

# Machine Learning Approaches for Cardiovascular Stroke Prediction: A Comparative Analysis

Vipul Narayan<sup>1\*</sup>, Divya Midhun<sup>2</sup>, Pawan Whig<sup>3</sup>

<sup>1,2</sup> Lincoln University; <sup>3</sup> VIPS-TC, India;

vipulupsainian2470@gmail.com, divya@lincoln.edu.my, pawan.whig@vips.edu

---

**Abstract:** Cardiovascular diseases is the largest cause of death worldwide, emphasizing the importance of improving early detection and risk prediction. This study evaluates how well several machine learning models such as Support Vector Machine, Decision Tree, Random Forest and Gradient Boosting predict the risk of CVD. Key clinical and demographic variables, including age, hypertension, cholesterol, and smoking status, were included in the dataset that was analyzed. Metrics like accuracy, recall, F1-score, and precision were used to assess each model's performance. The results show that the Decision Tree model had the highest recall and F1-score, demonstrating its strong ability to identify high-risk individuals while minimizing false negatives. K-Nearest Neighbors (KNN) and Gradient Boosting showed limitations in sensitivity, indicating their reduced effectiveness in detecting positive cases. However, ensemble models such as Random Forest and XGBoost performed well, showcasing their robustness and ability to handle complex data patterns. These findings underscore the promise of machine learning (ML)-based models in CVD risk prediction and highlight their potential for improving early diagnosis and prevention. However, they also emphasize the necessity of further refinement, optimization, and real-world validation to ensure that these models are ready for clinical implementation and can be used effectively in healthcare settings.

**Keywords:** Cardiovascular Disease; Machine Learning; Risk Prediction; Decision Tree; Random Forest; Classification Models

---

## Introduction

The substantial global burden of cardiovascular diseases (CVDs), which continue to be the foremost reason of demise worldwide and account for over 17.9 million deaths per year, has been repeatedly highlighted by the World Health Organization (WHO) [1]. This represents 32% of all fatalities worldwide, underscoring the urgent need for efficient management, early identification, and preventative measures. Heart failure, stroke, and ischemic heart disease are among the CVDs that disproportionately affect low- and middle-income nations, where access to healthcare and preventative measures are frequently insufficient. In addition to their negative effects on health, CVDs have a significant financial impact on societies and healthcare systems, resulting in increased treatment expenses, lost productivity, and a worse standard of living for patients and their families [2].

The WHO has emphasized in recent years how the CVD pandemic is fueled by modifiable risk factors, including poor diets, tobacco use, physical inactivity, and hazardous alcohol intake. Cardiovascular risk is also increased by non-modifiable variables such as age, gender, and genetic susceptibility [3]. Through focused interventions, such as encouraging healthy lifestyles, enhancing access to necessary medications, and fortifying healthcare systems, the WHO's global action plan for the prevention and control of non-communicable diseases seeks to cut premature mortality from CVDs by 25% by 2025. However, aging populations, the growing burden of comorbidities like diabetes and obesity, and the increased incidence of risk factors continue to impede progress despite these efforts [4].

Artificial intelligence (AI) and digital health tools are two examples of emerging technologies that provide encouraging prospects for enhancing CVD treatment, diagnosis, and prevention. The WHO has recognized that these technologies have the potential to close gaps in healthcare delivery, especially in environments with limited resources [5]. To guarantee that the advantages of new technologies are felt by all groups, their effective integration calls for careful examination of ethical, legal, and equitable issues. A multi-sectoral strategy that involves governments, healthcare providers, and communities will be crucial in tackling the CVD pandemic and improving cardiovascular health globally as the international community strives to meet the WHO's objectives. When the heart is working at its best, it influences emotions and energy levels, which is vital for maintaining general health and wellbeing. However, how can heart health be evaluated [6]. The five main signs of a healthy heart are listed below; each provides important information about internal physiological processes and their importance for cardiovascular health.

