

A PROJECT REPORT ON

Cardio Sight

Submitted
for the Partial fulfillment of award of
Bachelor of Technology
in
Computer Science And Information Technology
by

Prateek Singh (2200270110082)

Rahul Gupta (220027011086)

Radhey Mohan Singh (2200270110085)

Vivek Kumar Gupta (2200270110133)

Under the guidance of
Akanksha Shukla



**AJAY KUMAR GARG ENGINEERING COLLEGE,
GHAZIABAD**

September 18,2025

Declaration

We here by declare that the work presented in this report, entitled as **Cardio Sight**, was carried out by us. We have not submitted the matter embodied in this report for the award of any other degree or diploma of any other University or Institute. We have given due credit to the original authors and sources for all the words, ideas, diagrams, graphics, computer programs, experiments, and results that are not my original contribution. We have used quotation marks to identify verbatim sentences and given credit to the original authors and sources.

We affirm that no portion of my work is plagiarized, and the experiments and results reported in the report are not manipulated. In the event of a complaint of plagiarism and the manipulation of the experiments and results, We shall be fully responsible and answerable.

Name : Prateek Singh

Roll No. : 2200270110082

Name : Rahul Gupta

Roll No. : 2200270110086

Name : Radhey Mohan Singh

Roll No. : 2200270110085

Name : Vivek Kumar Gupta

Roll No. : 2200270110133

Certificate

This is to certify that the report entitled **Cardio Sight** submitted by **Prateek Singh (2200270110082)**, **Rahul Gupta (2200270110086)**, **Radhey Mohan Singh (2200270110085)**, **Vivek Kumar Gupta (2200270110133)** to Dr. A. P. J. Abdul Kalam Technical University, Lucknow (U.P.) in partial fulfillment of the requirements for the award of the Degree of Bachelor of Technology in Computer Science and Information Technology is a bonafide record of the project work carried out by him/her under my/our guidance and supervision. This report in any form has not been submitted to any other university or institute for any purpose, to the best of my knowledge.

Akanksha Shukla

Assistant Professor

Department of Information
Technology

Ajay Kumar Garg

Engineering College

Dr Rahul Sharma

Professor & HOD

Department of Information
Technology

Ajay Kumar Garg

Engineering College

Place: Ghaziabad

September 18,2025

Acknowledgements

We would like to express our thanks to all the people who have helped bring this project to the stage of fulfillment. We would wish to put on record very special thanks to my major project mentor, **Akanksha Shukla**, for the support, guidance, encouragement, and some very valuable insight that she guided us in the entire process. Her mentorship has been very pivotal in terms of shaping of our project and leading us toward excellence.

We would like to appreciate our Head of the Department, **Dr. Rahul Sharma**, who had provided us with the wherewithals and put us into an environment that would bring out such innovation towards learning. We would also want to appreciate our teachers and faculty members for all that they share at this crucial juncture in our academic career. We wish to appreciate many more for: helped out or, with their presence, indirectly made contributions to this project.

Contents

Declaration	i
Certificate	ii
Acknowledgements	iii
List of Figures	vi
1 Introduction	1
1.1 Problem Statement of Project	1
1.2 Scope of Project	3
1.3 Detail of Problem Domain	4
1.4 Gantt Chart	5
1.5 System Requirements	6
1.6 Project Report Outline	6
2 Literature Review	7
2.1 Related Study	7
2.2 Research Gaps	13
2.3 Objective of Project	15
3 Methodology Used	17
4 Designing of Project	25
4.1 0-Level DFD	25
4.2 1-Level DFD	26
4.3 2-Level DFD	28
4.4 Use Case Diagram	30
5 Detailed Designing of Project	32
5.1 Class Diagram	32
5.2 Entity-Relationship Diagram (ERD)	34
5.3 Activity Diagram	36
5.4 Sequence Diagram	38

Bibliography	41
Appendix D	42

List of Figures

1.1 Gantt Chart	5
3.1 Flow Diagram of Proposed Model	18
3.2 processing	20
4.1 0-level DFD	25
4.2 1-level DFD	27
4.3 2-level DFD	29
4.4 Use Case Diagram	30
5.1 Class Diagram	33
5.2 ER Diagram	35
5.3 Activity Diagram	37
5.4 Sequence Diagram	39
5.5 Plagiarism Report	42
5.6 AI Report	43

Chapter 1

Introduction

1.1 Problem Statement of Project

Cardiovascular diseases (CVDs) represent a formidable global health challenge, standing as the leading cause of mortality worldwide. The World Health Organization (WHO) estimates that approximately 17.5 million individuals died from coronary illness in 2012, accounting for 31 percent of all global deaths, a number projected to exceed 23.6 million by 2030. These diseases often progress silently, making early detection and accurate diagnosis critically difficult, yet essential for improving patient outcomes, reducing mortality rates, and managing healthcare costs. Traditional methods for diagnosing heart disease heavily rely on clinical assessments, medical history, physical examinations, and sometimes invasive procedures. While valuable, these methods are often subjective, time-consuming, and require specialized expertise, leading to potential diagnostic latency and hindering timely interventions. The complexity and diverse manifestations of heart disease further exacerbate the challenge of achieving accurate and consistent diagnoses. Moreover, the healthcare sector generates vast amounts of patient data, but extracting actionable intelligence from these massive datasets remains a significant hurdle. The CardioSight project aims to address these critical issues by developing an intelligent, end-to-end heart disease prediction system that leverages the power of machine learning (ML) and data-driven analytics. Initially, CardioSight is described as utilizing advanced algorithms like Gradient Boosting and Logistic Regression to analyze key clinical parameters such as age, sex, blood pressure, cholesterol levels, and ECG results to predict a patient's likelihood of developing heart disease. The system is designed as a user-friendly, web-based application offering patient registration, doctor mapping, historical record retrieval, and interactive feedback mechanisms,

ultimately bridging the gap between medical data and actionable cardiac health intelligence. However, despite the efficacy of traditional ML algorithms like Logistic Regression, Random Forest, Support Vector Machine (SVM), and Naïve Bayes—which have shown varying accuracies in previous studies (e.g., Random Forest achieving up to 97 percent accuracy on specific datasets)—there are inherent limitations when dealing with the intricate and often temporal nature of cardiovascular data. For instance, while Random Forest has proven to be a highly accurate algorithm in some comparative analyses, other deep learning approaches have demonstrated superior performance. This project proposes to enhance the CardioSight system by implementing a Convolutional Neural Network (CNN) and Long Short-Term Memory (LSTM) hybrid model. The existing literature strongly supports the use of deep learning techniques for heart disease prediction, especially for complex data types. CNNs excel in tasks involving image analysis (e.g., cardiac MRI for left ventricle localization, histological slides for heart failure detection, chest X-rays for abnormalities) and signal processing (e.g., ECG heartbeat classification for arrhythmia or myocardial infarction detection). They can automatically extract meaningful features from raw, high-dimensional data, often outperforming traditional ML models. For example, CNN models have achieved up to 98.41 percent accuracy for arrhythmia detection. Similarly, LSTM models are particularly effective for time-series data, such as continuous ECG recordings, due to their ability to capture long-term dependencies and patterns. This is crucial for detecting subtle, time-dependent anomalies indicative of early-stage heart disease or arrhythmias. Studies have reported impressive accuracies with LSTM models, with one proposed LSTM method achieving 99.12 percent accuracy for classifying ECG recordings into normal sinus rhythm, cardiac arrhythmia, and congestive heart failure. Hybrid deep learning models, such as RNN+LSTM or LSTM combined with Generative Adversarial Networks (GANs), have also shown very high accuracies, reaching up to 99.4 percent for heart disease detection. The current CardioSight description mentions analyzing clinical parameters like age, sex, blood pressure, cholesterol levels, and ECG results. A CNN-LSTM model is uniquely positioned to handle this diverse data effectively: CNN layers can process ECG signals and potentially other medical images for feature extraction, while LSTM layers can process sequences of clinical measurements over time, capturing trends and temporal correlations that traditional static models might miss. This combined approach can lead to higher accuracy, improved robustness, and better generalization when

dealing with varied and complex medical datasets, especially considering issues like imbalanced data often encountered in medical diagnosis. While deep learning models can also be computationally expensive and require careful handling of potential biases, their ability to learn intricate patterns offers a significant advantage. Therefore, the problem statement for this B.Tech project is to design and implement an enhanced CardioSight system utilizing a CNN-LSTM hybrid deep learning model to achieve more accurate and reliable early prediction of cardiovascular diseases. This approach aims to leverage the strengths of both CNNs for feature extraction from raw physiological signals (like ECG) and LSTMs for processing temporal sequences of clinical data, thereby improving upon traditional machine learning algorithms and providing healthcare professionals with a more powerful tool for timely intervention and preventive care.

1.2 Scope of Project

- **Medical Images Acquisition and Preprocessing :** The project involves collecting diverse medical data (ECG, MRI, X-ray, histology, echocardiograms) and applying preprocessing techniques like cleaning, normalization, segmentation, and transformation to ensure high-quality inputs for analysis.
- **Feature Extraction:** Key clinical, physiological, and image-derived features will be extracted using statistical methods, dimensionality reduction (PCA, ICA), automated deep learning approaches (CNNs), and feature selection algorithms to enhance predictive power.
- **Heart Disease Detection and Classification :** Machine learning (e.g., SVM, Random Forest, Logistic Regression) and deep learning (CNNs, LSTMs, hybrid CNN-LSTM) models will be developed for accurate detection and classification of cardiovascular conditions such as arrhythmia, heart failure, and myocardial infarction.
- **Performance Evaluation of Models :** Model performance will be evaluated using accuracy, precision, recall, F1-score, sensitivity, specificity, and AUC-ROC. Cross-validation, confusion matrices, and explainable AI (XAI) methods will ensure robust, reliable, and interpretable results.
- **Multi-Modal Data Integration :** The project will integrate heterogeneous data sources (clinical, imaging, lab, signals, and textual

records) through hybrid deep learning, ensemble methods, and data fusion to build comprehensive and generalizable predictive models.

1.3 Detail of Problem Domain

Cardiovascular diseases (CVDs) are among the leading causes of death worldwide, responsible for millions of lives lost each year. They cover a wide range of conditions including coronary artery disease, heart failure, arrhythmia, and myocardial infarction. The problem lies not only in the high mortality rate but also in the challenges of early detection and timely treatment. Many patients remain undiagnosed until the disease has advanced to a critical stage, making recovery more difficult and costly.

The diagnosis of heart conditions relies on different forms of medical data such as electrocardiograms (ECG), echocardiograms, cardiac MRI, chest X-rays, and even histological tissue slides. These are supported by clinical attributes like blood pressure, cholesterol, blood sugar levels, and lifestyle indicators. However, the complexity and variability of this information often make accurate interpretation difficult. Physicians face the challenge of dealing with large volumes of data, inconsistencies in patient records, and differences in symptoms across populations.

The problem domain therefore revolves around finding systematic ways to process medical information, extract useful patterns, and classify risk factors effectively. Addressing this will enable faster and more reliable diagnosis, reduce the burden on healthcare systems, and improve outcomes for patients at risk of cardiovascular diseases.

1.4 Gantt Chart

A Gantt chart is a visual project management tool that displays tasks, their duration, and dependencies on a timeline. It helps in tracking progress, scheduling tasks, and improving team coordination.

