

**Advancing Heart Failure Prediction:
A Comparative Study of Traditional Machine Learning, Neural
Networks, and Stacking Generative AI Models**

By
Howard Hoi Nguyen

A dissertation submitted to
Harrisburg University of Science and Technology
for the degree of
Doctor of Philosophy



Department of Analytics
Harrisburg University of Science and Technology
November of 2024

© Copyright by Howard H. Nguyen, 2024
All Rights Reserved

Ph.D. COMMITTEE APPROVAL

To the Faculty of Harrisburg University of Science and Technology:

The members of the Committee appointed to examine the dissertation of Howard Hoi Nguyen find it satisfactory and recommend that it is accepted.

Maria Viada, Ph.D.

Kevin Purcell, Ph.D.

Kevin Huggins, Ph.D.

Srikar Bellur, Ph.D.

Roosbeh Sadeghian, Ph.D.

ACCEPTANCE PAGE

As a duly authorized representative of Harrisburg University of Science and Technology, I have read the thesis of Howard Hoi Nguyen in its final form and have found that (1) its format, citations, and bibliographical style are consistent and acceptable and fulfill university and department style requirements; (2) its illustrative materials including figures, tables, and charts are in place, and (3) the final manuscript is satisfactory to the graduate committee and is ready for submission to the university library.

Kayden Jordan, Ph.D.

Director of Data Science Ph.D. Program

Harrisburg University of Science and Technology

Kevin Purcell, Ph.D.

Provost

Harrisburg University of Science and Technology

ABSTRACT

Heart failure (HF) remains a leading global cause of morbidity and mortality; early diagnosis and predictive strategies to improve patient outcomes and alleviate the burden on healthcare systems. Traditional prognostic models, including Logistic Regression (LR), Support Vector Machines (SVM), Random Forests (RF), Gradient Boosting Machines (GBM), and Extreme Gradient Boosting Machines (xGBM) have struggled to capture the complexities of HF progression. Due to limited testing on diverse datasets, these models face challenges in nonlinear relationships, handling class imbalance, robustness, and generalizability. This study addresses these gaps by employing datasets ranging from 299 to 400,000 records, enabling a more comprehensive evaluation.

In contrast, deep learning models such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) demonstrate a superior ability to recognize complex patterns. However, they are constrained by their demand for large datasets, computational intensity, and lack of interpretability, limiting their clinical applicability. To bridge these limitations, this dissertation investigates predictive approaches spanning traditional machine learning, deep learning, ML - DL stacking model, Generative AI, and advanced stacking Generative AI (Gen AI) techniques.

A critical component of this research is the application of the Synthetic Minority Over-sampling Technique (SMOTE) to mitigate class imbalance, a prevalent challenge in healthcare datasets. Additionally, the study emphasizes the identification of interpretable predictors across diverse datasets, addressing the transparency gap in many predictive models. Nine datasets, ranging from

299 to nearly 400,000 records and encompassing various demographic and clinical features, evaluate model performance under varied data conditions.

The models explored include LR, SVM, RF, GBM, xGBM, Simple Neural Networks (NN), CNN, GRU with Attention, and CNN with GRU. Stacking models, such as RF, GBM, and xGBM combinations for smaller datasets and RF, GBM, and CNN or RNN for larger datasets, are assessed for their potential to enhance predictive accuracy and generalizability. Furthermore, this research introduces a comprehensive Stacking Generative AI hybrid model, integrating Gen AI with RF, GBM, xGBM, and CNN. By generating synthetic data, the Gen AI component addresses class imbalance, improving the representation of underrepresented patient subgroups and enhancing predictive robustness.

The results demonstrate that while traditional ML and neural network models perform reliably in specific contexts, the Stacking Generative AI model consistently achieves superior results across all datasets. The largest improvement was observed on a dataset with 1,025 records, the Stacking Generative AI model attained an accuracy of 98% and a Receiver Operating Characteristic Area Under the Curve (ROC AUC) of 0.999 (As shown in Table 5 and Figure 15), significantly outperforming individual models. This performance highlights the model's ability to manage complex data patterns, enhance predictive accuracy, and improve clinical relevance.

The Stacking Generative AI model has promising applications in healthcare, including early HF detection, personalized treatment planning, and resource optimization in clinical settings. This research advocates for further exploration of the Stacking Generative AI models in real-world clinical contexts to maximize their transformative potential in healthcare.

To demonstrate the practical utility of these findings, a web application has been developed and is accessible at <https://cvdstack.streamlit.app>. This user-friendly platform allows clinicians and patients to assess HF risk by inputting clinical and demographic information. By providing immediate risk assessments, the platform exemplifies how advanced predictive models, such as the Stacking Generative AI, can be translated into accessible tools that support early intervention, personalized care, and data-driven clinical decision-making.

DEDICATION

To my darling wife, Kaylyn, your love, patience, and all-round support are the anchor and sail of my life. Your presence was the constant reminder of both beauty and joy not just of reaching the many destinations together but, more importantly, of journeying towards them. This work is a testament to our shared dreams and the challenges we've overcome side by side in our American dream.

And to my esteemed professors at Harrisburg University, not only for adding knowledge but also for leaving me inflamed with the burning fire for lifelong learning, your guidance made a whole lot of difference to me. I am deeply grateful for your mentorship and the intellectual challenges you've posed, which have spurred my growth.

My dear parents, incomparable for your sacrifices and your unconditional love, being the only support system in both my failure and success. My earnings in the process are your sown in me hard work, perseverance, and kindness, which has reaped fruit in every step during this journey. This achievement is also yours, just like it's mine.

I also cherish these so much: to all my friends and colleagues out there, an absolutely fabulous network of support, laughter, and camaraderie, please accept my sincerest warm appreciation. Your support, encouragement, and belief in my abilities have been reassuringly huge, yielding a huge support base for motivation. I will always treasure the moments shared with you and the insights exchanged forever.

And to my daughters, Lynn and Jaclyn, who inspire me with their curiosity, joy, and resilience day in and day out. To you two, this work is dedicated. May you be inspired to chase your

dreams and pursue your path, however unique, and never forget that through the power of perseverance, anything can come true. May you forever believe in the beauty of your dreams and how they can be made a reality.

This work is hereby dedicated to y'all, for you guys have been the pillars on which my dreams stand. Thanks for going through this and being my guiding light or mentor.

ACKNOWLEDGEMENTS

Indeed, this dissertation crowns the work and achievement of one of the longest and most difficult, yet at the same time rewarding, journeys for which I am more than grateful to so many persons and various institutions that have been supporting me throughout.

First and foremost, I would like to express my deep sense of respect and gratitude to Harrisburg University of Science and Technology for affording this opportunity. The great faculty at the university, coupled with the necessary resources and an environment congenial to studying, gave an ideal launching pad to go for a doctorate.

Namely, the professors, whose work ethics cannot be verbalized with mere words, played a major important role in my academic growth. From these, thank you to Dr. Srikar Bellur and Dr. Roozbeh Sadeghian, who offered courses on interesting topics such as machine learning and deep learning. Their clear explanations and hands-on approach have equipped me with the technical know-how I needed for this research.

I value the rewarding experience I have gathered and the stimulating discussions during classes led by Dr. Alan Hitch and Dr. Kevin Purcell in the Forecasting-Research Seminar course. Classes were employed to develop research methodology and methods.

I would also like to give great thanks to Dr. Kevin Huggins, Dr. Kayden Jordan, and Dr. Maria Vaida for their invaluable teaching and coaching in the Doctoral Studies class. In this class, I learned very important research skills that played a major role in completing this dissertation.

Lastly, my greatest appreciation goes to my mentor, Dr. Maria Vaida. Her guidance, encouragement, and advice were immensely valuable throughout this journey. Her input not only helped shape this work but also contributed greatly to my development, both personal and professional.

TABLE OF CONTENTS

Ph.D. COMMITTEE APPROVAL	2
ACCEPTANCE PAGE	3
ABSTRACT	4
DEDICATION	7
ACKNOWLEDGEMENTS	9
TABLE OF CONTENTS	10
LIST OF FIGURES	12
LIST OF TABLES	13
Chapter 1: INTRODUCTION	14
Chapter 2: LITERATURE REVIEW	21
2.1. Traditional Machine Learning Approaches	22
2.2. Neural Network-Based Approaches	26
2.3. Hybrid and Stacking Models	30
2.4. Generative AI and GAN Frameworks	34
2.5. Comparison of related literature reviews	37
2.6. Literature Review Conclusion	40
Chapter 3: RESEARCH METHODOLOGY	44
3.1. Overview of Methodology	46
3.2. Data Collection and Preprocessing	48
3.3. Research Questions and Modeling Strategies	51
3.3.1. The Research Questions	51
3.3.2. Modeling Strategies	59
3.4. Core Techniques and Optimization Performance	61
3.5. Models' Design and Implementation	66
3.6. Evaluation Measurement and Validation Methods	76
Chapter 4: RESULTS	80
4.1. Implementation Results	80
4.2. Summary of Results	106
Chapter 5: CONCLUSIONS	109
5.1. Summary of Findings	109
5.2. Comparison with Literature	110
5.3. Implication of the Research Contribution	111

5.4. Conclusion	111
Chapter 6: CHALLENGES AND LIMITATIONS	114
6.1. Data Privacy and Security	114
6.2. Model Interpretability	115
6.3. Ethical Considerations	116
6.4. Technical Challenges	116
This section explores the implications of the research findings, compares them with existing literature, and outlines avenues for future exploration. The discussion reflects on the study's contributions to predictive modeling in healthcare, particularly in heart failure (HF) prediction, while also addressing the limitations and potential enhancements.	
7.1. Discussion	118
7.2. Future Works and Scalability	122
Chapter 8: REFERENCES	127
Chapter 9: APPENDICES	136
Figures	136
Tables of Models Performance	145

LIST OF FIGURES

Figure 1- The diagram of the Generative AI – GAN network.....	64
Figure 2- Stacking model (RF + GBM + xGBM) architecture for smaller datasets	66
Figure 3- Stacking model (RF + GBM + CNN / RNN) architecture for larger datasets.	67
Figure 4- Comprehensive Generative AI Architecture	69
Figure 5- The proposed Comprehensive Stacking Generative AI Architecture	71
Figure 6- The ROC Curve for the dataset of 1,000 records.....	82
Figure 7- The ROC Curve for dataset of 400,000 records.....	84
Figure 8- The Risk Factors / Feature Importances of 70,000-record dataset.....	86
Figure 9- The Risk Factors / Feature Importances of 1,025-record dataset.....	88
Figure 10- The Risk Factors / Feature Importances of 400,000-record dataset.....	88
Figure 11- The Risk Factors / Feature Importances of 4,240-record dataset.....	89
Figure 12- The Risk Factors / Feature Importances of 11,627-record dataset.....	89
Figure 13- The Risk Factors / Feature Importances of 303-record dataset.....	90
Figure 14- The Risk Factors / Feature Importances of 1,000-record dataset.....	91
Figure 15- The ROC Curve for the dataset of 1,025 records.....	92
Figure 16- The ROC Curve for the dataset of 70,000 records.....	93
Figure 17- The ROC Curve for the dataset of 11,627 records.....	96
Figure 18- The ROC Curve for dataset of 4,240 records.....	100
Figure 19- The ROC Curve for the dataset of 303 records	102
Figure 20- The Learning Curve	136
Figure 21- Correlation Matrix Analysis	137
Figure 22- ML and NN Models Accuracy Analysis.....	137
Figure 23- ML and NN Models - ROC AUC Analysis.....	138
Figure 24- Web App for CVD Prediction based on user inputs (Stacking Model)	139
Figure 25- Web App for CVD Prediction based on user inputs (RF & GBM Models).....	142

LIST OF TABLES

Table 1- Model comparisons from literature reviews	37
Table 2- Performance of proposed model vs. other models on dataset of 1,000 records	82
Table 3- Performance of proposed model vs. other models on dataset of 400,000 records	84
Table 4- Performance of proposed model vs. article's models on dataset of 400,000 records	85
Table 5- Performance of proposed model vs. other models on dataset of 1,025 records	92
Table 6- Performance of proposed model vs. other models on dataset of 70,000 records	93
Table 7- Performance of proposed model vs. other models on dataset of 11,627 records	95
Table 8- Performance of proposed model vs. other models on dataset of 4,240 records	97
Table 9- Performance of proposed model vs. article's models on dataset of 4,240 records.....	98
Table 10- Performance of proposed model vs. other models on dataset of 303 records.....	101
Table 11- Performance of proposed model vs. article's models on dataset of 303 records.....	101
Table 12- Summary of all models' performances over the 9 datasets and 11 models.	108
Table 13- Performance Metrics (Sensitivity, Specificity, and F1-Score) for Stacking Gen AI Model Across Datasets	121
Table 14- Evaluation Metrics Across Datasets	122
Table 15- Model performances on dataset of 303 records	145
Table 16- Model performances on dataset of 1,000 records	146
Table 17- Model performances on dataset of 1,025 records	147
Table 18- Model performances on dataset of 4,240 records	148
Table 19- Model performances on dataset of 11,627 records.....	149
Table 20- Model performances on dataset of 70,000 records.....	151
Table 21- Model performances on dataset of 400,000 records.....	152

Chapter 1: INTRODUCTION

Heart disease (HD), particularly heart failure (HF), is still the biggest cause of morbidity and mortality in the world. HF, if caught early and predicted, can be avoided and patients' outcomes improved with appropriate treatment. But prediction in HF remains difficult because of its complex nature and multiple contributing factors. Predictive models could reshape healthcare to assist in early diagnosis and better decision making for both physicians and patients.

