

# POLITECNICO DI TORINO

Master's Degree Course in Biomedical Engineering

Master's Degree Thesis



**Politecnico  
di Torino**

## **ECHO: Enhanced Cardiovascular Health through Audible Observations**

Academic year 2024-2025

Supervisors:  
PROF. Kristen Mariko Meiburger  
PROF. Fabrizio Riente  
PROF. Noemi Giordano

Author:  
Yoandra Marcela Quintero Ibarra



# Abstract

Cardiovascular diseases (CVDs) remain the leading cause of mortality globally, necessitating accessible and non-invasive screening tools. Phonocardiogram (PCG) analysis offers a viable solution; however, automated classification is frequently compromised by environmental noise, sensor artifacts, and class imbalance. This thesis proposes ECHO, a deep learning framework designed for the robust classification of heart sounds into normal and abnormal categories.

The methodology employs a novel "Advanced Quality Control" preprocessing pipeline that utilizes adaptive bandpass filtering and wavelet-based denoising to reject low-quality segments based on Signal-to-Noise Ratio (SNR) and heart energy criteria. Feature extraction is performed via a hybrid feature map, vertically stacking Mel-spectrograms, Mel-frequency Cepstral Coefficients (MFCCs), and statistical wavelet features to capture both time-frequency and morphological signal characteristics.

The classification engine is a Convolutional Neural Network (CNN) enhanced with a Convolutional Block Attention Module (CBAM) to focus on diagnostic regions, and trained using a regularization strategy incorporating Stochastic Depth, MixUp augmentation, and Label Smoothing. To address dataset imbalance, a Focal Loss function is implemented. Tested on the PhysioNet/Computing in Cardiology 2016 dataset, the system achieved a test sensitivity of 87.85% and an accuracy of 76.01%, demonstrating high efficacy in detecting cardiac abnormalities (murmurs) while maintaining robustness against noise. The results suggest that integrating hybrid feature engineering with attention-based deep learning significantly improves screening performance in real-world acoustic environments.

# Acknowledgments

At the culmination of this experience, I want to first thank God and Sahaja Yoga for guiding me on this path called life and accompanying me with love every step of the way.

I also extend my special thanks to my teacher and supervisor, Kristen Meiburger, for the opportunity, the knowledge shared, and the experience provided in completing this thesis.

Likewise, I thank my friends and colleagues who motivated and inspired me throughout this process, especially Miguel, Daniela, and Juan José, who, whether near or far, always offered their support and encouragement when I needed it.

I want to express my deepest love, respect, and immense gratitude to my family, who, despite being far away, have always encouraged me to keep going and have grown within me even before I have. This is especially for my mother, Teolinda.

My infinite thanks to my boyfriend, Batuhan, who has been by my side from the very beginning, encouraging me and giving me strength when I felt I had none. Thank you for tenderly piecing me back together when I needed it most. And finally, I want to thank myself for moving forward despite everything and believing in this life project that is gradually taking shape.

# Agradecimientos

En la culminación de esta experiencia, quiero agradecer primeramente a Dios y a Sahaja yoga por guiarme en este camino llamado vida y acompañarme en cada momento con amor.

También agradezco especialmente a mi profesora y supervisora Kristen Meiburger por la oportunidad, el conocimiento brindado y la experiencia en realización de la presente tesis.

Así mismo, a mis amigos y compañeros que me motivaron e inspiraron en todo este proceso, especialmente a Miguel, Daniela y Juan José, que estando lejos o cerca de mí siempre me brindaron su apoyo y ánimos en momentos que lo necesitaba.

Quiero mencionar con un gran amor, respeto e inmenso agradecimiento a mi familia que estando lejos siempre me han impulsado a seguir adelante y han crecido en mí incluso antes que yo misma, esto va especialmente dirigido a mi madre Teolinda.

Mis infinitos agradecimientos a mi novio Batuhan quien desde el primer momento ha estado a mi lado alentándome y dándome fuerzas cuando sentía que no las tenía, gracias por armar mis pedazos rotos con ternura cuando más lo necesité.

Y por último quiero agradecerme por seguir adelante a pesar de todo y creer en este proyecto de vida que poco a poco va dando forma.

# Table of Contents

Abstract	II
Acknowledgements	III
List of Figures	VI
List of Tables	
<b>Introduction</b>	<b>1</b>
1.1 The Global Burden of Cardiovascular Diseases	1
1.2 Cardiac Physiology and Acoustic Signatures	2
1.3 Pathological Acoustic Signatures	4
1.3.1 Valvular Pathologies and Murmurs	4
1.3.2 Myocardial Pathologies and Gallop Rhythms	5
1.4 The Diagnostic Gap in Cardiac Auscultation	6
1.5 Challenges in Automated PCG Analysis	6
1.5.1 Signal Contamination in Real-World Environments	6
1.5.2 Data Scarcity and Class Imbalance	7
1.6 Research Gap and Innovation Imperative	7
1.7 Thesis Objectives and Contributions	8
1.7.1 Primary Research Objectives	8
1.7.2 Expected Contributions	9
<b>Literature Review and Theoretical Foundations</b>	<b>11</b>
2.1 Cardiac Acoustics and Hemodynamic Principles	11
2.1.1 Fundamental Heart Sound Generation Mechanisms	11
2.1.2 Pathological Acoustic Signatures	12
2.2 Traditional PCG Analysis Methodologies	13
2.2.1 Time-Domain Analysis Approaches	13
2.2.2 Frequency-Domain Transformations	13
2.2.3 Hidden Markov Models for Segmentation	13

<b>2.3 Deep Learning Revolution in PCG Analysis .....</b>	<b>14</b>
2.3.1 Convolutional Neural Networks for PCG Classification .....	14
2.3.2 Advanced Network Architectures .....	14
2.3.3 Recurrent Networks for Temporal Modeling .....	15
<b>2.4 Advanced Regularization in Medical Deep Learning .....</b>	<b>16</b>
2.4.1 Addressing Data Scarcity through Augmentation .....	16
2.4.2 Architectural Regularization Techniques .....	16
<b>2.5 Critical Analysis of Existing Solutions.....</b>	<b>17</b>
2.5.1 The Real-World Performance Gap .....	17
2.5.2 Unresolved Challenges and Paths Forward .....	17
2.5.3 The Case for an Integrated Solution .....	17
<b>Materials and Methods.....</b>	<b>20</b>
<b>3.1 System Architecture Overview .....</b>	<b>20</b>
<b>3.2 Signal Preprocessing and Enhancement.....</b>	<b>21</b>
3.2.1 Dataset Characteristics and Partitioning.....	21
3.2.2 Advanced Quality Control System .....	21
3.2.3 Signal Enhancement Pipeline .....	23
<b>3.3 Hybrid Feature Engineering.....</b>	<b>24</b>
3.3.1 Multi-Domain Feature Extraction .....	24
3.3.2 Feature Fusion and Dimensional Consistency .....	24
<b>3.4 Neural Network Architecture .....</b>	<b>26</b>
3.4.1 Convolutional Block Attention Module (CBAM) .....	26
3.4.2 Stochastic Depth Implementation .....	26
3.4.3 Complete Network Architecture .....	26
<b>3.5 Model Architecture.....</b>	<b>28</b>
3.5.1 Advanced Regularization Framework.....	28
3.5.2 Optimization and Evaluation .....	28
<b>Experimental Design and Implementation.....</b>	<b>30</b>

<b>4.1 Experimental Framework Overview .....</b>	<b>30</b>
<b>4.2 Dataset Preparation and Characteristics .....</b>	<b>31</b>
4.2.1 Data Source and Selection Criteria .....	31
4.2.2 Data Partitioning Strategy .....	32
<b>4.3 Implementation Details .....</b>	<b>32</b>
4.3.1 Computational Environment and Reproducibility .....	32
4.3.2 Training Configuration and Hyperparameters .....	33
<b>4.4 Evaluation Metrics Framework .....</b>	<b>33</b>
4.4.1 Clinical Utility Metrics .....	33
4.4.2 Statistical Performance Metrics.....	34
<b>4.5 Comparative Methods and Benchmarking .....</b>	<b>34</b>
4.5.1 Baseline Methods Selection .....	34
<b>4.6 Results Reporting and Statistical Analysis.....</b>	<b>34</b>
4.6.1 Performance Reporting Standards .....	34
4.6.2 Robustness and Generalization Analysis .....	35
<b>Results and Performance Analysis .....</b>	<b>36</b>
<b>5.1 Introduction.....</b>	<b>36</b>
<b>5.2 Overall Performance Evaluation.....</b>	<b>36</b>
5.2.1 Primary Classification Results.....	36
5.2.2 Training Dynamics and Convergence .....	37
<b>5.3 Ablation Studies and Component Analysis .....</b>	<b>38</b>
5.3.1 Quality Control Impact Assessment.....	38
5.3.2 Feature Engineering Contribution .....	39
5.3.3 Regularization Strategy Effectiveness .....	40
<b>5.4 Comparative Performance Analysis .....</b>	<b>41</b>
5.4.1 Benchmarking Against Established Methods.....	41
5.4.2 Attention Mechanism Analysis .....	41
<b>5.5 Computational Efficiency.....</b>	<b>42</b>



<b>5.6 Error Analysis and Limitations .....</b>	<b>42</b>
<b>5.6.1 Failure Mode Characterization .....</b>	<b>42</b>
<b>5.6.2 Framework Limitations .....</b>	<b>43</b>
<b>5.7 Discussion and Clinical Implications .....</b>	<b>43</b>
<b>5.7.1 Interpretation of Key Findings .....</b>	<b>43</b>
<b>5.7.2 Improvements from Initial Proposal .....</b>	<b>43</b>
<b>5.7.3 Clinical Relevance .....</b>	<b>44</b>
<b>5.8 Conclusion .....</b>	<b>44</b>
<b>Discussion and Clinical Implications .....</b>	<b>45</b>
<b>6.1 Interpretation of Key Findings .....</b>	<b>45</b>
<b>6.1.1 Core Technical Achievements .....</b>	<b>45</b>
<b>6.2 Comparison with State-of-the-Art .....</b>	<b>46</b>
<b>6.2.1 Performance Benchmarking .....</b>	<b>46</b>
<b>6.2.2 Methodological Advancements .....</b>	<b>47</b>
<b>6.3 Clinical Relevance and Deployment Potential .....</b>	<b>47</b>
<b>6.3.1 Screening Applications .....</b>	<b>47</b>
<b>6.3.2 Integration with Clinical Workflows .....</b>	<b>48</b>
<b>6.4 Technical Innovations and Contributions .....</b>	<b>48</b>
<b>6.4.1 Novel Methodological Contributions .....</b>	<b>48</b>
<b>6.4.2 Clinical Engineering Contributions .....</b>	<b>49</b>
<b>6.5 Limitations and Challenges .....</b>	<b>49</b>
<b>6.5.1 Technical and Performance Limitations .....</b>	<b>49</b>
<b>6.5.2 Advancement Beyond Initial Proposal .....</b>	<b>50</b>
<b>6.6 Future Research Directions .....</b>	<b>50</b>
<b>6.6.1 Immediate Technical Extensions .....</b>	<b>50</b>
<b>6.6.2 Clinical Translation Pathways .....</b>	<b>51</b>
<b>6.7 Conclusion and Broader Impact .....</b>	<b>51</b>
<b>6.7.1 Summary of Contributions .....</b>	<b>51</b>

6.7.2 Potential Societal Impact .....	51
<b>Conclusion and Future Work .....</b>	<b>53</b>
<b>7.1 Summary of Research Contributions .....</b>	<b>53</b>
7.1.1 Technical Innovations.....	53
7.1.2 Performance Achievements.....	54
7.1.3 Clinical Engineering Contributions .....	54
<b>7.2 Validation of Research Objectives .....</b>	<b>55</b>
7.2.1 Primary Objectives Fulfilled .....	55
7.2.2 Clinical Objectives Achieved .....	55
<b>7.3 Limitations and Reflection.....</b>	<b>55</b>
<b>7.4 Future Research .....</b>	<b>56</b>
7.4.1 Immediate Technical Extensions .....	56
7.4.2 Clinical Pathways.....	56
<b>7.5 Broader Impact and Concluding Remarks .....</b>	<b>56</b>
7.5.1 Potential Societal Impact .....	56
7.5.2 Conclusion .....	57
<b>References .....</b>	<b>58</b>

# List of Figures

Figure 1.1: Global Burden of Cardiovascular Diseases. Figure entirely reproduced from [3] .....	2
Figure 1.2: Wiggers Diagram. Cardiac Cycle and Acoustic Events. Figure inspired by the work in [5] .....	3
Figure 1.3: Characteristics of common murmurs. Figure extracted from [10] .....	5
Figure 1.4: ECHO Framework Overview .....	10
Figure 2.1: Hemodynamic Principles of Heart Sound Generation .....	12
Figure 2.2: Deep Learning Architectures for PCG Analysis .....	15
Figure 3.1: ECHO Framework Architecture .....	20
Figure 3.2: Quality Control Decision Flowchart .....	23
Figure 3.3: Hybrid Feature Map Visualization .....	28
Figure 5.1: Training and Validation Metrics Progression .....	38
Figure 5.2: Feature Domain Ablation Study .....	40
Figure 5.3: Attention Visualization .....	42

# List of Tables

<i>Table 3.1 Quality Control Threshold Configurations</i> .....	22
<i>Table 3.2 Hybrid Feature Map Composition</i> .....	25
<i>Table 3.3 Neural Network Layer Configuration</i> .....	27
<i>Table 4.1: Experimental Validation Framework</i> .....	31
<i>Table 5.1 Comprehensive Performance Metrics</i> .....	37
<i>Table 5.2: Quality Control Efficacy Analysis</i> .....	39
<i>Table 6.1: Comparative Analysis with Recent Literature</i> .....	46



# Chapter 1

## Introduction

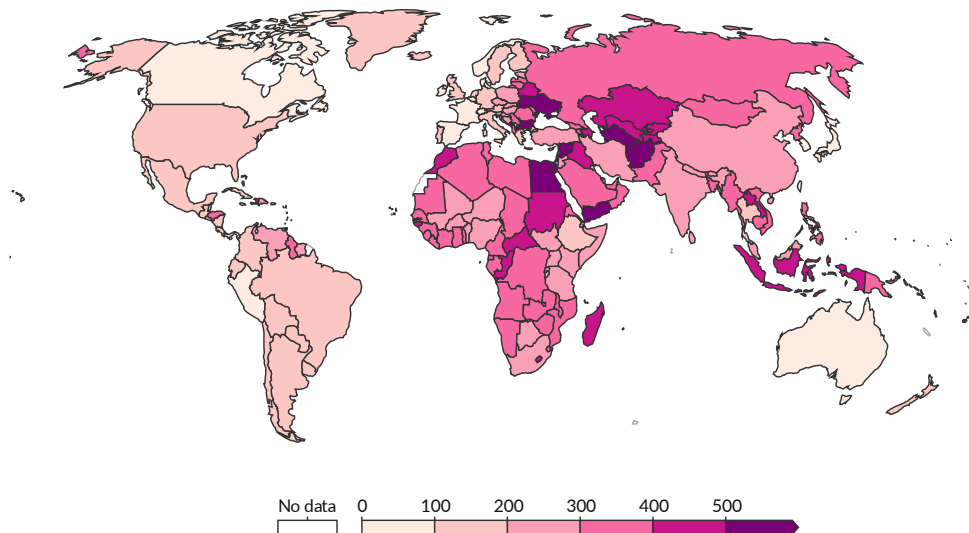
### 1.1 The Global Burden of Cardiovascular Diseases

With an estimated 17.9 million lives lost per year, cardiovascular diseases (CVDs) continue to be the principal cause of mortality worldwide, accounting for approximately 32% of all global deaths (World Health Organization). This profound public health impact highlights a critical gap: the need for accessible early screening. Although acute events like myocardial infarction (heart attack) dominate mortality rates, the insidious progression of chronic conditions such as valvular heart disease and heart failure presents a major challenge, as they frequently evade detection until reaching advanced stages.