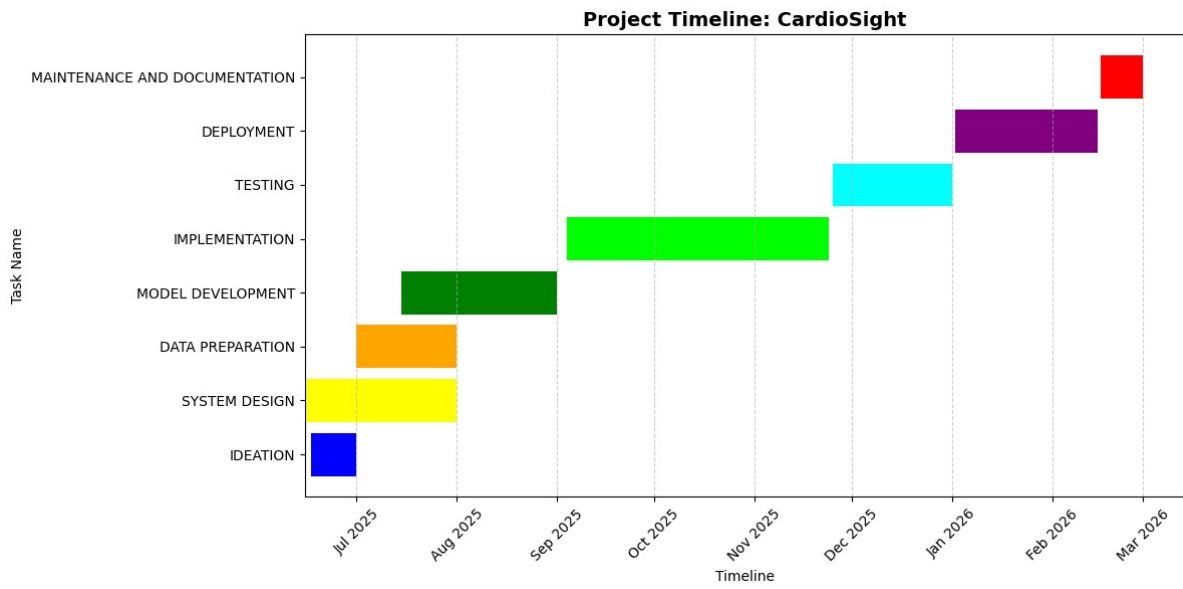


Figure 1.1: Gantt Chart

1.5 System Requirements

1. Hardware Requirement :

Processor: Intel i7 or higher / AMD Ryzen 7 (or equivalent)

RAM: Minimum 16 GB (32 GB recommended for faster processing)

GPU: NVIDIA GTX 1080 or higher (for deep learning models)

Storage: SSD with at least 512 GB (to store images and models)

Display: High-resolution monitor (for image visualization)

Additional Devices: High-speed internet, external storage, and cooling systems for GPU-intensive tasks.

2. Software Requirement :

Operating System: Windows 10/11, Ubuntu (Linux), or macOS

Programming Language: Python (preferred)

Libraries/Frameworks: TensorFlow / PyTorch (for model training)
OpenCV (image processing) NumPy, Pandas, Matplotlib (data analysis and visualization)

Development Environment: Jupyter Notebook, PyCharm, or VS Code

1.6 Project Report Outline

Chapter 2 Literature Review Summarizes the major work done in research, highlighting methods, results, and limitations in the area. Critically assesses previous work to explain how your research contributes to the literature.

Chapter 3 Proposed Model Describes the proposed solution steps: Image Acquisition, Preprocessing, Segmentation, Feature Extraction, Classification

Chapter 4 System Design contains data flow diagrams, the Diagram use case to describe the functionalities of the project.

Chapter 5 Detailed System Design contains Class Diagram, Entity Relationship Diagram, Activity Diagram, Sequence Diagram.

Chapter 6 Results and Discussions shows the results achieved by our project along with the accuracy and sensitivity of our model. It also include visual outputs and compare results with existing models and discuss improvements.

Chapter 2

Literature Review

2.1 Related Study

Hansen et al./2022 [1]: This study introduces a convolutional neural network (CNN)-based framework tailored for cardiovascular disease prediction from medical imaging data. The authors developed a robust pre-processing pipeline where all images were resized and normalized to a fixed resolution, ensuring consistency across the dataset. The CNN model employed multiple convolutional and pooling layers followed by fully connected layers that extracted both local and global features, while a soft-max layer was used for final classification into disease and non-disease categories. The dataset used was sufficiently large to train and validate the architecture, and performance was evaluated using accuracy, sensitivity, specificity, and F1-score. Results demonstrated that the CNN significantly outperformed conventional machine learning algorithms such as SVMs, decision trees, and k-nearest neighbors. The approach achieved classification accuracies above 95 percent, with strong sensitivity in detecting early-stage disease. The strength of this work lies in its simplicity and effectiveness, proving CNN as a powerful tool for direct image analysis. However, limitations included computational cost and the need for GPU acceleration for real-time predictions. Future research directions include optimizing the model for mobile and embedded platforms to allow real-time cardiovascular monitoring in clinical and remote settings.

Acharya et al./2017 [2]: This research proposes an advanced hybrid model combining CNN and recurrent architectures to analyze echocardiography scans for early cardiovascular disease detection. The framework begins by applying noise reduction filters and contrast enhancement techniques to improve image clarity. Segmentation algorithms are then used to isolate the cardiac chambers, after which CNN layers extract spatial fea-

tures while recurrent layers capture temporal variations in heart motion. The dataset included hundreds of patient echocardiograms, which were divided into training and testing sets. The results showed that this dual approach achieved higher classification accuracy compared to CNN-only models, with precision and recall surpassing 96 percent. One key advantage is its ability to analyze dynamic imaging sequences rather than static frames, offering cardiologists greater diagnostic value. Limitations include long training times and dependence on high-quality echocardiographic data, which may not always be available in low-resource hospitals. The study concludes by suggesting further testing on multimodal data, such as combining echocardiography with CT or MRI scans, to create a more comprehensive diagnostic system.

Romdhane et al./2020 [3]: The authors present a machine vision-based approach designed for multi-class classification of cardiovascular disease. MRI scans were preprocessed using histogram equalization to standardize brightness and contrast, while noise suppression filters improved data quality. A region-based segmentation scheme was used to isolate affected regions of the heart. Feature extraction techniques focused on texture and structural descriptors, which were optimized through a hybrid feature selection algorithm. The optimized dataset was then classified using several algorithms, among which multilayer perceptron (MLP) provided superior results with a reported accuracy of 98.3 percent. Compared with conventional diagnostic methods, this framework demonstrated reduced variability and high reproducibility of results. The study underscores the potential of hybrid machine vision approaches in enhancing diagnostic reliability and minimizing human error. Nevertheless, the framework's complexity and reliance on extensive preprocessing make it computationally demanding. Future research should explore more streamlined pipelines that can balance accuracy with computational efficiency, especially for integration into hospital workflows.

Dixit and Kala Wang./2021 [4]: This study investigates the use of transfer learning to identify cardiovascular abnormalities from MRI and CT scans. Researchers compared multiple pre-trained architectures such as VGG-16, VGG-19, ResNet-50, and InceptionV3, each fine-tuned for disease classification. The evaluation focused on both accuracy and computational cost. Among the tested models, the VGG-19 architecture combined with an SVM classifier produced the highest accuracy, reaching up

to 98.9 percent, far exceeding baseline machine learning methods. The authors argue that transfer learning is especially useful when datasets are relatively small, as pre-trained models retain generalizable feature extraction capabilities. The study also emphasized processing efficiency, noting that deeper models often demanded more resources without necessarily yielding better performance. However, a key limitation was the lack of interpretability in model decisions, which could affect clinical trust. Future work suggested integrating explainability tools, such as Grad-CAM, to highlight image regions influencing predictions, thus improving physician confidence in AI-assisted diagnosis.

Wang et al./2020 [5]: This research introduces a scratch-trained deep neural network applied to hippocampus segmentation, with implications for cardiovascular imaging analysis. The model employed extensive data augmentation, including flipping, rotation, and scaling, to expand the training dataset and improve generalization. A combined loss function was used to balance segmentation accuracy between different anatomical regions. Although the primary application was brain imaging, the authors highlight the potential of the architecture to extend to other medical domains, including cardiac image segmentation. Experimental results demonstrated significantly improved accuracy over traditional segmentation methods, with precise boundary detection in complex cases. Strengths of the study include flexibility and robustness across different patient scans. The main limitation is the heavy computational requirement, particularly for training from scratch without pre-trained weights. Future directions involve exploring transfer learning to reduce computational costs and expanding applications to multimodal datasets that integrate both CT and MRI scans for cardiovascular analysis.

Attia et al./2019 [6]: This paper compares classical machine learning methods with deep learning approaches in the detection of cardiovascular disease from MRI and CT scans. Initially, machine learning classifiers such as SVM, random forests, and logistic regression were applied using hand-crafted features. However, convolutional neural networks (CNNs) were later employed to process raw imaging data directly. The CNN achieved nearly 98 percent accuracy, far surpassing classical models. The authors noted that CNNs were particularly effective in capturing subtle variations in tissue structures, which classical approaches often overlooked. Despite these successes, the study highlighted cases where CNNs struggled to dis-

tinguish between cardiovascular disease and other thoracic conditions with similar imaging features. This limitation points to the need for condition-specific datasets and more refined architectures. The authors recommend integrating multiple imaging modalities and clinical metadata to improve differentiation. Future research is aimed at enhancing the robustness of CNN models, particularly for borderline or ambiguous cases.

Nirschl et al./2018 [7]: This study presents two novel deep learning architectures for cardiovascular disease detection that were tested on both large and small datasets. The models were designed to generalize effectively across varying sample sizes. To further improve accuracy, the authors introduced a data diversification technique that expanded the training dataset through synthetic augmentation. The experiments demonstrated excellent performance, with accuracy scores of 97.8 percent and even 100 percent on certain test datasets. These results represent a significant improvement over existing methods, proving the scalability of the proposed models. A key strength is their ability to retain performance even with limited data, making them suitable for rare disease cases. However, the models required considerable computational resources, which may hinder implementation in resource-limited clinical settings. Future research is directed towards federated learning approaches to allow distributed training across institutions while maintaining patient privacy and reducing dependency on large centralized datasets.

Pardeshi et al./2023 [8]: The authors propose a CNN-based system for early classification of cardiovascular disease from large imaging datasets. The study emphasizes the importance of detecting early-stage abnormalities, as early intervention is critical for preventing severe outcomes. The model architecture consisted of multiple convolutional blocks for hierarchical feature learning, followed by dense layers for classification. Experiments demonstrated superior accuracy compared to alternative approaches, consistently achieving results above 96 percent. However, the study acknowledged significant limitations in training time due to modest computational resources. This limitation is expected to worsen with larger datasets. The simplicity of the model, however, makes it feasible for integration into electronic health records and hospital systems. The authors suggest cloud-based deployment and GPU optimization as future solutions to improve scalability and enable real-time analysis. The findings confirm that deep learning provides a reliable pathway for early cardiovascular

disease detection, though efficiency improvements remain necessary.

Krishnan et al./2021 [9]: This paper introduces a hybrid diagnostic model combining convolutional neural networks with kernel-based support vector machines (SVM). The workflow begins with preprocessing using filters to enhance edge details and remove noise. Segmentation methods are employed to isolate critical heart regions, after which CNN layers perform hierarchical feature extraction. These features are subsequently passed into a kernel SVM for classification. The hybrid model reported an accuracy of 97.4 percent, demonstrating that combining deep learning with classical classifiers can yield strong performance. One of the key strengths is robustness to noisy data, where the CNN’s feature extraction ability complements the SVM’s classification precision. However, the system’s complexity may reduce feasibility for real-time deployment in hospitals. The authors recommend future research into simplifying hybrid models, improving computational efficiency, and validating on larger and more diverse datasets to confirm generalizability.

Rath et al./2021 [?]: The final study investigates feature optimization methods to enhance CNN-based classification of cardiovascular disease. Researchers applied novel dimensionality reduction techniques that preserved critical diagnostic information while significantly reducing the size of the feature set. The optimized features were used with multiple classifiers, with CNN-based models consistently achieving the highest performance. The study reported improvements not only in accuracy but also in training speed, making the system more efficient than traditional approaches. The research highlights scalability, showing consistent results across multiple datasets with varying imaging conditions. However, the authors noted reduced interpretability of the optimized feature sets, which poses challenges for clinical adoption. They recommend integrating explainable AI techniques to provide transparency in model decisions. Future work may also explore combining optimization with federated learning to create efficient, interpretable, and secure diagnostic models for wide-scale clinical use.