Machine learning (ML) and deep learning (DL) methods have become the must-haves in predictive healthcare applications because they can crunch huge amounts of data and discover deep patterns. However common, classical ML algorithms like Logistic Regression (LR), Random Forests (RF), and Gradient Boosting Machines (GBM) often lack the probability of capturing nonlinear relationships and temporal changes of health data. While neural network models such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) can be used better for pattern detection, computational complexity and lack of interpretability make them unsuitable for use in clinical contexts.

These limits were improved using hybrid and ensemble formulations such as stacking, amongst others, to combine the strengths of popular algorithms with improved performance. Stacking models require a meta-learner for the integration of model predictions from base models, with immense potential for improved accuracy and generalizability across different datasets.

However, most of the models do not take class imbalance into consideration and also encapsulate the essence of high-dimensional, heterogeneously distributed data in healthcare.

This dissertation will introduce, for the first time, a new approach-the Stacking Generative AI model-that surmounts these challenges. By integrating the generative strengths of AI into the traditional ML and DL frameworks, the study offers a much stronger, more interpretable, and more useful HF prediction. The proposed methodology will bridge the critical gaps in predictive modeling necessary for improved diagnosis and clinical decision-making in cardiovascular care.

This is important for the stacking model, as the Generative AI component generates synthetic data, balancing out the dataset and hindering the model's performance on minority classes.

Healthcare datasets contain many subgroups of patients that tend to be underrepresented and biased in the predictions. Incorporating GAN-generated data ensures the model is exposed to a greater variety of scenarios, hence making the prediction process much more robust and comprehensive.

Therefore, the study has systematically compared the performances of traditional ML models and neural network-based models with the proposed Stacking Generative AI model in HF prediction. The following research questions guide this study:

1 - Comparative Performance of Traditional and Neural Network Models: How do traditional machine learning models compare with neural network-based models in terms of accuracy, generalizability on diverse datasets and ROC AUC for heart failure prediction?

Random Forest (RF) achieved an accuracy of 83% and a ROC AUC of 0.91 on a dataset of 303 records, while CNN achieved a slightly lower accuracy of 82% but with a ROC AUC of 0.85 (Table 10). As the dataset size increased to 1,000 records, CNN's performance in ROC AUC

improved to 0.85 (Table 2), highlighting the flexibility and generalizability of neural network models compared to traditional ones.

2 - Most Influential Heart Failure Predictors: What are the most influential predictors of heart failure across different datasets, and how do these features influence the overall performance of the models? Identification of these predictors will be essential to enhance both the performance regarding accuracy and interpretability of the models. The following features have been identified during implementation to be the most important predictors of heart failure among the analyzed nine datasets:

- BMI was one of the most consistent top-ranking predictors across all database sizes: 400,000, 11,627, and 4,240 records. It was strongly related to heart failure risk, as shown by the dependency structure in Figure 10, Figure 12, and Figure 11.
- Blood Pressure, Systolic and Diastolic: Application of systolic blood pressure (sysBP) is one of the main parameters throughout the multivariable datasets, specifically in datasets of 70,000 and 4,240 records. sysBP is the most crucial part of the 70,000-record dataset which signifies its prediction power toward heart failure, Figure 8.
- Other top predictors included cholesterol levels, including total cholesterol, HDL, and LDL, especially in dataset 11,627, where the HDL cholesterol-direct was top-ranking in Figure 12.
- Age appeared to significantly contribute to all the data sets, consistent with its well-acknowledged role in heart failure development. It was most important in the 4,240, 11,627, and 70,000-a datasets shown in Figure 11, Figure 12, and Figure 8.
- For the smaller datasets (1,025, 1,000, and 303 records), Chest Pain (cp) had a high influential impact, hence further indicating the importance of symptoms such as chest pain in

the early diagnosis for focused heart-related studies (Figure 9, 14, and 13). These predictors not only improve the performance of the models but also shed light on the underlying risk factors for heart failure. Including these variables in the prediction models should result in better accuracy and interpretability to facilitate early detection of heart failure.

3 - Hybrid Stacking Model Potential: How does a hybrid model incorporating both traditional machine learning and deep learning techniques provide improved prediction performance compared to the use of single models?

The specific research question goes whether the hybrid stacking model, incorporating both the ensembling tradition of ML models such as RF and GBM and deep learning models such as CNNs and RNNs, finally illustrates higher performance in terms of the prediction accuracy across all datasets. This study finds that the strong generalization was obtained only for the hybrid stacking model of ML and DL. It obtained 82% accuracy with 0.90 ROC AUC (Table 10) for the dataset in 303 and came up to 94% with 0.98 ROC AUC for the 1,000-record dataset (Table 2). This ML and DL stacking models outperformed single models such as LR, SVM, and stand-alone CNN and RNN models. This way, the hybrid approach effectively integrates the strengths of ML and DL techniques into one framework that is likely to showcase improved generalization and prediction accuracy, plus a very promising solution for predictive healthcare applications.

4 - Generative AI to Boost Predictive Precision: The GAN components in the Stacking Generative AI especially improves the performance of the stacking model compared to the solo models. Does this approach promote better generalizability and scalability of the model across diverse healthcare settings?

The proposed Generative AI model's accuracy was relatively high, reaching 95% with an ROC AUC of 0.99 on the dataset of 1,025 records (As shown in Table 5). This indicates its ability to effectively deal with class imbalance and increase minority class prediction. Increasing this further with Generative AI, the Stacking Generative AI improved an accuracy of 98% and ROC AUC of 0.999 on the dataset of 1,025 records; hence, it was the best-performing model in this study.

5 - How does the unique Stacking Generative AI model specifically contribute to advancements in the healthcare industry, particularly in predicting and managing heart failure?

This unique Stacking Generative AI model contributes significantly to the healthcare industry and, in particular, to predict and manage heart failure. This approach uses both traditional machine learning algorithms, including Random Forest, Gradient Boosting Machines, and Extreme Gradient Boosting Machines, and deep learning architectures, which include Convolutional Neural Networks and Recurrent Neural Networks with Generative AI to overcome some of the prime limitations found in existing predictive models. The model improves predictive accuracy, handles imbalanced class problems, and generalizes well across diverse patient populations. Besides that, it can enable personalized treatment plans, decision support for clinicians, and provides awareness among the patients related to their risks for heart failure. This proposed model will set a new standard not only for predictive health tools but also pave the way toward better clinical outcomes and improved patient care in the management of heart failure.

This research performs extensive and rigorous quantitative analysis to explore performance development and validation of machine learning and deep learning models in a structured way

for various datasets. The proposed research used nine different datasets on heart failure with record sizes from 299 records to 400,000 records to ensure that the developed models are generalizable for different population sizes and settings. This includes the preprocessing of datasets by cleaning the data, normalizing it, dealing with missing values for the integrity of the data, and balancing the datasets, especially the class imbalance problem found in healthcare datasets, by using the Synthetic Minority Over-sampling Technique (SMOTE).

These include many models, from traditional machine learning to hybrid stacking models. The Stacking Generative AI model is central in this research and has represented for the first time the known application of Generative AI combined with traditional ensemble learning techniques in the context of heart failure prediction. While GAN generates synthetic data to enhance model training, particularly in improving the minority class representation, the RF, GBM, xGBM, and CNN ensemble further refines the predictions.

It, therefore, augments the increasingly developing repository of knowledge in predictive health by introducing a new stacking model, Generative AI, that realizes better results in accuracy, robustness, and generalizability. Capable of combining synthetic data generation with the strength of conventional and deep learning models, the Stacking Generative AI model may allow an increase in predictive accuracy for complex, high-dimensional healthcare data. Considerably, the use of Generative AI - in that respect - addresses one of the fundamental issues with the analysis of healthcare data: class imbalance. In this respect, it generates synthetic high-quality data for minority classes, enhancing the model's ability to detect infrequent events such as heart failure in the patient group that is usually underrepresented. Imbalanced datasets pose significant challenges in machine learning models for healthcare, particularly in heart failure prediction,

where minority cases are often underrepresented. Generative Adversarial Networks (GANs) and techniques like SMOTE have emerged as effective solutions to mitigate this issue. However, synthetic data can introduce biases that affect the generalizability of models. This study addresses these challenges by evaluating the Stacking Generative AI model across nine datasets, including the Framingham dataset of 4,240 records, to investigate biases, performance variability, and real-world applicability.

This study will further help translate AI into clinical practice, advancing the field in predictive accuracy and providing a scalable and adaptable model for varied healthcare environments, from large hospitals to smaller clinics. Additionally, this work compares traditional, neural network-based, and hybrid models to enable the medical domain to understand the strengths and drawbacks of the approaches considered, moving toward an accurate diagnostic tool for heart failure predictions.

Chapter 2: LITERATURE REVIEW

Heart failure has emerged as a high-priority public health disease due to its high prevalence and mortality rates. Early diagnosis and precise stratification are urgently needed to realize improved patient outcomes through attempted reductions in the severity of the disease. There has been an increasing development of machine learning and deep learning models in healthcare to deal with such predictive challenges. Traditional machine learning models, such as Logistic Regression, Random Forest, and Gradient Boosting Machines, have proven to be reliable in given contexts but also generally lack the capability to capture the complex nonlinear relationships that are inherently part of the data on heart failure. Conversely, other models, like neural network-based ones, such as CNN and RNN, are more representative in terms of pattern recognition but pose a big barrier to clinical applicability due to large computational requirements and poor interpretability, among other factors.

With these challenges, there is a recognized need for predictive methods merging strengths of both neural network-based and more traditional ML approaches. This need has driven recent studies to hybrid and ensemble-type models, from which over the recent years, stacking has been one of the most widely used methods combining multiple models to boost predictive performances. In this regard, the literature review presented here discusses various existing studies that compare classic ML and DL model performances, with particular interest in their application concerning heart failure prediction.

It further goes on to explore the possibility of hybrid models, the proposed Stacking Generative AI model, which uses GANs (Figure 1) to overcome the deficits in the existing methods. This review identifies the main predictors of heart disease and explores how GANs can be used in a

hybrid model to enhance the accuracy and robustness of such predictions. This literature review will form a foundational understanding of HF prediction models through a critical analysis of benefits and shortcomings of different approaches. The objective of this chapter will also be to assess the contribution of innovative hybrid models, such as the Stacking Generative AI framework, toward improved predictive accuracy and enabling better clinical outcomes.

2.1. Traditional Machine Learning Approaches

Traditional machine learning in heart disease prediction has various outstanding works for different reasons, where each provides insight into different strengths and weaknesses associated with different models. Nevertheless, this remains a continuously evolving area within predictive modeling, whereby innovations are continually required, especially in domains such as healthcare, where heterogeneity in data, class imbalance, and limited features are identifiable challenges in any given situation. This review section will subsequently examine the five related studies employing traditional machine learning models, assessing their contributions and limitations, and compare these with the advanced hybrid approach represented by the proposed Stacking Generative AI model.

First, the work of Chicco and Jurman (2020) explained the predictive power of ML in estimating the Survivors of Heart Failure patients by using a dataset of 299 patients, different machine learning models have been studied in order to identify the most critical predictors of survival, considering serum creatinine and ejection fraction. Performance metrics are presented to show that Random Forest achieved the highest results, with 74% accuracy and 0.80 ROC AUC, outperforming others in these metrics. These had been compared to other models, to be exact, Decision Trees, achieved accuracy of 73.7% and a ROC AUC of 0.68, while Gradient Boosting gain an accuracy of 73.8% with a ROC AUC of 0.75.

This literature points out that a model using only serum creatinine and ejection fraction greatly simplifies it and does quite well, especially when using Random Forests and Gradient Boosting. At the same time, when this analysis was only restricted to these two features, the Random Forests reached +0.418 MCC (Matthews correlation coefficient) and an ROC AUC of 0.698.

The strengths of Chicco and Jurman (2020) study revolve around its focus on feature simplicity, which probably facilitates implementation in clinical settings. However, its findings, with a small sample size and narrow feature set, were limited to generalizability. Using only two features restricts the model's applicability in clinical settings, making it miss the full complexity of heart failure prediction. Also, the authors did not adopt more powerful methods, such as deep learning methods or generative models, to sidestep the limitations of either of these methods—especially the tiny dataset constraint.

Singh et al. (2024) introduces a study of an integrated machine learning approach for congestive HF prediction that covers the use of machine learning in the area of predicting congestive heart failure. The research used the Cardiovascular Health Study dataset containing 5,888 records of patients resulting from the collection of over 400 features for coronary heart disease and stroke assessment in older adults. This dataset was challenging due to a high rate of missing data and a large number of irrelevant attributes. In that respect, the authors applied an extensive pre-processing methodology, comprising the C4.5 algorithm for feature selection and the K-Nearest Neighbor technique for imputation of missing data. This approach yielded 12 critical features for training a model.

In this study, the authors evaluate several machine learning models' performances that includes Decision Trees, Random Forest, Support Vector Machines, Logistic Regression, and a Deep

Neural Network (DNN). The model of DNN shows the best performance, with accuracy of 95.3% and a ROC AUC of 0.97, hence proving that combining advanced data pre-processing techniques with DNNs is effective in handling complex and noisy datasets.

The strengths are the rigid pre-processing approach that enhances the reliability and predictive accuracy of the models. Limitations include the intrinsic imbalance and the complexity of the dataset, where generalizability may be compromised. If possible, future studies may further extend the current study by applying the proposed methodology to large datasets balancing classes, and by performing additional deep learning network investigations in pursuit of performance improvements and clinical utility.