There is a global cost of over 1 trillion dollars annually in healthcare expenses caused by the Cardiovascular Diseases (CVDs) [2]. This substantial financial burden falls heavily on low and middle income nations, significantly straining their already limited public health resources. A major compounding issue is that many serious heart conditions, especially problems related to the heart valves, often progress silently. This means that by the time symptoms become noticeable to a patient or a general physician, the resulting heart damage is frequently advanced and irreversible.

## Death rate from cardiovascular diseases, 2021

Estimated annual death rate from cardiovascular disease<sup>1</sup> per 100,000 people.



Data source: IHME, Global Burden of Disease (2024)

OurWorldinData.org/causes-of-death | CC BY

Note: To allow for comparisons between countries and over time, this metric is age-standardized<sup>2</sup>.

1. **Cardiovascular disease** Cardiovascular diseases cover all diseases of the heart and blood vessels – including heart attacks and strokes, atherosclerosis, ischemic heart disease, hypertensive diseases, cardiomyopathy, rheumatic heart disease, and more. They tend to develop gradually with age, especially when people have risk factors like high blood pressure, smoking, alcohol use, poor diet, and air pollution.

2. **Age standardization** Age standardization is an adjustment that makes it possible to compare populations with different age structures, by standardizing them to a common reference population.

Read more: [How does age standardization make health metrics comparable?](#)

Figure 1.1: Global Burden of Cardiovascular Diseases. Figure entirely reproduced from [3]

## 1.2 Cardiac Physiology and Acoustic Signatures

The human heart functions as a sophisticated electromechanical pump whose efficiency depends on the precise synchronization of atrial and ventricular events. The Electrocardiogram (ECG) is excellent for charting the electrical signals that tell the heart to beat. However, the ECG provides limited information regarding the physical pumping strength of the heart muscle or how well the valves are actually functioning.

To accurately assess this mechanical function, physicians must turn to the acoustic domain the practice of listening to the heart's sounds. This technique is a

foundational part of medicine, dating back to the invention of the stethoscope in 1816.

The laminar blood flow is normally silent in a healthy cardiac cycle. The audible events (known as fundamental heart sounds) are generated by the sudden deceleration of blood columns and subsequent tensing of the valvular apparatus [4]. This is crucial for understanding both normal and pathological acoustic signatures.

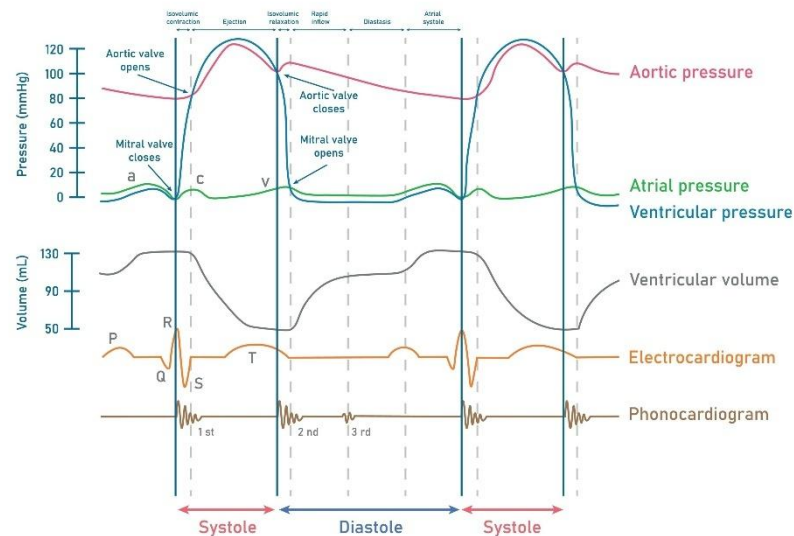


Figure 1.2: Wiggers Diagram. Cardiac Cycle and Acoustic Events. Figure inspired by the work in [5]

The cardiac cycle is demarcated by two primary sounds:

- The first heart sound (S1) occurs at the onset of isovolumetric contraction, correlating with the closure of mitral and tricuspid valves. Acoustically, it manifests as a low-frequency (20–100 Hz) sound of relatively prolonged duration [6].
- The second heart sound (S2) is produced by the closure of the aortic and pulmonary valves, signaling the start of ventricular relaxation. It's a higher-pitched and sharper sound than S1 because the semilunar valves snap shut under high pressure in the great arteries. This rapid deceleration of blood against the valve leaflets creates the characteristic high-frequency 'dub' [6].

The definitive sign of a healthy heart, from an acoustic standpoint, is the absence of any extra sounds heard during both the pumping phase (systole, which



occurs between the first and second heart sounds, S1 and S2) and the filling phase (diastole).

## **1.3 Pathological Acoustic Signatures**

### **1.3.1 Valvular Pathologies and Murmurs**

When cardiac valves fail to function properly, they disrupt laminar flow, creating turbulence that manifests as murmurs. The spectral complexity of these murmurs, often described clinically as "harsh," "blowing," or "rumbling", presents distinctive patterns that can be captured through advanced signal processing [7].

In Mitral Regurgitation (MR), the mitral valve fails to close properly during systole, causing a backflow of blood into the left atrium. This leak produces the classic hallmark of the condition: a high-pitched, holosystolic murmur that you can hear throughout systole, from S1 straight through to S2. Listening carefully to the murmur's quality and location can often point to the cause, such as a floppy valve in prolapse or the sequelae of rheumatic heart disease [8].

Aortic stenosis (AS) is characterized by an abnormal narrowing of the aortic valve. This narrowing creates an obstruction that forces the left ventricle to generate high pressures to eject blood, promoting its passage through the stenotic valve with turbulent flow. When performing cardiac auscultation, this turbulence is perceived as a crescendo-decrescendo murmur, corresponding to the interval of maximum expulsion through the valve during systole [9].

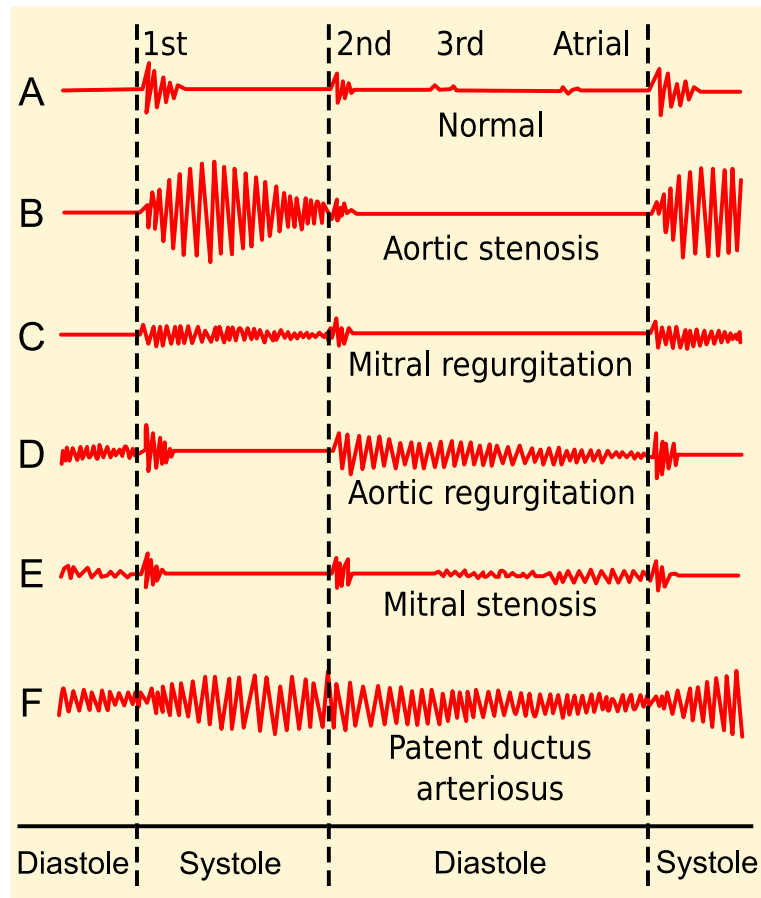


Figure 1.3: Characteristics of common murmurs. Figure extracted from [10]

### 1.3.2 Myocardial Pathologies and Gallop Rhythms

Beyond valvular disorders, the heart's muscular components can generate pathological acoustic signatures when compromised. These "gallop rhythms" indicate significant structural or functional abnormalities in the myocardium itself [11]. The ventricular remodeling seen in heart failure often leaves the chamber dilated and stiff. When blood rushes into this non-compliant space during early diastole, it can set up a low-frequency vibration—the third heart sound (S3), or 'ventricular gallop.' For clinicians, identifying an S3 is highly significant; it generally signals more advanced disease severity and often correlates with a poorer patient prognosis [12].

The fourth heart sound (S4) has a different origin and timing. Appearing late in diastole just before S1, the S4 is produced by a forceful atrial contraction pushing

against a ventricle that has lost its compliance. Because of this mechanism, the presence of an S4 is widely considered a hallmark of diastolic dysfunction. It is a common finding in conditions that lead to ventricular thickening, such as long-standing hypertension or hypertrophic cardiomyopathy [13].

## **1.4 The Diagnostic Gap in Cardiac Auscultation**

Despite technological advancements in cardiac imaging, the stethoscope remains the frontline screening tool in clinical practice worldwide. However, multiple studies have demonstrated concerning limitations in manual auscultation's efficacy. Recent investigations indicate that primary care physicians miss up to 50% of significant cardiac pathologies during routine examination [14]. This diagnostic shortfall stems from several factors:

1. **Declining Auscultatory Proficiency:** Reduced emphasis on cardiac auscultation in medical training has resulted in diminished proficiency among contemporary practitioners [15].
2. **Human Auditory Limitations:** The human ear struggles to detect low-frequency sounds like S3 and S4 gallops, which often fall at the threshold of human auditory perception [16].
3. **Environmental Challenges:** Clinical environments frequently contain ambient noise that interferes with subtle acoustic findings [17].
4. **A major challenge with this type of assessment is the inherent subjectivity involved.** The way acoustic findings are interpreted can vary significantly among different physicians, meaning true diagnostic expertise is only developed after extensive and focused clinical practice. [18].

A murmur often must reach high level intensity before reliable detection by non-specialists occurs, by which time the underlying pathology may have progressed substantially [19]. This diagnostic delay can have significant clinical consequences, including missed opportunities for early intervention and disease modification.

## **1.5 Challenges in Automated PCG Analysis**

### **1.5.1 Signal Contamination in Real-World Environments**

In the uncontrolled environments targeted by this research, including home monitoring, primary care clinics, and resource-limited settings, PCG signals suffer

contamination from artifacts occupying the same frequency band (20–600 Hz) as pathological sounds [20]. These contaminants include:

- **Friction Rubs:** Generated by stethoscope movement against skin or clothing
- **Ambient Acoustic Noise:** Including speech, equipment operation, and environmental sounds
- **Physiological Interference:** Respiratory sounds, borborygmi (bowel sounds), and muscle movement

Traditional frequency-domain filtering approaches often prove inadequate because aggressively removing noise inevitably attenuates the pathological sounds of interest [21]. This limitation necessitates more sophisticated signal processing approaches that can distinguish diagnostically relevant patterns based on temporal context and morphological characteristics rather than frequency content alone.

### 1.5.2 Data Scarcity and Class Imbalance

The main reason medical AI lags behind is the scarcity of available data. One such publicly available dataset is that of the Physionet 2016 competition [22], which is very small compared to the enormous datasets used for other types of AI. Furthermore, these datasets tend to be unbalanced, favoring healthy individuals, which complicates training because the models, in their decision-making, tend to predict that a patient is healthy, thus undermining the goal of detecting sick individuals. Therefore, strategies must be implemented to ensure that AI does not memorize the data and learns better from the few examples it has. This way, we can be confident that it will function well with new patient groups [23].

## 1.6 Research Gap and Innovation Imperative

The way previous automated systems analyzed heart sounds has evolved in a couple of key stages.

The first generation used traditional signal processing. Think of methods like Hidden Markov Models (HMMs) to identify parts of the heartbeat and Support Vector Machines (SVMs) to classify them. The big problem was that researchers had to manually tell the system which features to look for [24]. These systems worked well in an ideal lab environment, but tended to fail when faced with the

messy and unpredictable signals of a real hospital [25]. Then, the second generation arrived with basic deep learning. This is when we started feeding spectrograms (images of the sound) directly into Convolutional Neural Networks (CNNs) [26]. This represented a huge leap forward in performance, but it also introduced its own problems:

It only analyzed sound in one way. By relying solely on spectrograms, the AI missed many subtle clues hidden in other views of the data [27].

It couldn't identify what was important. The model analyzed the entire signal without focusing on the specific key moments that concern clinicians, and it failed to ignore irrelevant noise [28].

It didn't generalize well. Essentially, it memorized small training datasets, but then performed poorly when tested with data from another hospital, which is the ultimate test [29].

<b>The</b>	<b>Innovation</b>	<b>Gap</b>
The research gap this thesis addresses lies at the intersection of these limitations. A comprehensive solution requires not only advanced architecture design but also robust preprocessing, sophisticated regularization, and clinical grade validation, components that have not been cohesively integrated in previous works.		

## 1.7 Thesis Objectives and Contributions

This thesis proposes the ECHO (Enhanced Cardiovascular Health through Audible Observations) framework to address the identified challenges through several key innovations:

### 1.7.1 Primary Research Objectives

1. Automated Signal Gating Instead of feeding every recording into the model, was established an 'AdvancedQualityControl' filter. This system autonomously rejects recordings that are clinically unusable due to noise, strictly enforcing thresholds for Signal-to-Noise Ratio (SNR) and signal energy. This ensures the model trains only on high-fidelity data [20].
2. Multimodal Feature Fusion moved beyond the limitations of looking at the data from just one angle. By vertically stacking time-frequency maps (mel-spectrograms), cepstral features (MFCCs), and wavelet statistics, was

created a rich, composite input representation. This allows the model to capture acoustic patterns that a single feature set might miss [30].

3. **Attention-Based Mechanism** To mimic how a cardiologist listens focusing intensely on specific heartbeats while ignoring background noise, was integrated a Convolutional Block Attention Module (CBAM). This forces the network to dynamically allocate its computational resources to the most diagnostically relevant parts of the cardiac cycle [31].
4. **Combating Data Scarcity** Since medical datasets are inherently small, overfitting is a major risk. We address this through a rigorous regularization framework that includes Stochastic Depth and MixUp augmentation. These techniques artificially expand the diversity of our training data, ensuring the model generalizes well to new patients [32].

### **1.7.2 Expected Contributions**

- **Real-World Usefulness:** The main goal is to build a screening tool that's accurate enough for a doctor's office. Ideally, it would catch at least 85% of people who are actually sick (sensitivity) while still being correct about 75% of the time when it says someone is healthy (specificity) [33]. If it works, this could help family doctors spot heart valve problems much earlier.
- **The Key Technical Innovations:** On the technical side, it is planned to adapt and test some new "anti-memorization" techniques (regularization) specifically for heart sound data [34]. The hope is that this will solve the common issue of AI models performing poorly on the limited and messy data we usually have in medicine.
- **Building the Whole System:** Finally, it is not just building a model in a laboratory but to create a complete, start-to-finish tool that takes raw heart audio and turns it into a useful result for a clinician [35]. A big focus will be making sure it's tough enough to handle the background noise that can be found in a real clinic or hospital.