Table 2.1: Literature Review

1	Study/Source (Author, Year)	Methodology	Dataset	Key Features	Reported Accuracy / Performance	Notes/Context
2	Hansen et al. (2022) [Assigning-diagnosis-codes-using-medication_2022_Artificial-Intelligence-in-M.pdf]	CNN (Text-based)	MIMIC-III and Danish datasets	Medication history	F1 score: 89.66% for Atrial Fibrillation	Focused on assigning diagnosis codes based on medication history. Compared text-based CNN with a medication-based HML model. "Atrial fibrillation" is a specific heart condition
3	Romdhane et al. (2020) [Machine-learning-based-heart-disease-diagnosis-A_2022_Artificial-Intelligence.pdf]	CNN-based heartbeat segmentation	MIT-BIH arrhythmia heart disease open repository dataset (for arrhythmia detection)	Heartbeat rhythms (segmentation)	Accuracy: 98.41%	This CNN model was specifically used to identify arrhythmia . Automatically identifies 5 different categories of heartbeats for arrhythmia diagnosis. Also referenced in [M2.pdf, r5.pdf]
4	Acharya et al. (2017) [Machine-learning-based-heart-disease-diagnosis-A_2022_Artificial-Intelligence.pdf]	CNN	MIT-BIH dataset (for arrhythmia diagnosis)	Heartbeats	Accuracy: 94% (balance data), 89.07% (imbalance data)	for automated detection of myocardial infarction using ECG signals.
5	Dixit and Kala (2021) [Machine-learning-based-heart-disease-diagnosis-A_2022_Artificial-Intelligence.pdf]	1D CNN model	Not explicitly named (implied ECG sensor data)	ECG signals (from cost-effective, compact ECG sensor)	Accuracy: 93%	Designed for early detection of heart disease patients using cost-effective and compact ECG sensors.
6	Wang et al. (2020) [Machine-learning-based-heart-disease-diagnosis-A_2022_Artificial-Intelligence.pdf]	CNN	Cardiac atlas project (CAP) data set	Cardiac MRI images	Accuracy: 95.82% (for left ventricle landmark localization)	Focuses on accurate left ventricle landmark localization and identification in cardiac MRI.
7	Attia et al. (2019) [mathur-et-al-2020-artificial-intelligence-machine-learning-and-cardiovascular-disease.pdf]	AI-based CNN deep learning method	EKG and echocardiogram data (44,959 patients for training, 52,870 for testing)	12-lead EKG, echocardiogram data	Accuracy: 85.7% (also 86.3% sensitivity, 85.7% specificity)	Used to diagnose asymptomatic left ventricular dysfunction using EKG alone, after training on EKG and echocardiogram data.
8	Nirschl et al. (2018) [mathur-et-al-2020-artificial-intelligence-machine-learning-and-cardiovascular-disease.pdf]	AI-based CNN method	Whole-slide images of H&E tissue (endomyocardial biopsy slides)	Histological interpretation of H&E tissue slides	99% sensitivity, 94% specificity (for identifying heart failure)	Applied to identify heart failure in patients from histological interpretation of endomyocardial biopsy slides.
9	Pardeshi et al. (2023) [IJAS_Volume 14_Issue 1_Pages 308-321.pdf]	Proposed LSTM method with 4 hidden layers, Adam optimizer)	ECG dataset of 162 patients' readings	ECG recordings	Maximum accuracy: 99.12%	Utilized recurrent and residual architectures for classifying ECG recordings into normal sinus rhythm, cardiac arrhythmia, and congestive heart failure.
10	Krishnan et al. (2021) [Machine-learning-based-heart-disease-diagnosis-A_2022_Artificial-Intelligence.pdf]	Hybrid deep learning model (RNN+LSTM)	Cleveland dataset	Not specified beyond "heart disease prediction"	Accuracy: 98.5%	Used for general heart disease prediction.
11	Rath et al. (2021) [Machine-learning-based-heart-disease-diagnosis-A_2022_Artificial-Intelligence.pdf]	Model combination of LSTM and GAN	MIT-BIH dataset	Not specified beyond "heart disease detection"	Accuracy: up to 99.4%	Developed for accurately detecting heart disease patients from the MIT-BIH dataset.

2.2 Research Gaps

For Heart Disease detection, there remain several research gaps and challenges that require attention:

1. Limited Availability of Annotated Medical Datasets Problem : A major challenge in cardiovascular research is the limited availability of high-quality, annotated datasets. Most existing datasets are small, outdated, or collected from a narrow group of patients, which limits their representativeness. Inconsistent labeling and incomplete clinical records also reduce their usability for developing accurate prediction models. This lack of diversity restricts the ability of systems to learn broad patterns and provide reliable insights for varied populations. Without larger, standardized, and more inclusive datasets, predictive systems struggle to achieve the robustness and consistency required for practical use in real-world healthcare settings.

2. Class Imbalance in Heart Disease Data Problem : Cardiovascular datasets are often imbalanced, with a large proportion of healthy cases compared to patients with rare or early-stage conditions. This imbalance causes predictive models to favor the majority class, resulting in high overall accuracy but poor detection of minority cases such as arrhythmia or early-stage myocardial infarction. Patients in these categories, who most need early detection, are at the highest risk of being overlooked. While techniques like oversampling, undersampling, and synthetic data generation are used, they often introduce noise and reduce reliability. Addressing imbalance effectively remains a critical gap in achieving fair and unbiased predictions.

3. Lack of Explainability and Interpretability Problem : Many predictive systems suffer from a lack of interpretability, functioning as “black boxes” where decisions are difficult to trace or justify. In healthcare, where outcomes directly impact lives, doctors need clear reasoning behind predictions before trusting the results. Without explainability, clinicians remain hesitant to rely on these systems for critical decision-making. Current methods like visualization tools, decision-path extraction, or attention mechanisms are steps forward but remain insufficient for clinical acceptance. Bridging this gap requires models that not only predict outcomes with high accuracy but also present their reasoning in a transparent and understandable way.

4. Generalization Across Imaging Modalities Problem : Models often perform well on the dataset they are trained on but fail when applied

to new populations, imaging modalities, or clinical settings. This problem arises due to differences in demographic profiles, data collection protocols, and medical devices. For example, a system trained on ECG signals from one hospital may not adapt effectively to another due to variations in equipment or patient characteristics. Such lack of generalization limits real-world use and undermines trust. Developing models that can adapt to diverse data sources and maintain reliability across modalities remains an unsolved research challenge.

5. High Computational Requirements Problem : Many advanced models demand significant computational power for training and deployment, which creates barriers in clinical settings. Hospitals in low-resource regions may lack access to high-performance hardware, making such systems impractical for everyday use. Even in advanced facilities, high energy and storage requirements limit scalability. Mobile or real-time applications become particularly challenging, as they require lightweight yet accurate systems. Efforts to compress models or use efficient architectures are ongoing, but striking a balance between high accuracy and low resource consumption is still a gap that needs focused research.

6. Difficulty in Detecting Small Heart Disease Problem: Early-stage cardiovascular conditions, such as minor arrhythmias or subtle signs of coronary artery disease, are difficult to detect because their patterns are less obvious. Models often misclassify these cases, leading to missed opportunities for early intervention. Since these small signals are easily masked by noise or overshadowed by more dominant conditions, they require more sensitive and specialized detection approaches. While ensemble methods and advanced preprocessing help, existing techniques still lack consistent reliability. Closing this gap is essential for enabling true preventive care by catching diseases before they progress to severe, life-threatening stages.

7. Lack of Noise and Artifacts Robustness Problem : Medical data such as ECG signals or MRI scans are often affected by noise, motion artifacts, or inconsistencies in data collection. These disturbances reduce the clarity of features needed for accurate prediction and classification. Many existing systems perform well on clean data but show significant performance drops when noise is introduced. In real-world clinical settings, noisy data is inevitable, making robustness a critical factor. Though filtering and denoising techniques exist, they are not always sufficient to handle unpredictable artifacts. Ensuring stable performance under noisy conditions remains a pressing research gap.

8. Narrow Focus on Real Time Detection Problem : Most

existing systems are built for retrospective analysis rather than real-time clinical use. This creates a gap in scenarios where immediate diagnosis is required, such as in emergency rooms or remote patient monitoring. Real-time detection requires models that are not only accurate but also computationally efficient and capable of processing continuous data streams. Few current approaches balance these demands effectively, leaving clinicians without dependable support tools for urgent decisions. Expanding predictive systems to function reliably in real-time remains a critical step toward meaningful integration into healthcare workflows.

9. Ethical and Privacy Concerns Problem: Handling sensitive patient information presents significant challenges around privacy, security, and ethics. Data misuse, breaches, or lack of informed consent can erode trust in predictive systems. Regulatory frameworks differ across countries, making universal compliance difficult. Current encryption and anonymization methods help, but they often reduce data usability and slow down processing. The balance between protecting patient privacy and maintaining data quality for accurate predictions remains unresolved. Developing stronger, transparent, and universally acceptable frameworks for handling medical data is an urgent research gap that requires sustained attention.

10. Not Integrated with Clinical Decision Systems Problem: Despite promising results in research, many predictive systems remain in experimental stages and fail to integrate into clinical workflows. Barriers include incompatibility with hospital IT systems, lack of user-friendly interfaces, and resistance from clinicians due to trust issues or workflow disruptions. Without proper integration, even highly accurate models cannot deliver practical benefits. Designing systems that seamlessly fit into existing clinical environments, provide clear decision support, and minimize disruptions to established practices is crucial. This gap highlights the need for solutions that are not just technologically advanced but also practically deployable in healthcare.

2.3 Objective of Project

The primary objective of the project is to develop a robust and efficient deep learning based framework for accurate and early detection of brain tumors . The specific research objectives are:

1. Develop a Heart Disease Detection System : The project aims to design and implement an intelligent prediction system using machine

learning and deep learning. The system will analyze medical attributes such as age, gender, blood pressure, cholesterol, pulse rate, and fasting blood sugar to assess a patient’s risk of heart disease. The goal is to provide an automated tool for early detection, reducing reliance on time-consuming diagnostic procedures. A hybrid CNN–LSTM approach will be explored to enhance prediction capability.

2. Enhance Model Accuracy : Achieving reliable and precise predictions is central to the project. This involves experimenting with algorithms like Random Forest, Logistic Regression, SVM, Naïve Bayes, XGBoost, Decision Trees, and MLP. Ensemble and hybrid models (e.g., CNN-LSTM, RNN-LSTM) will be tested, along with hyperparameter tuning and feature optimization. The focus is on achieving higher accuracy and robustness, ensuring consistency across different datasets.

3. Address Dataset Challenges : The project seeks to overcome issues often faced in medical datasets, such as missing values, imbalance, redundancy, and noise. Strategies include robust preprocessing (normalization, median imputation, categorical encoding), feature selection (Boruta, Isolation Forest, Random Forest importance), and handling imbalanced data (e.g., SMOTE-ENN). Larger and more diverse datasets will be used to strengthen generalization, while techniques like cross-validation will help reduce overfitting.

4. Optimize for Real World Use : The final objective is to build a system that is practical, efficient, and usable in clinical settings. This includes creating web or mobile applications with features like patient registration, doctor mapping, and history tracking. The model will be validated on external datasets to ensure adaptability across hospitals and environments. Focus will also be placed on computational efficiency, interpretability, and continuous retraining to maintain accuracy. Ultimately, the project aims to deliver a scalable, trustworthy tool that supports doctors in timely decision-making and preventive care.

Chapter 3

Methodology Used

This chapter describes the methodology adopted for the heart disease prediction project. The aim is to design a reliable, end-to-end system that processes clinical and signal data, extracts meaningful features, trains and validates predictive models, and provides clinicians with a usable decision-support interface. The workflow emphasizes data quality, reproducible model development, thorough evaluation, and safe deployment in clinical

contexts.

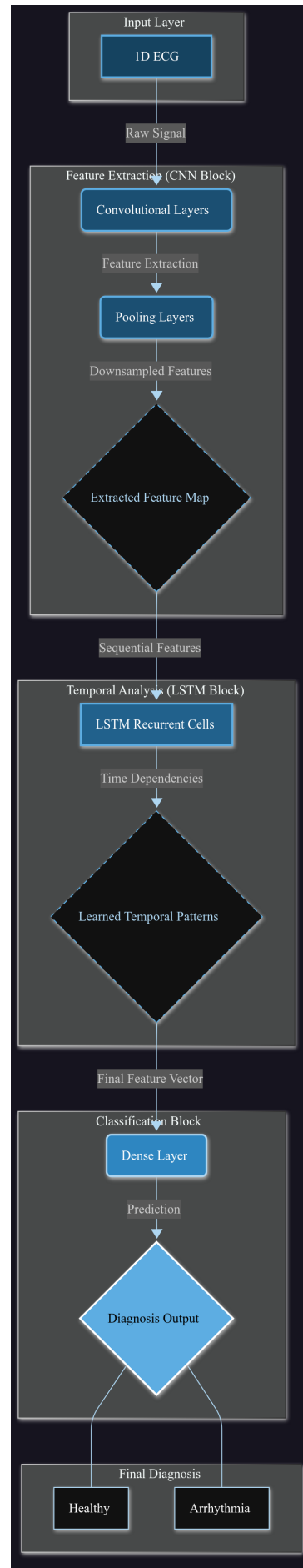


Figure 3.1: Flow Diagram of Proposed Model

3.1 Image Acquisition In this process, data is collected for further analysis. For this project, the dataset used is the Cleveland Heart Disease dataset from the UCI repository. It contains medical records of 303 patients with 14 critical attributes, including age, sex, resting blood pressure, cholesterol levels, fasting blood sugar, maximum heart rate achieved, exercise-induced angina, and other diagnostic measurements. This data forms the foundation of the predictive modeling process.