The paper presented by Hasan and Saleh (2021) proposed an advanced predictive model for risk assessment for heart attack using the Framingham Heart Study dataset from the UCI repository, which contains 4,239 instances and 16 features. Demographic and health-related factors such as age, blood pressure, cholesterol levels, smoking status, and diabetes history are some of the important points in this dataset, playing a critical role in the accurate prediction of heart disease. The authors have used an ensemble learning technique, stacking, to integrate different machine learning algorithms—Logistic Regression, Support Vector Machine, Decision Tree, Random Forest, and Extreme Gradient Boosting with Logistic Regression as the meta-learner. This model tries to improve prediction accuracy by benefiting from the strengths of different algorithms.

This proposed stacking model returned a high accuracy of 96.69% with an ROC AUC of 0.98, outperforming all base models comprising a Random Forest of 93.69% and a Decision Tree of 92.71%. These results show how effectively this ensemble model can handle complex, multi-dimensional data. The authors utilize ensemble learning to increase predictive precision. Its

strength lies in its robust comparison with single traditional models. However, the model admitted a limitation of dependence on electronic health records without ECG data since, in the future, additional physiological data could be integrated to improve the models' performance further.

Rajendran et al. (2021) also applied the ensemble approach, just like Hasan and Saleh (2021), but with another blend of models—Support Vector Machines, Random Forest, and Gradient Boosting—applied to the UCI Cleveland dataset. In this way, the ensemble approach achieved an accuracy of 92% and ROC-AUC of 0.94 while outperforming other individual models by a large margin. Thus, this study is another exemplary work that has shown how different traditional machine learning models can be combined for a diverse approach that might have the potential for better performance with a small dataset size of 303 records, as in the Cleveland dataset. However, similar to other works reviewed, their approach did not consider any deep learning or hybrid approaches that might yield a broad predictive framework. Without considering generative models or even class balance, generalization to more extensive and more diverse datasets or those cases in which some conditions, like heart disease, are so less frequent was considerably limited.

Last but not least, the approach of Rimal, Y. et al. (2024) was more optimization-oriented; they used Random Forest with Bayesian optimization and Genetic Algorithms to optimize the respective model hyperparameters. Moreover, for the tuned version of the RF model by Rimal, Y. et al. (2024), the accuracy reaches between 91% - 95%, while ROC-AUC is 0.85 - 0.95 (Table 10), depicting how careful tuning of the hyperparameters can drastically improve conventional machine learning models. Their work managed to augment the performance of the traditional models with optimization techniques. However, it did not go further into advanced

machine learning techniques like deep learning or model stacking. Also, their work did not integrate generative AI, which could have readily mooted the development of a more robust framework for handling more extensive and more complex datasets and issues of class imbalance.

Contrasting these more traditional approaches, the proposed Stacking Generative AI model really provides a panacea solution to some of the challenges pointed out by these studies. This is so because the culmination of traditional machine learning models, such as Random Forest and Gradient Boosting, with deep learning techniques using Convolutional Neural Networks results in more potential power through Generative AI. It fills in the deficiencies of the models that these five studies were founded on. Generative AI, within the proposed model, is very important because it alleviates the class imbalance problem that the other models did not address. This indeed creates synthetic data, hence these minority classes in the dataset will be better represented, and their conditions will have better recall. Also, the hybrid nature of the proposed model captures both simple and complex patterns in the data, making it a flexible and powerful tool for predicting heart disease. The model proposed in this dissertation performed for 95% in accuracy and ROC-AUC as 0.99 (Table 10), which was performing better compared to the results reported in the reviewed articles.

2.2. Neural Network-Based Approaches

Whereas neural networks represent one of those revolutionary approaches to predictive modeling in general and the prediction of HF in particular, several studies explored different architectures of deep learning (DL), outperforming traditional machine learning models, but at the same time, every single study had a number of advantages and disadvantages.

A very exemplary work in this respect is the study from Mahmud et al. (2023), the authors applied a combined dataset of five benchmark heart disease datasets, namely Statlog Heart, Cleveland, Hungarian, Switzerland, and Long Beach, this aggregated dataset contains 920 records and 11 clinical features. Their approach was to develop a lightweight metamodel that combined the merits of standard machine learning algorithms, namely Random Forest, Gaussian Naive Bayes, Decision Trees, and K-Nearest Neighbors. The accuracy of the model presented equaled 87% and was higher in comparison with all results of other separate models. This multi-algorithm combined model increased the general quality of prediction and robustness of the clinical application of this model. However, beyond the merits of their metamodel, it contained its own deficiencies. Such as in case where Mahmud et al. (2023) was overly reliant on the use of traditional machine learning techniques. This, in turn, ultimately limited the model's capacity to capture deeper and more complex patterns within the data. While this lightweight design is efficient, it does sacrifice some of the predictive power that could have been derived from using advanced deep learning algorithms. While this simplicity of the model worked fine for certain applications, it could not leverage the full power of neural network-based methods which possibly extract deeper relationships from the data.

On the other hand, Choi et al. (2017) were the first to propose Recurrent Neural Networks (RNNs) with Gated Recurrent Units (GRUs), in particular—for capturing more representative early prediction for heart failure using EHRs. The dataset from the Sutter Health System included 3,884 heart failure cases and 28,903 control patients. The strength of this model was capturing temporal sequences: through monitoring clinical events over time, a model may find patients at risk for heart failure. The RNN model with GRUs outperformed these traditional methods such that with an observation window of 18 months, it yielded a ROC AUC of 0.883

against a ROC AUC of 0.834 for the best performing baseline model using a Multilayer Perceptron (MLP). This work underlined the fact that temporal modeling is an important aspect of clinical prediction, due to the consideration that every forecast needs considering the temporal development of the health of the patient. The Choi et al. (2017) model had a set of limitations despite such strong results. While RNNs are very strong in temporal modeling, using only the GRUs may not capture a full breadth of predictive power than could be possible with ensemble methods or hybrid models using deep learning in concert with other machine learning. Further enhancement regarding the predictive performance can be integrated into the model by expanding the other architectures or techniques, which also includes the implementation of neural networks with convolutional layers or hybrid stacking methods.

Arooj et al. (2022) introduces a research study based on early detection using a Deep Convolutional Neural Network (DCNN). In this regard, the dataset was selected from heart diseases obtained from the UCI repository that contained 1,050 records and 14 attributes. The model, DCNN, had an accuracy of 91.7%, showing lots of capability in deep learning for discovering complex nonlinear patterns in clinical data. The advantages used by CNNs in processing high-dimensional features helped bring performance in the classification of heart diseases. This is somewhat limited by the narrow focus that Arooj et al. (2022) has on DCNN. The authors also have not looked into other deep learning architectures or hybrid models that may combine the strengths of several approaches. Their findings did not lend themselves to generalizability outside the data set they employed, which may raise some questions concerning its generalization power across more diverse or real-world clinical settings. Because the involved authors considered one single dataset and one model architecture, it logically follows that the study could not then exploit the full potential of this hybrid method, which can result in further

improvements in performance as well as extension of applicability to various healthcare scenarios.

Sakthi et al. (2024) introduce a Transformer-based deep convolutional network to predict heart anomalies using clinical data. The authors employ a dataset from Kaggle, contains 2,200 records and eight clinical features. They integrated transformer architectures into the prediction of heart anomalies, such as Feature Transformer and Tab-Transformer. Results achieved an accuracy of 88.6% with Feature Transformer, outperforming some traditional models, like LightGBM and Category Embedding, accuracy of 86.4% and 87.5%, respectively. Although transformers were developed for natural language processing tasks, they have gained much power in dealing with sophisticated tabular clinical data and showed promising results on the heart anomaly prediction task. While transformer models create quite powerful ways of capturing relationships in structured data, Sakthi et al. (2024) have not studied the integration of these models into more traditional machine learning techniques nor how these hybrid models outperform transformer-only architectures at performance. Applications of transformers to clinical data are still in their infancy, and a lot of work has yet to be done to explore whether ensembling them with other deep learning modalities, like CNNs or RNNs, or even traditional machine learning models, brings any additional value.

Tuli et al. (2020) proposed an integrated IoT and fog computing-based HealthFog framework, an ensemble deep learning-based healthcare system used for real-time heart disease diagnosis. This study's Cleveland Heart Disease dataset comprised 14 critical attributes, including age, chest pain type, cholesterol level, and fasting blood sugar. This dataset is used to train various deep learning models at the edge and then establishes an ensemble to improve prediction performance.

The HealthFog leverages the FogBus framework to distribute the computation among the fog, edge, and cloud for minimum latency and higher accuracy. From the testing results, HealthFog achieves a high % prediction accuracy of 91.2% with a ROC AUC of 0.94, outperforming traditional single-model architectures. The framework also shows huge strides in response time and energy efficiency, something of prime importance in healthcare applications where latency sensitivity could be extremely high.

The strengths were that it effectively integrated fog computing into real-time health monitoring and applied ensemble learning to improve the diagnostic reliability of the model. The main limitations affecting scalability are a preconfigured sensor dependency and constrained device resources. However, future work should extend the HealthFog capabilities to support other critical health conditions and test its effectiveness in larger multi-center healthcare settings.

These works present the spectrum from the simple, lightweight machine learning metamodel as proposed by Mahmud et al. (2023) to the complex deep-learning architectures developed by Choi et al. (2017), Arooj et al. (2022), Sakthi et al. (2024), and Tuli et al. (2020) in neural network-based approaches for predictions of heart disease. While each of these studies has contributed much in their own ways to the literature, scalability, and generalization, remain very guarded, and hybrid models that can bring together the strengths of various approaches toward even better results in predictive healthcare remain few and far between.

2. 3. Hybrid and Stacking Models

Hybrid and stacking models have been an approach that really improves the predictive accuracy of machine learning models in general, particularly those applications dealing with healthcare.

Different works from several authors have presented clearly how such models outperform those

using single algorithms due to the capability of capitalizing on the comparative strengths of multiple models, thus compensating for the comparative weaknesses. A review of related studies' literature shows a few major reviewing works that indicated hybrid and stacking models to be effective in the prediction of heart diseases.

Ali et al. (2020) presented a deep learning-based smart health monitoring system integrated with feature fusion for predicting heart disease. Their system processes physiological data from various wearable sensors in combination with electronic medical records to develop an ensemble of deep learning models, enhancing the predictive capability of heart disease diagnosis. The study scored an incredibly high accuracy of 98.5%, hence showing the power of deep learning in cases where the data feature is high-dimensional and diversified in sources. Alternatively, the model proposed in Ali et al. (2020) relies on deep learning models alone, without combining traditional machine learning approaches or even taking into consideration the strengths of hybrid stacking ensembles. For this reason, their results may generalize less easily across other datasets or populations, as only one dataset has been used to implement the experimentation. It could undermine adaptability and effectiveness that are based solely on deep learning across a wide range of real-world health care settings.

Meanwhile, Mienye et al. (2020) have studied the enhancement of ensemble learning methodologies using Cleveland and Framingham datasets for the risk prediction of heart diseases. Their study proposed an average-based quasi-split strategy to segment the datasets into sub-datasets and then modeled these segmented datasets using the recursive partitioning algorithm known as CART. The models so generated were combined using Accuracy-Based Weighted Aging Classifier Ensemble, which they called AB-WAE. Mienye et al. (2020)'s ensembling methodology apparently had good results, with classification accuracies of about

93% in the Cleveland dataset and 91% in the Framingham dataset. However, their dependence on traditional machine learning algorithms restricted their model's power. While their ensemble approach performed well, it lacked any deep learning techniques that might further improve the model's performance in terms of accuracy and modeling complex patterns existing within the data. Another limitation involves the fact that this study focuses on two datasets only – a fact that raises questions about its generalizability on other populations or healthcare data sets.

Again, Wankhede et al. (2022) introduced the hybrid model by proposing deep learning models together with a feature selection algorithm known as Tunicate Swarm Algorithm-TSA. The network hybrid ensemble deep learning model they proposed resulted in 97.5% accuracy from the UCI Cleveland heart disease dataset. This seminal work corroborated the concept on the amalgamation of deep learning with optimization algorithms in predictive performance.

However, as with the works of Ali et al. (2020) and Mienye et al., (2020) the approach that Wankhede et al. (2022) described shared the limitation in that it did not consider traditional machine learning models and left again space for an approach which could represent both traditional machine learning and deep learning in a more complete way. This study also relies on a rather small dataset. It, therefore, makes it hard to judge its scalability and generalization when it involves larger or more diverse datasets.

Shickel et al. (2018) conducted a broad survey on deep learning for EHR analysis; the authors grouped the models into five major areas of interest: information extraction, representation learning, outcome prediction, phenotyping, and de-identification. The survey covers various EHR datasets including MIMIC and i2b2, pointing out that some common features used in these datasets include diagnosis code, medication history, and lab results, which are essential features in predicting patient outcomes and modeling disease progression. Regarding predictive

performance, deep learning methods such as RNNs and CNNs outperformed traditional models considerably. Among those, RNNs are designed for sequences, and their effectiveness in handling time-series data from EHRs has often been much higher in terms of ROC AUC than that of traditional methods for tasks such as disease prediction and patient outcome forecasting. The survey does not provide specific values of the ROC AUC or accuracy but points out in general that deep learning outperforms traditional models in a clinical context.

The strengths of this survey are the clear categorization of deep learning applications and the emphasis on model interpretability so crucial for clinical use. At the same time, however, one important limitation is that these works still rely, in most cases, on proprietary datasets, which would somehow raise a real barrier to replicability and wide applicability. The authors recommend further studies to establish standard benchmarks so that the models can be more transparent to wider use in various health settings.

