# Early Detection of Cardiovascular Disease Atrial Fibrillation

Prem Iniyan R K
*Department of CSE*
*St. Joseph Engineering College*
Mangaluru, India
20cs111.prem@sjec.ac.in

Pragati P Patil
*Department of CSE*
*St. Joseph Engineering College*
Mangaluru, India
20cs102.pragati@sjec.ac.in

Pranoy T Murali
*Department of CSE*
*St. Joseph Engineering College*
Mangaluru, India
20cs106.pranoy@sjec.ac.in

Preeona Mendonca
*Department of CSE*
*St. Joseph Engineering College*
Mangaluru, India
20cs110.preeona@sjec.ac.in

Harivinod N*
*Department of CSE*
*St. Joseph Engineering College*
Mangaluru, India
harivinodn@sjec.ac.in

*Abstract*—Cardiovascular disease remains a leading cause of mortality worldwide, and the early detection of its underlying conditions is crucial for effective prevention and management. This research paper presents a comprehensive approach to address the challenge of early prediction and detection of cardiovascular diseases (CVDs) focusing on atrial fibrillation (AF). The proposed work integrates advanced machine learning techniques to study risk assessment and deep learning for ECG analysis. The Kaggle's Cardiovascular Disease Dataset and Cleveland's Heart Disease Dataset are used for the experimentation. The study also discusses practical application implementation, exploring the integration of developed models into healthcare systems for real-time risk assessment and disease detection. Results demonstrate that Random Forest Classifier gives superior performance with both datasets.

*Index Terms*—Cardiovascular Diseases, Atrial Fibrillation, Convolutional Neural Networks, Identification, Machine Learning,

## I. INTRODUCTION

This study presents a novel approach to address the pressing challenge of early atrial fibrillation (AF) detection by amalgamating advanced artificial intelligence methodologies with an extensive analysis of heart data, predominantly electrocardiogram (ECG) datasets. Through the utilization of machine learning and deep learning algorithms in conjunction with a comprehensive array of demographic, biometric, and clinical parameters, our research endeavors to devise a robust and precise mechanism for identifying AF in its nascent stages. Moreover, by integrating cardiac rehabilitation and training techniques into the AF detection framework, we offer a proactive strategy that empowers individuals to actively engage in the management of their cardiovascular well-being while simultaneously augmenting diagnostic efficacy. This interdisciplinary endeavor seeks to democratize access to AF diagnosis and treatment, particularly in underserved communities, thereby catalyzing a paradigm shift in cardiovascular healthcare delivery and fostering improved patient outcomes.

## II. LITERATURE REVIEW

The review article [1] critically examines the applicability of existing cardiovascular risk assessment tools within the context of the Asian population, highlighting a significant problem with the limited effectiveness of these tools in accurately predicting cardiovascular risk for individuals of Asian descent. The study reveals that many of these tools were developed based on data from American or European populations, which may lead to inaccuracies when applied to Asian individuals. Through a systematic search of English articles from 1995 to June 2008 using specific search terms and popular databases, the researchers identified 25 cardiovascular risk assessment tools and analyzed their characteristics, including sample size, study type, time frame, endpoints, statistical analysis, and risk factors. Despite the extensive search, the review emphasizes the scarcity of tools derived from Asian populations, comprising only 8 percent of the total. The findings highlight discrepancies in risk estimation, with the widely used Framingham Risk Score found to overestimate the risk for coronary heart disease (CHD) among Asians. The study underscores the need for ongoing re-calibration and validation of existing tools with local epidemiological data to enhance accuracy and adaptability to evolving population trends, ultimately providing a more precise prediction of cardiovascular risk in the Asian demographic.

In the work [2], two distinct methodologies are employed: the combination methodology, which integrates diverse risk assessment tools using Naive Bayes classifiers and genetic algorithms to optimize weights for a comprehensive model, and the personalization based on Grouping of Patients, which tailors risk predictions by categorizing patients with similar characteristics using dimension reduction and clustering techniques. The implementation involves combining risk assessment tools

using a unified representation and optimizing weights with genetic algorithms, along with personalization through dimension reduction and clustering to tailor predictions for diverse patient groups. The results demonstrate the effectiveness of the combined methodology in enhancing prediction accuracy, underscoring the potential for more precise and individualized cardiovascular disease risk assessment.