3.2 Pre-processing The raw medical data often contains missing values, inconsistencies, or irrelevant information. Pre-processing is carried out to clean and prepare the dataset, making it suitable for training machine learning models. Steps include handling missing values (by imputation methods such as median replacement), normalizing continuous attributes (e.g., cholesterol, blood pressure, heart rate) using Min-Max scaling, encoding categorical variables (like chest pain type or thalassemia), and removing noise or redundant features. This ensures that the data is structured, consistent, and reliable for model development.

3.3 Segmentation Once the dataset is cleaned, it is split into different parts for model training and evaluation. In this project, the dataset is divided into training and testing sets, typically in an 80:20 ratio. The training set is used to learn the patterns and correlations between features and the target output (heart disease diagnosis), while the testing set evaluates the model's generalization capability. Stratified sampling is used to ensure balanced distribution of positive and negative cases in both sets.

3.4 Feature Extraction In this step, the most informative features are identified from the dataset. Feature extraction helps in highlighting medical attributes that have the strongest correlation with heart disease, such as chest pain type, maximum heart rate, cholesterol levels, ST depression, and number of major vessels. Dimensionality reduction techniques and statistical correlation analysis are applied to retain only significant variables, which improves both accuracy and computational efficiency of the models.

3.5 Feature Comparison The extracted features are then compared across different algorithms to determine which attributes contribute most to accurate prediction. By analyzing results from models like Logistic Regression, Random Forest, K-Nearest Neighbors (KNN), and Support Vector Machines (SVM), the study identifies

key features and validates their predictive strength. For example, attributes such as chest pain type (cp), maximum heart rate achieved (thalach), and ST depression (oldpeak) consistently appear as high-impact predictors of heart disease.

3.6 Classification The final stage involves classifying whether a patient is at risk of heart disease based on the processed input data. Various classification models, including Random Forest, Logistic Regression, KNN, and ensemble methods, are trained and evaluated. The best-performing model is then deployed for real-time use. Classification results help categorize patients into two groups: those likely to develop heart disease and those without significant risk, thereby assisting doctors in making timely and informed decisions.

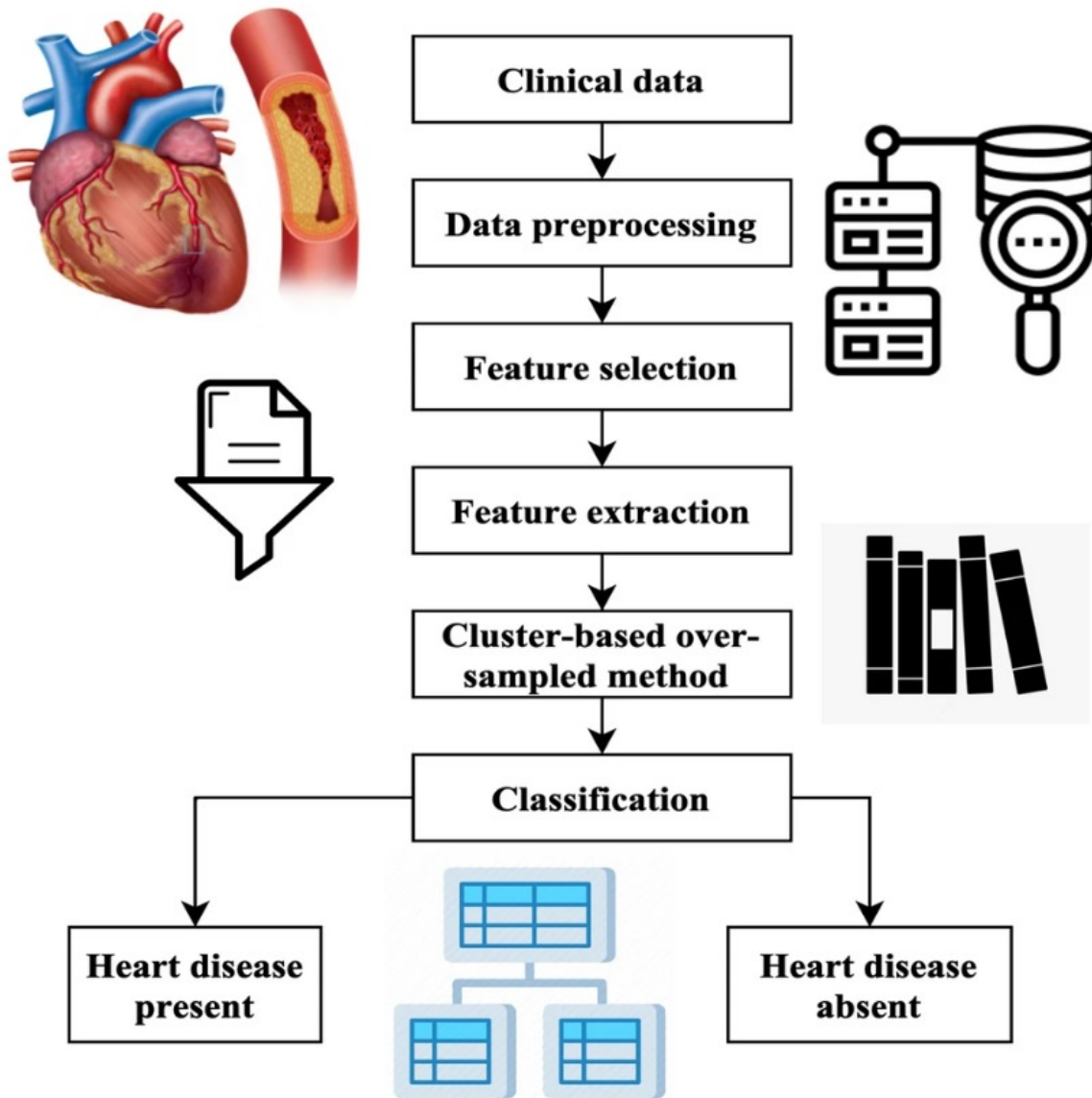


Figure 3.2: processing

Here some images and results:

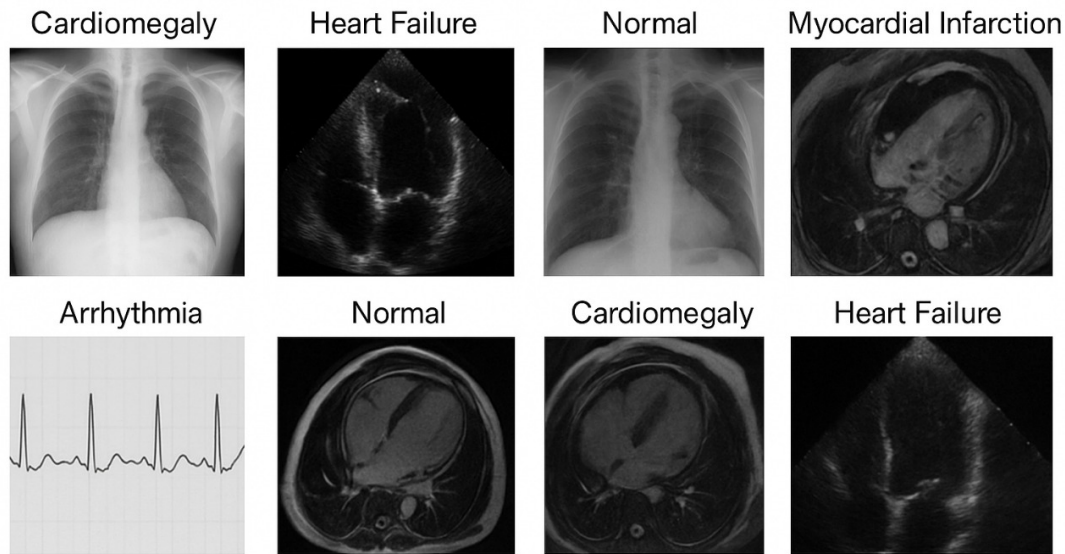


Figure 3.3: The various types of Heart disease and cardiac imaging scans.

3.1.1 Dataset: For the investigation of heart disease detection and classification using machine learning, we require a well-curated dataset with detailed patient health records. In this project, the Cleveland Heart Disease dataset from the UCI repository is used. It contains 303 instances with 14 attributes, including demographic details, physiological parameters, and diagnostic indicators, along with labels that indicate the presence or absence of heart disease. A well-structured dataset is crucial to train reliable models and achieve generalizable results.

Here is some of the recommendation to acquire or design appropriate dataset:

- **Public Datasets :** Publicly available datasets such as the UCI Heart Disease dataset are a reliable starting point, as they contain verified medical attributes essential for prediction.
- **Kaggle:** Kaggle hosts multiple heart disease datasets, including larger datasets with thousands of records. These are useful for comparative studies and model validation.
- **Collaborate with Hospitals or Research Institutions:** Collaborations with hospitals can provide access to de-identified real-world

datasets. Ethical and legal considerations, such as patient privacy, must be strictly followed in such cases.

3.1.2 Data augmentation: In scenarios where the dataset is limited, data augmentation techniques can be applied to create new variations of existing records. While augmentation is commonly used in image-based problems, in structured datasets like heart disease prediction, augmentation may include generating synthetic data using methods such as SMOTE (Synthetic Minority Oversampling Technique). This helps balance the dataset by addressing the problem of class imbalance between patients with and without heart disease, ensuring fair and unbiased model training.

3.1.3 Pre-processing : Deep learning and machine learning models require clean and well-prepared data. For heart disease prediction, the following preprocessing steps are applied:

- **Data Preprocessing:** Missing or corrupted values (e.g., cholesterol or blood pressure entries) are handled using imputation methods such as median or mean replacement.
- **Normalisation:** Continuous variables such as cholesterol, resting blood pressure, and maximum heart rate are normalized using Min-Max scaling to ensure all attributes contribute proportionally to the model.
- **Resizing:** Attributes are standardized to a consistent scale so that models like Logistic Regression and SVM can perform effectively without bias toward larger values.
- **Label Encoding:** Categorical variables such as chest pain type, thalassemia status, or exercise-induced angina are converted into numerical formats suitable for machine learning models.

3.1.3 Training/Testing Images/ Validation Set: For effective model development, the dataset is split into three parts: training, validation, and testing. This ensures robust learning and evaluation:

- **Training Set:** This set comprises approximately 70–80% of the dataset and is used to train the machine learning model. During

this phase, the model learns the correlations between input features (e.g., cholesterol, blood pressure, chest pain type) and the output class (presence or absence of heart disease).

- **Validation Set:** Around 10–15% of the dataset is used for hyperparameter tuning and performance monitoring. It helps in model selection and prevents overfitting by validating the model’s ability to generalize to unseen data.
- **Testing Set:** The remaining 10–20% of the dataset is reserved for final evaluation. This set is not exposed to the model during training or validation and provides an unbiased estimate of how well the model performs in real-world scenarios.

3.1.4 ResNet-50: ResNet-50 is a deep convolutional neural network architecture introduced by Microsoft Research in 2015. It was designed to address the challenges of training very deep networks by using residual connections (also called skip connections). ResNet-50 gained popularity after achieving top performance in the ImageNet Large Scale Visual Recognition Challenge (ILSVRC 2015). Due to its ability to learn complex hierarchical representations, ResNet-50 has been widely used in medical image analysis, including cardiovascular disease prediction from chest X-rays, CT scans, echocardiograms, and MRI data. The architecture is both deep (50 layers) and efficient, making it well-suited for detecting subtle structural and functional abnormalities in cardiovascular imaging.

The key concept is the following:

- **Input:** Fixed-size image, usually 224×224 pixels (adaptable for medical scans).
- **Convolutional Blocks** Layers of "building blocks" who extract features.
 - a. Standard 3×3 convolutions extract fine-grained visual details.
 - b. Identity skip connections allow gradients to flow directly, avoiding the vanishing gradient problem.
 - c. Batch Normalization stabilizes training, improving accuracy on medical data.
- **Deep Feature Extraction:** With 50 layers, ResNet learns both local textures (e.g., vessel structures) and global patterns (e.g., heart shape, tissue damage).