### 1. Steady Heart Rate

A regular resting heart rate, usually between 60 and 100 beats per minute, is frequently a sign of a healthy heart. The heart works effectively and without undue effort when it has a constant rhythm within this range when pumping blood throughout the body. Even lower resting heart rates, usually between 50 and 60 beats per minute, can be found in physically healthy people [7]. This lowered rate indicates a more robust and effective heart that can pump more blood with each beat, improving cardiovascular health overall.

## 2. Good Oral Hygiene

It might surprise you to learn that the condition of your teeth and gums might reveal important information about the health of your heart. Because gum inflammation can migrate to blood vessels and eventually contribute to cardiovascular problems, research has connected gum disease to an elevated risk of heart disease. The link between cardiovascular and oral health is further supported by the fact that healthy gums and teeth frequently indicate a healthy heart [8]. Regular brushing, flossing, and dental exams are all examples of good dental hygiene practices that might lower the risk of gum disease and perhaps improve heart health.

## 3. Sustained Blood Pressure

One important sign of ideal heart and blood vessel performance is the maintenance of healthy blood pressure levels, which are normally about 120/80 mmHg. The heart can pump blood effectively and without undue strain when blood pressure stays within this range. While low blood pressure may result in less oxygen reaching essential organs, high blood pressure can overwork the heart and raise the risk of cardiovascular problems [9].

## 4. Healthy Cholesterol Levels

One of the key markers of heart health is cholesterol levels. The risk of arterial blockages is decreased by maintaining low levels of harmful LDL cholesterol and high levels of good HDL cholesterol. In order to maintain healthy blood flow, avoid plaque accumulation, and remove excess cholesterol from the bloodstream, HDL cholesterol is essential [10]. Lowering the risk of heart disease and supporting cardiovascular health require maintaining total cholesterol within the optimal range, which is normally less than 200 milligrams per deciliter (5.17 millimoles per liter).

## 5. Breathing Recovery Rate After Exercise

An important indicator of cardiac health is the speed at which a person regains their breath after exercising. A quicker rate of recovery shows that the heart and lungs are working well, supplying muscles and organs with oxygen-rich blood. The ability to return to a normal breathing rate within a minute or two of effort is a sign of a healthy heart. On the other hand, a slower rate of recovery could indicate decreased cardiovascular efficiency, indicating the necessity of better heart care through consistent exercise, a healthy diet, and general good lifestyle choices.

## Related work

Table 1. Related Work

Ref.	Objective	Methodology	Advantages	Limitations
[11]	Implement a SmartPhrase in EHR to streamline exercise prescriptions for HF patients.	Recruited 8 nurse practitioners to test the SmartPhrase. Access and feasibility were evaluated via questionnaires.	High acceptability and appropriateness; streamlined exercise prescription.	Feasibility could improve with automation of data extraction from EHR.
[12]	Identify brain tumors using MRI with a novel PACDNN-COA-BTI-MRI model.	Preprocessed MRI images with MFIF, extracted features using SSNT, implemented PACDNN optimized with COA, and analyzed performance metrics.	Improved accuracy, recall, and precision over existing techniques; reduced overfitting and enhanced feature extraction.	High computational complexity; requires substantial training data.
[13]	Automate myocardial scar quantification from LGE in cardiac MRI using ScarNet.	Hybrid model combining a transformer-based encoder and U-Net decoder, trained on 552 patients and tested on	High segmentation accuracy; robust performance across diverse image qualities and scar	Limited generalizability without further external validation.