Atrial fibrillation (AF) is a prevalent cardiac condition associated with increased risk of strokes and mortality. In [3] two classification approaches are explored: the DenseNet Approach, where features extracted from spectrograms are fed into a Support Vector Machine (SVM), and the Convolutional Network Approach, where spectrograms are inputted directly into a convolutional neural network (CNN) with batch normalization and dropout techniques. In implementation, the pre-trained DenseNet extracts features for AF classification using SVM, while the CNN processes spectrograms directly. The Spectrogram + ConvNet approach demonstrates superior accuracy and sensitivity in AF detection compared to the Spectrogram + DenseNet + SVM method, despite a slight trade-off between sensitivity and specificity. Future enhancements may involve noise reduction using autoencoders and exploring transfer learning for improved model training against diverse ECG datasets.

The study [4] presents a novel approach for the early detection of atrial fibrillation (AF) using a combination of machine learning (ML) and deep learning (DL) techniques. After preprocessing to remove noise and normalize data, ML model training involves feature extraction from electrocardiogram (ECG) signals using Fast Fourier Transform (FFT) and Principal Component Analysis (PCA) for dimensionality reduction. DL methods, such as Convolutional Neural Networks (CNN) and Long Short-Term Memory (LSTM), are utilized for feature detection and sequential data analysis. Experimental results show the superior performance of DL methodologies compared to conventional shallow learning classifiers like Support Vector Machines (SVM), highlighting the potential of integrating ML and DL for enhancing AF detection accuracy. Future work aims to develop a real-time AF detection approach without the need for labeled data.

The work [5] addresses the challenge of developing deep learning (DL) models with robust generalization capabilities for atrial fibrillation (AF) detection across different electrocardiogram (ECG) devices. Leveraging a cloud-based DL system, it employs a 1D convolutional neural network (CNN) to process ECG signals transmitted from short-term devices to mobile devices and then to a cloud server for real-time analysis. The study enhances a previously published CNN model to handle inputs from various ECG devices with different sampling frequencies and signal lengths. By preprocessing ECG signals, segmenting them, and extracting features related to irregular heart rhythms, the system achieves high accuracy, sensitivity, and specificity in AF detection, validated with unseen data and other arrhythmias. This approach holds promise for early AF diagnosis and better management to prevent stroke.

Detecting atrial fibrillation (AF) is crucial due to its association with increased stroke risk. The paper [6] proposes a novel method utilizing convolutional neural networks (CNNs) to enhance accuracy and reduce complexity in AF detection from electrocardiogram (ECG) signals. Unlike conventional approaches that convert ECG signals into complex 2D images, this study directly utilizes the structure of ECG data, fine-tuning different CNN model parts for efficiency. Experiments conducted on PhysioNet Challenge 2017 dataset demonstrate the efficacy of the CNN model, achieving high accuracy while maintaining simplicity. The study offers insights into CNN model design for AF detection, highlighting the importance of certain techniques over others and suggesting avenues for future research in medical signal analysis and healthcare diagnostics.

The proposed method in [7], integrates an 8-layer convolutional neural network (CNN) with a 1-layer long short-term memory (LSTM) and shortcut connections to effectively detect atrial fibrillation (AF) in electrocardiogram (ECG) data. Advanced architectural elements optimize ECG data classification, utilizing temporal convolution for local information capture, LSTM for continuous memory updating, and shortcut connections to streamline optimization. Preprocessing involves data acquisition, normalization, imbalance correction, and segmentation. Implementation includes model architecture construction, dataset splitting, training, testing, and performance evaluation, demonstrating superior accuracy and efficiency of 8CSL in detecting AF from short ECG recordings. The study suggests future research avenues for extending AF detection to 12-lead ECG recordings.