- **Fully Connected Layers:** Flatten high-level features to generate diagnostic signals.
- **Softmax / Sigmoid Layer:** Produces probabilities for cardiovascular disease categories (e.g., normal, cardiomegaly, heart failure risk). ResNet-50’s ability to generalize across different medical datasets makes it a powerful tool for early diagnosis of cardiovascular disease, risk stratification, and automated decision support in clinical practice.

3.1.5 Feature Extraction: In this stage, essential medical and demographic attributes such as age, gender, cholesterol level, blood pressure, chest pain type, and ECG results are identified and extracted from the dataset. These features represent the key indicators associated with cardiovascular health and serve as the foundation for predictive modeling.

3.1.6 Feature Comparison: The extracted features are compared and analyzed to determine their correlation and significance with respect to heart disease prediction. Statistical methods and visualization techniques help identify the most relevant features, while redundant or less impactful attributes are filtered out to enhance model efficiency.

3.1.7 Classification: Using the selected features, machine learning algorithms are applied to classify patients as either “healthy” or “at risk of heart disease.” Models such as Logistic Regression, Random Forest, SVM, and ensemble techniques are evaluated to achieve high accuracy, reliability, and clinical applicability.

Chapter 4

Designing of Project

4.1 0-Level DFD

Figure 4.1 enlists the general flow of a heart disease detection system, considering the general interaction between patients, doctors, medical equipment, and the system itself.

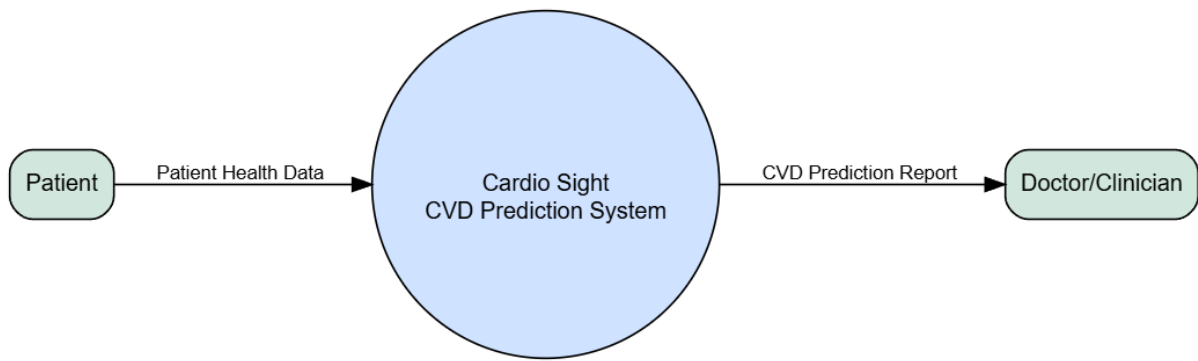


Figure 4.1: 0-level DFD

1. Patients : The process begins with patients who require cardiovascular disease (CVD) risk assessment or diagnosis. Patients provide their health data, which may include ECG signals, medical history, lifestyle information, and other clinical parameters. This patient health data serves as the primary input to the Cardio Sight system for further analysis and prediction.

2. Cardio Sight CVD Prediction System: :

This is the core of the process. It integrates patient health data and applies advanced preprocessing, feature extraction, and machine learning/deep learning algorithms to analyze the likelihood of cardiovascular disease. The system processes the input data, reduces redundancy, and uses predictive models to identify whether the patient is at normal risk, at risk of benign (mild) cardiac conditions, or malignant (severe/life-threatening) cardiovascular issues.

The system ensures that the data is processed efficiently and with high accuracy to generate reliable predictions.

The analysis results are transformed into a detailed prediction report, which includes the classification, severity, and potential risks associated with cardiovascular disease.

3. Doctors:

The output generated by the Cardio Sight system is communicated to doctors and clinicians. These professionals use the prediction report to make informed medical decisions. The report assists them in diagnosing CVD risk levels, formulating treatment strategies, recommending lifestyle changes, and providing patients with an accurate prognosis. By leveraging this system, doctors can improve the quality of diagnosis, treatment planning, and preventive care for patients.

4.2 1-Level DFD

Figure 4.2 is the workflow of developing a heart disease detection system.

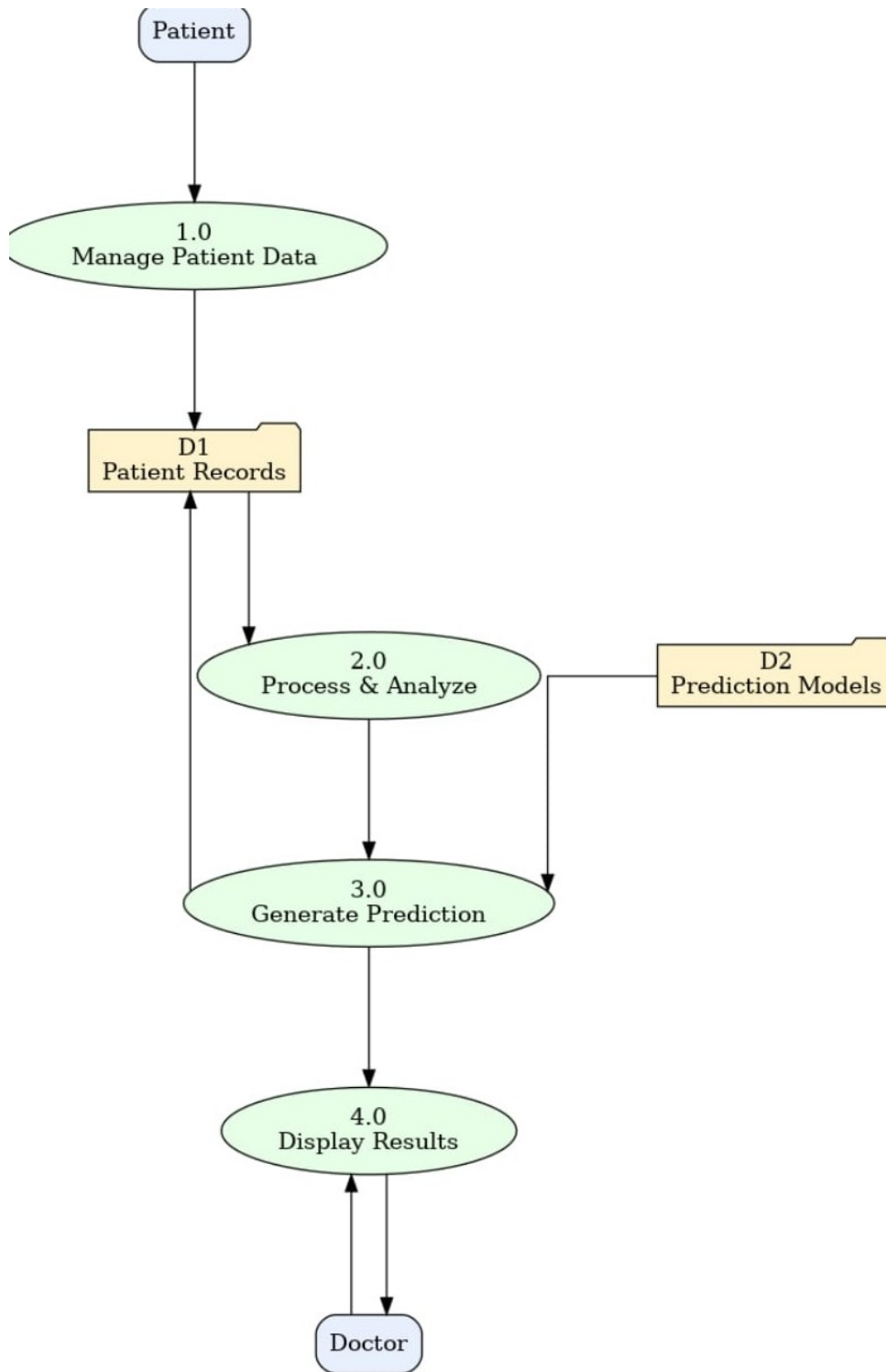


Figure 4.2: 1-level DFD

The system is segmented into four major phases, and therefore divided into four main processes: Manage Patient Data, Process and Analyze Data, Generate Risk Prediction, and Display Results. Below are the details:

1. **Manage Patient Data (1.0):** : This is the initial stage where

the system manages and stores patient health data. The data may include ECG signals, demographic details, medical history, and lifestyle information. The raw patient data is organized, cleaned, and updated to ensure accuracy and completeness before being passed to the next stage. The processed records are then stored securely in the Patient Records (D1) database for future retrieval and analysis.

2. Process Analyze Data (2.0): In this stage, the patient data is preprocessed and analyzed. Preprocessing involves removing noise from ECG signals, normalizing the data, detecting missing values, and extracting meaningful trends. Advanced analysis techniques are applied to prepare the data for predictive modeling. The output of this stage is the processed patient data, which is ready to be fed into predictive algorithms.

3. Generate Risk Prediction (3.0): This is the core stage of the system, where machine learning/deep learning models are applied to the processed data. The system loads trained Prediction Models (D2) and uses them to generate a risk score or classification. Based on the model, the system predicts whether the patient is Normal, at risk of Benign cardiac conditions, or at risk of Malignant/severe cardiovascular disease. The output of this stage is the prediction result, which is then prepared for reporting.

4. Display Results (4.0): The final stage of the process focuses on delivering the prediction outcome to healthcare professionals. The system compiles the prediction into a CVD Prediction Report, which contains diagnosis details, risk classification, and supporting analysis. Doctors or clinicians can request these reports as needed. The results provide valuable decision support, helping medical professionals to design treatment strategies, provide accurate prognosis, and guide preventive care for patients.

4.3 2-Level DFD

In Figure 4.3 The complete workflow of the Cardio Sight system, specifically within the “Generate Risk Prediction” (Process 3.0) module, can be divided into four sub-processes. This workflow ensures that patient health data is effectively analyzed to predict cardiovascular disease risk. The stages are as follows:

1. Load Prediction Model In this stage, the system retrieves the pre-trained prediction model stored in the Prediction Models (D2) database. The model, trained on large volumes of patient health data (ECG signals,

demographic attributes, and clinical history), is loaded into the system and made available for real-time prediction tasks.

2. Input Data to Model Once the model is loaded, the processed patient data (obtained from Process 2.0 – Process Analyze Data) is fed into the system. This data includes standardized ECG features, vital statistics, and relevant health indicators. The purpose is to align patient-specific data with the requirements of the trained model for prediction.

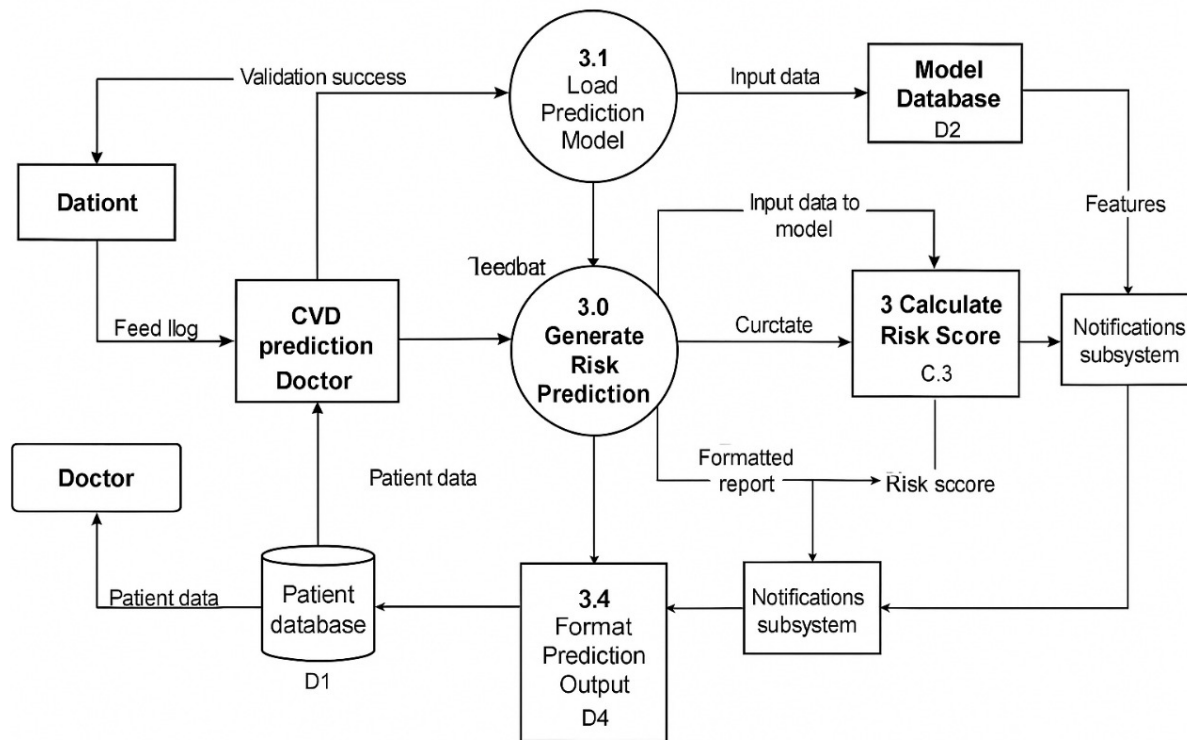


Figure 4.3: 2-level DFD

3. Calculate Risk Score

At this stage, the prediction model performs the analysis and calculates a raw risk score or classification output. The score reflects the likelihood of the patient being in one of the categories:

Normal (no CVD risk)

Benign (mild/moderate cardiac condition)

Malignant (severe/life-threatening CVD risk)

This step constitutes the core decision-making process of the system, where machine learning algorithms analyze input patterns against learned knowledge.