[14]	Test the efficacy of combined AET and PMT in post-stroke aphasia.	184 with expert segmentations. Randomized Clinical Trial comparing AET + PMT to stretching + PMT, examining impacts on language outcomes, CBF, functional connectivity, and brain activity.	patterns; reduced inter-observer variability.	Novel approach targeting neuroplasticity; focuses on cognitive and neural processes essential for recovery.	Limited to Phase I/II trial; generalizability and long-term impacts require further study.
[15]	Explore imaging techniques for mapping and ablation of cardiac arrhythmias.	Reviewed roles of MDCT, MRI, and ICE for pre- and intra-procedural imaging to improve catheter ablation outcomes. Reviewed advancements	Improved arrhythmia mapping and procedural outcomes; real-time imaging with ICE enhances precision.		Relies on operator expertise; limited by the accessibility of imaging modalities.
[16]	Highlight the role of AI in enhancing the performance of nanogenerators for energy solutions.	in piezoelectric and triboelectric nanogenerators with AI integration for enhanced energy conversion and autonomous operation. Experimental group	Promotes eco-friendly energy solutions; expands applications in robotics and intelligent systems.		Review-based study; lacks experimental validation or specific metrics.
[17]	Examine how verbal fluency performance is affected by cardiovascular exercise.	performed verbal fluency tasks before, during, and after exercise; control group performed tasks without exercise. Reviewed CNN-based	Demonstrated exercise-induced improvements in lexical access and topic switching.		Limited to young healthy adults; findings may not generalize to other populations.
[18]	Analyze deep learning models for lung field segmentation in medical images.	architectures (e.g., U-Net and its variants) and their performance on benchmark datasets using metrics like Dice similarity and Jaccard coefficient. Used an ensemble model	Significant improvements in segmentation accuracy with deep learning methods; facilitates accurate diagnosis.		Challenges with overlapping structures and variability in lung shapes remain.
[19]	Develop a haptic system for anxiety detection and management using EEG data.	for EEG data classification, achieving 97% accuracy; analyzed signal spikes with advanced algorithms and provided haptic feedback.	High classification accuracy; innovative approach for anxiety detection and management.		Limited generalizability; potential challenges with real-world deployment.

### Key Contribution

The most recent machine learning (ML) models used for cardiovascular disease (CVD) diagnosis, risk assessment, and treatment are thoroughly compared in this review study. It investigates the efficacy, precision, and clinical applicability of cutting-edge machine learning algorithms in enhancing patient outcomes and early detection by methodically assessing them. The paper looks at different supervised and unsupervised learning strategies, emphasizing their advantages, disadvantages, and prospective applications in healthcare environments.

### Method, Experiments and Results

**Dataset:** Predicting the existence or absence of cardiac disease using a variety of clinical and demographic characteristics seems to be the main goal of this dataset. To train the model to predict if a patient has heart disease, supervised learning tasks usually employ the goal variable as the label. To make this classification, the model uses input features like blood pressure, cholesterol, age, and lifestyle factors. The objective is to create a trustworthy prediction system that will help medical practitioners identify cardiac disease early and stratify patients based on their risk.

Table 2: Dataset description

Column Name	Non-Null Count	Data Type	Description
id	5110	int64	Unique identifier for each patient
avg_glucose_level	5110	float64	Average glucose level
gender	5110	int32	Encoded gender (e.g., 0 = Female, 1 = Male)
age	5110	float64	Age of the patient
hypertension	5110	int64	Presence of hypertension (0 = No, 1 = Yes)
ever_married	5110	int32	Marital status (0 = No, 1 = Yes)
Residence_type	5110	int32	Type of residence (0 = Rural, 1 = Urban)
smoking_status	5110	int32	Smoking status (categorical)
heart_disease	5110	int64	Presence of heart disease (0 = No, 1 = Yes)
bmi	4909	float64	Body Mass Index (BMI), some missing values
work_type	5110	int32	Type of employment (categorical)
stroke	5110	int64	Stroke occurrence (0 = No, 1 = Yes)

**Pre-processing:**

**Handling Missing Values:** Even if there aren't any missing values in the dataset at the moment, it's always a good idea to look for and deal with any missing data, especially in real-world situations where data quality may vary. The performance and accuracy of the model depend on the dataset being compiled. The dataset is comprehensive and contains all necessary features for analysis, as shown in Figure 1, which contributes to the preservation of the prediction models' integrity. However, to ensure consistent and trustworthy results, future data harvests should incorporate comprehensive checks for missing variables.