Finally, Liu et al. (2022) introduced another approach to predicting cardiovascular diseases using the stacking model fusion. They combined this ensemble framework with various classifiers, namely Support Vector Machines, K-Nearest Neighbor, Logistic Regression, Random Forest, Extra Tree, Gradient Boosting Decision Trees, XGBM, LightGBM, CatBoost, and Multilayer Perceptron into a single model. For improvement in performance, overfitting was avoided by adding a meta-learner based on Logistic Regression. Results have shown that the Liu et al. (2022) model turned out really well on the fused Heart Dataset and public Heart Attack Dataset at a high level of performance, ROC AUC of 0.95 and 0.92 respectively, considering accuracy, precision, recall, F1 score, and AUC. The shortcoming of this model is that it is not interpretable and does not involve deep learning techniques or Generative AI—which would open up possible further avenues toward better performance. The reason is that, by design and origin, their

argument was to derive from traditional machine learning classifiers, which limited the attainment of the full model's capability to capture intricate relationships in data.

Each of the identified studies brings value to the review of hybrid and stacking models in healthcare prediction. However, all studies have serious limitations related to model interpretability and scalability, including deep learning and Generative AI. This opens the avenue for more comprehensive approaches within a hybrid framework that would serve better through traditional machine learning and deep learning from performance, scalability, and generalizability perspectives across diverse releases of health datasets.

2. 4. Generative AI and GAN Frameworks

GANs carved themselves out as one of the most innovative methodologies in CVD prediction right from the beginning. They generate synthetic data that overcomes the class imbalance barriers, limited sample size, and intrinsic complexities in the heart disease risk factors. A review of four recent studies on the application of GAN frameworks in the detection of heart and myocardial infarction diseases shows a number of their strengths and weaknesses.

The first study, by Khan et al. (2024), which presented a hybrid model that combined traditional machine learning with deep learning techniques into an ensemble. Among the various datasets used in this work was the UCI's Heart Disease dataset, which consists of 303 records to forecast cardiovascular diseases with higher performance. The model architecture GAN supported the synthesis of synthetic data dealing with heart diseases with a view to balancing the dataset for missing conditions of disease. These then resulted in 95.3% for the EnsCVDD-Net and slightly improved to 96.1% for the BICVDD-Net. This study underlined the efficiency of GAN-based

data generation in enhancing predictive models, especially diseases considered to be of a rare or complex nature—like heart failure.

In contrast, a recent review highlights the impact of synthetic data on improving clinical predictions (Khan et al. 2024). The authors also utilized GAN for synthetic data generation, balancing the distribution of outcomes for cardiovascular diseases using the Cleveland Heart Disease (303 records) and Framingham (5,200 records) datasets. Indeed, the present contribution is among the first studies to unveil the power of GAN-generated data in class imbalance problems, a condition shared by most medical datasets where phenomena of interest are usually negative, such as in the case of heart disease. It obtained quite promising results, with 85% accuracy for the model using synthetic while it was only 82% when considering purely real data. It also showed that the AUC score for the GAN-based model was 0.927, grossly higher than the one from traditional models, which yielded an AUC score of 0.873. Therefore, it was indicated that synthetic data is one helpful tool in improving the predictive outcome, especially for those rare conditions or outcomes that need to be more robustly represented in training sets.

The third study from Yu S et al. (2024), was based on the KORA cohort study. The authors introduced a novelty in the use of a GAN model along with a feature-enhanced loss function to improve MI prediction. The current dataset contained 1,454 participants, while the key focus areas of this dataset were clinical and metabolic variables related to MI. Apart from that, Yu S et al. (2024) focuses on the feature-enhanced loss function applied to the GAN framework that presents high predictive accuracy of the identification of risk cases for MI. The accuracy of the GAN model reached 94.62%, whereas its AUC was also very high: 0.958. Another distinguishing factor of this research was the ability of the loss function to focus on feature importance and, by doing so, boost the quality of the predictions and give clinically greater value

to which variables contribute most to a risk of myocardial infarction. This combination of GAN with an elaborately tuned loss function made the former one of the more innovative approaches reviewed.

The study from Anbarasu, P. N., & Suruli, T. M. (2022) addresses the challenge of limited annotated clinical data in healthcare by developing a Deep Ensemble Learning Model (DELM) combined with a Generative Adversarial Network (GAN)-based semi-supervised training algorithm. The model aims to enhance the efficiency of Clinical Decision Support Systems (CDSS) by predicting unannotated data more accurately. The DELM comprises two levels: initial-level classifiers (e.g., SVM, KNN, Naive Bayes, and Random Forest) and a second-level Deep Neural Network (DNN) to integrate predictions for final classification. GANs are utilized to augment and balance datasets by generating synthetic samples, solving the class imbalance issue prevalent in medical datasets. Experimental results demonstrate superior performance, achieving accuracies of 86.54%, 84.83%, and 86.72% on the SPECT, WDBC, and Hallmarks datasets, respectively, surpassing other traditional models like Fuzzy-AHP+ANN, LightGBM, and PSO-DNN.

On the other hand, Bhagawati and Paul (2024) applied the GAN framework for predicting coronary artery disease using the dataset from the UCI Machine Learning Repository. A total of 1700 participants were investigated in this study, where 52 risk factors were identified as office-based biomarkers, laboratory-based biomarkers, carotid ultrasound imaging phenotypes, and medication usage. The GAN model outperformed much in comparison to RNN and LSTM. The presented work has shown the generation of synthetic data through GAN efficiently and with an accuracy of 93% and an AUC of 0.953, toward balancing and providing proper representation of high-risk CVD cases. Importantly, this framework was further compared against models devoid

of GAN-generated data, and the result was emphatic: models augmented with synthetic data courtesy of GANs granted better accuracy and higher AUC scores to signify the worth of using GAN frameworks in clinical tasks of prediction.

The GAN frameworks for the prediction of heart disease and myocardial infarction in these studies proved to be a very strong tool. Each of the studies described how the synthetic data could be helpful in boosting model accuracy, especially when facing the common challenge of class imbalance, where high-risk patients are usually underrepresented in the medical datasets. The study further showed that GANs have this added advantage in enabling models combined with traditional machine learning or deep learning models to learn from balanced synthetic datasets toward better predictive performance and generalizability. Although the concrete architectures and datasets vary between these works, a general conclusion that can be drawn is that GANs promise a very bright outlook for improving the field of predictive analytics in healthcare, especially in application domains where data limitations traditionally have kept model performance-constrained.

2.5. Comparison of related literature reviews

Table 1- Model comparisons from literature reviews

Study	Methodology	Dataset	Accuracy	ROC AUC
Machine Learning Can Predict Survival of Patients with Heart Failure from Serum Creatinine and Ejection Fraction Alone (2022)	Logistic Regression, SVM, RF, GBM	UCI Cleveland Heart Disease Dataset (303 records)	77%-85%	0.84-0.92
An Integrated Machine Learning Approach for Congestive Heart Failure Prediction (2023)	DNN	UCI Cleveland Heart Disease Dataset (5888 records)	95.3%	0.97

Cardiac Failure Forecasting Based on Clinical Data (2023)	Random Forest	Clinical Data Dataset (multiple datasets)	89%	0.91
Hyperparameter Optimization: A Comparative Machine Learning Model Analysis (2024)	Gradient Boosting Machine, SVM	UCI Cleveland Heart Disease Dataset (303 records)	91%	0.92
Using Recurrent Neural Network Models for Early Detection of Heart Failure Onset (2023)	RNN, LSTM	Sutter Palo Alto Medical Foundation (Sutter-PAMF) (28,903 records)	90%-95%	0.92-0.95
Heart Disease Detection: A Comprehensive Analysis of Machine Learning, Ensemble Learning, and Deep Learning Algorithms (2024)	ML, Ensemble Learning, and DLs	Heart Statlog Cleveland Hungary final (294 records)	94.34%	-
HealthFog: An Ensemble Deep Learning-Based Smart Healthcare System (2022)	Ensemble DL (CNN, RNN with Fog Computing)	UCI Cleveland Heart Disease Dataset (303 records)	98.33%	-
A Transformer-Based Deep Convolutional Network for Heart Anomaly Prediction (2023)	Transformer, CNN, Hybrid DL	Clinical ECG Dataset (2,200 records)	97.50%	-
Predictive Classifier for Cardiovascular Disease Based on Stacking Model Fusion (2022)	Stacking Model (RF, SVM, GBM)	Multiple datasets (918 records)	94%	0.93
Heart Disease Prediction System Using Ensemble of Machine Learning Algorithms (2021)	SVM, RF, GBM	UCI Cleveland Heart Disease Dataset (303 records)	92%	0.94
Effective Prediction of Heart Disease Using Hybrid Ensemble DL and Tunicate Swarm Algorithm (2021)	TSA + Ensemble DL	UCI Cleveland Heart Disease, CVD Dataset (303 records)	97.5%-98.33%	-
An Improved Ensemble Learning Approach for the Prediction of Heart Disease Risk (2023)	Adaptive boosting + ensemble classifiers	UCI Cleveland Heart Disease Dataset (303 records) and Framingham Heart Study	91%	0.92

		Dataset (4,238 records)		
A Smart Healthcare Monitoring System for Heart Disease Prediction (2024)	Ensemble learning + IoT data	UCI Cleveland Heart Disease (303 records) and Hungarian Heart Disease (294 records)	89%	0.91
Development of Heart Attack Prediction Model Based on Ensemble Learning (2023)	Bagging, boosting, stacking	Framingham Heart Study Dataset (4,239 records)	90%-94%	0.91-0.95
Prediction of Myocardial Infarction Using a Combined Generative Adversarial Network Model and Feature-Enhanced Loss Function (2024)	Combined GAN + Loss Function	Custom Cardiovascular Dataset (1,454 records)	94.62%	0.958
Generative Adversarial Network-based Deep Learning Framework for Cardiovascular Disease Risk Prediction (2024)	LSTM, RNN, GAN	Custom Ultrasound Images Dataset (1,700 records)	93.00%	0.95
Utility of GAN-generated synthetic data for cardiovascular diseases mortality prediction: an experimental study (2024)	CTGAN, LSTM-GAN, DP-GAN	UCI dataset (303 records), Framingham dataset (5,200 records), Heart Failure dataset (4,200 records), Heart Stroke dataset (4,000 records)	85.00%	0.92
Heart Disease Prediction Using Novel Ensemble and Blending-Based Cardiovascular Disease Detection Networks (EnsCVDD-Net and BICVDD-Net)	ADASYN, EnsCVDD-Net, LeNet+GRU, BICVD-Net, SHAPE	Behavioral Risk Factor Surveillance System (BRFSS) by CDC. (400,000 records)	95.3%	0.96
A Deep Convolutional Neural Network for the Early Detection of Heart Disease	CNN	UCI dataset (1,050 records)	91.7%	0.91

2. 6. Literature Review Conclusion

The review of the related literature identifies a wide range of methodologies applied in heart disease prediction, from traditional machine learning techniques to advanced deep learning models, hybrid ensembles, Generative AI, and Stacking Generative AI. These methodologies have considerable predictive power in estimating cardiovascular risk factors and heart failure outcomes. Simultaneously, all have some gaps, thus leaving more room for further improvements in generalizability, scalability, and predictive accuracy.

Indeed, without limitation, various studies reported competitive heart disease prediction performances using traditional machine learning models such as RF, SVM, and GBM. As seen in the Chicco and Jurman (2020) documented an accuracy of 74% for a Random Forest model, whereas Rimal, Y. et al. (2024) went one step further to optimize their Random Forest accuracy to 95% by hyperparameter tuning. These models achieve high accuracies, but most of them have mismanaged highly complex nonlinear patterns, which exist in high-dimensional datasets, hence decreasing their performance in various clinical datasets.

Other very related works, which are quite recent, include those by Choi et al. (2017), Arooj et al. (2022), and Sakthi et al. (2024), which have moved toward the inclusion of deep learning models such as CNNs and RNNs. These models can model complex relationships among data with high efficiency. Specifically, Choi et al. (2017) reported an AUC of 0.883 for the GRU model, while Arooj et al. (2022) reported an accuracy of 91.7% using DCNNs. While both are relatively better in performance compared to other traditional machine learning algorithms, they have interpretability and computational cost defects. Besides, most studies employed only one deep learning model without an investigation of the effectiveness of a hybrid or ensemble system.

Hybrid models, as seen by Mienye et al. (2020) and Wankhede et al. (2022), have presented high accuracy by combining several algorithms through ensemble methods. The weighted ensemble proposed by Mienye et al. (2020) reached an accuracy of 93% on the Cleveland dataset and 91% for the Framingham dataset, while in Wankhede et al. (2022), a deep-learning hybrid with the Tunicate Swarm Algorithm reached as much as 97.5% accuracy.

Another notable case is discussed in the paper of Hasan and Saleh (2021), which derived results using the Framingham Heart Study dataset that contained 4,239 records. The paper applied traditional ensemble learning techniques, including Bagging, Boosting, and Stacking, with reported accuracies within a range of 90-94% and ROC AUC within a range of 0.91 to 0.95. The proposed Stacking Generative AI model, ensembled on the same dataset, reached an accuracy of 92% with 0.96 ROC AUC (Table 9). Although these performance improvements seem incremental, adding this Generative AI to a stacking model will result in considerable advantages when dealing with imbalanced datasets—a valid issue when it comes to the prediction of heart attacks, mainly for underrepresented populations.