4. Format Prediction Output The raw prediction results are then

structured into a comprehensive and interpretable format. This includes generating a Prediction Report that summarizes the risk level, highlights contributing factors, and provides insights for clinicians. The formatted output is both stored in the Patient Records (D1) database and forwarded to the next stage (Display Results – 4.0) for doctors and clinicians to access.

4.4 Use Case Diagram

Below Figure 4.4 depicts a use case for a system of heart disease detection and classification involving two main actors, namely the Developer and the User .

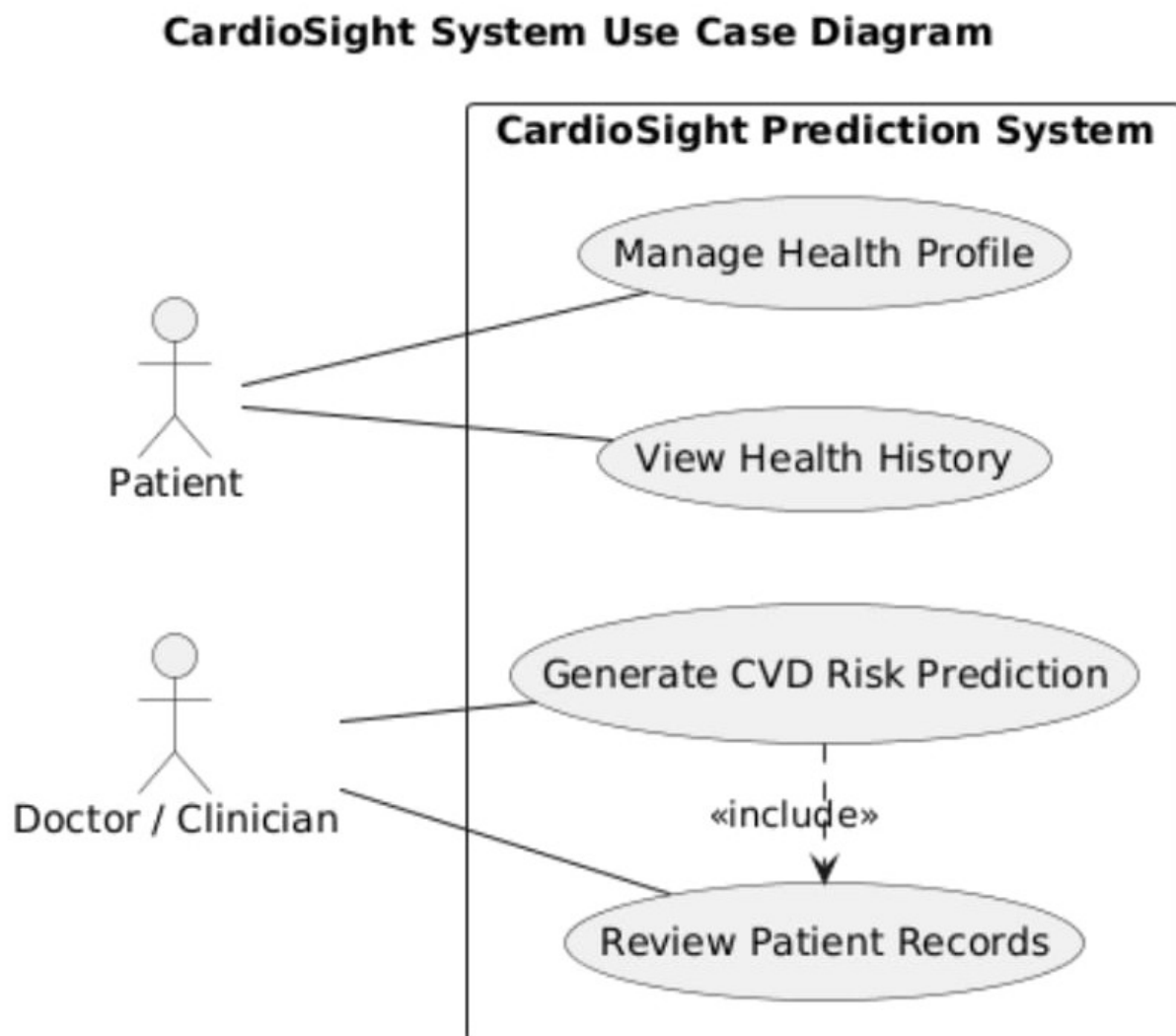


Figure 4.4: Use Case Diagram

1. Role of Doctor

The Doctor primarily works with patient records and provides professional consultation based on system predictions. Their major responsibilities include:

- **Access Patient Records:** Retrieve and review patient health data stored in the system..
- **View Predictions:** Analyze the system-generated predictions to support clinical decision-making.
- **Provide Consultation:** Offer expert advice and treatment recommendations to patients based on predictions and medical history.

2. Role of Patient

The Patient interacts with the Cardio Sight system through a user-friendly interface. Their critical tasks include:

- **Register / Login :** Gain secure access to the system.
- **Upload Health Data :** Provide clinical or diagnostic data such as ECG readings, lifestyle metrics, or blood reports for system analysis..
- **View Prediction :** Receive the system's prediction regarding the likelihood of cardiovascular disease, empowering them with timely insights for early diagnosis and prevention.

Chapter 5

Detailed Designing of Project

5.1 Class Diagram

The following Figure 5.1 represents the overall structure of the system, showing different classes, their attributes, operations, and relationships among them. The diagram defines how data flows between the application layers, including authentication, data preprocessing, model inference, and result management. Class of interface allows the user to access the system. The major methods include:

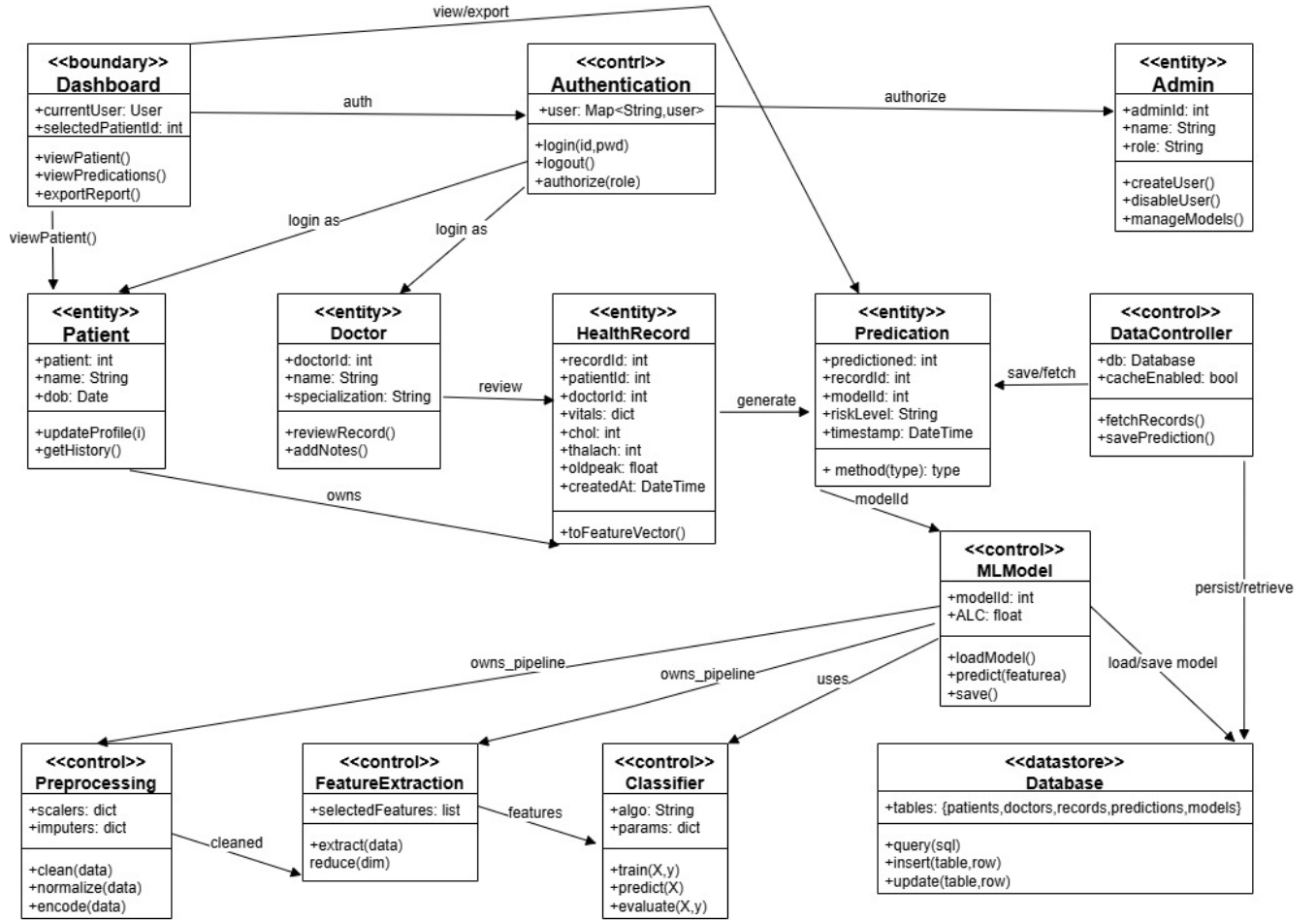


Figure 5.1: Class Diagram

The Dashboard class acts as the system's user interface, allowing users to view patient details, access predictions, and export reports. The Authentication class handles secure login, logout, and role-based authorization, ensuring that only verified users can access specific functionalities. The Patient class stores personal information such as patient ID, name, and date of birth, allowing users to update profiles and access medical history. The Doctor class represents healthcare professionals who can review health records, add clinical notes, and manage patient cases.

The HealthRecord class stores patient diagnostic information, including vital signs, cholesterol levels, and ECG-related data. It acts as a bridge between patients, doctors, and the prediction module. The Predication class contains the results generated by machine learning models, recording the risk level, timestamp, and associated record and model IDs.

To ensure data readiness, the Preprocessing class cleans, normalizes, and encodes raw data, while the FeatureExtraction class extracts and reduces key clinical features. The Classifier class performs the main predic-

tive task using machine learning algorithms for training, prediction, and evaluation. The MLModel class represents the trained model, handling loading, saving, and prediction generation processes. The DataController class manages the transfer of information between system components and the Database, ensuring smooth data retrieval and storage.

Finally, the Admin class supervises the overall system by managing users, roles, and model configurations. The Database class acts as the system's persistent storage, containing structured tables for patients, doctors, records, predictions, and models.

This class diagram provides a clear view of how CardioSight integrates data processing, authentication, and machine learning workflows to deliver secure, interpretable, and efficient cardiovascular disease predictions.

5.2 Entity-Relationship Diagram (ERD)

This figure 5.2 attached image is the Entity-Relationship (ER) Diagram of a system that detects and classifies Cardio Sight.