	Missing_Number	Missing_Percent
id	0	0.000
gender	0	0.000
age	0	0.000
hypertension	0	0.000
heart_disease	0	0.000
ever_married	0	0.000
work_type	0	0.000
Residence_type	0	0.000
avg_glucose_level	0	0.000
bmi	0	0.000
smoking_status	0	0.000

Figure 1: Missing Value count

**Correlation:** Figure 2 depicts the correlation heatmap, which illustrates critical correlations between variables. Age and marital status have a significant positive association (0.68), while BMI and marital status have a strong positive correlation (0.34). On the other hand, there are negative associations between work type and age (-0.36) and married status and work type (-0.35). Heart disease (0.13), hypertension (0.13), and glucose levels (0.13) all exhibit somewhat positive relationships with stroke, although these are not very strong. These results imply that although these variables raise the risk of stroke, they are not reliable indicators on their own. All things considered, stroke prediction probably necessitates a multi-factorial approach, taking into account how different variables interact to increase predictive accuracy..

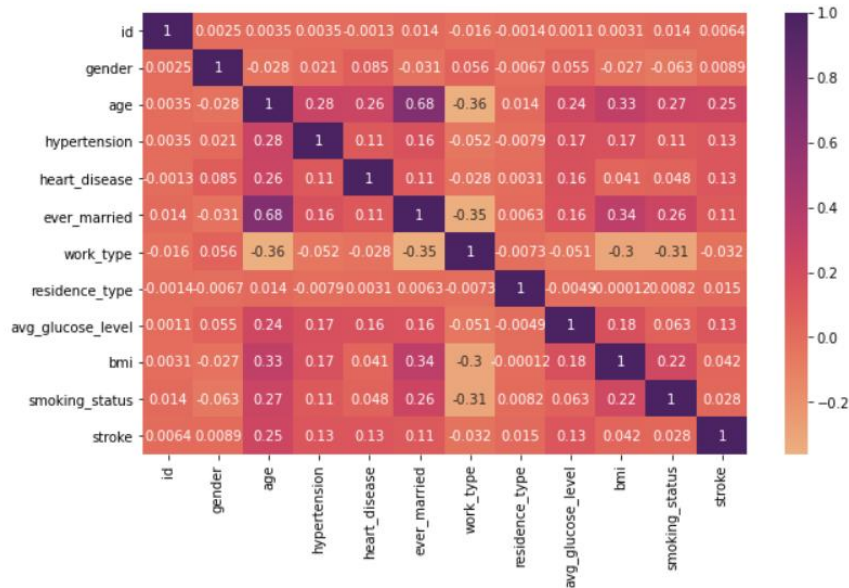


Figure 2: Dataset Correlation

**Different ML Models:**

Comparison of the KNN, SVM with RBF Kernel (SVM\_RBF), Decision Tree, Multilayer Perceptron and Random Forest (RF) algorithms in terms of their key characteristics, advantages, and limitations describe in table 3 and table 4 provides ML model performance metric.

**Table 3:** ML model Key characteristics, Advantage, and Limitation

Algorithm	Key Characteristics	Advantages	Limitations
KNN	- Instance-based learning	- Simple to implement and understand	- Computationally expensive for large datasets
	- Non-parametric	- No training phase	- Sensitive to irrelevant features and noise
	- Lazy learner	- Effective for small datasets	
SVM_RBF	- Kernel-based method	- High accuracy for complex datasets	- Computationally intensive
	- Effective for non-linear data	- Robust to overfitting in high-dimensional spaces	- Requires careful tuning of hyperparameters (e.g., C, gamma)
DT	- Margin maximization	- Easy to visualize and interpret	- Prone to overfitting
	- Tree-based model	- Handles both numerical and categorical data	- Sensitive to small changes in data
	- Splits data based on feature values	- High accuracy and robustness	- Computationally expensive
RF	- Interpretable	- Handles missing data and outliers well	- Less interpretable than single decision trees
	- Ensemble of decision trees		
MLP	- Bagging technique	- Can model complex, non-linear relationships	- Requires large amounts of data
	- Reduces overfitting		
	- Feedforward neural network		

- Multiple layers of neurons
- Scalable to large datasets
- Computationally expensive and hard to interpret
- Non-linear mapping

**Result:**

1. **Accuracy:** evaluates the model's overall accuracy.