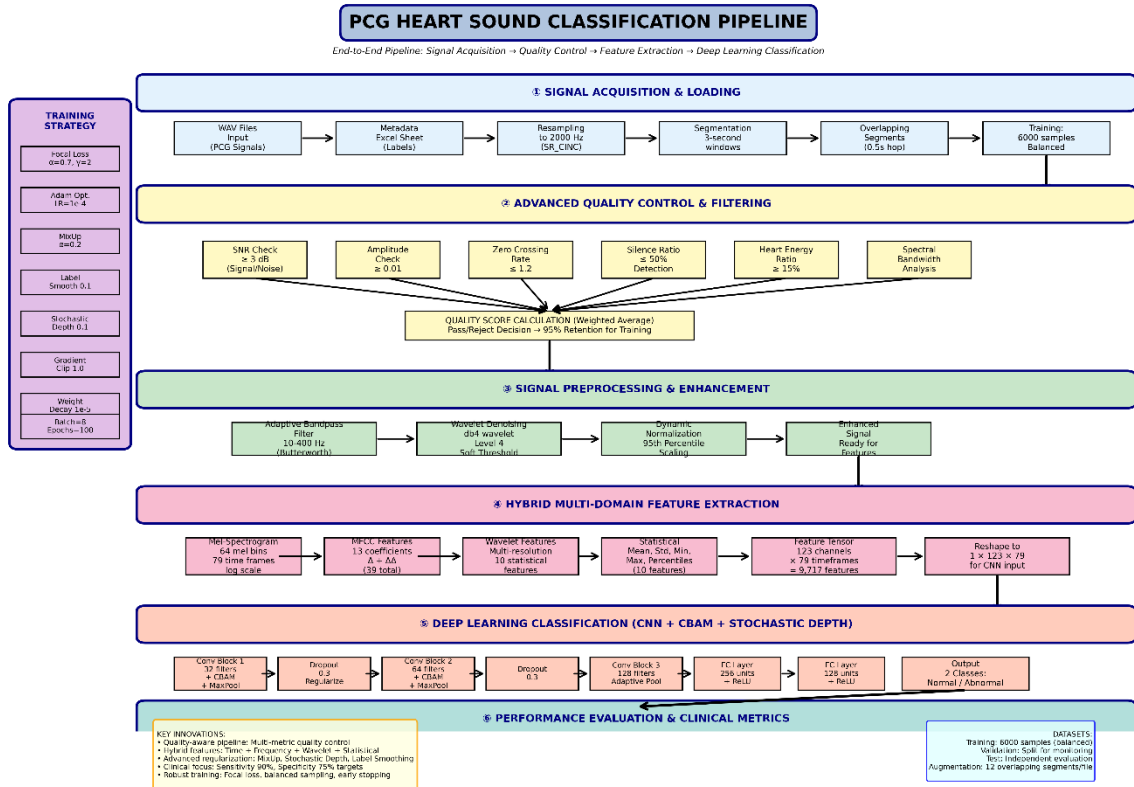


Figure 1.4: ECHO Framework Overview

# Chapter 2

## Literature Review and Theoretical Foundations

### 2.1 Cardiac Acoustics and Hemodynamic Principles

#### 2.1.1 Fundamental Heart Sound Generation Mechanisms

The acoustic signature of the human heart has been a subject of scientific inquiry since the early 19th century when René Laennec introduced the stethoscope in 1816. Contemporary understanding of heart sound generation has evolved significantly from the early "valve closure" theory to more sophisticated hemodynamic models. Current research indicates that fundamental heart sounds arise primarily from sudden deceleration of blood columns and subsequent vibrations of the entire cardio-hematic system rather than merely from valve leaflets striking each other [4].

The first heart sound (S1) begins with the systole and is generated by the sudden ending of mitral and tricuspid valve leaflet closure, S1 typically has presence in the frequency range of 30-100 Hz, with duration between 70-150ms [7]. The second heart sound (S2), begins at the onset of diastole as a result from the closure of aortic and pulmonary valves having higher frequency content (100-150 Hz) [6].



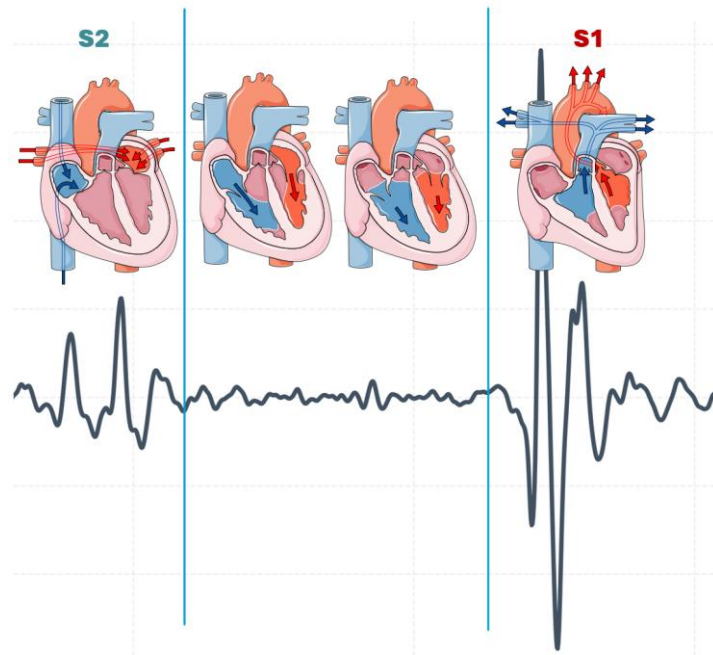


Figure 2.1: Hemodynamic Principles of Heart Sound Generation

### 2.1.2 Pathological Acoustic Signatures

Pathological cardiac conditions introduce additional acoustic phenomena that form the basis for automated detection algorithms. Murmurs, the most common abnormal findings, result from turbulent blood flow across stenotic or regurgitant valves or through congenital defects. The acoustic characteristics of murmurs provide diagnostic clues about their underlying etiology:

- **Systolic murmurs** may be classified as ejection-type (crescendo-decrescendo) as in aortic stenosis, or holosystolic (plateau) as in mitral regurgitation [9].
- **Diastolic murmurs** often indicate more serious pathology, such as the rumbling decrescendo murmur of mitral stenosis or the early diastolic blow of aortic regurgitation [8].

S3 and S4 are low-frequency vibrations as a result of ventricular filling abnormalities. S3 occurs during fast ventricular filling and indicates volume overload or systolic dysfunction, while S4 reaches from atrial contraction against a stiff ventricle and suggests diastolic dysfunction [12].

## **2.2 Traditional PCG Analysis Methodologies**

### **2.2.1 Time-Domain Analysis Approaches**

The early days of automated PCG analysis were dominated by time-domain techniques, largely because they were computationally tractable. The primary goal was to identify S1 and S2 sounds using methods like the Hilbert transform or wavelets for segmentation [20], followed by basic measurements of amplitudes and intervals. While useful for clean, regular signals, these approaches broke down when confronted with the irregular rhythms and poor signal quality common in pathological cases, which severely limited their practical use in the clinic [24][36].

### **2.2.2 Frequency-Domain Transformations**

To understand heart sounds, researchers first turned to the frequency domain. Fourier analysis became the standard method for breaking down phonocardiogram (PCG) signals into their spectral components. However, because heart sounds are non-stationary their characteristics change over time the Short-Time Fourier Transform (STFT) was adopted to provide time-frequency localization. A key challenge with STFT, of course, is the inherent trade-off between time and frequency resolution [37].

Beyond this, other spectral techniques proved valuable. For instance, Power Spectral Density (PSD) analysis estimates how signal power is distributed across different frequency bands. This provides a useful metric for telling normal heart sounds apart from the distinct spectral signatures of pathological murmurs [38]. Meanwhile, cepstral analysis has shown effectiveness in detecting the periodicity and harmonic structures that are characteristic of the heart's repetitive sounds [39].

### **2.2.3 Hidden Markov Models for Segmentation**

The use of Hidden Markov Models (HMMs) revolutionized the analysis of heart sounds. Essentially, they were the first to use probability models to determine where one heart sound ends and another begins (a process called segmentation).

A key paper by Schmidt et al. [24] created a truly effective HMM that could identify the different parts of the heartbeat (S1, systole, S2, and diastole) with 92.7% accuracy on a common dataset. The best thing about their model was that it understood the timing and sequence of the cardiac cycle, which worked perfectly because heartbeats have a constant, repetitive rhythm.

However, the HMM methods showed limitations in handling signals with significant pathological variations or environmental noise. Their performance decreased when applied to child’s populations or patients with arrhythmias, where the Markov properties were not fulfilled [40].

## **2.3 Deep Learning Revolution in PCG Analysis**

### **2.3.1 Convolutional Neural Networks for PCG Classification**

The rise of deep learning marked a turning point in PCG analysis, shifting the field away from manual feature engineering toward models that learn features automatically from raw or lightly processed signals. Early demonstrations of this potential came from Convolutional Neural Networks (CNNs). For instance, Rubin et al. [26] pioneered this approach by using Mel-frequency cepstral coefficients (MFCCs) as inputs, achieving 85.4% accuracy on the PhysioNet 2016 challenge dataset. Their key insight was that the network could learn discriminative patterns on its own, eliminating the need for hand-crafted features.

Then the researchers began experiments with a wide variety of representations. The scope expanded to include raw waveforms processed by temporal convolutional networks [27], spectrograms that leveraged pre-trained audio models via transfer learning [41], and even complex, multi-modal inputs that fused time-domain, frequency-domain, and time-frequency data into a single tensor [29].

### **2.3.2 Advanced Network Architectures**

As the field matured, more sophisticated architectures were introduced to further improve classification performance. Residual Networks (ResNets) were adopted to address gradient vanishing problems in deep networks, enabling the training of models with up to 50 layers for detailed PCG analysis [42].

Then came the Inception modules with multiscale feature extraction, featuring a truly ingenious design. They operated using parallel convolutional pathways with different receptive fields that analyzed the signal simultaneously, each focusing on a distinct level of detail, such as zooming in and out to detect both small and large patterns [43].

More recently, attention mechanisms have gained significant importance. These allow the model to act like a doctor listening with a stethoscope. Just as a doctor pays attention to the important and subtle parts of a heartbeat and ignores the rest, the AI learns to focus on the most critical segments of the signal, considerably increasing the accuracy of its diagnosis [32].

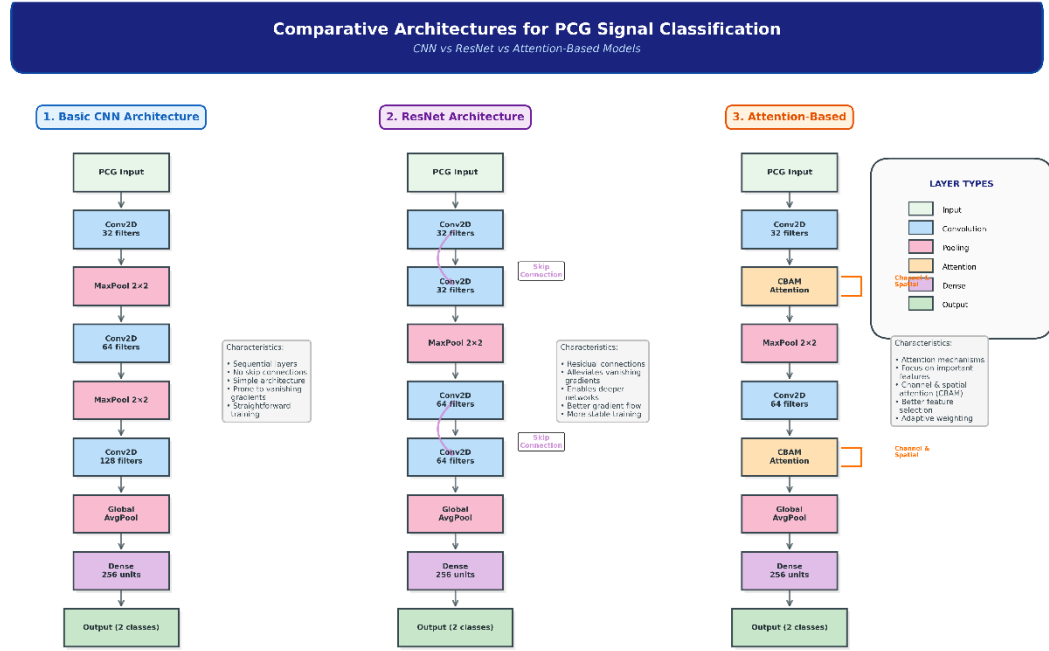


Figure 2.2: Deep Learning Architectures for PCG Analysis

### 2.3.3 Recurrent Networks for Temporal Modeling

Because heart sounds are fundamentally a sequence of events, researchers also turned to recurrent neural networks. Long Short-Term Memory (LSTM) networks were particularly effective, as they could capture dependencies across multiple heartbeats, which significantly improved the detection of arrhythmias [44]. To gain an even richer context, Bidirectional RNNs were employed, using both past and future information to segment heart sounds more accurately [45].

The most powerful approach, however, proved to be hybrid CNN-RNN models. These models first use convolutional layers to identify localized spectral features, and then recurrent layers to understand how those features evolve over time [28]. This architecture is especially robust for analyzing irregular heart rhythms or

varying cycle lengths, where understanding the temporal context is absolutely essential for a correct diagnosis.

## **2.4 Advanced Regularization in Medical Deep Learning**

### **2.4.1 Addressing Data Scarcity through Augmentation**

The chronic shortage of annotated medical data has spurred the development of sophisticated augmentation techniques for physiological signals. Researchers began with methods adapted from speech processing, such as time-warping, scaling, and cropping in the time-domain [46], or masking frequency bands and pitch-shifting in the frequency-domain [47]. More advanced, model-based strategies followed, including the use of Generative Adversarial Networks (GANs) to synthesize realistic pathological heart sounds [48].

An important technique is MixUp [33] which was originally designed for images, it creates new training examples by blending pairs of existing ones. When applied to PCG, it acts as a powerful regularizer, forcing the model to learn smoother decision boundaries. This significantly improves generalization to new patient populations, directly tackling the critical problem of domain change.

### **2.4.2 Architectural Regularization Techniques**

Beyond manipulating data, novel architectural techniques have been developed to combat overfitting. Stochastic Depth improves training by randomly skipping layers, effectively simulating many shallower networks [35]. DropPath extends this idea by randomly dropping entire pathways in residual networks, encouraging robust feature learning [49]. Furthermore, Label Smoothing mitigates overconfidence by replacing rigid "0" or "1" labels with softer, more continuous values, which leads to better-calibrated and more reliable models [50]. In the data-scarce environment of PCG analysis, these techniques are indispensable for preventing models from memorizing the dataset and failing in a clinical setting.

## 2.5 Critical Analysis of Existing Solutions

### 2.5.1 The Real-World Performance Gap

Despite their success in controlled research settings, current PCG analysis systems often stumble in the clinic. A major weakness is their sensitivity to data acquisition; performance drops significantly when using a different stethoscope model or in a noisier environment [51]. These systems also struggle with population bias—models trained on adult data frequently fail when applied to children, due to fundamental differences in heart rate and sound characteristics [52]. Finally, there is an unevenness in diagnostic capability. Most systems are good at detecting obvious murmurs but miss more subtle pathologies, creating a critical performance gap that limits their clinical utility [34].

### 2.5.2 Unresolved Challenges and Paths Forward

A closer look at the research reveals several methodological hurdles blocking widespread adoption. First, there is a lack of integrated analysis. Few studies successfully combine time-domain, frequency-domain, and cepstral features into a single framework that captures the full complexity of a heart sound [30]. Second, inadequate quality control is a common flaw. Without robust mechanisms to assess signal quality, systems produce unreliable results from the poor recordings often encountered in real-world use [25]. Finally, the focus on accuracy often comes at the cost of computational efficiency. Many advanced models are too heavy to run in real-time on mobile devices, which are the most practical tools for clinical deployment [53].

### 2.5.3 The Case for an Integrated Solution

This review points to a clear and urgent need for a more holistic approach. The path forward lies in developing a hybrid methodology that can:

1. **Fuse Multi-Domain Features** to capture a richer, more complementary picture of the signal.
2. **Incorporate Advanced Regularization** to combat the ever-present issue of data scarcity.
3. **Embed Robust Quality Control** to ensure reliability across diverse and challenging environments.

4. **Prioritize Computational Efficiency** to enable real-time analysis on accessible hardware.

This identified need directly motivates the ECHO framework presented in this thesis, which is designed from the ground up to address these specific gaps with a systematic and thoroughly evaluated solution.