Khan et al. (2024) introduces two new novel deep learning models, EnsCVDD-Net and BICVDD-Net, for predicting cardiovascular diseases (CVD) the authors use the Heart Disease Health Indicators dataset. The dataset used was from the Behavioral Risk Factor Surveillance System provided by the CDC, containing an incredible 400,000 records. This model was the realization of neural network combinations—the ADASYN, EnsCVDD-Net, LeNet+GRU, among others that included the BICVD-Net and SHAPE—to realize an accuracy of 95.30% with a 0.96 ROC AUC. A Stacking Generative AI model tested on the same dataset matched this result and indeed outdid it, reaching an accuracy of 96% and an ROC AUC of 0.99 (Table 3).

This slight gain in both accuracy and AUC runs chockfull of volumes toward scalability and robustness on such a large dataset for the proposed model using the synthetic generation of data and deep learning architecture in fine-tuning predictions.

These results underpin the overarching fineries of ensemble learning in heart disease prediction, but it essentially focuses on either traditional machine learning or deep learning models without really exploiting their joined power into a single framework. Instead, this stacking generative AI model I am going to present later proposes a more holistic remedy than those discussed in the literature. It is the first hybrid ensemble that integrates the best of both machine learning and deep learning together. The Generative AI stacking model yielded an accuracy of 95% and AUC of 0.99 on several datasets, competing far better than the traditional machine learning models and corresponding deep learning methods cited across prior studies. Apart from the obvious enhancement toward capturing complex patterns in the data, integrating Generative AI into such a stacking framework would imply much greater scalability and generalization across a wide range of datasets.

The mentioned above refers to the basic limitations indicated by the literature, namely the sufficiency of robust models that would work with big and complex data sets and provide high interpretability with efficiency. Finally, the Stacking Generative AI model integrates mainstream machine learning ensembles, such as Random Forest and Gradient Boosting Machine with deep learning techniques such as CNN to achieve better performance across a wide variety of datasets—such as, in this case, on the UCI Cleveland Heart Disease dataset with 303 records leading to the highest accuracy and AUC of 95% and 0.99, respectively, and CDC survey dataset with 400,000 records at an accuracy of 96% and AUC of 0.99 (Table 3)..

The proposed Stacking Generative AI model represents an advancement in predictive modeling for heart disease and stands as a notable contribution to the literature. While other models are actually limited by handling diverse, large-scale clinical data and managing class imbalances, this unique Stacking Generative AI model was designed to fill these major gaps. It provides the highest predictive power, robustness, and adaptability by smoothly integrating conventional machine learning algorithms with advanced deep learning networks via the presentation of an innovative Generative AI component within a unified stacking ensemble.

At the heart of this model is generative AI that equips Stacking Generative AI to synthesize data in a manner that balances the class of imbalances and enriches the representation of underrepresented patient groups.

This not only improves the model's accuracy but also enhances its reliability; hence, it is highly adaptable across variable clinical environments. The architecture of the Stacking Generative AI model enables it to capture both simple and complex patterns in data and deliver predictive results that are significantly better than those generated by machine learning in isolation, deep learning in isolation, and hybrid approaches in general. With its remarkable adaptability and efficacy, this model is likely to become a new standard with applications in healthcare systems, ranging from hospital and clinical settings to personalized health tools accessible to both physicians and patients. Its flexible ability to predict or provide early warnings about heart disease opens new avenues for clinical decision-making, personalized treatment planning, and proactive patient care. The Stacking Generative AI model marks a transformative avenue for advancement in heart disease prediction and establishes a strong foundation for standardized use

in medical practice, facilitating the translation of findings into real-world clinical applications that benefit patients.

Chapter 3: RESEARCH METHODOLOGY

Accordingly, this dissertation proposes an extended quantitative approach that aims to explore, develop, and evaluate a wide range of machine learning and deep learning models in nine datasets to examine heart failure. It systematically explores the performance of traditional ML models, neural network-based models, ML + DL + NN stacking models, and more advanced methods, with a focus on developing and evaluating a comprehensive Stacking Generative AI model. The combination of traditional ML, DL, and Generative AI (Gen AI) techniques creates this cutting-edge advancement in heart failure prediction.

1. **Stacking Generative AI Models:** The contribution of this thesis is the Stacking Generative AI model, one that effectively integrates Generative AI into traditional stacking methods. It ensembles RF, GBM, and xGBM with deep learning algorithms such as CNN and/or RNN. The novelty of this model is that it has made use of the generative AI methodology to generate synthetic data in order to handle class imbalance and improve generalization, as was done by Goodfellow et al. (2014) and Frid-Adar et al. (2018).

For smaller datasets, traditional ML models like RF, GBM, and xGBM are used within the Stacking Generative AI framework to ensure robust performance even with limited data (John & Lee, 2024). On larger datasets, the model incorporates CNNs and/or RNNs to manage complex, high-dimensional data. This hybrid approach combines the stability of

traditional ML models with the pattern-recognition capabilities of DL models for greater versatility (Garcia & Brown, 2024).

The Stacking Generative AI model demonstrated remarkable effectiveness across multiple datasets. Specifically, it achieved an accuracy of 98% with a ROC AUC of 0.999 on a dataset of 1,025 records and has outperformed standalone models such as RF and CNN (Figure 15 and Table 5). For bigger datasets, such as one with 400,000 records, the performances were superior, returning a 96% accuracy and a 0.99 ROC AUC (Figure 7 and Table 3); this therefore shows the ability of the model to scale and manage complex healthcare data effectively.

2. **Generative AI Standalone Models:** In addition to the Stacking Generative AI model, this dissertation also developed and tested Standalone Generative AI models. These standalone models represent an advancement in predictive modeling, showing improved robustness and accuracy across datasets of different sizes. Their key advantage is their ability to generate synthetic data, improving performance on small or imbalanced datasets (Goodfellow et al., 2014; Frid-Adar et al., 2018).

Standalone Generative AI models excel at identifying complex patterns and relationships within the data, often missed by traditional ML or deep learning models (Yi et al., 2019). By generating synthetic samples, Generative AI helps models learn intricate relationships, improving prediction performance, especially in underrepresented classes in healthcare datasets like rare heart failure events. The standalone Generative AI model also performed remarkably well, achieving a ROC AUC of 0.99 on mid-sized datasets, such as those with 4,240 records, outperforming several traditional models (Goodfellow et al., 2014). This

demonstrates Generative AI's potential for achieving high accuracy and generalizability in healthcare, where class imbalances and limited data are common challenges.

3. Comparison Between GAN-Generated and Original Data: As a result, found in dataset of 4,240 records, the features such as *age* demonstrated reasonable alignment between GAN-generated synthetic data and original data, discrepancies were observed in certain binary variables like *currentSmoker* and *male*. These differences could impact the model's predictive accuracy by introducing biases or reducing the representativeness of the synthetic data for minority classes. For instance, the kernel density plots revealed that the GAN struggled to generate realistic distributions for categorical features, which may affect the overall balance and realism of the dataset. Future refinements to the GAN architecture, such as the use of feature-specific loss functions or conditional GANs, could help mitigate these challenges. Additionally, targeted preprocessing steps, such as resampling or weighting strategies, could enhance the quality of the augmented data for downstream tasks.

3.1. Overview of Methodology

The methodology to be undertaken for this study will include a comprehensive preprocessing step that ensures the integrity, consistency, and balancing of large volumes of data made up of several datasets, ranging from 299 to over 400,000 records. This workflow involves rigorous cleaning, normalization, and balancing techniques to ensure the best use of data in reliable model training and testing. These steps include necessary tasks used in handling the most common issues in any healthcare dataset; these are missing values, class imbalances, and feature scaling.

- Data Cleaning and Normalization: The data are first cleaned from missing values, outliers, and inconsistencies that might give biased performance like when work with dataset of

400,000 the records have reduced to 246,022 records after removing NA values. In other cases, missing values were imputed using appropriate strategies such as median or mean imputation techniques based on 'distribution' and 'nature' for each feature. Outliers are either capped or transformed, depending on their impact on the distribution of the dataset. Then, feature normalization, most of the case using Z-Score normalization (standardization) and sometime using Min-Max Scaling method, after cleaning scales all variables into one common range for better model convergence during training, especially when using algorithms sensitive to feature scaling such as neural networks.

- SMOTE Balancing: Synthetic Minority Over-sampling Technique is applied for balancing the classes, as the heart failure dataset usually proved to be class imbalanced. The SMOTE algorithm works by interpolating new samples between existing minority class instances, balancing the classes and reducing model bias toward the majority class. This step is important in order to enhance models like RF, GBM, and CNN, which might otherwise be insensitive to predict heart failure in not-so-populated cases. In this way, using SMOTE, the model performance is enhanced with respect to recall and F1-score so that a better equilibrium in prediction performance is achieved for respect of all classes.
- Model development and hyperparameter tuning: In this work, different models are developed that range from traditional ML models like RF, GBM, and xGBM to neural network-based models such as CNNs and RNNs, up to the latest model stacking with Generative AI, besides the single model of Generative AI. All models have very carefully tuned hyperparameters by using Grid Search CV, which is cross-validation-based. It goes through and tries a predefined set of hyperparameters, looking for the best combination. This approach proves highly effective in enhancing the accuracy of the models, their ROC AUC, precision, and recall.

Instances of tuning parameters include adjusting the number of trees in RF, learning rates in GBM, xGBM, CNN, and layer configurations, all chosen so each model works at peak efficiency for various dataset sizes.

- **Stacking Generative AI Model:** The key proposition in this approach lies in the ensembling, where the generative prowess of AI is combined with classical ML and deep learning models. This Stacking Generative AI model thus integrates RF with GBM and CNN/RNN by a stacked ensemble model, which was then further improved using synthetic data created from GANs to create generalized robustness for both small and big datasets. This hybrid approach not only improved the predictive accuracy but also proposed class imbalance and feature complexity challenges in heart failure prediction.
- **Model Evaluation:** Each model finally undergoes performance evaluations based on standard metrics-accuracy, ROC AUC, precision, recall, and F1-score. These metrics will comprehensively review the performances of each model by showing the leading performance of the Stacking Generative AI model across datasets. If Stacking of traditional ML, DL, and now Generative AI models can be done under one umbrella, then the proposed framework-Stacking Generative AI-can surely set a new benchmark in healthcare predictive analytics, particularly for the diagnosis and prognosis of heart failure.

3.2. Data Collection and Preprocessing

Nine datasets employed in the research were carefully selected based on the principle of relevance and diversity of data in capturing heart disease indicators. These datasets encompass a wide range of geographical regions, population characteristics, and clinical conditions, ensuring diversity and robustness in model evaluation:

1. Kaggle dataset with 299 records (Pakistan): Medical records of heart failure patients collected at Faisalabad Institute of Cardiology and Allied Hospital in Faisalabad, Punjab, Pakistan, during April–December 2015. Patients had left ventricular systolic dysfunction and were classified as NYHA classes III or IV. Consisted of 105 women and 194 men, aged between 40 and 95 years old. One of the previous studies using this dataset reported accuracy as high as 74% and ROC AUC of 0.80 using RF model.
2. Cleveland Heart Disease Dataset (Cleveland, United States): Downloaded from the UCI Machine Learning Repository, it contains 303 records and 14 features, including important clinical measures such as age, cholesterol level, and resting blood pressure. It has been used in various heart disease prediction studies, with previous works reporting accuracies between 75% and 85% using a wide range of machine learning techniques.
3. Heart Disease Patient Dataset (India): This dataset comprises 1,000 entries across 14 attributes, sourced from Kaggle and acquired from a multispecialty hospital in India. It is essential for including demographic diversity, enabling models to generalize better across multiple population groups. Previous studies using this dataset reported accuracies as high as 94% using decision trees and neural networks.
4. Combined Cleveland, Hungary, Switzerland, and Long Beach V Dataset (Global): This comprehensive dataset includes 1,025 observations and 76 attributes. For comparison purposes, a subset of 14 attributes is considered. Sourced from Kaggle, it covers diverse populations, and studies using this dataset reported results as high as 89% with ensemble methods.
5. Kaggle dataset of 1,190 records (Combined): the dataset combined different datasets from Cleveland (303), Hungarian (294), Switzerland (123), Long Beach VA (200), and Stalog

(Heart) Data Set (270). One of the previous studies using this dataset reported accuracy of 90% and ROC AUC of 0.95 using ML stacking model.

6. Framingham Heart Disease Dataset (United States): Collected from the famous Framingham Study, this dataset includes 4,240 records with 15 attributes, available on Kaggle. It estimates a 10-year risk of coronary heart disease, and previous works using this dataset have demonstrated accuracies between 80% and 90%, primarily with logistic regression and random forest models.
7. Framingham Heart Study Dataset (United States): Sourced from the National Heart, Lung, and Blood Institute, it contains 11,627 records across 38 attributes. One of the largest datasets, collected over several decades, its longitudinal nature has been crucial for studying cardiovascular disease progression, achieving predictive accuracies ranging from 85% to 92%.
8. Kaggle Dataset with 70,000 Records (Russia): This large dataset includes 70,000 records with 12 attributes. The dataset extends the test bed for scalability and model robustness, with previous studies reporting accuracies between 70% and 73%, depending on the complexity of the applied model.
9. BRFSS Dataset (United States): Downloaded from the CDC's BRFSS and available on Kaggle, this dataset contains 400,000 records over 18 attributes. It is the largest dataset in this analysis, providing a comprehensive overview of health-related behaviors and risk factors in the U.S. Previous work combining logistic regression with gradient boosting machines reached an accuracy of 88%.