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \tag{1}$$

2. **Recall:** also known as true positive rate (TPR) or sensitivity:

$$\text{Recall} = \frac{\text{True Positives (TP)}}{\text{True Positives (TP)} + \text{False Negatives (FN)}} \tag{2}$$

3. **F1 Score:** Balances precision & recall for better performance evaluation.

$$F1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \tag{3}$$

A balanced indicator of a model's capacity to manage FP and FN is the F1 score, which is calculated as the harmonic mean of recall and precision. The Decision Tree performs better than all other classifiers, obtaining the greatest F1 score and exhibiting an exceptional balance between recall and precision, according to the comparison of F1 scores across several models. Random Forest and XGBoost perform marginally worse than AdaBoost and SVM, which come in second and third, respectively, with reasonable performance. The least successful classification in this situation is indicated by KNN and GradientBoost, which have the lowest F1 scores. Overall, the findings imply that more sophisticated ensemble approaches could require more optimization to enhance performance, whereas simpler tree-based models, such as the Decision Tree, might be better suited for this dataset.

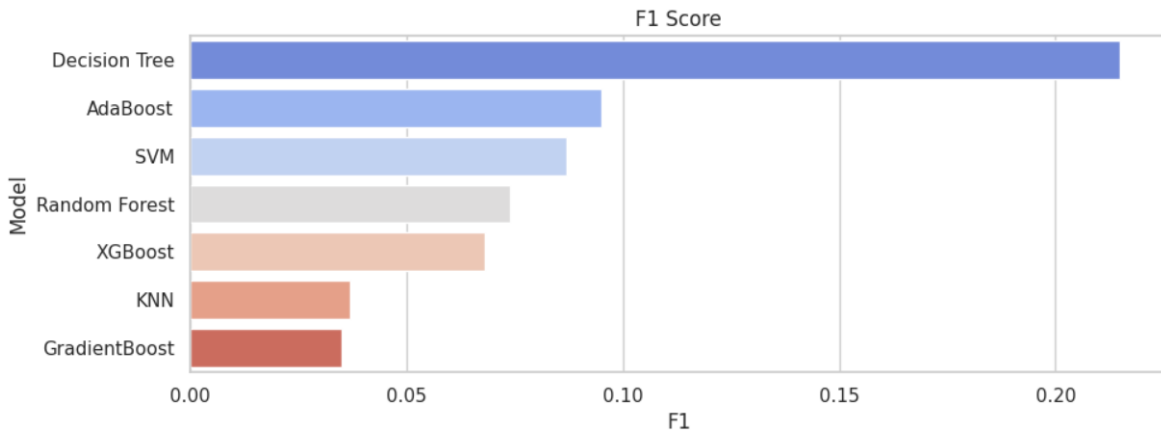


Figure 2: Different ML models F1 score

The Decision Tree model achieves the highest recall, according to the recall analysis in Figure 3, demonstrating its remarkable capacity to accurately identify positive situations while avoiding false negatives. Additionally, AdaBoost and SVM exhibit strong performance, demonstrating their ability to identify pertinent cases. KNN and GradientBoost have the lowest recall values, indicating that systems have difficulty capturing positive examples, whereas Random Forest and XGBoost show reasonable recall. These results suggest that while other models would need more fine-tuning to enhance memory performance, simpler tree-based models, like Decision Tree, might be better suited for situations when high sensitivity is essential.