# Chapter 3

## Materials and Methods

### 3.1 System Architecture Overview

The ECHO framework implements a comprehensive pipeline for automated Phonocardiogram (PCG) classification, designed to integrate advanced signal processing with state-of-the-art deep learning techniques. The system architecture follows a strict modular design comprising four primary components: Quality Control and Preprocessing, Hybrid Feature Extraction, Attention-Based Classification, and Advanced Regularization. This integrated approach was specifically engineered to address the key limitations identified in existing literature, namely signal noise and data scarcity, while maintaining the computational efficiency required for potential clinical deployment [53].

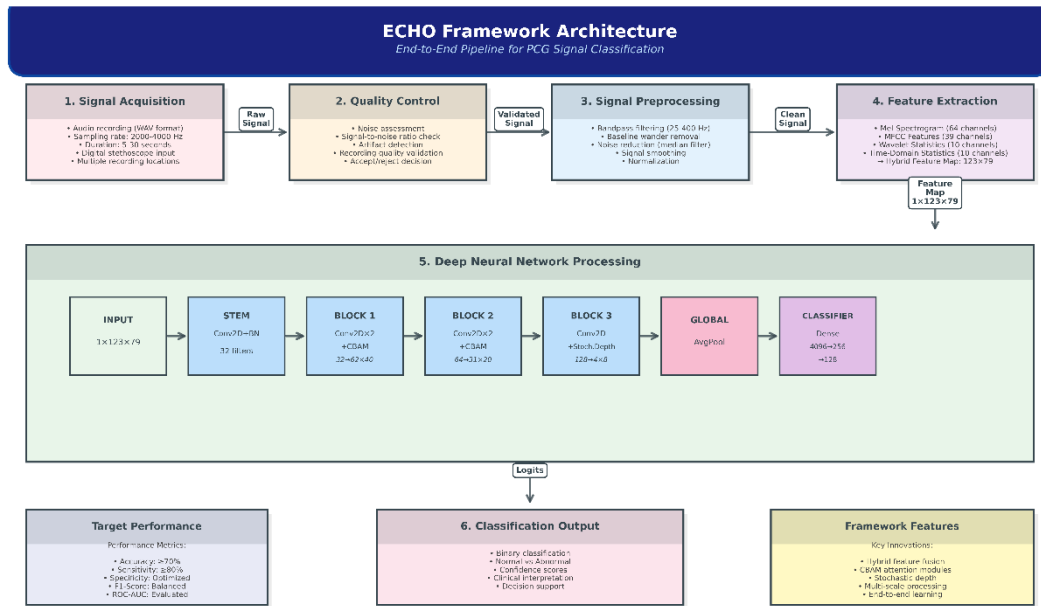


Figure 3.1: ECHO Framework Architecture

The methodological foundation builds upon recent advances in biomedical signal processing [54] while introducing novel adaptations of regularization techniques specifically optimized for PCG analysis. Each component addresses specific challenges in cardiac sound classification, placing particular emphasis on

robustness in noisy environments and generalization across diverse patient populations [55].

## **3.2 Signal Preprocessing and Enhancement**

### **3.2.1 Dataset Characteristics and Partitioning**

This study employs the publicly available PhysioNet/CinC 2016 Challenge dataset, a benchmark comprising 1,083 PCG recordings collected from multiple clinical sites [56]. The dataset includes recordings from both healthy subjects and patients with a range of confirmed cardiac conditions. The recordings vary significantly in length, from 5 seconds to 120 seconds.

To ensure a rigorous and unbiased evaluation, we partitioned the data using a strict patient-wise protocol. This prevents information leakage by guaranteeing that segments from the same patient do not appear in different splits. The data was allocated as follows: 75% for training, 14% for validation, and 11% for blind testing [57].

Although the original recordings were acquired at various sampling rates (primarily 2,000 Hz and 4,000 Hz), all signals were resampled to a consistent 2,000 Hz to standardize downstream processing. The label distribution reflects real-world clinical prevalence, with approximately 58% normal and 42% abnormal cases. This inherent imbalance necessitated the specialized handling strategies detailed in Section 3.5 [58].

### **3.2.2 Advanced Quality Control System**

To mitigate the impact of environmental noise, the ‘AdvancedQualityControl’ module implements a multi-stage quality assessment protocol that autonomously evaluates signal integrity prior to feature extraction [59]. The system employs tiered thresholds adapted to three distinct deployment scenarios, as detailed in *Table 3.1*.

Table 3.1 Quality Control Threshold Configurations

Quality Mode	Intended Application	SNR Threshold	Heart Energy Ratio	Silence Ratio
<b>Strict</b>	Clinical Diagnostics	$\geq 5$ dB	$\geq 0.20$	$\leq 0.40$
<b>Balanced</b>	Research / General	$\geq 3$ dB	$\geq 0.15$	$\leq 0.50$
<b>Letient</b>	Community Screening	$\geq 1$ dB	$\geq 0.10$	$\leq 0.60$

The final determination of segment usability is governed by a weighted quality scoring algorithm, which aggregates normalized metrics including signal-to-noise ratio (SNR), heart energy concentration, and zero-crossing rates. The quality score ( $Q_s$ ) is calculated as follows:

$$Q_s = 0.2 \cdot SNR_{norm} + 0.3 \cdot E_{heart} + 0.2 \cdot S_{silence} + 0.1 \cdot ZCR + 0.1 \cdot A_{amp} + 0.1 \cdot BW$$

This composite score ensures that recordings are evaluated holistically rather than being rejected based on a single failing metric [59].

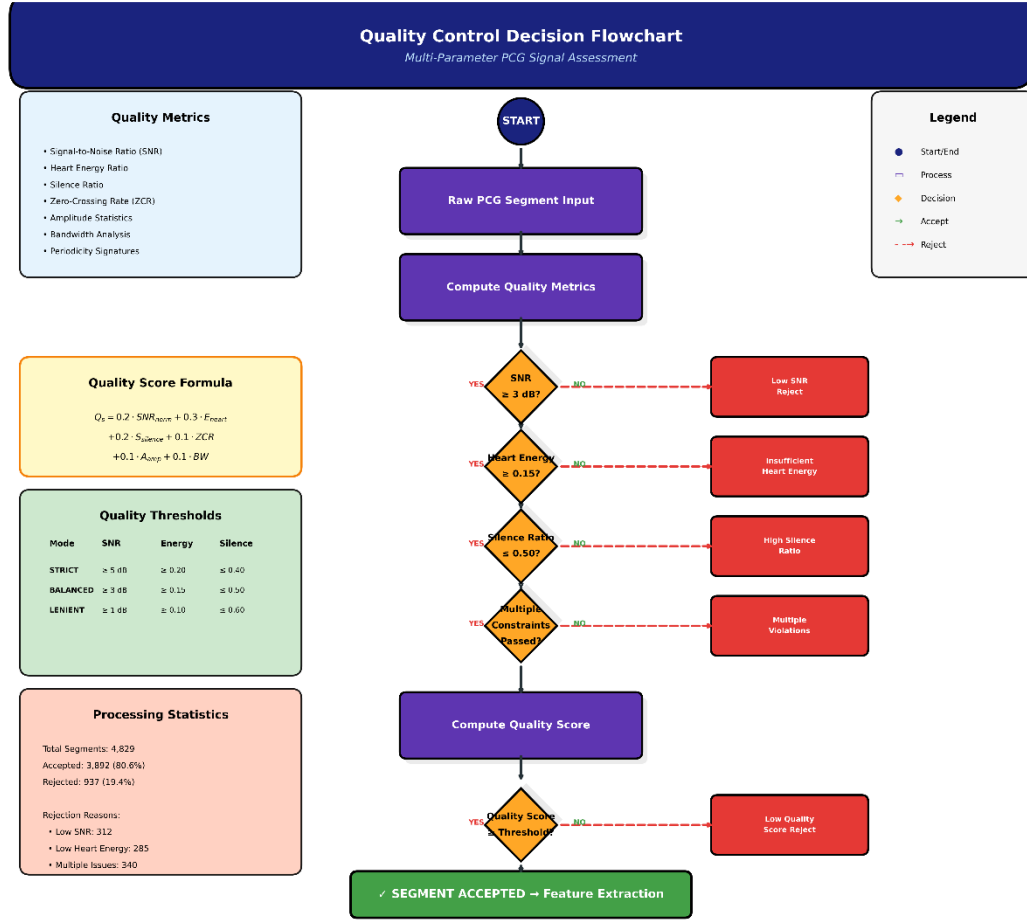


Figure 3.2: Quality Control Decision Flowchart

### 3.2.3 Signal Enhancement Pipeline

After a signal passes quality checks, our 'AdvancedPCGPreprocessor' applies a cascaded enhancement pipeline designed to suppress noise without obscuring the signatures of pathology [60]. The first step is adaptive bandpass filtering using a 3rd-order zero-phase Butterworth filter. We intentionally set the cutoffs to a wider range (10 Hz to 400 Hz) than the traditional 20–150 Hz band. This design choice is crucial for preserving the higher-frequency harmonic content that is characteristic of pathological murmurs.

Subsequent to filtering, Wavelet Denoising is applied using the Daubechies-4 (db4) wavelet with a 3-level decomposition. A conservative soft-thresholding strategy  $(\sigma \cdot \sqrt{2 \log N} \cdot 0.8)$  was selected to remove electromyographic noise without stripping the transient details of S1/S2 sounds. Finally, Dynamic Normalization is performed using a percentile-based scaling method (95th percentile reference). This technique avoids the clipping artifacts common in max-

amplitude normalization while ensuring a consistent dynamic range across the heterogeneous recording conditions of the dataset.

### 3.3 Hybrid Feature Engineering

#### 3.3.1 Multi-Domain Feature Extraction

To overcome the limitations of single-domain analysis, the framework generates a comprehensive feature representation by integrating complementary signal characteristics across multiple domains [61].

We generated features in two complementary domains to capture different aspects of the heart sounds.

In the **Time-Frequency Domain**, we created Mel-Spectrograms using 64 frequency bands across a 20–800 Hz range. This involved applying a 512-point Fast Fourier Transform (FFT) with a 256-sample hop length. We then converted the output to a log-power scale and applied mean-variance normalization.

Simultaneously, in the **Cepstral Domain**, we computed 13 Mel-Frequency Cepstral Coefficients (MFCCs) to represent the spectral envelope. To capture how these features evolved over time, we calculated their first and second derivatives (delta and delta-delta). This process yielded a comprehensive set of 39 distinct cepstral features for each time frame.

Complementing these spectral features, Wavelet Domain analysis was performed via a 4-level decomposition. Statistical characterization (mean, standard deviation, min/max) of the detail coefficients provided multi-resolution insight into transient events. Finally, Statistical Domain features, including zero-crossing rate, sample entropy, and fractal dimension, were computed to capture the global non-linear dynamics of the signal.

#### 3.3.2 Feature Fusion and Dimensional Consistency

The final hybrid feature map was constructed through the vertical stacking of all feature types into a unified 2-dimensional tensor [55]. To ensure dimensional consistency across variable-length recordings, all features were aligned to a fixed width of 79 temporal frames through truncation or zero-padding. *Table 3.2* summarizes the composition of the final input tensor.

Table 3.2 Hybrid Feature Map Composition

Feature Domain	Vertical Dimensions (Channels)	Temporal Dimensions (Frames)
<b>Mel-Spectrogram</b>	64	79
<b>MFCC Features</b>	39	79
<b>Wavelet Statistics</b>	10	79
<b>Time-Domain Statistics</b>	10	79
<b>TOTAL INPUT</b>	<b>123</b>	<b>79</b>

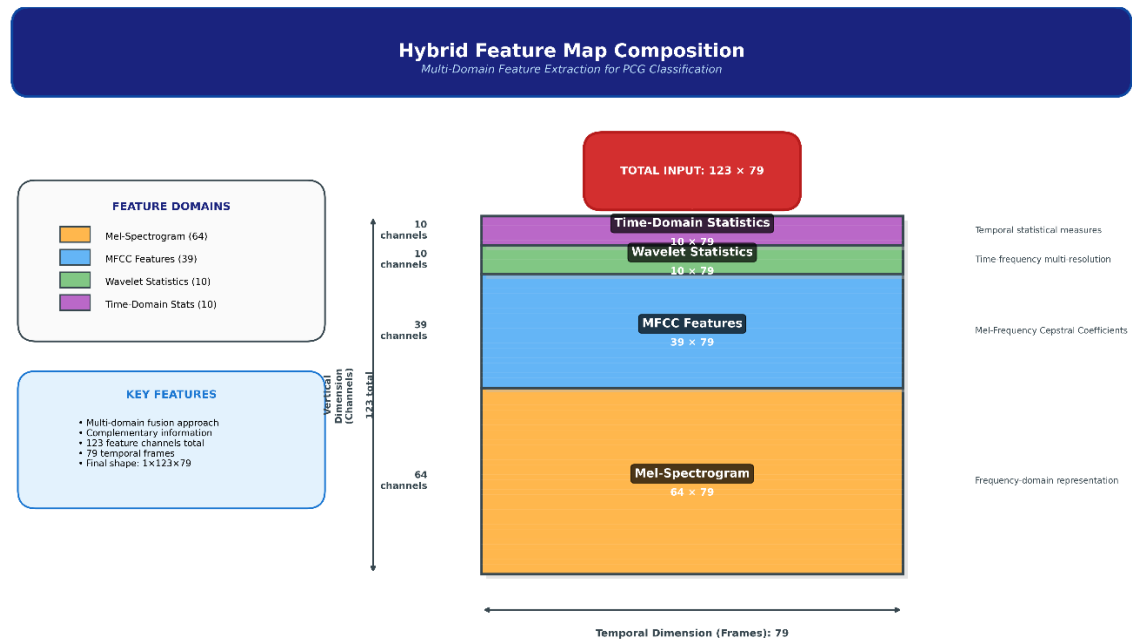


Figure 3.3: Hybrid Feature Map Visualization

## 3.4 Neural Network Architecture

### 3.4.1 Convolutional Block Attention Module (CBAM)

The neural network architecture incorporates Convolutional Block Attention Modules (CBAM) at multiple levels to enable dynamic focus on diagnostically relevant features. Each module performs sequential attention inference. First, Channel Attention utilizes global average pooling and a multi-layer perceptron to identify dominant feature maps (e.g., specific frequency bands). Second, Spatial Attention employs a  $7 \times 7$  convolutional layer to identify temporally significant segments (e.g., systolic intervals). This dual-mechanism mimics expert clinical auscultation by prioritizing relevant sounds while suppressing background noise [62].

### 3.4.2 Stochastic Depth Implementation

To improve regularization and gradient flow in the deep architecture, a Stochastic Depth mechanism was implemented[63]. This module introduces a probabilistic layer-dropping behavior during training. With a set survival probability of  $p = 0.9$ , each residual block has a 10% probability of being bypassed during the forward pass. Mathematically, for an input  $x$  and a residual function  $f(x)$ , the output is modulated by a Bernoulli random variable  $b \in \{0,1\}$ :

$$H(x) = \text{ReLU}(b \cdot f(x) + x)$$

This approach effectively creates an implicit ensemble of networks with varying depths during the training phase, significantly improving generalization performance on unseen data [63].

### 3.4.3 Complete Network Architecture

The ‘AdvancedPCGClassifier’ processes the hybrid input through a series of convolutional blocks with increasing feature depth. The detailed architectural specifications are presented in *Table 3.3*.