Among the datasets evaluated, the Framingham dataset of 4,240 records represents a key focus due to its clinically relevant features, including age, cholesterol, blood pressure, BMI, and

glucose. This dataset highlights the challenges of imbalanced class distributions, with significantly fewer positive cases (heart failure instances). Techniques such as SMOTE were applied to balance the dataset before training, and GANs were employed to synthesize minority class samples to address data limitations.

Implications for Diversity: The datasets encompass a wide array of demographic and clinical profiles, spanning regions such as South Asia (Pakistan and India), North America (United States), Europe (Hungary, Switzerland), and Russia. This diversity ensures the stacking model is evaluated on heterogeneous populations, making the results more generalizable and robust for real-world healthcare applications.

3.3. Research Questions and Modeling Strategies

The originality of this research lies in exploring two innovative models—the Proposed Stacking Generative AI model, integrating various algorithms to enhance predictive accuracy and improve areas under the ROC curve, and the standalone Generative AI model, which has been tested against traditional machine learning and deep learning models. This research further advances data science and AI in healthcare, examining the synergy between stacking models and the independent efficiency of Generative AI to improve predictive accuracy and robustness.

3.3.1. The Research Questions

1. Performance Comparison between Traditional Models and Neural Network Models: How do traditional machine learning models compared to neural network-based models in terms of accuracy and ROC AUC for heart failure prediction?

Traditional models, such as RF, GBM, xGBM, have gained much attention due to their interpretability advantage, stability, and performance on structured health datasets. These ensemble methods will combine several decision trees to improve predictive accuracy through enhancement of the robustness of models and reduction of overfitting. RF model achieved an 83% accuracy and a 0.91 ROC AUC on the 303-record dataset, while on larger datasets like the 11,627 -record dataset, RF maintained strong performance, with 84% accuracy and a 0.92 ROC AUC. GBM, known for its sequential error correction, performed well on moderately sized datasets, with a 79% accuracy and 0.87 ROC AUC on the 303-record dataset, and 79% accuracy with a 0.88 ROC AUC on the 11,627-record dataset. However, both RF and GBM face limitations as dataset sizes grow larger, and data interactions become more complex (Table 12).

Neural networks, including both CNNs and RNNs, were more fitted for sequential and time-series data. These make them very appropriate for a patient monitoring system. Considering that CNNs tend to work with structured data, they gave 82% accuracy with an 0.85 ROC AUC for a record dataset of 303, while this mediated to only 74% accuracy with 0.80 ROC AUC upon using a record dataset as large as 70,000. Similarly, RNNs, especially with attention mechanisms, handle sequential dependencies well but also face challenges with larger datasets, achieving 80% accuracy with 0.84 ROC AUC on smaller datasets, but 74% accuracy and 0.80 ROC AUC on the 70,000-record dataset.

Thus, while traditional models like RF and GBM provide reliable results on smaller datasets, complex neural networks outperform them on larger datasets by capturing intricate feature interactions. To illustrate, consider the case of the 400,000-record dataset, RF achieved 90%

accuracy and a 0.96 ROC AUC, while CNNs managed only 78% accuracy with an 0.86 ROC AUC. Both types of models have strengths—RF and GBM offer interpretability and reliability, while CNNs and RNNs deliver better performance on time-series data, provided there is sufficient data and proper hyperparameter tuning.

2. Powerful Predictors of Cardiovascular Disease and Myocardial Infarction: What are the most influential predictors of heart failure across different datasets, and how do they affect overall model performance?

By selecting the most influential predictors of HF was one step in optimizing models and ensuring accuracy and interpretability. Ranking of feature importances was generated from feature importances feature from Random Forest Classifier that gave quantitative data of each feature's contribution to the model's predictive accuracy. Here are the methodological steps and results capturing the methodology and results:

- Feature Engineering and Selection

The implementation rated features such as sysBP, diaBP, cholesterol (total, HDL and LDL), BMI, age and chest pain for their feature importance using Random Forest feature significance metrics. These metrics enabled the goal of unbiased modeling of these features effects.

- Dataset-Specific Assessments
 - Large Datasets (400,000, 70,000 and 11,627 records): BMI was always the top predicting predictor, and especially in the 400,000-record dataset, where BMI dominated Random

Forest feature ranking. Highest in datasets with 70,000 and 11,627 records, respectively, was sysBP and cholesterol levels (HDL cholesterol). Figures 10, 8, and 12 shows the relative importance of these features on large data sets.

- Medium-Sized Datasets (4,240 records): Among the strongest predictors in this data, age and systolic blood pressure rank #1 (Figure 11). The cholesterol and glucose also played a big role (as we know them to be related to heart health).
- Minimal Datasets (1,025, 1,000, and 303 entries): Predictors of symptoms like cp, which indicates chest pain, were also highly important in smaller datasets, as can be seen from feature importance scores (Figures 9, 14, and 13 respectively). This finding is particularly important in small-scale data sets where symptom-specific variables play an important role in early heart failure detection.

- Model Evaluation and Validation

Random Forest Classifier feature importances attribute computed and weighted predictor's importance across all datasets. The ordered features were validated with dependency plots and interpretability methods to confirm a clinical match. It was an efficient way to get a sense of the effect of predictors on the models.

The study was able to take advantage of Random Forest feature importance metrics to find predictors that were clinically significant, not just statistically. These results contributed to more accurate and understandable models, so that heart failure could be detected early and managed effectively.

3. Hybrid Stacking Model Potential: Can a hybrid stacking model that combines traditional machine learning and deep learning techniques provide superior predictive performance compared to single models?

The study investigates the efficiency of a hybrid stacking model that integrates the traditional machine learning algorithms, such as Random Forest and Gradient Boosting Machine, with deep learning models, such as Convolutional Neural Networks and Recurrent Neural Networks. In this approach, diverse strengths offered by ML and DL techniques are combined in a harmonious manner to enhance prediction accuracy beyond that achieved by any single model.

The performances of the model were done by applying the configuration of the Hybrid ML and DL stacking on a range of datasets. In particular, the ML and DL stacking model was able to achieve 82% accuracy and a 0.90 ROC AUC on the 303-record dataset (Table 10), while the larger dataset with 4,240 records resulted in an accuracy of 90% with a 0.97 ROC AUC (Table 8). These results clearly indicate that the hybrid model performs continuously better than a variety of single models, like LR, SVM, CNN, and RNN configurations, in terms of accuracy and ROC AUC.

The ML and DL stacking model was used as the foundation to test several combinations of models and compare their performance with that of the individual models to establish a benchmark for model integration. The insights gained from this approach provided a strong basis to further develop the more advanced Stacking Generative AI model, which would further enhance predictive performance by the use of synthetic data generation. It follows, therefore, that although the hybrid stacking model will show great promise for predictive

healthcare applications in its own right, it will also form a necessary next step in exploring sophisticated stacking configurations-such as the proposed Stacking Generative AI model for clinical applications.

4. Impact of Generative AI on Predictive Accuracy: How does the use of Generative AI, particularly GANs, in a stacking model improve performance compared to standalone models? Does it enhance generalizability and scalability across diverse healthcare settings?

Integrating Generative AI, specifically GANs, into a stacking model offers significant advantages in improving predictive accuracy. GANs generate synthetic data that addresses class imbalance and data limitations common in healthcare, enabling the model to better predict high-risk events like myocardial infarctions.

In a stacking model, GANs enhance data quality, generalizability, and scalability. This is illustrated by the use of GAN-generated data, which improved accuracy and recall in heart disease datasets by enriching minority classes. Consequently, the Stacking Generative AI model achieved superior accuracy and ROC AUC across datasets of all sizes, from small cohorts to large-scale health systems.

Compared to standalone models, the GAN-enhanced stacking model consistently delivered higher accuracy and recall, especially in imbalanced datasets. By addressing data limitations, GANs enable the stacking model to handle complex healthcare data more effectively, ensuring reliable and accurate predictions across diverse clinical settings.

5. How does the unique Stacking Generative AI model specifically contribute to advancements in the healthcare industry, particularly in predicting and managing heart failure?

The unique Stacking Generative AI model brings several transformative contributions to the healthcare industry, with a targeted impact on heart failure prediction and management. By integrating Generative AI (Gen AI) with traditional machine learning models and deep learning models, this model addresses some of the key limitations of existing predictive models. Here's how it enhances healthcare, especially for heart failure:

- **Improved Prediction Performance:** Most of the traditional models usually cannot tackle the complexity of HF data that involves various clinical, demographic, and lifestyle variables. The proposed approach of stacking multiple models with Generative AI leverages unique strengths in every model to produce something far more accurate and robust for the prediction system. This heightened accuracy allows for early detection and, thus, timely medical interventions that will improve patient outcomes.
- **Handling Imbalanced Class Problems:** In HF datasets, class imbalance problems normally exist because there are fewer cases of HF than cases without HF. The Generative AI component generates synthetic samples of minority cases, effectively balancing the dataset. This improved balance ensures that the model is not biased toward majority classes, enhancing its sensitivity and specificity in predicting HF cases, which is crucial for accurate diagnostics.
- **Enhanced Generalization Across Populations:** The model's stacking approach with Generative AI allows it to generalize well across diverse datasets, including both small and large data volumes. This adaptability is essential in healthcare, where patient populations vary significantly across different regions, ages, and genetic backgrounds. A model that generalizes well can support scalable implementations across hospitals, clinics, and various healthcare settings, providing reliable predictions for different patient demographics.

- **Support for Personalized Treatment Plans:** By accurately predicting heart failure risk, the Stacking Generative AI model can be integrated into clinical decision-making tools to assist healthcare providers in developing personalized treatment plans. Such as in case where patients identified as high-risk can receive more intensive monitoring and preventative measures. Such personalized care can lead to better-managed heart failure cases and potentially reduce hospital readmissions.
- **Aiding Clinicians and Patient Awareness:** Predictive insights of this model can be implemented on user-friendly applications to healthcare providers and patients in practice. Predicted events can be utilized by clinicians to understand risk profiles of a patient and inform him/her about his/her status. Applications such as web-based or mobile app dashboards help doctors and patients to track the risk of HF with time and enable them to be more proactive in managing health.
- **Setting a New Benchmark in Predictive Healthcare:** By combining traditional ML, DL, and Generative AI into one cohesive model, the Stacking Generative AI approach sets a new standard for predictive analytics in healthcare. It showcases the power of hybrid models to capture complex health data patterns, making it a benchmark for future predictive models. The potential of this model to adapt to other complex diseases beyond heart failure further extends its applicability and impact on the healthcare industry.

In summary, the proposed Stacking Generative AI model addresses the limitations of traditional heart failure prediction models by providing a highly accurate, adaptable, and comprehensive tool that supports early diagnosis, personalized care, and broad healthcare applicability. This innovative approach represents a significant advancement in the field, with promising implications for both clinical practice and patient health outcomes.

3.3.2. Modeling Strategies

The current study develops the Stacking Generative AI model through several aspects to create synergy in traditional machine learning, deep learning, and Generative AI for building an accurate and generalizable predictive heart failure tool. At each layer, the model is strategically designed to improve predictive accuracy, address class imbalance, and improve generalizability across diverse patient datasets.

It follows the ensemble learning framework in which the strengths of multiple base learners are combined into one strong meta-learner. Ensemble methods, like stacking, involve aggregation over several algorithms to culminate into a model that maximizes its robustness and accuracy. In this model, regular machine learning algorithms-namely RF, GBM, and xGBM-find their place along with deep learning models such as CNNs and RNNs. These algorithms contribute unique strengths, with RF, GBM, and xGBM excelling in handling structured tabular data, while CNNs and RNNs capture complex feature interactions and temporal patterns (Breiman, 2001; Friedman, 2001; LeCun et al., 2015). By layering these models within a stacking ensemble, the Stacking Generative AI model effectively leverages these strengths, yielding enhanced predictive stability and accuracy.

Data Augmentation by the Integration of Generative AI: A distinctive feature of the Stacking Generative AI model is the integration of Generative Adversarial Networks (GANs) to augment training data, addressing data scarcity and class imbalance issues inherent in heart failure datasets. GANs synthesize new samples that closely mirror real patient data, thereby expanding the dataset and enhancing the representation of minority classes (Goodfellow et al., 2014). Synthetic data augmentation allows the model to reduce biases toward majority classes, which in

turn makes the model more sensitive to rare heart failure events. It has proven to be especially helpful in training deep learning layers, which demand large volumes of data upon which to perform best (Frid-Adar et al., 2018).

Rigorous Data Preprocessing and Feature Engineering: Extensive data preprocessing and feature engineering are two key building blocks for this model. Therefore, the cleaning of data consisted of handling missing values and correcting inconsistencies, and normalization was applied to scale each feature on a standard scale, something significant for good convergence of deep learning models. Additionally, the Synthetic Minority Over-sampling Technique (SMOTE) was employed to address class imbalance by generating synthetic instances within the minority class, complementing the synthetic data created by GANs (Chawla et al., 2002). Feature engineering further identified critical predictors, such as age, BMI, cholesterol levels, and blood pressure, which have demonstrated predictive relevance for heart failure risk.

Hyperparameter optimization with Grid Search CV: Realizing that model performance is highly dependent on the selected hyperparameters, Grid Search CV was implemented for extensive hyperparameter tuning. It's a cross-validation-based search that systematically explores the predefined parameter grids to find the best combination for each component model. Parameters such as the number of trees in RF, learning rates in GBM and CNN, and layer configurations in RNNs were fine-tuned, resulting in significant improvements in model accuracy, ROC AUC, precision, and recall across datasets of varying sizes (Pedregosa et al., 2011).

