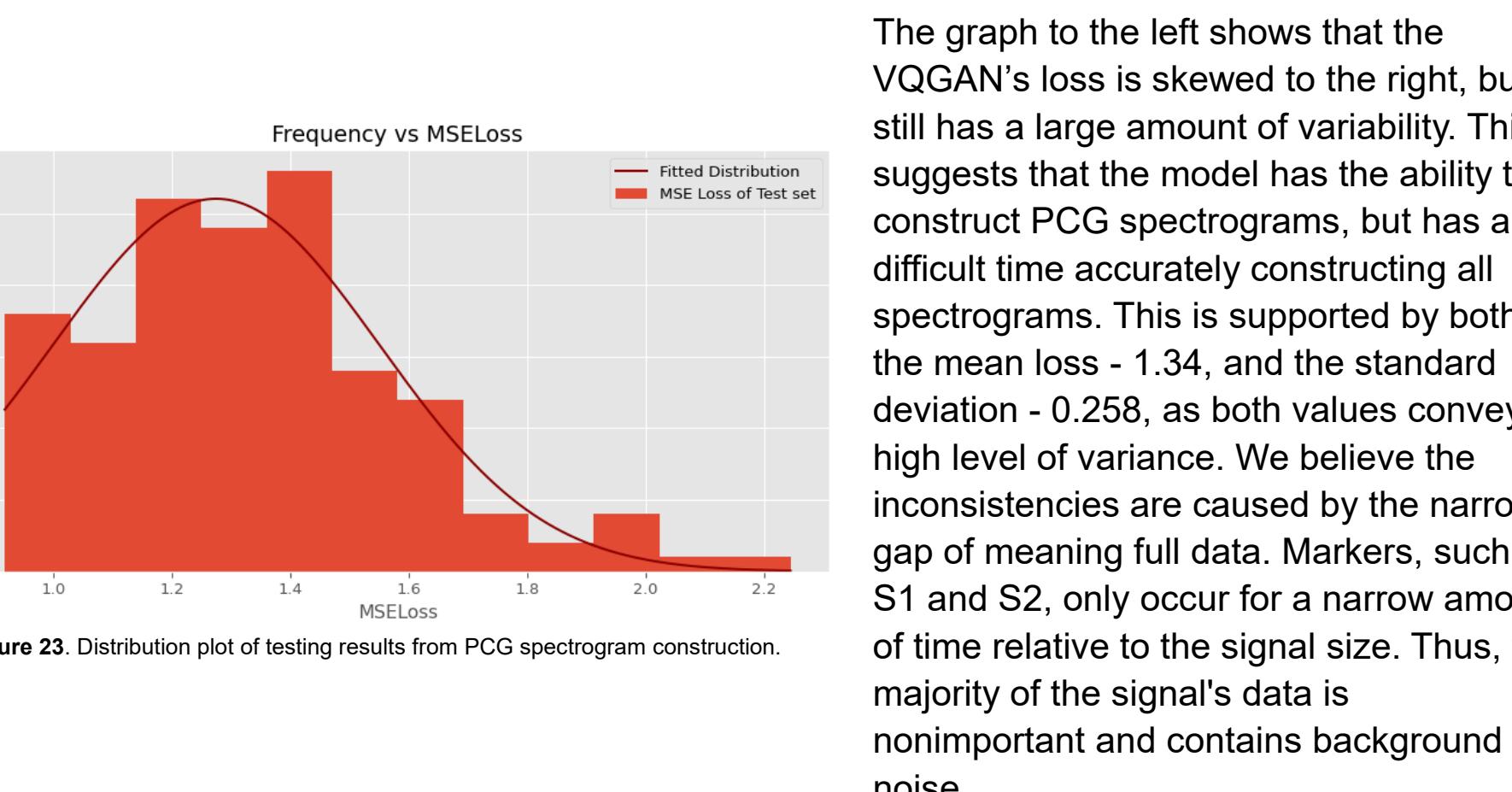


Figure 21-22. Plot of PCG spectrograms over epochs and PCG construction accuracy.

The plots above show the progression of heart sound construction from ECGs over epochs of the validation set. Ideally, he would want the reconstruction of the PCG spectrogram to identical to that of the ground truth. In practice, this doesn't occur, some features may be lost in the latent representation of the ECG spectrogram. These missing features will cause a spatial anti-aliasing effect, as the latent space doesn't have the dimensionality to extract pixel-to-pixel information. The series of spectrograms show the development of features in the latent space through the epochs. For example, the frequency of the first peak (S1) varies significantly until the ~86th epoch. This feature is important because it determines the rate of the S1 or "lub" sound; thus, creating the illusion that the sound is occurring faster relative to the ground truth. Furthermore, the S2 marker is severely softened, this is due to the rapid change in frequencies surrounding the marker and the light vertical bars to the right in each spectrogram. This suggests that much of the information regarding S2 will be lost when converting the spectrogram into a wave signal.

## VQGAN: Testing Sample Distributions



The graph to the left shows that the VQGAN's loss is skewed to the right, but still has a large amount of variability. This suggests that the model has the ability to construct PCG spectrograms, but has a difficult time accurately constructing all spectrograms. This is supported by both the mean loss - 1.34, and the standard deviation - 0.258, as both values convey a high level of variance. We believe the inconsistencies are caused by the narrow gap of meaning full data. Markers, such as S1 and S2, only occur for a narrow amount of time relative to the signal size. Thus, the majority of the signal's data is nonimportant and contains background noise.

## Further Exploration and Application

- Deploying the model with an app that is available to 3rd world countries that can't afford to conduct in-depth testing regularly
  - Integrate and serve the model using FastAPI
- Use abnormal heart sound unsupervised datasets as a basis of categorical arrhythmia classification
  - Using low dimensional visualization techniques like t-SNE or UMAP
  - Cluster data using methods like K means and hierarchical clustering
- Create a classifiable latent representation of PCG signal biomarkers that can be represented with accuracy and precision
  - Created by VAE that are fed the PCG signals directly instead of a spectrogram
  - Investigate training a heart sound discriminator from generated PCG data from ECG datasets
  - Reconstructing speech (wav) recordings from the human auditory cortex (EEG) using techniques used for PCG construction