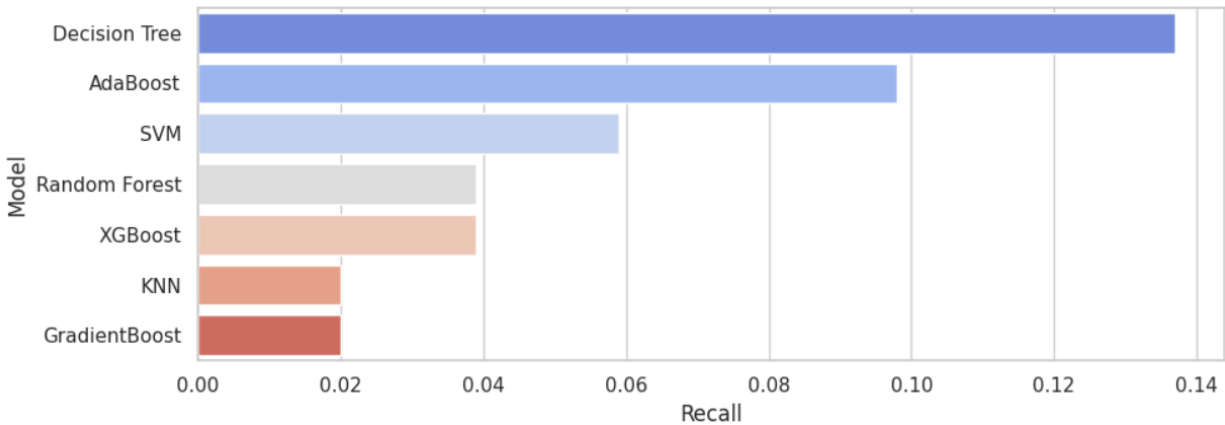


Figure 3: Different ML models Recall

Figure 4's accuracy analysis demonstrates that all models function similarly, achieving high accuracy ratings overall. There aren't many differences between Decision Tree, AdaBoost, SVM, Random Forest, XGBoost, KNN, and GradientBoost, indicating that accuracy might not be the best way to assess model performance in this situation. Metrics like F1-score and recall are essential for evaluating how well the models manage class imbalances and false negatives, even with their high accuracy. To ensure excellent classification performance, more research is required to determine which model provides the best balance between precision and recall.

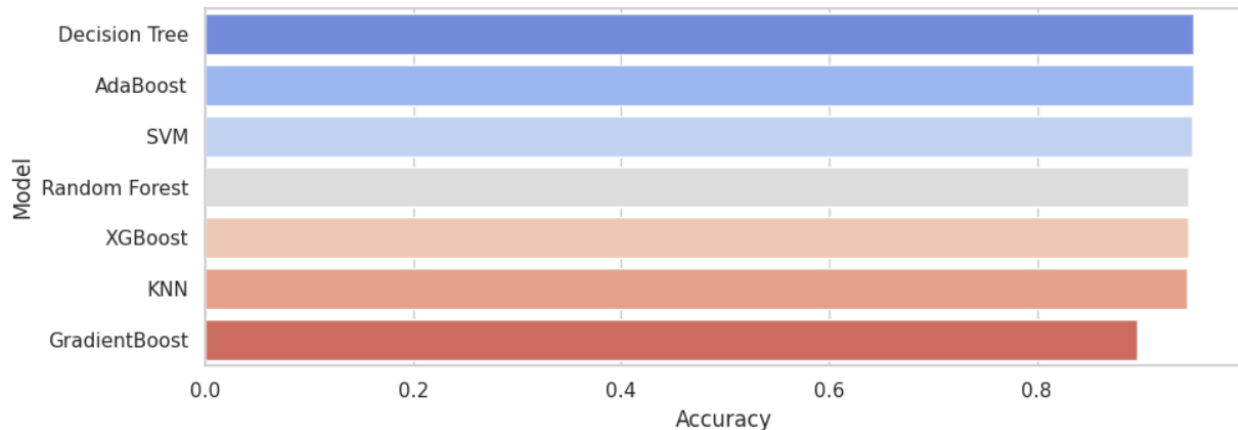


Figure 4: Different ML model Accuracy

### Discussions

The study's findings show that by using clinical and demographic data, machine learning approaches may significantly enhance the prognosis of cardiovascular illnesses. With the greatest recall and F1-score among the models tested, the Decision Tree classifier was particularly good at detecting positive instances while reducing false negatives. In medical applications, where failing to notice a high-risk patient might have serious repercussions, this is crucial. Additionally, ensemble techniques like XGBoost and Random Forest shown excellent classification accuracy and resilience. They are less effective than simpler models, though, since they demand more processing power.