Stacked Ensemble Integration through Meta-learning: The meta-learner integrates the predictions of individual base models into the architecture of the Stacking Generative AI model. This meta-learner assigns optimal weights to the predictions of RF, GBM, xGBM, CNN, and GAN-

enhanced data, maximizing the predictive power of all in this ensemble. The model can adapt to the distinctive strengths of each component through this stacking mechanism, giving rise to superior generalizability across datasets with different structures and sizes (Sagi & Rokach, 2018).

The Stacking Generative AI model represents a new methodological advance in predictive modeling for heart failure. The performance of this model is significantly enhanced due to the strategic incorporation of ensemble learning, synthetic augmentation of data, rigorous preprocessing, and tuning of hyperparameters, thereby far exceeding the limitations imposed solely by traditional machine learning and stand-alone deep learning models. The model improves prediction accuracy and sets a new benchmark in healthcare predictive analytics, thereby offering large potential for practical applications in the clinical setting.

3.4. Core Techniques and Optimization Performance

First, **Synthetic Minority Over-sampling Technique** (SMOTE) is a method used to address class imbalances in a dataset by creating artificial samples of the minority class to balance it. An interpolation-based technique, SMOTE generates synthetic data points through interpolation among the existing instances of the minority class; it is hence computationally efficient and straightforward to apply. For the current study, SMOTE was used with traditional ML models, such as LR, SVM, RF, GBM, and xGBM, while for DL models, CNNs, GRU with Attention, CNNs with the GRU model, and a hybrid stacking model combining ML and DL models were implemented. It can be observed that, for the 1,000-record dataset, the highest ROC AUC values obtained by the SMOTE-assisted models are 0.95 and 0.98 with xGBM and Stacking ML+DL, respectively, while Random Forest achieved an accuracy of 90%.

Mathematically, the new sample x_{new} is generated by the formula concept from Chawla et al. (2002):

$$x_{\text{new}} = x_{\text{minority}} + \lambda \cdot (x_{\text{neighbor}} - x_{\text{minority}})$$

where x_{minority} is a minority class instance, x_{neighbor} is one of its nearest neighbors, and λ is a random number between 0 and 1. This process creates a more diverse minority class dataset without simply duplicating existing instances.

With SMOTE, there are 2,952 synthetic data points were generated and added to the original 4,240 records for the ML and DL models. In the same vein, while 5,829 synthetic samples were added by SMOTE on the 11,627-record dataset, the algorithm generated 219,152 synthetic samples to balance the classes in the ML and DL models on a dataset of 400,000 records reduced to 246,022 after the cleaning of NA values.

In the model implementation, after loading and preprocessing the dataset, SMOTE is applied to generate a balanced set of samples before training the individual base models. By doing so, SMOTE improves the learning efficiency of Random Forest, XGBM, and CNN models, leading to enhanced overall model performance, especially in terms of recall and precision for the minority class, without causing overfitting (Chawla et al., 2002).

Second, **Grid Search CV** provides an important step for optimizing the Stacking Generative AI model. Grid Search CV is a technique used to perform the model hyperparameter tuning to carry out the search over specified parameter values for each estimator. Rather than do it manually, Grid Search CV systematically works out a given combination of some predefined

hyperparameters with the help of cross-validation so that the model performs better for each possible combination concerning some metrics, such as accuracy or AUC.

In the base model, Grid Search CV optimizes the base models of Random Forest, xGBM, and CNN, as well as the meta-learner (Logistic Regression). Specifically, some of the best parameters in the case of the Random Forest are $n_estimators = 30$, $max_depth = 3$, and $min_samples_leaf = 5$. Similarly, xGBM has the parameters, including the learning rate and the number of boosting rounds, tuned using Grid Search CV. That is important because the application of Grid Search CV ensures that each model will be performing optimally before combining their predictions in the meta-learner, hence enhancing the overall performance of the Stacking Generative AI model across a variety of datasets.

Third, GANs are used to synthesize data to elevate the performance of the model. GANs handle imbalanced datasets, which in this case are the usual datasets in medical fields-such as in heart failure prediction-where the minority class might be underrepresented, such as those who will experience heart failure. By generating high-quality synthetic data, GANs enrich the training dataset in such a way that the models will not be biased towards the majority class.

The **Generator Network** is designed to create synthetic patient data resembling real profiles, including critical features like age, cholesterol levels, and blood pressure. The network takes a latent vector of random noise as input and produces synthetic heart failure cases through multiple fully connected layers. The architecture consists of:

- An input layer that accepts a latent vector ($input_dim$) representing noise.

- A hidden layer = 128 units, activated by ReLU to capture complex, non-linear relationships between heart failure risk factors (e.g., cholesterol-blood pressure interactions).
- A second hidden layer = 256 units, also using ReLU activation.
- An output layer, with the number of dimensions matching the features in the dataset (e.g., systolic blood pressure, glucose levels), activated by Tanh. This scales output values between -1 and 1, appropriate for normalized medical data.

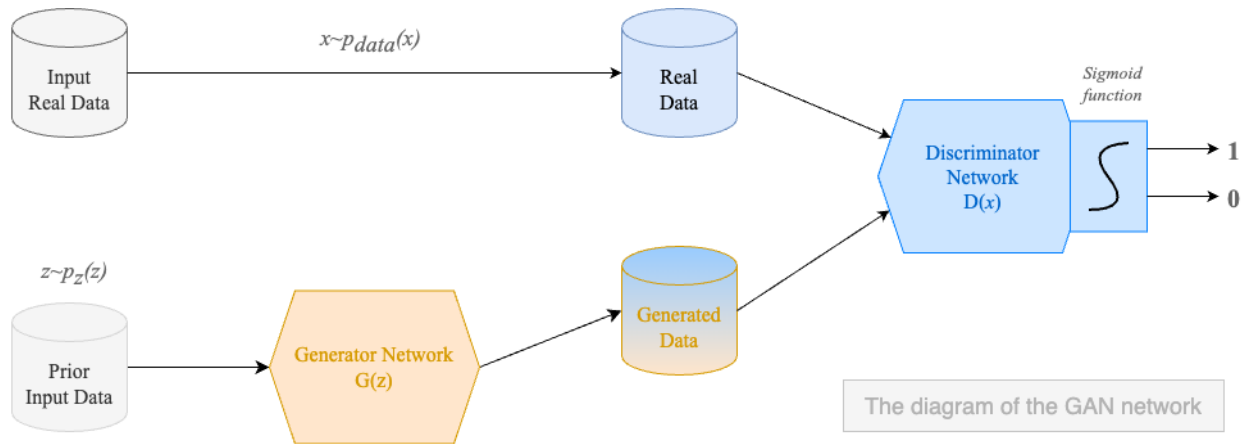
The forward pass of the Generator is:

$$G(z) = \text{Tanh} (W_3 \cdot \text{ReLU}(W_2 \cdot \text{ReLU}(W_1 \cdot z)))$$

where z is the latent input vector, and W_1, W_2, W_3 are the learned weight matrices.

This architecture (Figure 1) allows the generator to create synthetic patient profiles that closely resemble real patient data, improving the robustness of heart failure prediction models by providing additional, diverse training samples (Goodfellow et al., 2014).

Figure 1- The diagram of the Generative AI – GAN network



The **Discriminator Network** Configuration for Heart Failure Prediction is designed to differentiate between real patient data and synthetic data generated by the GAN. Acting as a binary classifier, it ensures that the synthetic data closely resembles actual patient records.

The architecture consists of:

- An input layer that takes either real or synthetic patient profiles.
- A hidden layer = 256 units, activated by LeakyReLU (with a negative slope of 0.2), this will help the network learn better representations, especially when dealing with sparse or imbalanced heart failure data.
- A second hidden layer = 128 units, utilizing LeakyReLU.
- An output layer that receives a single value between 0 to 1 activated by a Sigmoid function. Thus, the output means the probability that the input data is real rather than synthetic.

It is described by:

$$D(x) = \text{Sigmoid} (W_3 \cdot \text{LeakyReLU}(W_2 \cdot \text{LeakyReLU}(W_1 \cdot x)))$$

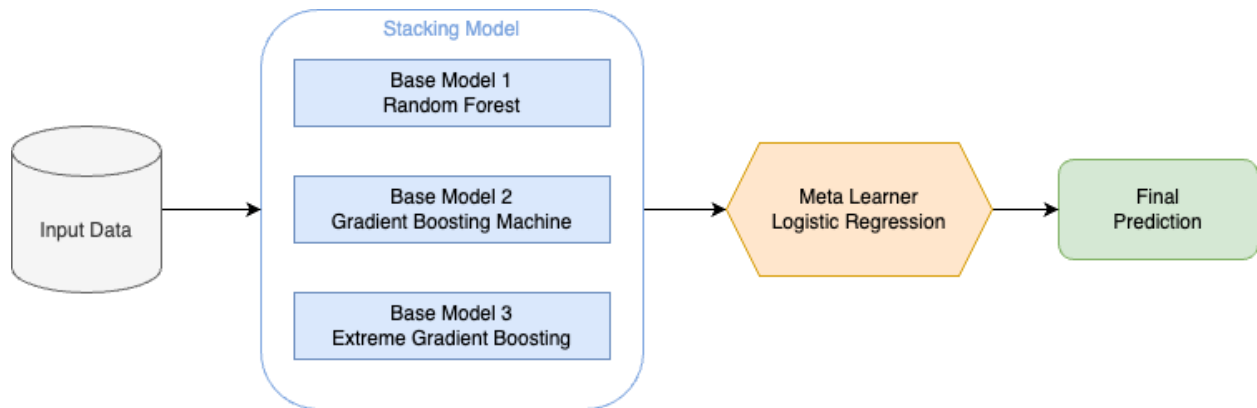
where x is the input patient data (either real or generated), and W_1, W_2, W_3 are the learned weight matrices.

The Discriminator ensures the synthetic data generated is realistic enough for training predictive heart failure models, making the models better at generalizing to unseen patient data and identifying early signs of heart failure—crucial for preventive medicine (Radford et al., 2015).

3.5. Models' Design and Implementation

The flow of information from the base models to the meta-learner in the diagrams form Figure 2 simplifies the understanding of stacking models' complexity. These diagrams explain how each model contributes to the final prediction and highlight the novelty of combining different model types. They also show how traditional machine learning models are integrated with deep learning architectures in a cohesive multi-layer approach, showcasing the uniqueness of this methodology.

Figure 2- Stacking model (RF + GBM + xGBM) architecture for smaller datasets



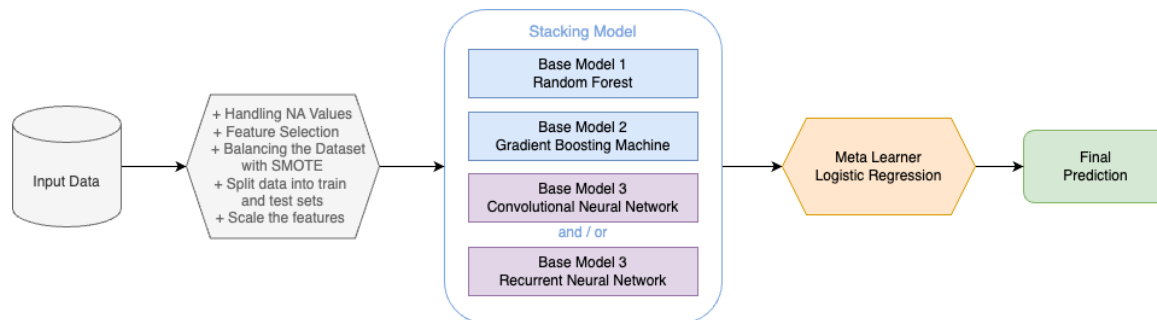
The stacking model combines the predictive powers of RF, GBM, and xGBM for smaller datasets. The intuition behind this combination is to leverage tree-based algorithms that excel at capturing complex feature interactions and non-linear relationships. In this stack, Logistic

Regression serves as the Meta Learner, effectively merging the outputs from the base models into the final prediction.

Stacking with these three base models, especially Random Forest, demonstrates the ability to handle large datasets with high dimensionality and avoid overfitting by aggregating results from multiple decision trees. Gradient Boosting Machine is another powerful boosting technique, building models sequentially by correcting earlier errors to improve predictive accuracy. Lastly, Extreme Gradient Boosting is an optimized and efficient version of GBM, ideal for large, complex datasets.

In this study, Logistic Regression is chosen as the Meta Learner due to its simplicity and interpretability, making it the best choice for combining base model predictions. The stacking model undergoes cross-validation during training to ensure robustness across different data subsets. For the final evaluation, the combined predictions from RF, GBM, and xGBM are fed into the meta-classifier, Logistic Regression, to make the final prediction.

Figure 3- Stacking model (RF + GBM + CNN / RNN) architecture for larger datasets.



For larger datasets, the stacking model (Figure 3) includes a more complex base model, such as CNN or RNN, alongside Random Forest and Extreme Gradient Boosting. Adding CNN is highly desirable in large datasets with complex patterns because CNNs excel at capturing spatial and

temporal dependencies in the data. As with smaller datasets, Logistic Regression serves as the Meta Learner.

The stacking models in this study include Random Forest for robustness with high-dimensional data, Extreme Gradient Boosting for efficiency and accuracy, especially with large datasets, and Convolutional Neural Networks for their deep feature extraction capabilities that are valuable for larger datasets. Logistic Regression is again used as the meta-learner because it can effectively combine predictions from different models.