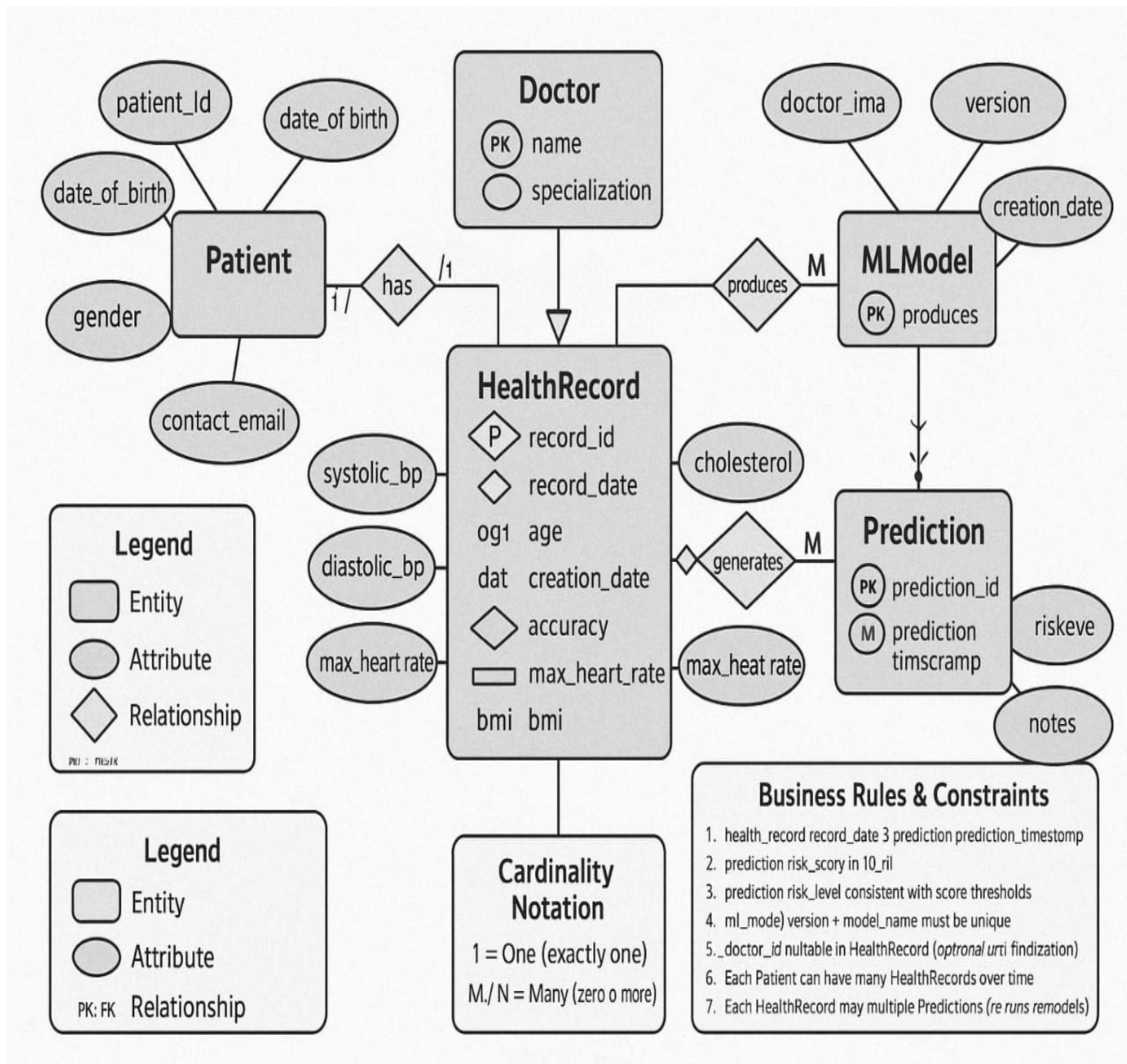


Figure 5.2: ER Diagram

It represents the relationships among different entities in the process, including patients, doctors, health records, machine learning models, and prediction outputs.

1. Patient Attributes: Patient ID, Name, Date of Birth, Gender, Contact Information, Family History of Cardiovascular Disease (CVD).

The patient entity stores essential demographic and medical background information. Family history of CVD is captured to assist in risk factor analysis. Each patient may have multiple health records over time, which are used for monitoring and prediction.

2. Doctor Attribute: Doctor ID, Name, Specialization, Hospital Affiliation.

This entity contains the professional details of the healthcare provider associated with the patient. Doctors can create or verify health records,

review prediction outcomes, and provide medical consultation based on system-generated results. Each doctor can be linked with several patient health records.

3. Health Record : Record ID, Patient ID (FK), Doctor ID (FK), Record Date, Age, Systolic Blood Pressure, Diastolic Blood Pressure, Cholesterol, Fasting Blood Sugar, Resting ECG, Maximum Heart Rate, BMI, Smoking Status, Chest Pain Type, Exercise Induced Angina, Number of Major Vessels.

The HealthRecord entity represents a collection of clinical parameters and test results for a patient at a specific point in time. It acts as the central entity connecting patients, doctors, and prediction data. Each health record serves as the primary input for the model to assess cardiovascular risk.

4. Prediction : Prediction ID, Record ID (FK), Model ID (FK), Prediction Timestamp, Risk Score, Risk Level.

The Prediction entity stores the outputs generated by the trained machine learning models. It provides both numerical and categorical insights, such as the risk score and risk category (e.g., low, moderate, or high risk). Each prediction corresponds to a particular health record and is linked to the ML model that produced it.

5. ML Model : Model ID, Model Name, Version, Creation Date, Accuracy, Model File Path.

The MLModel entity holds information about the trained machine learning models used in the prediction process. It records model metadata such as version, accuracy, and creation date. Each model must have a unique version and name combination, and it produces multiple predictions based on patient health data.

5.3 Activity Diagram

In this figure 5.3 schematic explains the process by which a multi-step machine learning technique performs the task of detection and classification of Cardio Sight . It has different significant steps on its workflow:

1. ECG/Heart Signal Acquisition: It first acquires ECG/heart signals from the patient, which are basically provided as input to the system. These signals are the raw data that form the foundation for further analysis.

2. Pre-processing: The acquired signals are preprocessed in order to

enhance quality. This includes noise removal, baseline wandering correction, normalization of data, and filtering unwanted frequency components. Preprocessing ensures that the signals are clean and suitable for analysis.

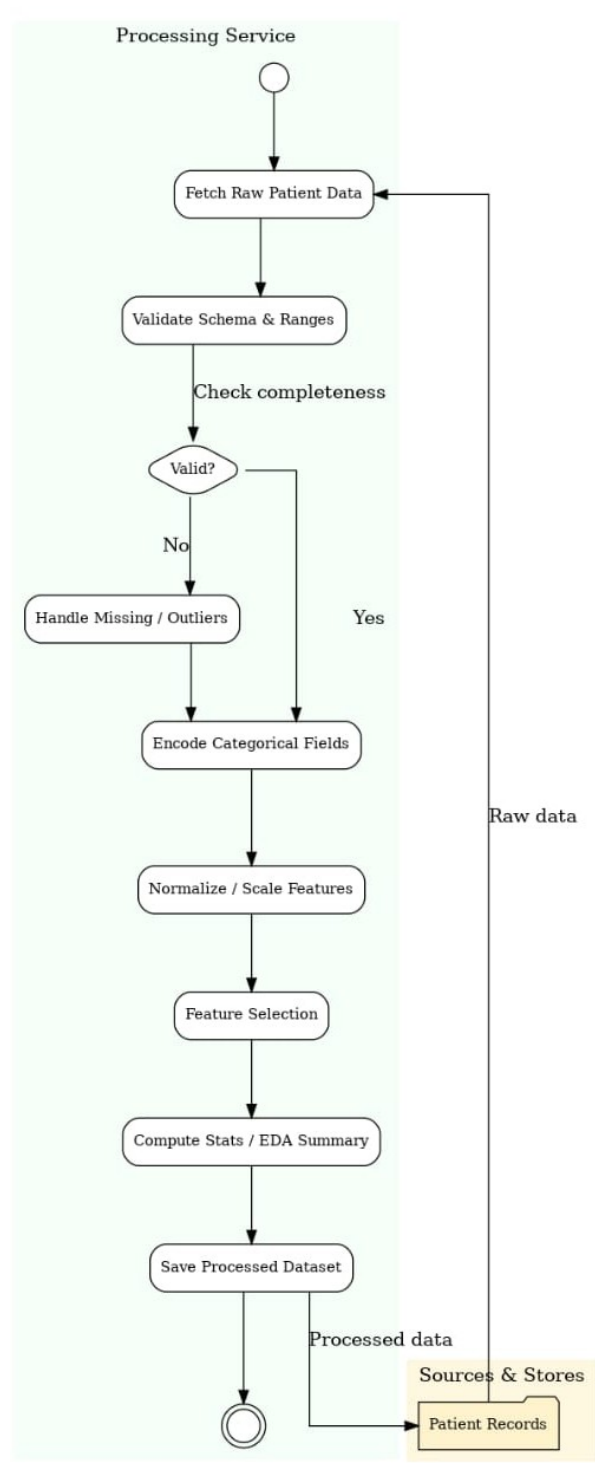


Figure 5.3: Activity Diagram

3. Feature Extraction : In this phase, relevant features such as heart rate variability (HRV), QRS complex characteristics, P-wave and T-wave

morphology, and frequency-domain measures are extracted from the ECG signals. These features capture the important characteristics needed to differentiate between normal and abnormal cardiac conditions.

4. Feature Reduction : In this step, redundant and irrelevant features are eliminated. This dimensionality reduction process retains only the most significant features, ensuring faster computation and improved accuracy during classification.

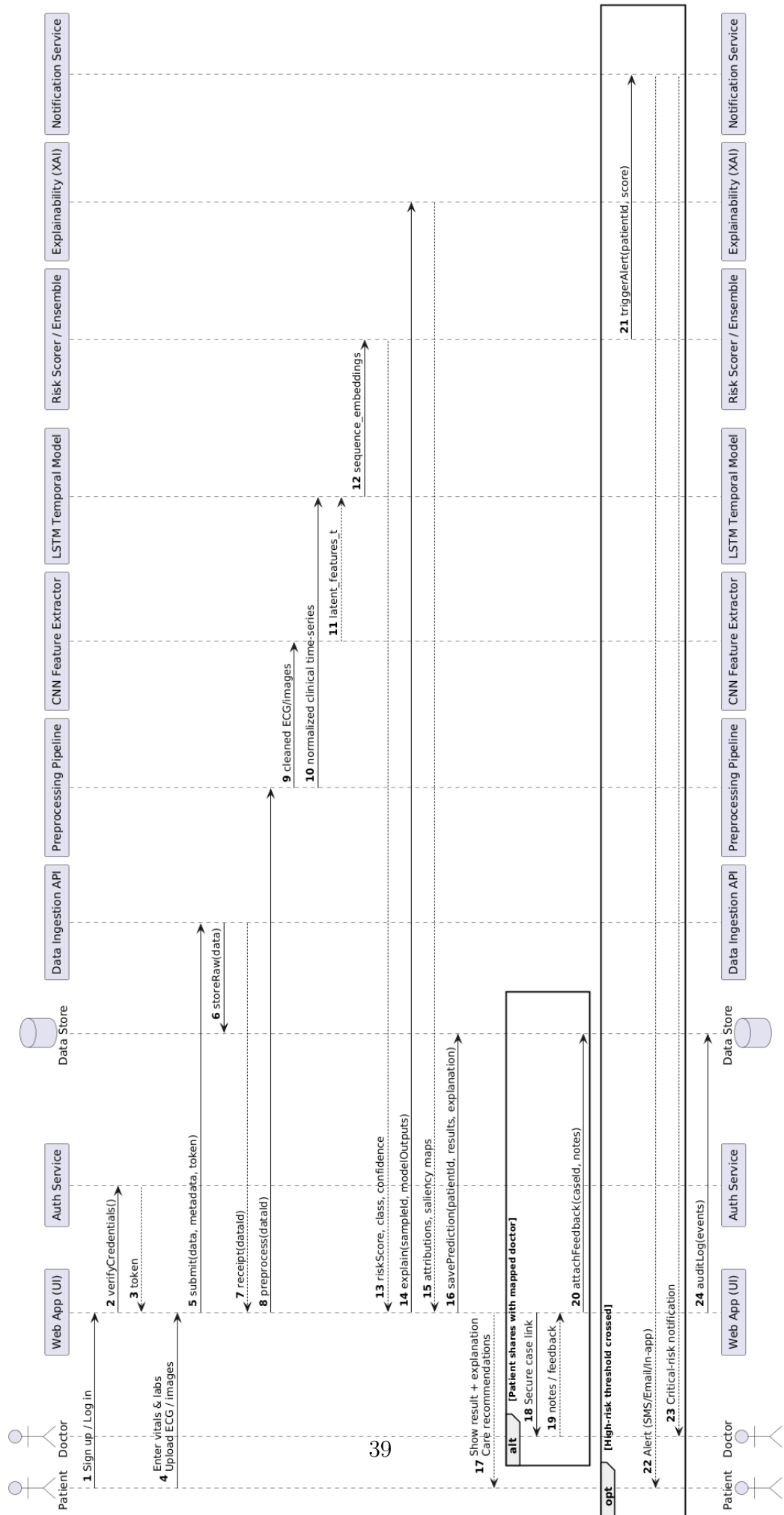
5. Training and Classification – Stage 1 : The selected features are fed into a pre-trained machine learning or deep learning model that classifies the heart condition as either “Normal” or “Abnormal”. If the system classifies the signal as “Normal,” then no further process is required. However, if the condition is classified as “Abnormal,” it proceeds to the next stage.