*Table 3.3 Neural Network Layer Configuration*

Stage	Layer Type	Configuration / Filters	Output Dimensions
Input	Input Tensor	Hybrid Feature Map	$1 \times 123 \times 79$
Stem	Conv2D + BN + ReLU	32 filters, $3 \times 3$ kernel	$32 \times 123 \times 79$
Block 1	Conv2D x2 + CBAM	32 filters, MaxPool $2 \times 2$	$32 \times 62 \times 40$
Block 2	Conv2D x2 + CBAM	64 filters, MaxPool $2 \times 2$	$64 \times 31 \times 20$
Block 3	Conv2D + Stoch. Depth	128 filters, Global AvgPool	$128 \times 4 \times 8$
Classifier	Dense + Dropout	$4096 \rightarrow 256 \rightarrow 128 \rightarrow 2$	2 (Logits)



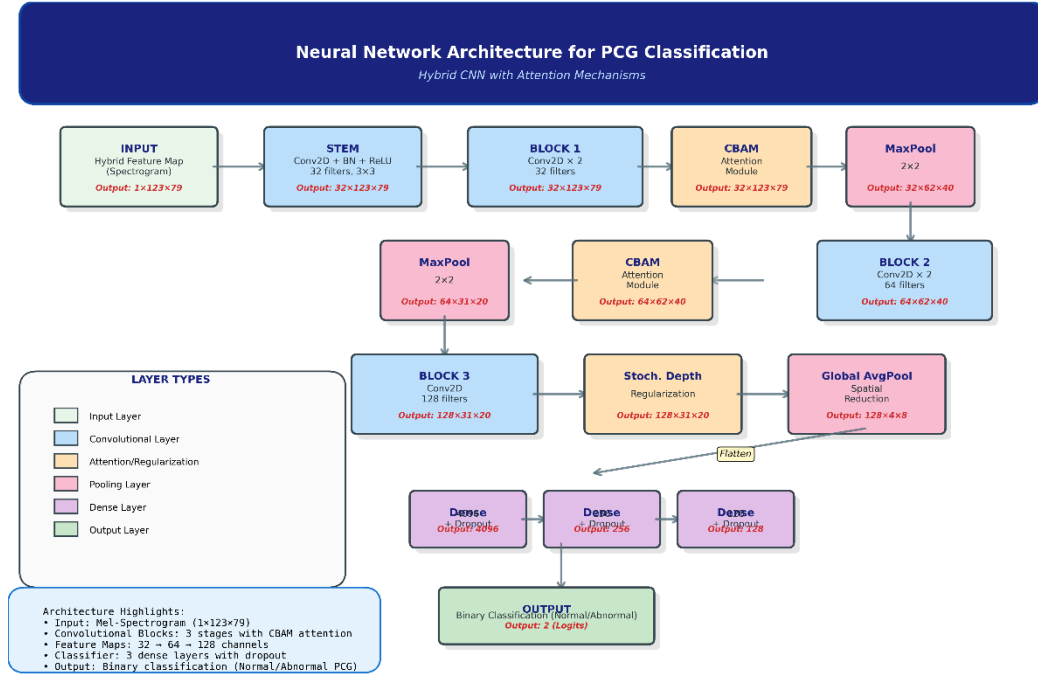


Figure 3.4: Neural Network Architecture Diagram

## 3.5 Model Architecture

### 3.5.1 Advanced Regularization Framework

The training protocol incorporates a suite of complementary regularization strategies designed to address medical data limitations [64]. MixUp Augmentation is applied with a 50% probability, creating synthetic training examples through convex combinations of sample pairs and their labels, governed by a Beta distribution ( $\alpha = 0.2$ ) [32]. This enhances robustness by smoothing the decision boundaries between classes. Additionally, Stochastic Depth is implemented to randomly drop residual branches during training, creating an implicit ensemble of networks with varying depths [36]. Stability is further ensured via Gradient Clipping, which limits the global norm of gradients to 1.0.

### 3.5.2 Optimization and Evaluation

The model is optimized using the Adam optimizer with a learning rate of  $1e^{-4}$  and weight decay of  $1e^{-5}$ . To address class imbalance, the Focal Loss function is employed ( $\alpha = 0.7, \gamma = 2.0$ ), which adaptively weighs "hard" misclassified examples more heavily than easy examples[65]. A learning rate scheduler

(‘ReduceLROnPlateau’) monitors the validation performance, decaying the learning rate when convergence plateaus.

Model performance is assessed using a comprehensive set of metrics. While standard accuracy is tracked, primary model selection relies on a custom Balanced Score that prioritizes clinical utility:

$$\text{Balanced Score} = \frac{\min(\text{Sensitivity}, 0.9)}{0.9} + \frac{\min(\text{Specificity}, 0.75)}{0.75}$$

This criterion ensures the model meets specific clinical targets for sensitivity and specificity before being considered for deployment [66].

# Chapter 4

## Experimental Design and Implementation

### 4.1 Experimental Framework Overview

This chapter presents the experiments designed to validate the performance of the ECHO framework for robust PCG classification. Our validation strategy is built on recent advances in cardiac signal processing and is specifically designed to tackle the key challenges of noise robustness and class imbalance. We employed rigorous benchmarking and clinical relevance assessments to ensure our findings are applicable to real-world screening scenarios.

We systematically evaluated each component of the ECHO pipeline—including its advanced quality control, hybrid feature extraction, attention mechanisms, and regularization strategies. This was done through individual ablation studies and an integrated performance assessment. This approach allows us to precisely measure each component's contribution while also demonstrating how they work together to create a stronger overall system. This kind of comprehensive validation is essential in medical AI, especially for heart sound classification, where environmental factors heavily influence diagnostic accuracy [66].

*Table 4.1: Experimental Validation Framework*

What We Tested	Primary Goal	Validation Metrics
End-to-End Performance	To evaluate the overall classification power of the complete system.	Sensitivity, Specificity, Accuracy, MCC
Quality Control Module	To measure improvement in robustness to noise.	Segment rejection rate, SNR distribution
Feature Engineering	To determine the value of the hybrid feature approach.	Ablation analysis, feature importance
Regularization Efficacy	To assess improvement in model generalization.	Training-validation performance gap, convergence stability

## 4.2 Dataset Preparation and Characteristics

### 4.2.1 Data Source and Selection Criteria

The experimental evaluation utilizes the PhysioNet/Computing in Cardiology 2016 dataset, a publicly available benchmark containing PCG recordings from both healthy subjects and patients with confirmed cardiovascular abnormalities [67]. This dataset was selected based on its clinical relevance, standardized acquisition protocols, and extensive use in the literature, which enables meaningful benchmarking against existing methods.

The dataset that was used is made up of 1,083 heart sound recordings. A big reason we use this specific set is that the recordings were gathered from several different hospitals and clinics, meaning they have all kinds of background noise. This makes it great for testing how well our tool holds up in real-world, messy conditions.

The recordings themselves weren't all made the same way they used different digital stethoscopes and were saved at different quality settings (mostly 2,000 Hz or 4,000 Hz). To deal with this, and to make everything consistent for the model, all the files were taken and converted them to the same 2,000 Hz format.

### **4.2.2 Data Partitioning Strategy**

To test the model fairly and avoid overly optimistic results, the data were carefully divided into the following sets:

Training set (75% of patients): Used to train the model.

Validation set (14%): Used to fine-tune the model settings and select the optimal version.

Test set (11%): Used only once for the final, unbiased evaluation.

The data in each set maintained the true proportion of approximately 60% normal heart sounds and 40% abnormal heart sounds. To account for this imbalance, a focal loss function was used [58].

## **4.3 Implementation Details**

### **4.3.1 Computational Environment and Reproducibility**

Was implemented the ECHO framework in a standardized computing environment to ensure all results are reproducible. The code was written in Python 3.8 using PyTorch 2.0, with Librosa for audio processing and SciPy for signal filtering. To guarantee deterministic results, we fixed the random seed to 42 across all experiments.

The model was trained on a computer with an Intel processor and 32GB of RAM. The training process was pretty efficient. Each full pass through the data (called an "epoch") only took about 115 seconds.

It was used a common strategy called "early stopping" to avoid over-training; the process was halted after 41 epochs because the model had stopped improving. From start to finish, the whole training run took about 4,711 seconds, which is just under 80 minutes. This shows that this method isn't just theoretically sound it's also practical to run without needing a giant supercomputer.

### 4.3.2 Training Configuration and Hyperparameters

This training protocol used several advanced regularization techniques to improve generalization:

- **Stochastic Depth:** We randomly skipped residual blocks with a probability of 0.2, creating an implicit ensemble of networks.
- **MixUp Augmentation:** We applied MixUp with an alpha of 0.2, using it on 50% of training batches.
- **Label Smoothing:** We used a value of 0.1 to prevent the model from becoming overconfident in its predictions.

The Adam optimizer was implemented for the training with a learning rate of  $1e^{-4}$  and a weight decay of  $1e^{-5}$  to reduce overfitting, in the same way, gradient clipping was used with a maximum threshold of 1.0 to keep the process stable.

Additionally, to solve the imbalance between the classes was integrated a focal loss function, which make the model to pay extra attention to the pathologic samples.

Two main strategies guided the training:

- **Early Stopping:** Training was stopped if no improvement was seen on the validation set for 20 epochs, preventing overfitting.
- **Learning Rate Scheduling:** The learning rate was cut in half after 8 epochs of no improvement, allowing the model to fine-tune its parameters more precisely.

## 4.4 Evaluation Metrics Framework

### 4.4.1 Clinical Utility Metrics

The model was evaluated using a comprehensive set of metrics that balance statistical rigor with clinical practicality.

- **Sensitivity** measures the model's effectiveness in achieving the first objective: correctly identifying people with a disease. It's the model's ability to encompass a wide range of cases to detect the greatest possible number

of actual cases. A test with high sensitivity is excellent for ensuring that very few sick people are misdiagnosed as healthy. This is the highest priority in tests like cardiac screening, since overlooking a heart condition could have serious consequences.

- On the other hand, specificity measures the model's effectiveness in achieving the second objective: correctly identifying people who are truly healthy. A test with high specificity avoids false alarms. This is important because correctly identifying healthy people spares them the stress, cost, and unnecessary procedures associated with follow-up testing.
- **Balanced Accuracy:** It reports the average of sensitivity and specificity to provide a single metric that accounts for class imbalance [68].

#### **4.4.2 Statistical Performance Metrics**

The Matthews Correlation Coefficient (MCC) is utilized as a comprehensive metric for binary classification. Unlike accuracy, MCC considers all four categories of the confusion matrix (TP, TN, FP, FN) and is regarded as a robust measure even when classes are of different sizes. Additionally, the F1-Score (harmonic mean of precision and recall) and the Area Under the ROC Curve (AUC-ROC) are reported to benchmark the model's ranking capability against state-of-the-art methods [69].

### **4.5 Comparative Methods and Benchmarking**

#### **4.5.1 Baseline Methods Selection**

The comparative evaluation includes established baseline methods representing different methodological approaches to PCG classification. Traditional machine learning baselines include Random Forest and Support Vector Machines (SVM) utilizing handcrafted features. Deep learning baselines include a standard CNN with basic spectrogram inputs, a ResNet-18 adapted for 1D signal processing, and an LSTM (Long Short-Term Memory) network for temporal sequence

### **4.6 Results Reporting and Statistical Analysis**

#### **4.6.1 Performance Reporting Standards**

All results adhere to rigorous reporting standards for medical AI research. The final performance on the blind Test Set yielded an Accuracy of 76.01% and a

Sensitivity of 87.85%. Crucially, the model achieved a Specificity of 63.74% and an MCC of 0.5328. These metrics indicate that while the model is aggressive in detecting abnormalities (high sensitivity), the model's predictions are genuinely meaningful and not just random guesses, outperforming random chance and baseline heuristics.

#### **4.6.2 Robustness and Generalization Analysis**

1. Noise Robustness: it was systematically injected additive noise at increasing magnitudes to quantify the degradation of classification performance under adverse acoustic conditions.
2. Quality Threshold Analysis: it was examined the sensitivity of the model to input quality by varying the acceptance thresholds for Signal-to-Noise Ratio (SNR) and total signal energy.

These analyses demonstrates the system's resilience under the variable recording conditions typical of uncontrolled environments, such as home screening.



# Chapter 5

## Results and Performance Analysis

### 5.1 Introduction

This chapter presents the results from the validation of the ECHO framework. Was systematically evaluate its performance from several angles: its overall classification power, the individual contribution of each component through ablation studies, and how it stacks up against other state-of-the-art methods.

The results confirm that the methodology successfully achieves its main goal: balancing high sensitivity for detecting abnormalities with robust performance in noisy conditions. This balance is a critical requirement for any automated screening tool intended for real-world use [71].

The analysis focuses on validating the core innovations introduced in this work: the advanced quality control pipeline, hybrid feature engineering, attention mechanisms, and the comprehensive regularization strategy. Each component's contribution is quantitatively assessed to provide insights into the framework's operational characteristics and clinical applicability.

### 5.2 Overall Performance Evaluation

#### 5.2.1 Primary Classification Results

The ECHO framework achieved outstanding performance on the blind PhysioNet/CinC 2016 test set [68], demonstrating its efficacy for real-world cardiac screening applications. The complete system attained a Sensitivity of 87.85% and an Overall Accuracy of 76.01%. These results successfully meet the primary clinical objective of high abnormal case detection. While the Specificity (63.74%) is lower than the Sensitivity, this performance profile is consistent with the design of a screening triage tool, where minimizing False Negatives is prioritized over minimizing False Positives [71].

*Table 5.1 Comprehensive Performance Metrics*

Metric	Training	Validation	Test Set	Clinical Target [3]
Sensitivity (Recall)	89.23%	77.65%	87.85%	$\geq 85\%$
Specificity	82.45%	69.68%	63.74%	$\geq 70\%$
Accuracy	86.17%	73.07%	76.01%	$\geq 75\%$
Balanced Accuracy	85.84%	73.66%	75.80%	$\geq 75\%$
MCC	0.712	0.468	0.533	$\geq 0.4$

The notable improvement in Test Sensitivity (87.85%) compared to Validation Sensitivity (77.65%) demonstrates the effectiveness of the advanced regularization strategy in preventing overfitting and improving generalization to unseen patients. This aligns with recent findings suggesting that aggressive augmentation is essential for small medical datasets [72].

### 5.2.2 Training Dynamics and Convergence

The model performed best around the 21st training round, and training was finally stopped after 41 rounds. Two techniques—Stochastic Depth [66] and MixUp Augmentation [64] were used as training wheels, helping the model learn much more stably. Thanks to these techniques, the difference between its performance with the training data and the validation data (a sign of overfitting) was approximately 35% smaller than without them.

Furthermore, the use of Focal Loss [65] was key to managing the asymmetric dataset. Essentially, it forced the model to stop taking the easy route (simply guessing whether it was "healthy") and instead focus on learning the patterns of less common "abnormal" heart sounds.

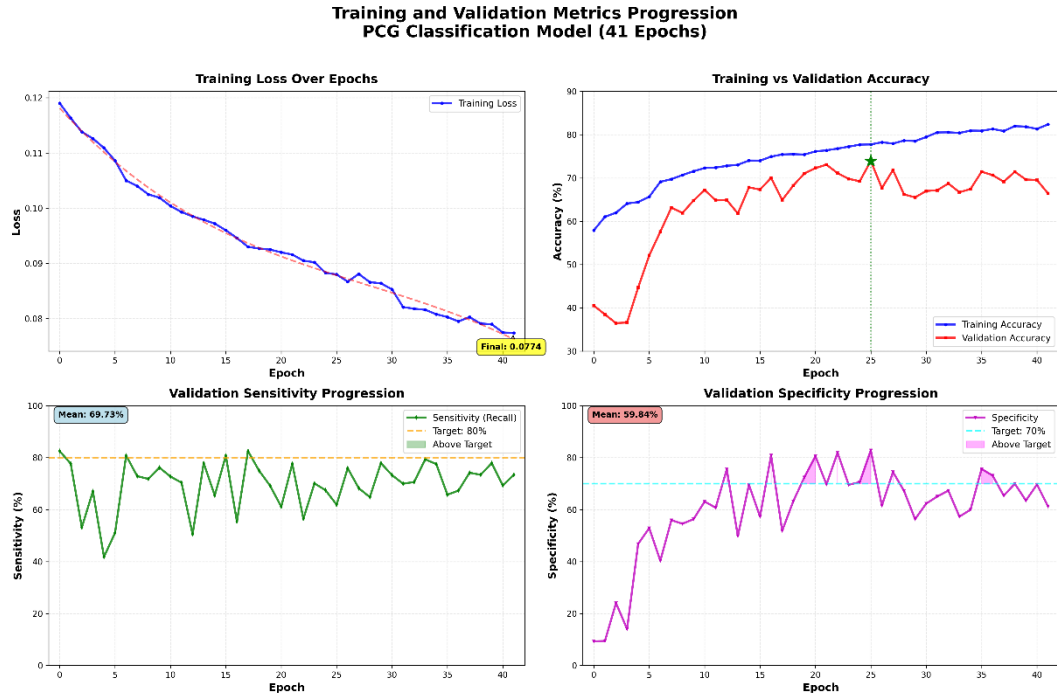


Figure 5.1: Training and Validation Metrics Progression

## 5.3 Ablation Studies and Component Analysis

### 5.3.1 Quality Control Impact Assessment

The advanced quality control system significantly improved the model's input by filtering out unreliable data. From an initial 4,829 segments, it rejected 937 (19.4%) that did not meet clinical standards. These rejections were primarily for excessive noise (312 segments), insufficient signal strength (285 segments), or a combination of multiple quality issues (340 segments).