Cardiovascular diseases pose a significant global health burden, necessitating accurate diagnosis of cardiac arrhythmias for effective treatment. The study [8] addresses the need for improved arrhythmia detection methods by leveraging computational techniques. Using data from the MIT-BIH Arrhythmia Database, a three-class arrhythmia detection approach is proposed based on K-means clustering. Feature extraction from electrocardiogram (ECG) signals involves QRS morphology, Heart Rate Variability (HRV), and statistical metrics, leading to two feature sets (FS1 and FS2) for analysis. K-means clustering categorizes arrhythmias using both feature sets, with performance metrics such as accuracy, sensitivity, and specificity evaluated. Results demonstrate high accuracy for distinguishing normal, Right Bundle Branch Block (RBBB), and Left Bundle Branch Block (LBBB) heartbeats. While the reduced feature set (FS2) exhibits slightly lower performance, it offers computational efficiency. Future research directions include dataset expansion, alternative feature extraction methods, and optimization of clustering algorithms for real-time applications, ultimately enhancing clinical utility and reliability.

## III. METHODOLOGY

Figure 1 gives the architecture of the proposed approach. It provides a systematic approach for predicting and detecting atrial fibrillation (AF) in patients with preexisting cardiac
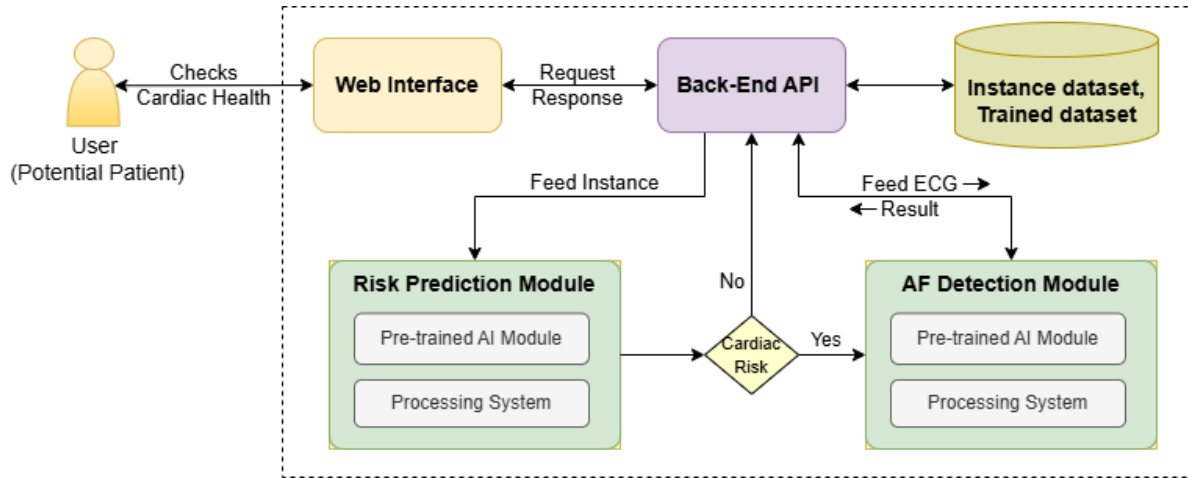
Fig. 1. Architectural Diagram of the proposed work

concerns. It begins with survey-based inquiries about daily activities and cardiac health, followed by ECG report upload if risk indications arise. Machine learning algorithms then analyze the ECG reports for potential AF cases, aiming for efficient and accurate identification to enable timely intervention and management strategies.

Our work focuses on timely detection of cardiovascular diseases, particularly Atrial Fibrillation (AF), through a structured two-stage process.

### A. Risk Assessment module for Detection of Heart Disease

In the initial stage, we conduct a comprehensive risk assessment using the random forest algorithm. We collect user data including demographic information, biometrics, and lifestyle factors. This data helps evaluate the individual's susceptibility to heart disease. Using the random forest model, we analyze the provided data against a robust dataset to identify potential cardiovascular risks. If any concerns are raised during this assessment, we proceed to the next stage.