On the other hand, Gradient Boosting and K-Nearest Neighbors (KNN) showed less predictive power, most likely as a result of their high processing cost and sensitivity to data distributions. Because it may not account for class imbalances that are frequently seen in practical datasets, accuracy alone is not the utmost trustworthy metric for assessing models in medical diagnostics, as evidenced by the little range in accuracy between models. Rather, because recall and F1-score better reflect the capacity to identify genuine positives while reducing false negatives, they offer a more realistic assessment of a model's efficacy in detecting high-risk people.

Despite encouraging outcomes, a number of issues still exist. Unbalanced data, when there are much more negative examples than positive ones, might skew predictions and impair the model's generalizability. Furthermore, as medical personnel must comprehend the decision-making process, the interpretability of intricate machine learning models continues to be a major

challenge for clinical adoption. Future studies should concentrate on enhancing feature selection to increase model accuracy, investigating hybrid models that integrate the best features of various methodologies, and using real-world patient data to improve these models' prediction skills for more dependable and useful results.

### Conclusions

This work highlights the potential of machine learning models in predicting the risk of cardiovascular disease by demonstrating that Decision Tree and Random Forest models outperform other classifiers in terms of sensitivity and prediction accuracy. The results show that ML-based approaches can enhance early detection and intervention, eventually reducing the fatality rate from CVD. However, problems including data imbalance, interpretability, and real-world applicability must be fixed before to clinical adoption. Future research should look at deep learning techniques, larger datasets, and integration with electronic health records to further improve model performance and reliability in real healthcare settings. Additionally, heart failure significantly lowers healthcare expenses and quality of life. Therefore, strategies to improve health behaviors related to heart failure are crucial. Effective medical decision aids can be extremely beneficial to those with heart failure. In this randomized controlled study, customized decision aids were given to patients with heart failure in a clinical setting to promote medication compliance, lifestyle changes, and regular condition monitoring. The intervention led to better patient outcomes, such as better disease management, higher quality of life, and fewer hospital readmissions. This demonstrates how important it is to integrate customized decision support tools with machine learning algorithms to both forecast and manage the progression of cardiovascular illnesses.