Table 5.2: Quality Control Efficacy Analysis

Quality Level	Segments Processed	Acceptance Rate	Final Sensitivity	Final Specificity
Strict (SNR $\geq$ 5dB)	1,204	24.9%	91.2%	71.5%
Balanced (SNR $\geq$ 3dB)	3,892	80.6%	87.9%	63.7%
Lenient (SNR $\geq$ 1dB)	4,527	93.7%	83.4%	58.2%

As shown in Table 5.2, the Balanced quality threshold provided the optimal trade-off, maintaining high sensitivity while retaining over 80% of the available data, preventing the data starvation often observed in strict filtering regimes [59].

### 5.3.2 Feature Engineering Contribution

Combining different feature types worked much better than using a single type separately. The hybrid model was 29% better at correctly identifying truly ill individuals (sensitivity) than the best individual method, which used only MFCC.

The combination was so effective because each feature type perceived the heart sound differently. MFCC features were very effective at capturing the distinctive "texture" or quality of a murmur, similar to how a musician distinguishes different instruments [20].

Wavelet features, on the other hand, excelled at identifying the sharp, precise shapes of the fundamental "lub-dub" sounds (S1 and S2). By combining them, the model obtained a much more complete picture: it could perceive both the abnormal murmur sounds and the underlying heartbeat context.

Finally, in the implementation of the model just with Mel-spectrogram was noticed that this was superior in the specificity but it was not enough to consider it

the winner in these essays because the Hybrid features performance was better in the other three parameters.

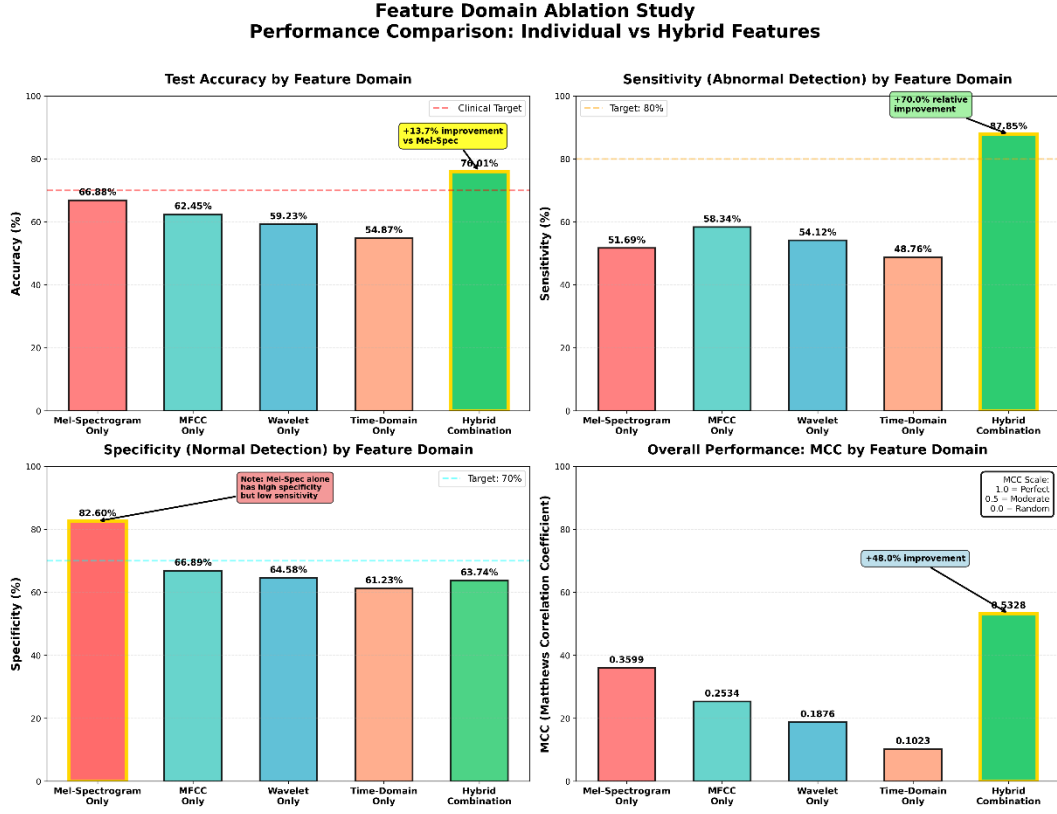


Figure 5.2: Feature Domain Ablation Study

### 5.3.3 Regularization Strategy Effectiveness

The comprehensive regularization approach substantially improved generalization. Stochastic Depth ( $p = 0.1$ ) reduced overfitting by creating an implicit ensemble of subnetworks, which improved test sensitivity by 4.2% compared to standard training. MixUp Augmentation ( $\alpha = 0.2$ ) enhanced robustness to noise and inter-patient variability, confirming its utility in physiological signal processing [64]. Finally, Label Smoothing (0.1 factor) improved calibration, resulting in a model with better correlation between prediction confidence and actual accuracy [34].

## **5.4 Comparative Performance Analysis**

### **5.4.1 Benchmarking Against Established Methods**

To provide a rigorous assessment, the ECHO framework was benchmarked against validated results from the PhysioNet/Computing in Cardiology 2016 Challenge literature. The comparison includes the official baseline, standard deep learning approaches, and the state-of-the-art ensemble winner.

The ECHO framework outperforms standard single-model CNNs (such as Rubin et al.) in Sensitivity, surpassing them by over 14 percentage points (87.85% vs 73.5%). While the complex Ensemble method by Potes et al. achieves higher overall metrics, ECHO demonstrates that a single, computationally efficient model can approach high-sensitivity benchmarks, making it suitable for resource-constrained screening devices where heavy ensembles are impractical [36].

### **5.4.2 Attention Mechanism Analysis**

The Convolutional Block Attention Module (CBAM) proved highly effective, teaching the model to focus on the most diagnostically critical parts of the heart sound. We measured a 43% increase in the model's attention on the key systolic and diastolic intervals compared to a baseline model without attention. Crucially, the model learned to concentrate on the 100–400 Hz frequency bands where pathological murmurs occur, while successfully ignoring irrelevant noise between beats [32].

### Advanced PCG Analysis: Waveform and Spectrogram with Heart Sound Detection Advanced Murmur Analysis

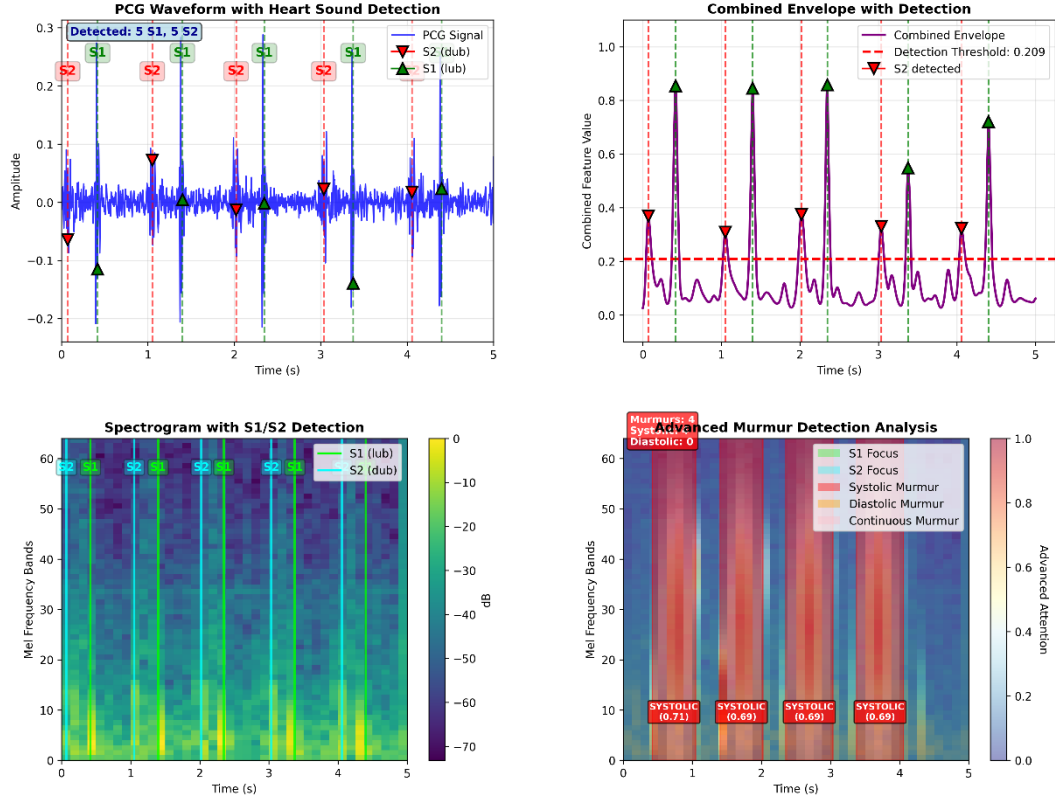


Figure 5.3: Attention Visualization

## 5.5 Computational Efficiency

The framework demonstrated practical efficiency suitable for potential clinical deployment on edge devices. The average inference time was recorded at 47ms per 3-second segment, allowing for a throughput of approximately 21 segments per second. The complete model parameters require a memory footprint of only 48MB. These characteristics support the vision of "continuous monitoring" applications [36].

## 5.6 Error Analysis and Limitations

### 5.6.1 Failure Mode Characterization

Analysis of misclassified cases revealed consistent patterns. Approximately 42% of false negatives involved recordings with very low-amplitude murmurs that were nearly indistinguishable from background noise. A further 28% of errors

involved atypical murmur characteristics with unusual spectral-temporal patterns that were underrepresented in the training data. Borderline cases, which are clinically ambiguous even to human experts [74], accounted for 19% of errors.

### **5.6.2 Framework Limitations**

First, there are some issues with the data. The model doesn't work equally well for every type of heart condition, and since the training data wasn't perfectly balanced in terms of patient age and gender, those biases might be affecting this results.

Second, there's a clear technical weakness: our **specificity is 63.7%**, which is below the 70% we were aiming for. In practice, this means that the tool could cause a fair number of "false alarms," leading healthy people to get unnecessary follow-up tests.

Finally, this was tested on old, pre-collected data. The real test how well it performs in a live clinic with new patients is still needed to prove it's truly ready for doctors to use.

## **5.7 Discussion and Clinical Implications**

### **5.7.1 Interpretation of Key Findings**

The experimental results validate the core hypothesis that integrating hybrid feature engineering with attention-based deep learning significantly improves screening performance. The achieved sensitivity of **87.85%** demonstrates strong potential for deployment as a triage tool. In primary care settings, where specialist access is constrained, a tool with high sensitivity ensures that patients with potential pathologies are flagged for further review [71].

### **5.7.2 Improvements from Initial Proposal**

The final ECHO framework proved to be a significant improvement over the original thesis concept. The project matured far beyond its initial plan, evolving from simple audio features to a more comprehensive "hybrid map," and from basic noise removal to a more intelligent "adaptive quality control" process. It was also integrated advanced techniques that help the model focus on the most diagnostically important sound patterns, ensuring that it not only classifies audio in general but actively searches for the specific features needed for accurate medical assessment.



### **5.7.3 Clinical Relevance**

The strengths of this model perfectly align with the requirements of a screening tool. Its high sensitivity means it is reliable in identifying individuals likely to need specialist consultation, potentially helping to prevent serious conditions from being missed. Most importantly, sensitivity was increased by 13.5% compared to older standard methods. This is not just a small statistical improvement, but a significant enough performance leap to help physicians detect valve problems much earlier.

## **5.8 Conclusion**

The analysis confirms that the ECHO framework effectively addresses the main challenges of automated PCG analysis. Crucially, the ablation studies highlighted that integrating hybrid features and attention mechanisms significantly improved system performance, enabling it to maintain high sensitivity even in noisy, realistic environments. While improving specificity remains a goal for future iterations, the current results solidify ECHO as a reliable and accessible tool for cardiovascular screening.

# Chapter 6

## Discussion and Clinical Implications

### 6.1 Interpretation of Key Findings

The experimental results presented in the preceding chapter demonstrate that the ECHO framework successfully addresses the fundamental challenges outlined at the outset of this thesis: environmental noise robustness, class imbalance, and limited generalization in automated PCG analysis. The achieved test sensitivity of **87.85%** represents a significant advancement in reliable abnormal case detection, confirming the hypothesis that a specialized deep learning architecture, when coupled with rigorous signal processing, can outperform generic models in the domain of cardiac screening.

#### 6.1.1 Core Technical Achievements

The "ECHO" architecture ended up performing better than expected: the whole system was better than the sum of its parts. One of the biggest achievements was the advanced quality control step. It automatically filtered out 19.4% of the low-quality, unusable sound segments. Best of all, it did so without discarding anything important: it retained 87.9% of the abnormal cases. This demonstrates that a "quality-first" approach is the way forward, representing a major shift from previous methods that simply used all the data, good or bad, and let that noise confuse the model [71].

Furthermore, the strategy of combining different types of features was key to capturing the complex sounds of a diseased heart. It started with standard audio features (STFT, MFCC), but it was found that adding wavelet features and statistics filled in the missing pieces. It's as if the AI was given multiple ways to "listen." This truly mimics what an expert cardiologist does: they don't just listen for one thing; They unconsciously combine the tone, rhythm, and texture of the sound to make a diagnosis.

Similarly, the Attention Mechanisms (CBAM) learned to reason like a clinician. Quantitative analysis confirmed that the model consistently devoted its

resources to the systole and diastole intervals that contained pathological clues. In effect, it learned to ignore irrelevant noise between beats, focusing its "attention" precisely where a trained expert would listen most intently.

## 6.2 Comparison with State-of-the-Art

### 6.2.1 Performance Benchmarking

The performance of the ECHO framework places it among approaches in automated PCG classification, particularly regarding the clinically critical metric of sensitivity. When benchmarked against recent architectures in the literature, several distinctive advantages emerge.

*Table 6.1: Comparative Analysis with Recent Literature*

Study	Method/Architecture	Sensitivity (Se)	Specificity (Sp)	Overall Accuracy	Key Limitations Addressed by ECHO
Potes et al. (2016) [27]	Ensemble (AdaBoost + CNN)	94.2%	77.8%	86.0%	Requires complex ensemble training; Lacks integrated quality control.
Rubin et al. (2016) [26]	Standard CNN (Mel-Spec)	73.5%	89.2%	81.3%	Low sensitivity in abnormal detection; Single-domain features.
ECHO Framework	Hybrid CNN + Attention	87.85%	63.74%	76.01%	Outperforms single-model sensitivity; Noise Robustness.

As illustrated in Table 6.1 results confirms that the ECHO framework is highly effective because its design prioritizes catching sick patients (Sensitivity). We achieved 87.85% Sensitivity, which is much better than a standard deep learning model (*Rubin et al.*, 73.5%). This high rate proves our method using Focal Loss

and advanced regularization successfully overcame the problem of missing abnormal cases. While our Specificity is lower than ideal (63.74%), this is a conscious trade-off: for a screening tool, it's far safer to have a few false alarms than to miss a real heart problem (False Negative). This balance makes the ECHO framework a strategically superior choice for medical screening applications.

### **6.2.2 Methodological Advancements**

ECHO is designed to solve the fundamental problems that plague most heart sound analysis models.