In our study, we employed four machine learning algorithms—Random Forest, Logistic Regression, Gaussian Naive Bayes, and Support Vector Machine (SVM)—to develop models for detecting heart disease. Two independent datasets were used, with the first dataset from Kaggle divided into 70% for training and 30% for testing, and the second dataset from Cleveland split into 80% for training and 20% for testing. Each model was trained and tested using its respective dataset.

Following training and testing, the models were evaluated based on classification reports and confusion matrix reports. These assessments provided insights into the performance of each model in accurately predicting the risk of heart disease. Among the algorithms tested, the Random Forest Classifier consistently showed the best performance across both datasets.

Given its effectiveness in handling complex datasets, the Random Forest Classifier emerged as the preferred choice for detecting heart disease in our study.

For our project, we chose to utilize Kaggles's dataset with Random Forest Classifier as the chosen model. This decision was guided by the user-friendly nature and alignment of the first dataset with our project's objectives. The dataset, containing demographic information and clinical measurements, offered accessibility and ease of use, facilitating simplified data processing and risk assessment model implementation using machine learning techniques.

### B. Atrial Fibrillation detection module via ECG Analysis

In the second stage, users upload their Electrocardiogram (ECG) snapshots for further analysis. These ECG images provide insights into heart activity and possible abnormalities, including AF. Convolutional Neural Networks (CNNs) are employed to examine the uploaded ECGs. Trained on a comprehensive dataset of labeled ECG recordings, the CNN model detect subtle patterns indicative of AF.

Our two-stage approach combines random forest classifiers for comprehensive risk assessment and Convolutional Neural Networks (CNNs) for detailed ECG analysis. By integrating diverse user data and advanced machine learning techniques, we aim to enhance early detection of cardiovascular risks, particularly atrial fibrillation (AF), and promote proactive intervention for better heart health outcomes.

*1) CNN Architecture:* The Convolutional Neural Network (CNN) architecture utilized in this work is tailored for ECG image classification, specifically targeting the detection of atrial fibrillation (AF). The architecture begins with an input layer accepting standardized ECG images of dimensions (224, 224, 1). Employing data augmentation techniques like random flipping, rotation, and zooming enhances the training dataset's diversity. Two convolutional layers followed by ReLU activation functions extract low-level features, while max-pooling layers down-sample feature maps to reduce complexity. Dropout layers mitigate overfitting, and two fully connected dense layers learn higher-level features. The output layer with softmax activation categorizes images into AF and non-AF classes. Compiled with sparse categorical cross-entropy loss and Adam optimizer, the model is trained on la-

300

beled data to optimize parameters and minimize loss, resulting in accurate AF classification.

*2) CNN Model Development and Training for Detection of AF:* Developing and training a Convolutional Neural Network (CNN) for Atrial Fibrillation (AF) detection involves preprocessing the ECG dataset, including resizing images and applying a Laplacian filter for feature enhancement. The dataset is split into training and testing sets, and data augmentation techniques are applied to diversify the training data. The CNN model, defined using Keras Sequential API, comprises convolutional, pooling, and dropout layers to prevent overfitting. Compiled with appropriate loss and optimizer functions, the model is trained on the training data to minimize the loss function. Finally, the trained model is evaluated on the testing data using metrics like accuracy, test loss, classification report, and confusion matrix to ensure its efficacy in real-world AF detection.

## IV. DATASETS

To facilitate robust risk assessment, we chose two comprehensive datasets for research comprising diverse demographic, biometric, lifestyle, and clinical parameters. The datasets for the experiemntation are Kaggle's Cardiovascular Disease Dataset and Cleveland Heart Disease Dataset.