6. Classification – Stage 2 : The abnormal conditions are further sub-classified into “Benign” (non-life-threatening arrhythmias or mild cardiac issues) and “Malignant” (serious or life-threatening cardiac abnormalities such as ventricular fibrillation or severe arrhythmias). This classification is crucial in guiding diagnosis and treatment planning.

5.4 Sequence Diagram

The CardioSight — Inference Sequence diagram 5.4 illustrates the end-to-end workflow of the heart disease prediction and explainability system, showing how various system components and actors interact from data submission to result interpretation and alert generation. The process involves multiple services, including data ingestion, preprocessing, model inference, explainability (XAI), and notification handling.

CardioSight — Inference Sequence



- 1. Sign-up / Log-in:** The Patient or Doctor accesses the system via the Web App (UI). Authentication is handled by the Auth Service, which verifies credentials and issues a secure token for session validation.
- 2. Data Submission:** The user uploads clinical data, ECG signals, or medical images through the web interface. The UI sends the submission (data, metadata, token) to the Data Ingestion API, which stores the raw data securely in the Data Store for further processing.
- 3. Preprocessing :** The Preprocessing Pipeline retrieves the stored data and performs cleaning, normalization, and transformation.
- 4. Feature Extraction :** The preprocessed data is passed to the CNN Feature Extractor, which identifies spatial features (e.g., waveform morphology, pixel intensity). These features are then fed into the LSTM Temporal Model, which captures sequential patterns over time. The result is a set of sequence embeddings representing latent temporal and spatial correlations in the patient's cardiac data..
- 5. Risk Scoring :** The Risk Scorer / Ensemble Model aggregates the outputs from the CNN-LSTM pipeline to predict a risk score, classify the result (e.g., Low, Moderate, or High risk), and compute a confidence level for the prediction.
- 6. Explainability Generation :** To ensure clinical transparency, the Explainability (XAI) module generates interpretability maps such as attributions and saliency visualizations, helping clinicians understand which data patterns influenced the decision.
- 7. Result Storage and Display:** The Web App stores both prediction results and explainability outputs in the Data Store. The UI then presents the risk score, explanation, and care recommendations to the user. The Patient can share this result securely with the mapped Doctor through a unique case link.
- 8. Doctor Feedback:** When shared, the Doctor accesses the case report, reviews predictions, and provides clinical notes or feedback, which are stored back into the Data Store for continuous model improvement and audit.

Bibliography

- [1] Hansen et al. Assigning diagnosis codes using medication history. *Artificial Intelligence in Medicine*, 2022. F1 score: 89.66
- [2] Romdhane et al. Cnn-based heartbeat segmentation for arrhythmia detection. *Artificial Intelligence in Medicine*, 2020. Accuracy: 98.41
- [3] Acharya et al. Convolutional neural network for arrhythmia diagnosis. *Artificial Intelligence in Medicine*, 2017. Accuracy: 94
- [4] Dixit and Kala. 1d cnn model for heart disease detection. *Artificial Intelligence in Medicine*, 2021. Accuracy: 93
- [5] Wang et al. Cnn-based left ventricle landmark localization in cardiac mri. *Artificial Intelligence in Medicine*, 2020. Accuracy: 95.82
- [6] Attia et al. Ai-based cnn for asymptomatic left ventricular dysfunction. *Circulation*, 2019. Accuracy: 85.7
- [7] Nirschl et al. Ai-based cnn on histological tissue slides for heart failure. *Artificial Intelligence in Medicine*, 2018. 99
- [8] Pardeshi et al. Efficient approach for detecting cardiovascular disease using lstm. *International Journal of Aquatic Science*, 2023. Maximum accuracy: 99.12
- [9] Krishnan et al. Hybrid deep learning model (rnn+lstm) for heart disease prediction. *Artificial Intelligence in Medicine*, 2021. Accuracy: 98.5

Appendix A

Figure A.1 shows the plagiarism report of our project, which is 8 percent-age

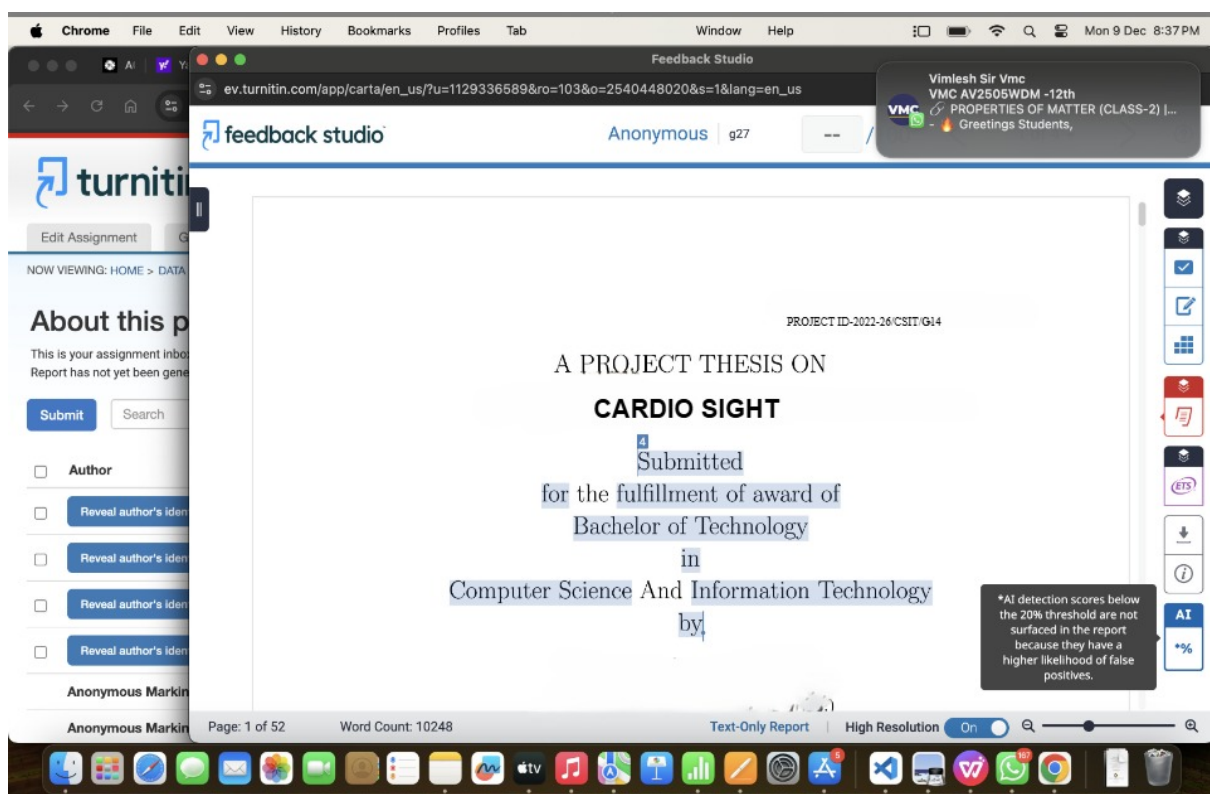


Figure 5.5: Plagiarism Report

Figure A.2 shows the AI report of our project, which is less than 20 percentage

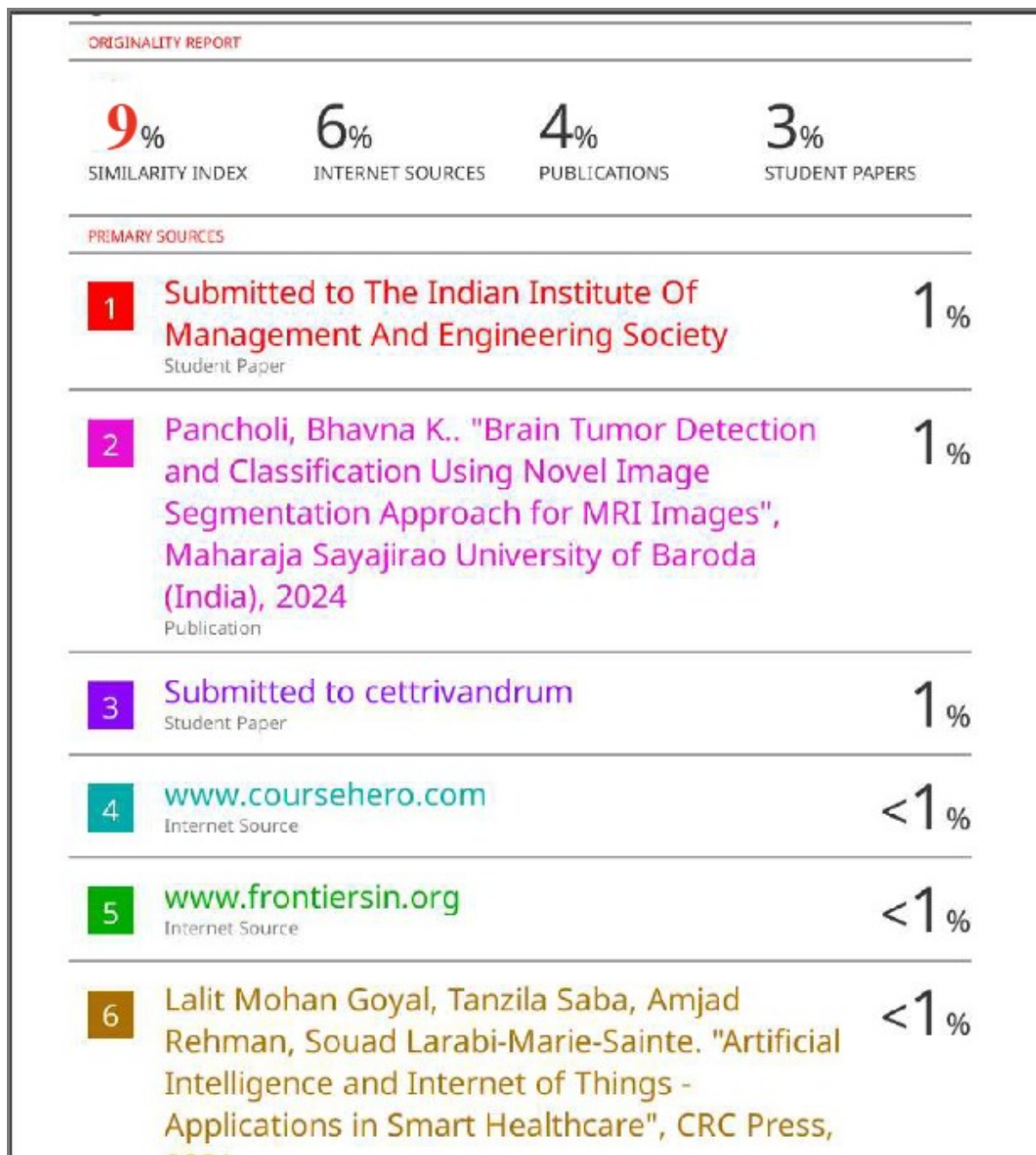


Figure 5.6: AI Report