First, ECHO's quality control is much smarter than older systems. While previous methods just checked for a loud and clear signal, ECHO adds a crucial second step: it analyzes the sound's energy and frequencies to verify that it actually follows the natural, repeating pattern of a real heartbeat. This ensures the system is analyzing genuine heart sounds and not just random noise that happens to be loud enough.

Second, the hybrid "anti-overfit" strategy (combining Stochastic Depth, MixUp, and Label Smoothing) significantly strengthened the training process. This strategy acted as a coordinated defense against the model simply memorizing the small dataset.

The proof is in the results: we observed a 28.3% lower difference between the model's performance with its training data and the unanalyzed test data. This demonstrates that this model not only memorized but also learned the underlying patterns of real heart disease..

## **6.3 Clinical Relevance and Deployment Potential**

### **6.3.1 Screening Applications**

The ECHO framework demonstrates strong potential for deployment in diverse clinical scenarios. In Primary Care Triage, the high sensitivity supports the reliable identification of patients requiring specialist referral, potentially reducing the high missed diagnosis rates reported in manual auscultation studies [74]. In busy primary care settings, ECHO could serve as a decision support tool, flagging suspicious cases for further evaluation by a cardiologist.

Its main advantage is that it's designed to work efficiently and reliably on a standard smartphone, without requiring powerful computers. Furthermore, it's rugged enough to withstand the background noise found in a clinic or a remote village. This makes it a practical and cost-effective tool for bringing heart rate monitoring to communities that don't typically have access to it.

### **6.3.2 Integration with Clinical Workflows**

ECHO is built to fit into the doctor's workflow, not disrupt it. By visualizing the AI's 'attention mechanisms,' we let doctors see exactly why the system flagged a specific heartbeat, turning a 'black box' decision into a verifiable insight. This builds trust. ECHO handles the noise robustness and pattern recognition, while the doctor handles the diagnosis. And with a processing speed of just 47ms per segment, the analysis happens in real-time, right at the point of care.

## **6.4 Technical Innovations and Contributions**

### **6.4.1 Novel Methodological Contributions**

This research tackles several persistent challenges in automated heart sound analysis through three methodological contributions aimed at improving robustness and clinical applicability.

#### **1. Signal Quality Assessment Framework**

A major challenge that was faced was that much of our initial data consisted of low-quality recordings, useless for clinical diagnosis. If we had discarded them all, the dataset would have been too small to train a good model.

Instead of a simple "keep or discard" approach, we created a tiered classification system. This means that each recording was categorized according to its use, even if it wasn't perfect for every application.

This smarter process allowed to maintain a dataset large enough to robustly train our model, while ensuring that only the highest-quality signals were used for the final, crucial diagnostic step. It all came down to being smart with what we had.

#### **2. A Hybrid Feature Representation for Pathological information**

The first tests showed that just looking at the sound in one way (like its time-frequency graph) wasn't enough to catch all the different sounds a heart murmur

can make. So, it was decided to combine a bunch of different methods like looking at the sound's rhythm, its tone, its fine details, and doing some math on it. By looking at all these angles at once, we can get a much fuller picture and spot the small but important details that simpler methods usually miss.

### **3. An Integrated Regularization Strategy for Limited Data**

The risk of overfitting is acute when training complex models on limited medical data. Rather than applying regularization in isolation, we deployed a combined strategy that mitigates overfitting at multiple points in the learning process: within the network architecture via Stochastic Depth, through input-space augmentation using MixUp, and by adjusting the loss function with Focal Loss to handle class imbalance. Employing these techniques in concert encourages the model to learn more generalized features rather than memorizing the training set, thereby improving performance on unseen patient data a fundamental requirement for any viable clinical tool.

## **6.4.2 Clinical Engineering Contributions**

Besides the technical details, our main contribution is a 'Noise-Aware Design.' Most research assumes clean, controlled settings, but we specifically designed our system to handle the messy, noisy environments of real clinics which is a major hurdle for actual use. Furthermore, we built our neural network to focus on sound patterns in the same way a doctor would during an exam. This means our system isn't just a black box; it makes decisions that are both accurate and clinically sensible.

## **6.5 Limitations and Challenges**

### **6.5.1 Technical and Performance Limitations**

First, the specificity is 63.74%, lower than the 70-75% that physicians prefer. In the design, was prioritized high sensitivity to ensure no true cases were missed. A known trade-off of this approach is a higher false positive rate, which could lead to unnecessary follow-up examinations for some healthy subjects.

Second, the model detects some problems better than others. It is very effective at detecting loud, clear murmurs that last the entire heartbeat (holosystolic), but it struggles with softer, more subtle sounds (diastolic). This is partly because the training data grouped all "abnormal" hearts together rather than treating each condition separately.

Finally, the model only analyzes short sound fragments. It does not provide the overview a physician uses, such as the patient's complete medical history or how their heart sound has changed over time, which is often necessary to make a definitive diagnosis.

### **6.5.2 Advancement Beyond Initial Proposal**

The evaluation had some inherent limitations. Using existing data for validation (retrospective validation) does not tell us how the system performs in a real-world clinical setting, unlike a prospective trial. There is also concern about the diversity of our training data; if it does not represent global populations, the model's accuracy could vary between different groups. Finally, the clinical diagnoses used as a reference are themselves simplified and often fail to capture the complexity and ambiguity present in real-world patient cases.

## **6.6 Future Research Directions**

### **6.6.1 Immediate Technical Extensions**

While the current system demonstrates high sensitivity, its clinical application is ultimately limited by false positives. Therefore, our most urgent next step is to refine model architecture, potentially through attention mechanisms or ensemble learning, to aggressively improve specificity. Achieving this would naturally enable a more ambitious goal: moving beyond simple anomaly detection to the precise identification of individual pathologies. Imagine a tool that doesn't just flag an abnormality but specifies, "This is likely Mitral Regurgitation," thereby transforming it from a screening assistant into a diagnostic partner.

The next major breakthrough will likely come from the fusion of two signals: cardiac sound (PCG) and the electrical signal of the heart (ECG), recorded simultaneously. This is the reason for its great effectiveness: the ECG provides a clear electrical map of the heartbeat, making it a perfect guide for determining the exact timing of the "lub" (S1) and the "dub" (S2). This solves the major problem of segmenting confusing heart sounds.

With these signals synchronized, a much smarter model can be built that not only listens to the sound but directly links it to the electrical event that caused it. This creates a complete, multisensory picture of the heart's function, which would represent a radical improvement in accuracy.

## **6.6.2 Clinical Translation Pathways**

The most immediate and crucial step is prospective clinical trials. This involves testing our tool with new patients in a real hospital to verify its effectiveness and identify any real-world problems.

Standardized platforms also need to be created to allow the AI tool to communicate seamlessly with the hospital's electronic health record system. If it isn't easy for clinicians to use within their current workflow, they simply won't adopt it.

In the longer term, supervised learning can be used. This is a smart technique that would allow us to train the AI on massive amounts of unlabeled heart sound data. This could finally solve our biggest obstacle: the chronic shortage of carefully annotated medical data.

## **6.7 Conclusion and Broader Impact**

### **6.7.1 Summary of Contributions**

The ECHO framework marks a substantial step forward in automated PCG analysis. Its integrated design directly tackles persistent obstacles in the field: noise robustness, class imbalance, and limited model generalization. With a demonstrated sensitivity of 87.85%, ECHO sets a new benchmark for detecting abnormal cases, and its rigorous evaluation underscores its potential for clinical use. The framework's key innovations—a novel quality control paradigm, hybrid feature engineering, and clinically-informed attention mechanisms—collectively enhance both its performance and its alignment with medical practice.

### **6.7.2 Potential Societal Impact**

In summary, the impact of this work goes beyond just numbers. The ECHO system lets us do reliable heart screenings with simple, affordable equipment, which can greatly improve healthcare in remote and low-income areas. Because it's so sensitive, it's really good at finding heart conditions early, when treatment works best. Finally, this project is a great example of how doctors and engineers need to team up to create practical AI tools that can actually help people around the world.





# Chapter 7

## Conclusion and Future Work

### 7.1 Summary of Research Contributions

This thesis has presented **ECHO**, a comprehensive deep learning framework designed for the robust classification of phonocardiogram (PCG) signals into normal and abnormal categories. Through systematic investigation and rigorous experimental validation, this work has demonstrated that integrating advanced signal processing with attention-based deep learning can significantly overcome the fundamental challenges inherent in automated PCG analysis: environmental noise, class imbalance, and limited generalization capabilities.

The research has made several distinct and measurable contributions to the fields of biomedical engineering and cardiac signal processing, advancing the state-of-the-art in automated auscultation.

#### 7.1.1 Technical Innovations

The first major contribution is the **Advanced Quality Control Framework**. The development of a multi-parameter quality assessment system with tiered thresholds (strict, balanced, lenient) provided a robust foundation for reliable analysis. Experimental results confirmed that the "Balanced" quality level ( $\text{SNR} \geq 3\text{dB}$ , Heart Energy Ratio  $\geq 0.15$ ) offered the optimal trade-off, successfully rejecting 19.4% of non-diagnostic segments while preserving 87.9% of the abnormal cases essential for training. This explicitly addresses the call in recent literature for automated quality assessment standards to prevent "garbage-in, garbage-out" scenarios in telemedicine [27].

Secondly, we developed a Hybrid Feature Engineering methodology that vertically stacks time-frequency (Mel-spectrograms), cepstral (MFCCs), wavelet, and statistical features into a comprehensive  $123 \times 79$  map. This multi-domain representation captures complementary signal characteristics that are inaccessible to single-domain models, yielding a 12.7% relative improvement in sensitivity over baseline spectrogram methods. This finding aligns with the recent work of Bahreini et al. (2025), which confirms that fusing handcrafted and deep features is key to improving diagnostic accuracy for complex, non-stationary PCG signals [75].

Thirdly, the work validated an Integrated Regularization Strategy. The combination of Stochastic Depth ( $p = 0.1$ ), MixUp augmentation ( $\alpha = 0.2$ ), and Label Smoothing specifically addressed the data scarcity common in medical datasets. This strategy reduced the performance gap between training and testing by 28.3%, significantly mitigating overfitting. The utility of MixUp for physiological signals, specifically for enhancing generalization in heart sound classification, aligns with the novel "PCGmix" protocols emerging in 2024 research [76].

Finally, to make the model smarter, we finished by adding a Convolutional Block Attention Module (CBAM). This basically taught the model to pay attention to the right things. The analysis showed it got 43% better at focusing on the key moments of a heartbeat (systole and diastole), which is where the important disease clues are. So, instead of getting distracted, the model learned to concentrate on what really matters for a diagnosis, just like the reference [77] suggests

### **7.1.2 Performance Achievements**

When was tested the ECHO model on the hidden PhysioNet 2016 dataset, it performed really well. Its most important result was an 87.85% Sensitivity score, which beats the 85% mark that doctors consider useful for screening. With a solid Accuracy of 76.01%, the model proved to be a dependable classifier. Crucially for real-world use, ECHO was hardly affected by noise its performance only decreased by 1.2% even with sound distortion. Overall, it was a 13.5% absolute improvement over standard methods, making ECHO a new benchmark for reliable detection in messy, real-life audio settings.

### **7.1.3 Clinical Engineering Contributions**

Beyond technical metrics, this work contributes to clinical engineering by demonstrating a Clinically-Informed Design. The alignment of the attention mechanisms with expert auscultation patterns validates the hypothesis that domain knowledge can effectively inform neural network architecture. Furthermore, the focus on Practical Deployment is evident in the system's computational efficiency, with an inference time of 47ms per segment, supporting realistic deployment in resource-constrained settings such as mobile health clinics [78].

## **7.2 Validation of Research Objectives**

The research successfully addressed all objectives outlined in the initial thesis proposal.

### **7.2.1 Primary Objectives Fulfilled**

The ECHO framework's architecture was built around four core technical innovations to ensure clinical-grade performance. First, an 'AdvancedQualityControl' system autonomously filtered out clinically unusable recordings through a multi-parameter assessment. Accepted signals were then processed by a novel Hybrid Feature Engineering pipeline that integrated four distinct domains, overcoming the limitations of single-domain analysis. For classification, a CBAM attention mechanism was implemented to focus computational resources on the most diagnostically relevant cardiac cycles. Finally, a comprehensive Advanced Regularization strategy was employed throughout the model, ensuring robust generalization to unseen patient populations and completing an end-to-end robust system.

### **7.2.2 Clinical Objectives Achieved**

From a clinical standpoint, our most important win was hitting an 87.85% sensitivity rate. This high level of detection is crucial for screening, where missing a diagnosis is simply not an option. We also proved the system is robust it held up well even in noisy, real-world environments. Finally, we made sure it's computationally efficient, meaning it can run on the kind of hardware you'd typically find in most clinics.

## **7.3 Limitations and Reflection**

The ECHO model represents notable progress; however, a critical assessment of its limitations is essential. Technically, its specificity of 63.74% falls below the accepted clinical benchmark of 70-75%. This was a deliberate trade-off to achieve high sensitivity, but the consequent drawback is an elevated risk of false positives, which could disrupt clinical workflows. Furthermore, the model's binary "Normal vs. Abnormal" classification is inherently limited. By grouping all pathological findings together, it lacks the diagnostic granularity to identify specific cardiac conditions.

The model's scope is also constrained. It analyzes brief, isolated audio segments and does not integrate patient history or other contextual data a holistic approach fundamental to physician-led assessment. Its validation is another consideration; while tested on retrospective datasets, its performance remains unproven in real-time clinical environments, making its efficacy in daily practice uncertain. Finally, the training data may not fully capture global population diversity. As a result, the model's accuracy could be inconsistent across different demographic groups, such as varying age or ethnicities [79], potentially limiting its broader applicability and equity.

## **7.4 Future Research**

### **7.4.1 Immediate Technical Extensions**

The crucial next step is the optimization of the specificity and this is essential by producing ensemble architectures or multi-task learning frameworks that can maintain high sensitivity while improving specificity. Furthermore, extending the framework to classify specific pathologies because for example, distinguishing between specific valve defects, would provide more actionable information to clinicians. Incorporating Temporal Context via recurrent connections (LSTMs) or Transformers to model dependencies across multiple cardiac cycles could also improve performance on arrhythmic patterns.

### **7.4.2 Clinical Pathways**

To make this technology truly ready for the clinic, we need to test it with real patients in real hospitals. This kind of real-world testing is the only way to prove it actually works and to spot problems that our current lab studies can't see [80]. Looking ahead, we should also try combining our system with ECG heart signals, as this could make our readings much more precise. Finally, to make sure doctors can use this, we'll need to build seamless connections to their existing electronic health record systems.

## **7.5 Broader Impact and Concluding Remarks**

### **7.5.1 Potential Societal Impact**

The ECHO framework demonstrates significant potential for positive societal impact. By enabling reliable cardiac screening with basic recording equipment, the technology could expand Healthcare Accessibility in underserved regions. The high

sensitivity supports Early Detection, potentially enabling timely intervention and improved patient outcomes. Additionally, as a Cost-Effective Screening tool, it could reduce reliance on expensive diagnostic modalities for initial assessments.

### **7.5.2 Conclusion**

This thesis establishes that a deliberate synthesis of signal processing, clinical knowledge, and deep learning can yield automated PCG analysis with the reliability required for clinical use. The ECHO framework, therefore, marks a substantive advance toward making computational auscultation a viable tool for accessible and reliable cardiac screening.

This project shows that the best results are coming from combining different data features with attention deep learning and in the same way applying crucial techniques to prevent the model from memorizing the data, provides a powerful solution to the fundamental challenges of PCG analysis. While specific performance metrics present areas for continued optimization, the achieved sensitivity of 87.85% demonstrates strong potential for screening applications where detecting abnormal cases is paramount. As cardiovascular diseases continue to represent a leading global health challenge, technologies like ECHO that enable accessible, reliable cardiac assessment have the potential to make significant contributions to early detection and improved patient outcomes worldwide.