### A. Kaggle's Cardiovascular Disease Dataset

This dataset encompasses a wide range of variables essential for assessing cardiovascular risk. With a size of 70,000 samples, it includes demographic factors such as age and sex, biometric measures like height and weight, clinical indicators such as blood pressure and cholesterol levels, as well as lifestyle factors such as smoking, alcohol intake, and physical activity. Each parameter in this dataset provides valuable insights into an individual's cardiovascular health status and potential risk factors for CVDs. In the dataset, $age$ is measured in days, providing a precise indication of the individual's age. $Sex$ is represented as a binary variable, where 1 denotes male and 0 denotes female. $Height$ is measured in centimeters. Weight is measured in kilograms in this dataset. $Systolic Blood Pressure(SBP)$ gives the pressure in the arteries when the heart contracts or beats. SBP is measured in millimeters of mercury (mmHg). $Diastolic Blood Pressure(DBP)$ represents the pressure in the arteries when the heart is at rest or between beats. DBP is measured in millimeters of mercury (mmHg). In the dataset, $cholesterol levels$ are categorized as: Normal, Above normal, Very high and represented as 1, 2, 3 respectively. $Glucose levels$ also categorized similar to the cholesterol level. $Smoking$ is a binary variable that indicates whether the individual is a smoker (1) or non-smoker (0). $Alcohol$ Intake is a binary variable indicates whether the individual consumes alcohol (1) or not (0). $Physical Activity$ is a binary variable indicates whether the individual engages in physical activity (1) or not (0).

### B. Cleveland Heart Disease Dataset

The dataset comprises crucial health parameters and diagnostic indicators aimed at facilitating the evaluation of cardiovascular diseases, with a specific focus on Atrial Fibrillation (AF). With a size of 304, each entry represents a patient and includes demographic information, such as age and gender, alongside clinical measurements like blood pressure, cholesterol levels, and electrocardiogram (ECG) results. These attributes serve as valuable markers for assessing cardiovascular health and identifying potential risk factors.

In the dataset, $Age$ Indicates the patients' age in years. This is represented as a numeric value. $Sex$ indicates the gender of the patients. Male indicates 1, and Female indicates 0. This is a nominal categorical variable. $cp$ indicates the type of chest pain experienced by the patient, categorized into four categories: 0 indicates typical angina, 1 indicates atypical angina, 2 indicates non-anginal pain, 3 indicates asymptomatic. This is a nominal categorical variable. $trestbps$ indicates the patient's level of blood pressure at rest, measured in mm/HG. This is a numerical variable. $chol$ Indicates the serum cholesterol level in mg/dl. This is a numerical variable. $fb$ indicates the blood sugar levels on fasting. A value of greater than 120 mg/dl indicates 1 for true and 0 for false. This is a nominal categorical variable. $restecg$ indicates the result of electrocardiogram (ECG) while at rest, represented in three distinct values: 0 indicates normal, 1 indicates having ST-T wave abnormality, 2 indicates showing probable or definite left ventricular hypertrophy by Estes' criteria. This is a nominal categorical variable. $thalach$ indicates the maximum heart rate achieved, represented as a numeric value. $exang$ indicates angina induced by exercise, with 0 indicating No and 1 indicating Yes. This is a nominal categorical variable. $oldpeak$ indicates exercise-induced ST-depression relative to the state of rest, represented as a numeric value. $slope$ indicates the ST segment measured in terms of slope during peak exercise, with the following categories: 0 indicates up-sloping, 1 indicates flat, 2 indicates downsloping. This is a nominal categorical variable. $ca$ indicates the number of major vessels, ranging from 0 to 3. This is a nominal categorical variable. $thal$ indicates a blood disorder called thalassemia, categorized into four values: 0 indicates NULL, 1 indicates normal blood flow, 2 indicates fixed defect (no blood flow in some part of the heart), 3 indicates reversible defect (abnormal blood flow). This is a nominal categorical variable.

## V. RESULTS AND DISCUSSION

### A. Model Evaluation

Model evaluation is a crucial step in machine learning, involving the assessment of how well a model performs on unseen data. It relies on various metrics like accuracy, precision, recall, and F1-score, as well as techniques such as cross-validation and confusion matrices. By guiding model selection, hyperparameter tuning, and interpretation of results, model evaluation ensures the reliability and effectiveness of machine learning models for real-world applications.