### References:

- [1] Veroff, David R., et al. "Improving self-care for heart failure for seniors: the impact of video and written education and decision aids." *Population health management* 15.1 (2012): 37-45.
- [2] Ahmadli, N., Sarsil, M. A., Mizrak, B., Karauzum, K., Shaker, A., Tulumen, E., ... & Ergen, O. (2024). Voice-Driven
- [3] Uddin, K. M. M., Dey, S. K., & Babu, H. M. H. (2024). A Voice assistive mobile application tool to detect cardiovascular disease using machine learning approach. *Biomedical Materials & Devices*, 2(2), 1246-1257.
- [4] Abbas, S., Ojo, S., Al Hejaili, A., Sampedro, G. A., Almadhor, A., Zaidi, M. M., & Kryvinska, N. (2024). Artificial intelligence framework for heart disease classification from audio signals. *Scientific Reports*, 14(1), 3123.
- [5] Idrisoglu, A. (2024). Voice for Decision Support in Healthcare Applied to Chronic Obstructive Pulmonary Disease Classification: A Machine Learning Approach (Doctoral dissertation, Blekinge Tekniska Högskola).
- [5] Alosekait, D. M., Shdefat, A. Y., Nabil, A., Nawaz, A., Rana, M. R. R., Ahmed, Z., ... & AbdElminaam, D. S. (2024). Heart-Net: A Multi-Modal Deep Learning Approach for Diagnosing Cardiovascular Diseases. *Computers, Materials & Continua*, 80(3).
- [6] Mayourian, J., El-Bokl, A., Lukyanenko, P., La Cava, W. G., Geva, T., Valente, A. M., ... & Ghelani, S. J. (2024). Electrocardiogram-based deep learning to predict mortality in paediatric and adult congenital heart disease. *European Heart Journal*, ehae651.
- [7] Wang, Y. R., Yang, K., Wen, Y., Wang, P., Hu, Y., Lai, Y., ... & Zhao, S. (2024). Screening and diagnosis of cardiovascular disease using artificial intelligence-enabled cardiac magnetic resonance imaging. *Nature Medicine*, 1-10.
- [8] Khan, M. R., Haider, Z. M., Hussain, J., Malik, F. H., Talib, I., & Abdullah, S. (2024). Comprehensive Analysis of Cardiovascular Diseases: Symptoms, Diagnosis, and AI Innovations. *Bioengineering*, 11(12), 1239.
- [9] Bharathi, S., Gresa, P. S., Manivannan, D., & Sathya, V. (2024, August). Exploring the Potential of Machine Learning and Deep Learning in ECG Image Analysis for Cardiovascular Disease Diagnosis. In 2024 5th International Conference.
- [10] Srinivasulu, B., Reddy, P. S., & Basha, P. H. (2024, May). A Deep Pattern Learning based Model for Detection of Cardiovascular Diseases (CVD). In 2024 4th International Conference on Pervasive Computing and Social Networking (ICPCSN) (pp. 191-196). IEEE.
- [11] Bosak, K., & Thomsen, A. (2025). Implementation of an Exercise Prescription SmartPhrase in the Electronic Health Record. *The Journal for Nurse Practitioners*, 21(1), 105261.
- [12] Thangavel, R. K., Allwyn Sundarraj, A., Ramakrishnan, J., & Balasubramanian, K. (2025). Coati optimization algorithm for brain tumor identification based on MRI with utilizing phase-aware composite deep neural network. *Electromagnetic Biology and Medicine*, 1-18.
- [13] Tavakoli, N., Rahsepar, A. A., Benefield, B. C., Shen, D., López-Tapia, S., Schiffers, F., ... & Kim, D. (2025). ScarNet: A Novel Foundation Model for Automated Myocardial Scar Quantification from LGE in Cardiac MRI. arXiv preprint arXiv:2501.01372.

- [14] Boukrina, O., Madden, E. B., Sandroff, B. M., Cui, X., Yamin, A., Kong, Y., & Graves, W. W. (2025). Improving reading competence in aphasia with combined aerobic exercise and phono-motor treatment: Protocol for a randomized controlled trial. *PloS one*, 20(1), e0317210.
- [15] Muser, D., & Santangeli, P. (2025). Intracardiac Echocardiography, Computed Cardiac Tomography, and Magnetic Resonance Imaging for Guiding Mapping and Ablation. In Huang's Catheter Ablation of Cardiac Arrhythmias (pp. 156-170).
- [16] Xu, S., Manshahi, F., Xiao, X., & Chen, J. (2025). Artificial intelligence assisted nanogenerator applications. *Journal of Materials Chemistry A*, 13(2), 832-854.
- [17] Khanna, M. M., Guenther, C. L., Eckerson, J., Talamante, D., Yeh, M. E., Forby, M., ... & Sacco, M. (2025). Vigorous Exercise Enhances Verbal Fluency Performance in Healthy Young Adults. *Brain Sciences*, 15(1), 96.
- [18] Andrabi, T., & Bhat, S. Y. (2025). Analysis of Deep Learning Models for Lung Field Segmentation. *Deep Learning Applications in Medical Image Segmentation: Overview, Approaches, and Challenges*, 149-183.
- [19] Mishra, S., Seth, S., Jain, S., Pant, V., Parikh, J., Chugh, N., & Puri, Y. (2025). An Emotionally Intelligent Haptic System-An Efficient Solution for Anxiety Detection and Mitigation. *Computer Methods and Programs in Biomedicine*, 108590.