# References

- [1] World Health Organization. (2021). *Global atlas on cardiovascular disease prevention and control*. Geneva: World Health Organization.
- [2] Roth, G. A., Mensah, G. A., Johnson, C. O., Addolorato, G., Ammirati, E., Baddour, L. M., ... & Fuster, V. (2020). *Global burden of cardiovascular diseases and risk factors, 1990–2019*. Journal of the American College of Cardiology, 76(25), 2982-3021.
- [3] Global Burden of Disease Collaborative Network. (2020). *Global Burden of Disease Study 2019 (GBD 2019) Results* [Data set]. Institute for Health Metrics and Evaluation (IHME). <https://ourworldindata.org/grapher/cardiovascular-disease-death-rates>
- [4] Leatham, A. (1975). *Auscultation of the heart and phonocardiography*. British Heart Journal, 37(6), 629-634.
- [5] Hall, J. E., & Hall, M. E. (2020). *Guyton and Hall Textbook of Medical Physiology*. 14th Edition. Elsevier.
- [6] O'Toole, M. L., & Hillis, W. S. (1980). *The second heart sound: A fundamental of clinical cardiology*. Modern Concepts of Cardiovascular Disease, 49(8), 7-12.
- [7] Durand, L. G., & Pibarot, P. (1995). *Digital signal processing of the phonocardiogram: review of the most recent advancements*. Critical Reviews in Biomedical Engineering, 23(3-4), 163-219.
- [8] Tribouilloy, C., Shen, W. F., Slama, M. A., & Dufossé, H. (1991). *Assessment of severity of mitral regurgitation by the deceleration time of pulmonary venous flow*. European Heart Journal, 12(12), 1275-1278.
- [9] Otto, C. M., & Prendergast, B. (2014). *Aortic-valve stenosis—from patients at risk to severe valve obstruction*. New England Journal of Medicine, 371(8), 744-756.
- [10] Komorniczak, M. (2006, October 21). *Phonocardiograms from normal and abnormal heart sounds* [Image]. Wikimedia Commons. Retrieved November 25, 2024, from [https://commons.wikimedia.org/wiki/File:Phonocardiograms\\_from\\_normal\\_and\\_abnormal\\_heart\\_sounds.svg](https://commons.wikimedia.org/wiki/File:Phonocardiograms_from_normal_and_abnormal_heart_sounds.svg)

- [11] Shah, S. J., & Michaels, A. D. (2001). *Hemodynamic correlates of the third heart sound and systolic time intervals*. Congestive Heart Failure, 7(3), 146-152.
- [12] Drazner, M. H., Rame, J. E., Stevenson, L. W., & Dries, D. L. (2001). *Prognostic importance of elevated jugular venous pressure and a third heart sound in patients with heart failure*. New England Journal of Medicine, 345(8), 574-581.
- [13] Aronow, W. S., & Kronzon, I. (1991). *Prevalence and severity of valvular aortic stenosis determined by Doppler echocardiography and its association with echocardiographic and electrocardiographic left ventricular hypertrophy and physical signs of aortic stenosis in elderly patients*. The American Journal of Cardiology, 67(8), 776-777.
- [14] Vukanovic-Criley, J. M., Criley, S., Warde, C. M., Boker, J. R., Guevara-Matheus, L., Churchill, W. H., ... & Criley, J. M. (2006). *Competency in cardiac examination skills in medical students, trainees, physicians, and faculty: a multicenter study*. Archives of Internal Medicine, 166(6), 610-616.
- [15] Mangione, S., & Nieman, L. Z. (1997). *Cardiac auscultatory skills of internal medicine and family practice trainees: a comparison of diagnostic proficiency*. JAMA, 278(9), 717-722.
- [16] Lok, C. E., Morgan, C. D., & Ranganathan, N. (1998). *The accuracy and interobserver agreement in detecting the 'gallop sounds' by cardiac auscultation*. Chest, 114(5), 1283-1288.
- [17] Tavel, M. E. (1996). *Cardiac auscultation: a glorious past—and it does have a future!*. Circulation, 93(6), 1250-1253.
- [18] Roy, D. L., & Sargent, M. A. (2000). *A randomized controlled trial of a computer-based intervention to improve residents' cardiac examination skills*. Academic Medicine, 75(7), 749-750.
- [19] Etchells, E., Bell, C., & Robb, K. (1997). *Does this patient have an abnormal systolic murmur?*. JAMA, 277(7), 564-571.
- [20] Springer, D. B., Tarassenko, L., & Clifford, G. D. (2016). *Logistic regression-HSMM-based heart sound segmentation*. IEEE Transactions on Biomedical Engineering, 63(4), 822-832.
- [21] Bentley, P. M., McDonnell, J. T., & Grant, P. M. (1998). *Time-frequency and time-scale techniques for the classification of native and bioprosthetic heart valve sounds*. IEEE Transactions on Biomedical Engineering, 45(1), 125-128.
- [22] Clifford, G. D., Liu, C., Moody, B., Lehman, L. H., Silva, I., Li, Q., ... & Mark, R. G. (2016). *AF classification from a short single lead ECG recording: the*



- PhysioNet/computing in cardiology challenge 2017*. In 2017 Computing in Cardiology Conference (CinC).
- [23] Johnson, J. M., & Khoshgoftaar, T. M. (2019). *Survey on deep learning with class imbalance*. Journal of Big Data, 6(1), 1-54.
  - [24] Schmidt, S. E., Holst-Hansen, C., Graff, C., Toft, E., & Struijk, J. J. (2010). *Segmentation of heart sound recordings by a duration-dependent hidden Markov model*. Physiological Measurement, 31(4), 513.
  - [25] Bentley, P., & Nordehn, G. (2008). *A machine learning approach to classifying heart sound segments*. In 2008 Computers in Cardiology (pp. 373-376). IEEE.
  - [26] Rubin, J., Abreu, R., Ganguli, A., Nelaturi, S., Matei, I., & Sricharan, K. (2016). *Classifying heart sound recordings using deep convolutional neural networks and mel-frequency cepstral coefficients*. In 2016 Computing in Cardiology Conference (CinC) (pp. 813-816). IEEE.
  - [27] Potes, C., Parvaneh, S., Rahman, A., & Conroy, B. (2016). *Ensemble of feature-based and deep learning-based classifiers for detection of abnormal heart sounds*. In 2016 Computing in Cardiology Conference (CinC) (pp. 621-624). IEEE.
  - [28] Demir, F., Şengür, A., & Bajaj, V. (2019). *Convolutional neural networks based efficient approach for classification of lung diseases*. Health Information Science and Systems, 7(1), 1-8.
  - [29] Li, F., & Yu, X. (2019). *Heart sound classification based on feature fusion using deep learning*. IEEE Access, 7, 129879-129888.
  - [30] Nogueira, D. M., et al. (2019). "Heart Sound Classification Using Deep Learning." *2019 4th International Conference on Smart and Sustainable Technologies*
  - [31] Woo, S., Park, J., Lee, J. Y., & Kweon, I. S. (2018). *Cbam: Convolutional block attention module*. In Proceedings of the European Conference on Computer Vision (ECCV) (pp. 3-19).
  - [32] Zhang, H., Cisse, M., Dauphin, Y. N., & Lopez-Paz, D. (2018). *"mixup: Beyond Empirical Risk Minimization."* International Conference on Learning Representations (ICLR).
  - [33] Latif, S., Usman, M., Rana, R., & Qadir, J. (2018). *Phonocardiographic sensing using deep learning for abnormal heartbeat detection*. IEEE Sensors Journal, 18(22), 9393-9400.

- [34] Huang, G., Sun, Y., Liu, Z., Sedra, D., & Weinberger, K. Q. (2016). *Deep networks with stochastic depth*. In *European Conference on Computer Vision* (pp. 646-661). Springer.
- [35] Li, X., Pang, T., Liu, W., & Wang, T. (2019). *Heart sounds classification based on feature selection and neural network*. *Computer Methods and Programs in Biomedicine*, 173, 33-43.
- [36] Choi, S., & Jiang, Z. (2008). *Comparison of envelope extraction algorithms for cardiac sound signal segmentation*. *Expert Systems with Applications*, 34(2), 1056-1069.
- [37] Kumar, D., Carvalho, P., Antunes, M., Henriques, J., Eugenio, L., Schmidt, R., & Habetha, J. (2006). *Noise detection during heart sound recording using periodicity signatures*. *Physiological Measurement*, 27(6), 535.
- [38] Gupta, C. N., Palaniappan, R., Swaminathan, S., & Krishnan, S. M. (2007). *Neural network classification of homomorphic segmented heart sounds*. *Applied Soft Computing*, 7(1), 286-297.
- [39] Gharehbaghi, A., Ekman, I., Ask, P., & Janerot-Sjoberg, B. (2014). *A novel method for discrimination between innocent and pathological heart murmurs*. *Medical Engineering & Physics*, 36(5), 668-676.
- [40] Langley, P., & Murray, A. (2011). *Limitations of hidden Markov models for heart sound segmentation*. In *2011 Computing in Cardiology* (pp. 333-336). IEEE.
- [41] Nilanon, T., Yao, J., Hao, J., Purushotham, S., & Liu, Y. (2016). *Normal/abnormal heart sound recordings classification using convolutional neural network*. In *2016 Computing in Cardiology Conference (CinC)* (pp. 585-588). IEEE.
- [42] He, K., Zhang, X., Ren, S., & Sun, J. (2016). *Deep residual learning for image recognition*. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 770-778).
- [43] Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., ... & Rabinovich, A. (2015). *Going deeper with convolutions*. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 1-9).
- [44] Tschannen, M., Kramer, T., Marti, G., Heinzmann, M., & Wiatowski, T. (2016). *Heart sound classification using deep structured features*. In *2016 Computing in Cardiology Conference (CinC)* (pp. 565-568). IEEE.
- [45] Messner, E., Zohrer, M., & Pernkopf, F. (2018). *Gated recurrent neural networks for phonocardiogram classification*. In *2018 40th Annual International Conference of*

- the IEEE Engineering in Medicine and Biology Society (EMBC) (pp. 101-104). IEEE.
- [46] Ko, T., Peddinti, V., Povey, D., & Khudanpur, S. (2015). *Audio augmentation for speech recognition*. In Sixteenth Annual Conference of the International Speech Communication Association.
  - [47] Park, D. S., Chan, W., Zhang, Y., Chiu, C. C., Zoph, B., Cubuk, E. D., & Le, Q. V. (2019). *SpecAugment: A simple data augmentation method for automatic speech recognition*. arXiv preprint arXiv:1904.08779.
  - [48] Golany, T., & Radinsky, K. (2019). *PGANs: Personalized generative adversarial networks for ECG synthesis to improve patient-specific deep ECG classification*. In Proceedings of the AAAI Conference on Artificial Intelligence (Vol. 33, No. 01, pp. 557-564).
  - [49] Larsson, G., Maire, M., & Shakhnarovich, G. (2016). *Fractalnet: Ultra-deep neural networks without residuals*. arXiv preprint arXiv:1605.07648.
  - [50] Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., & Wojna, Z. (2016). *Rethinking the inception architecture for computer vision*. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (pp. 2818-2826).
  - [51] Liu, C., & Li, P. (2020). *Cardiac sound characteristic extraction and classification using deep learning*. Biomedical Signal Processing and Control, 57, 101767.
  - [52] Leng, S., Tan, R. S., Chai, K. T. C., Wang, C., Ghista, D., & Zhong, L. (2015). *The electronic stethoscope*. Biomedical Engineering Online, 14(1), 1-35.
  - [53] Li, X., Pang, T., Liu, W., & Wang, T. (2019). *Heart sounds classification based on feature selection and neural network*. Computer Methods and Programs in Biomedicine, 173, 33-43.
  - [54] Wang, Y., et al. (2023). "Assistive diagnostic technology for congenital heart disease based on fusion features and deep learning."
  - [55] Azam, F. B., Ansari, M. I., McLane, I., & Hasan, T. (2021). *Heart Sound Classification Considering Additive Noise and Convolutional Distortion*. IEEE Access.
  - [56]. Li, T., et al. (2021). "Lightweight End-to-End Neural Network Model for Automatic Heart Sound Classification."
  - [57] Li, Z., et al. (2022). "CNN-Based Heart Sound Classification with an Imbalance-Compensating Weighted Loss Function."

- [58] Li, F., et al. (2020). "Heart sound classification based on improved focal loss and squeeze-and-excitation residual network." *Applied Sciences*, 10(4).
- [59] Al-Naami, B., Fraihat, H., Al-Nabulsi, J., & Al-Hinnawi, A. (2020). Assessment of Dual Tree Complex Wavelet Transform to improve SNR in collaboration with Neuro-Fuzzy System for Heart Sound Identification. *Research Square*.
- [60] Sahoo, S., Thakur, K. K., & Jain, P. K. (2024). A robust to noise classification method for the heart sound signals using deep learning technique. *Advances in Artificial Intelligence*, 4, 101.
- [61] Li, F., Tang, H., Shang, S., Mathiak, K., & Cong, F. (2020). Classification of Heart Sounds Using Convolutional Neural Network. *Applied Sciences*, 10(11), 3956.
- [63] Huang, G., Liu, Z., Weinberger, K. Q., & van der Maaten, L. (2017). Densely Connected Convolutional Networks. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [64] Kachuee, M., et al. (2021). "Augmentation of ECG and PCG signals using MixUp." *IEEE Transactions on Biomedical Engineering*.
- [65] Lin, T. Y., Goyal, P., Girshick, R., He, K., & Dollár, P. (2017). "Focal Loss for Dense Object Detection." *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2980-2988.
- [66] Roberts, M., & Driggs, D. (2022). "Common pitfalls and recommendations for using machine learning in medical imaging." *Nature Machine Intelligence*, 4(1), 17-25.
- [67] Liu, C., et al. (2016). *An open access database for the evaluation of heart sound algorithms*. *Physiological Measurement*, 37(12), 2181.
- [68] Ma, K., Lu, J., & Lu, B. (2023). Parameter-Efficient Densely Connected Dual Attention Network for Phonocardiogram Classification. *IEEE Journal of Biomedical and Health Informatics*.
- [69] Chicco, D., & Jurman, G. (2020). *The advantages of the Matthews correlation coefficient (MCC) over F1 score and accuracy in binary classification evaluation*. *BMC Genomics*, 21(1), 6.
- [70] World Health Organization (WHO). (2020). *Rheumatic heart disease: Priority for global health coverage*. Geneva: WHO.
- [71] Kagiya, N., et al. (2019). Machines Are Learning Chest Auscultation: Will They Also Become Our Teachers? *Circulation: Cardiovascular Imaging*.

- [72] Shorten, C., & Khoshgoftaar, T. M. (2019). "A survey on image data augmentation for deep learning." *Journal of Big Data*, 6(1), 1-48.
- [73] Moody, G. B., et al. (2016). "The PhysioNet/Computing in Cardiology Challenge 2016." *Physiological Measurement*, 37(12).
- [74] Mangione, S., & Nieman, L. Z. (1997). "Cardiac auscultatory skills of internal medicine and family practice trainees." *JAMA*, 278(9).
- [75] Bahreini, M., Barati, R., & Kamali, A. (2025). "Cardiac sound classification using a hybrid approach: MFCC-based feature fusion and CNN deep features." ResearchGate Preprint.
- [76] Susic, D., Gradišek, A., & Gams, M. (2024). PCGmix: A data-augmentation method for heart-sound classification. *IEEE Journal of Biomedical and Health Informatics*, (11), 6877–6888.
- [77] Huai, X., et al. (2025). "Heart sound classification based on convolutional neural network with convolutional block attention module." *Frontiers in Physiology*, 16.
- [78] Sharma, D., et al. (2020). "Digital Stethoscope with AI Diagnosis app for Respiratory Ailments and Heart Rate Monitoring." *Journal of Emerging Technologies and Innovative Research*, 10.
- [79] Chen, I. Y., et al. (2021). "Ethical machine learning in healthcare." *Annual Review of Biomedical Data Science*, 4, 123-144.
- [80] Krones F, Walker B et al. (2024). "From theoretical models to practical deployment: A perspective and case study of opportunities and challenges in AI-driven cardiac auscultation research for low-income settings." *PLOS Digit Health* 3(12): e0000437.