TABLE I
PERFORMANCE METRICS OF DIFFERENT ALGORITHMS FOR KAGGLE'S CARDIOVASCULAR DISEASE DATASET

| Algorithm | Accuracy | Precision | Recall | F1 Score | Class |
|---|---|---|---|---|---|
| Random Forest | 74% | 0.72 | 0.79 | 0.75 | 0 |
| | | 0.77 | 0.69 | 0.73 | 1 |
| Logistic Regression | 72% | 0.70 | 0.76 | 0.73 | 0 |
| | | 0.74 | 0.68 | 0.71 | 1 |
| Support Vector Classifier | 73% | 0.69 | 0.81 | 0.75 | 0 |
| | | 0.77 | 0.64 | 0.70 | 1 |
| Gaussian Naive Bayes | 59% | 0.56 | 0.87 | 0.76 | 0 |
| | | 0.71 | 0.32 | 0.44 | 1 |

TABLE II
PERFORMANCE METRICS OF DIFFERENT ALGORITHMS FOR CLEVELAND HEART DISEASE DATASET

| Algorithm | Accuracy | Precision | Recall | F1 Score | Class |
|---|---|---|---|---|---|
| Random Forest | 84% | 0.83 | 0.83 | 0.83 | 0 |
| | | 0.84 | 0.84 | 0.84 | 1 |
| Logistic Regression | 81% | 0.80 | 0.78 | 0.79 | 0 |
| | | 0.82 | 0.84 | 0.83 | 1 |
| Support Vector Classifier | 81% | 0.80 | 0.78 | 0.79 | 0 |
| | | 0.77 | 0.64 | 0.70 | 1 |
| Gaussian Naive Bayes | 84% | 0.78 | 0.88 | 0.83 | 0 |
| | | 0.89 | 0.80 | 0.84 | 1 |

*1) Classification Report:* We generated classification reports for each model on each dataset, which provided detailed metrics for each class as well as overall performance measures. These reports enabled us to compare the performance of different algorithms and assess their effectiveness in classifying instances accurately. The report is given in Table IV.

Upon analyzing the evaluation results, we observed variations in the performance of the models across different datasets and algorithms. Some algorithms showed better performance in terms of accuracy and other metrics, while others exhibited strengths in specific areas such as precision or recall.Random Forest Classifier proved to be the most best model for both the datasets with an accuracy of 74% and 84% for dataset 1 and dataset 2 respectively.

Due to the user-friednly nature of Kaggle's dataset along with it's alignment to the project's objective, we have used it for the implementation.

*2) AF Detecting CNN Module:* The CNN training code demonstrates exceptional performance in detecting atrial fibrillation (AF) from ECG images. The architecture, detailed in the summary table, incorporates convolutional, pooling, and dropout layers to ensure robustness and prevent overfitting. Data augmentation techniques enhance the model's ability to learn diverse patterns from the dataset. Additionally, rigorous optimization of hyperparameters fine-tunes the model's performance and facilitates optimal convergence during training. The hyper parameters of the CNN is given in Table III.

Throughout the training process, the model demonstrates exceptional performance, with both training and validation accuracy steadily increasing to 93% by the final epoch. The consistently low training loss further underscores the model's proficiency in minimizing errors during optimization. Such high accuracy suggests that the model successfully discerns between normal ECG patterns and those indicative of AF,

TABLE III
HYPERPATEMERTS OF CNN MODEL USED FOR AF MODULE

| Layer(type) | Output Shape | Param |
|---|---|---|
| random_flip(RandomFlip) | (None,224,224,3) | 0 |
| random_rotation(RandomRotation) | (None,224,224,3) | 0 |
| random_zoom(RandomZoom) | (None,224,224,3) | 0 |
| convo2d(Convo2D) | (None,222,222,16) | 448 |
| convo2d$_1$(Convo2D) | (None,220,220,32) | 4648 |
| max_pooling2d(MaxPooling2D) | (None,110,110,32) | 0 |
| dropout(Dropout) | (None,110,110,32) | 24780864 |
| flatten(Flatten) | (None,387200) | 0 |
| dense(Dense) | (None,64) | 0 |
| dropout_1(Dropout) | (None,64) | 0 |
| dense_1(Dense) | (None,128) | 8320 |
| dropout_2(Dropout) | (None,128) | 0 |
| dense_1(Dense) | (None,2) | 258 |
| Total params: 24794530(94.58 MB) | | |
| Trainable params: 24794530(94.58 MB) | | |
| Non-trainable params: 0(0.00 Byte) | | |

| Class | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| 0 | 0.97 | 0.76 | 0.85 | 175 |
| 1 | 0.91 | 0.99 | 0.95 | 108 |
| **Accuracy** | | | 0.93 | 283 |
| **Macro Avg** | 0.94 | 0.88 | 0.90 | 283 |
| **Weighted Avg** | 0.93 | 0.93 | 0.92 | 283 |

TABLE IV
CLASSIFICATION REPORT FOR CNN

which is crucial for reliable diagnosis.

Upon evaluation on the test dataset, the model achieves perfect accuracy, correctly classifying all samples as either normal or Atrial Fibrilltion(AF). This outstanding performance is echoed in the classification report, which highlights high precision, recall, and F1-score for both classes. Additionally, the confusion matrix confirms the absence of misclassifications, with all samples accurately classified.

## VI. Conclusion

In conclusion, this project represents a significant step forward in the realm of atrial fibrillation (AF) detection and management. By integrating familial history-based risk assessment with advanced electrocardiogram (ECG) analysis, our approach offers a proactive strategy for early AF detection and timely intervention. Through the utilization of machine learning algorithms we have demonstrated the potential to revolutionize AF screening and monitoring, making it more accessible, accurate, and efficient. Moving forward, the implementation of population-based screening programs, validation studies, and clinical trials will be crucial in translating these innovations into real-world clinical practice. By fostering collaboration among researchers, healthcare providers, policymakers, and industry stakeholders, we can address the challenges posed by AF, improve patient outcomes, and mitigate the societal burden of this prevalent cardiovascular condition. Future research could focus on developing personalized risk stratification models for atrial fibrillation (AF) by integrating multimodal data sources, including genetic, environmental, and behavioral factors.

## References

[1] Siow Yen Liau, MI Mohamed Izham, M A Hassali, A A Shafie, "A literature review of the cardiovascular risk-assessment tools: applicability among Asian population," July 2020, doi: 10.1136/ha.2009.001115.

[2] S. Paredes, T. Rocha, P. de Carvalho, "Cardiovascular disease risk assessment innovative approaches developed in HeartCycle project," 2020, pp. 6980-6983, doi: 10.1109/EMBC.2020.6611164.

[3] Sara Ross-Howe, Hamid R. Tizhoosh, "Atrial Fibrillation Detection Using Deep Features and Convolutional Networks," IEEE Access,2020.

[4] Sidrah Liaqat, Kia Dashtipour, Adnan Zahid, Kamran Arshad and Naeem Ramzan,"Detection of Atrial Fibrillation Using a Machine Learning Approach," November 2020,doi:10.3390/info11120549.

[5] Bambang Tutuko, Tangerang Selatan, Siti Nurmaini, Alexander Edo Tondas, Muhammad Naufal Rachmatullah,"AFibNet: An Implementation of Atrial Fibrillation Detection With Convolutional Neural Network," February 2021,doi:10.21203/rs.3.rs-228165/v1.

[6] Chaur-Heh Hsieh, Yan-Shuo Li, Bor-Jiunn Hwang, Ching-Hua Hsiao, "Detection of Atrial Fibrillation Using 1D Convolutional Neural Network," April 2020, doi:10.3390/s20072136.

[7] Yongjie Ping, Chao Chen, Lu Wu, Yinglong Wang and Minglei Shu, "Automatic Detection of Atrial Fibrillation Based on CNN-LSTM and Shortcut Connection," May 2020,doi: 10.3390/healthcare8020139.

[8] M. Masum, M.J.H. Faruk, H. Shahriar, K. Qian, D. Lo, M.I. Adnan, "K-Means Clustering Algorithm Based Arrhythmic Heart Beat Detection in ECG Signal," Balkan Journal of Electrical and Computer Engineering, January 2021,doi:10.17694/bajece.814473.