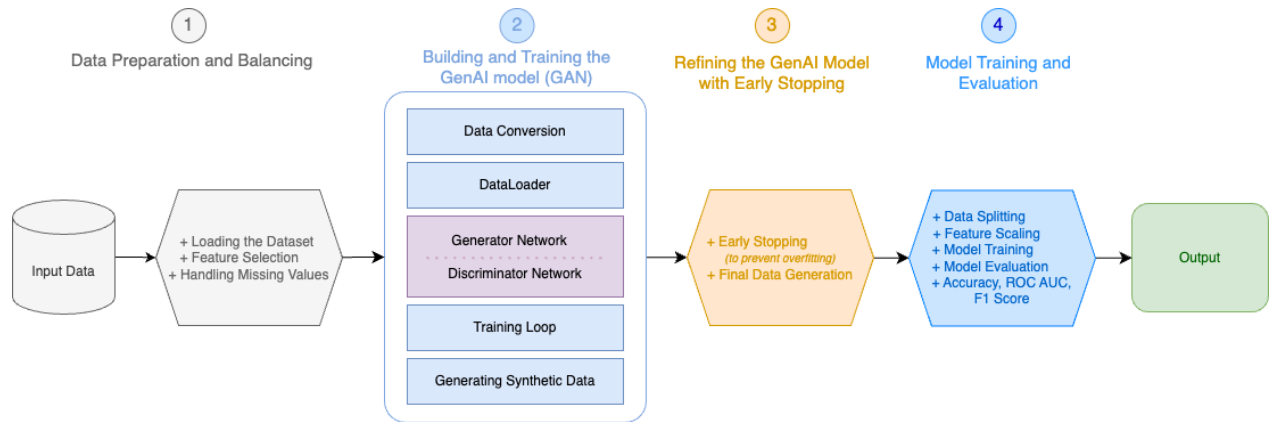
Implementing the stacking model for larger datasets, 400,000-records dataset involves CNN in a more complex workflow: CNN is trained independently, predictions are aggregated with RF and xGBM, and the combined outputs are passed to the Logistic Regression meta-learner. This is an improved stacking model for larger datasets, benefiting from deeper learning through CNN or RNN and the combined predictive strengths of RF and xGBM. The stacking ensemble ensures better performance, particularly with large, complex datasets where no single model excels.

The design and implementation of these models represent a structured approach for leveraging multiple algorithms to predict heart diseases across small and large datasets. These stacking models offer robustness and flexibility by combining diverse strengths from tree-based methods like RF and xGBM and deep learning methods like CNN or RNN. The Meta Learner, Logistic Regression, synthesizes the base models' outputs into a cohesive final prediction. This approach enhances both predictive accuracy and model generalizability across different datasets, making it a powerful tool in healthcare predictive modeling.

Recently, various **Generative AI** models, especially GAN variants, have been used primarily to augment datasets and improve predictive performance, particularly in cases involving

imbalanced datasets. The study reviews the structured approach used to develop and refine a Generative AI model (Figure 4) for heart failure prediction, using a dataset featuring cardiovascular health-related attributes.

Figure 4- Comprehensive Generative AI Architecture



Step 1: Data Preparation and Balancing implement on the dataset of 1,025 records – Relevant features for cardiovascular conditions were selected from the dataset for heart failure prediction. These included age, sex, chest pain type (cp), resting blood pressure (trestbps), cholesterol (chol), fasting blood sugar (fbs), resting electrocardiographic results (restecg), maximum heart rate achieved (thalach), exercise-induced angina (exang), ST depression by exercise (oldpeak), peak exercise ST segment slope, number of vessels colored by fluoroscopy (ca), and thalassemia (thal). The target variable was cardiovascular disease (cvd), indicating heart failure. Missing values were addressed by replacing them with the column mean, ensuring data completeness without dropping any rows. Balancing was critical, especially considering potential class imbalances where heart failure cases were fewer than non-heart failure cases. This thorough data preparation laid a solid foundation for the modeling stages.

Step 2: Creation and Training of the Generative AI Model using GAN – With the dataset ready, the next step was developing a Generative AI model to enhance heart failure prediction using a GAN. Features and targets were converted to PyTorch tensors for neural network processing. A DataLoader was used to batch the data efficiently during training. The GAN comprised two neural networks: a generator, which created artificial data starting from random noise and converting it into patient-like data points, and a discriminator, which classified data points as either real or synthetic. The GAN training alternated between these networks for 5,000 epochs, gradually improving the generator's ability to produce synthetic data that became increasingly difficult for the discriminator to distinguish from real data. The synthetic data generated by the GAN was added to the original dataset, augmenting it for further model training.

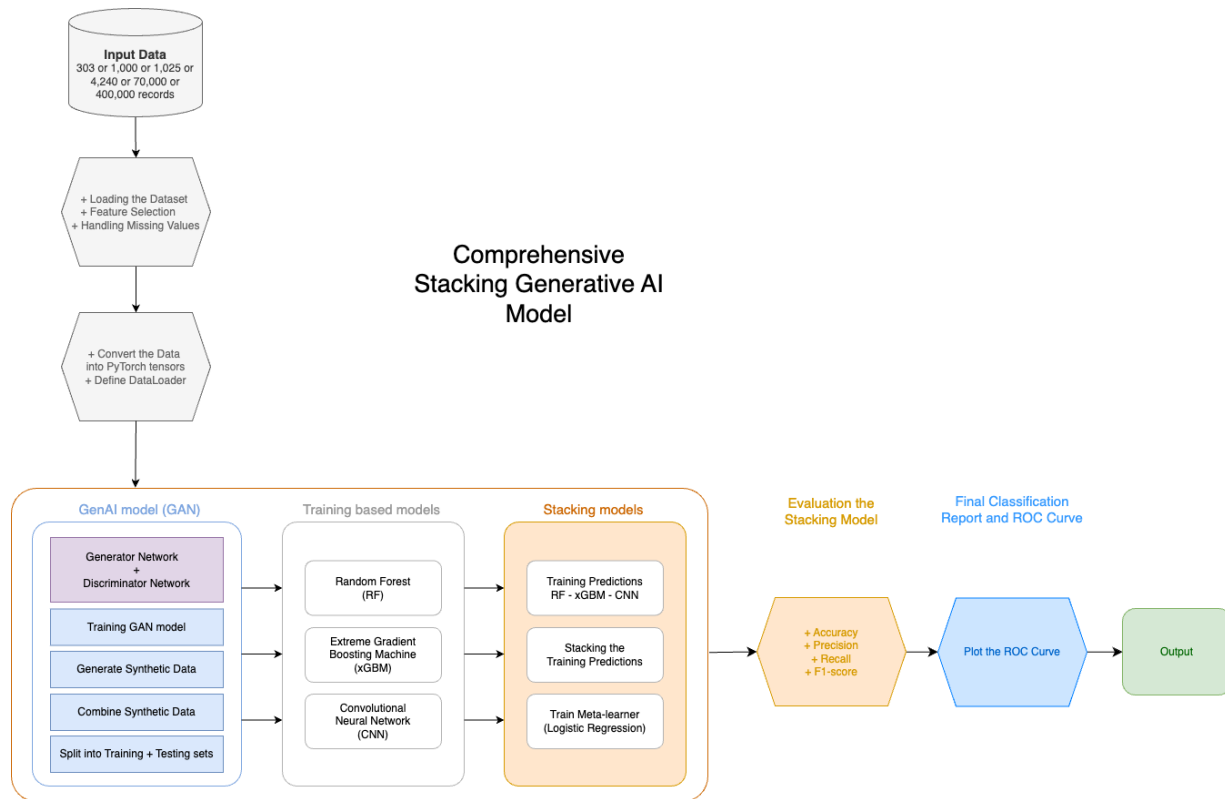
Step 3: Fine Tuning of Generative AI Model with Early Stopping – Early stopping was implemented to prevent overfitting and optimize the training process. This involved monitoring the discriminator loss, and if it failed to improve after a certain number of epochs, the training was stopped. Early stopping not only conserved computational resources but also protected the model from overfitting to the training data. Once the GAN was trained using early stopping, additional synthetic data representing heart failure cases (the positive class) was generated. When applying this to a dataset of 4,240 records, the GAN produced 2,952 synthetic data points to balance the majority class population of 3,596, which originally had a minority population of 644 records. This newly generated data was integrated with the original dataset, shuffled to eliminate order bias, and subsequently used for the final model training and evaluation.

Step 4: Training and Evaluation – The final phase was the training of a machine learning model on the augmented data, and the subsequent evaluation of that model. First, the combined data of real and synthetic were divided into a training set and a test set to validate the model. Feature

scaling using Standard Scaler was applied to standardize all features, ensuring they contributed equally during training. The Random Forest Classifier model, with setting of `n_estimators` was 100, `max_depth` equal 5, and `random_state` for 42, was chosen for its robustness and suitability for large datasets with complex feature interactions. The model trains and tests on the test set as described below; important metrics include Accuracy, ROC AUC, and a detailed Classification Report. The ROC AUC is plotted, which helps in visualizing the model's ability to differentiate between heart failure and non-heart failure patients, with the AUC indicative of overall performance.

This approach to developing the heart failure predictor utilized GAN-based data augmentation followed by training a Random Forest Classifier, demonstrating the model's potential in handling imbalanced data. The structured process involved data preparation, synthetic data generation using GAN, early stopping during training, and final evaluation using traditional machine learning techniques, resulting in a robust model capable of predicting heart failure accurately. This method highlights the importance of each step in producing a reliable predictive model in healthcare, where precision and accuracy are crucial for patient outcomes.

Figure 5- The proposed Comprehensive Stacking Generative AI Architecture



The **Comprehensive Stacking Generative AI model** for heart failure prediction integrates multiple machine learning techniques, combining traditional models like Random Forest (RF), Gradient Boosting Machine (GBM), and Extreme Gradient Boosting Machine (xGBM) with Convolutional Neural Networks (CNN) and Generative Adversarial Networks (GANs). This approach allows the model to handle class imbalance, generate synthetic data to improve learning, and combine multiple model predictions to achieve better performance.

Step 1: Data Preparation, Balancing, and Processing – The heart failure dataset, which includes key cardiovascular features such as age, cholesterol levels, and resting blood pressure, is first loaded. The target variable indicates whether a patient experienced heart failure. Initially, missing values are handled by applying appropriate imputation techniques to maintain data integrity. Standard Scaler is used to scale the balanced dataset, ensuring that all features are scaled consistently-a very important factor in training neural networks (Pedregosa et al., 2011).

Step 2: Defining Generator and Discriminator Networks for GAN - This stage defines the Generator and Discriminator networks for GAN. The generator generates synthesized data resembling real heart failure patient data while the discriminator differentiates whether that data is real or generated as shown in Figure 1. The Generator network receives a latent vector (random noise) as input, which passes through multiple fully connected layers. Each layer is activated using ReLU functions, with 128 units in the first hidden layer and 256 units in the second hidden layer. The final output is generated using a Tanh activation function, which scales the generated data between -1 and 1, appropriate for normalized medical data (Radford et al., 2015). This synthetic data can then be added to the real dataset to enhance the diversity of the training data.

Sample of the generator network is structured in Python code as follows:

```
class Generator(nn.Module):
    def __init__(self, input_dim, output_dim):
        super(Generator, self).__init__()
        self.network = nn.Sequential(
            nn.Linear(input_dim, 128),
            nn.ReLU(),
            nn.Linear(128, 256),
            nn.ReLU(),
            nn.Linear(256, output_dim),
            nn.Tanh()
        )
    def forward(self, x):
        return self.network(x)
```

The Discriminator network acts as a binary classifier, determining whether a heart failure record is real or synthetic. This input is processed through fully connected layers activated by LeakyReLU. The final output is given as a probability score, output from the Sigmoid function, indicating whether the input is real or fake. The adversarial training ensures that over time, the generator produces increasingly realistic synthetic data (Goodfellow et al., 2014).

Sample of the discriminator network is structured in Python code as follows:

```
class Discriminator(nn.Module):
    def __init__(self, input_dim):
        super(Discriminator, self).__init__()
        self.network = nn.Sequential(
            nn.Linear(input_dim, 256),
            nn.LeakyReLU(0.2),
            nn.Linear(256, 128),
            nn.LeakyReLU(0.2),
            nn.Linear(128, 1),
            nn.Sigmoid()
        )
    def forward(self, x):
        return self.network(x)
```

Step 3: Generating the Model – Then the GAN is trained over 5000 epochs with Adam optimizers (0.00005). Generator and discriminator are trained separately so that generator becomes proficient at producing authentic synthetic heart failure data and discriminator becomes skilled in distinguishing real from fake data. This training ensures the quality of the synthetic data that will be used later on in the actual dataset for modeling improvement.

Step 4: Synthetic Data Generating with GAN – Synthetic data is computed after the GAN is fully trained, by-passing random noise vectors into the generator. The artificial data are merged with the original heart failure data to generate a large training set consisting of real and artificial patient samples. That's a way of teaching models from a wider set of samples and generalizing to new unseen data.

Step 5: Divisible the Data into Training and Test Sets (80/20) – Once the combined data has been created with synthetic data, the overall dataset is divided into training and test sets so that the model gets run against the unseen data to determine how well it performs in the real world. Data is then normalized with a Standard Scaler so that all the input features are also scaled,

which is very important for neural networks like CNNs, where feature scaling is very important to learning.

Step 6: Training the Base Models (RF, xGBM, CNN) – Now it's time to train the individual models – Random Forest (RF), Extreme Gradient Boosting (xGBM), and Convolutional Neural Network (CNN)- shown in the diagram of Figure 5. The Random Forest has 100 trees, maximum depth 10; min samples split 10; random state 42. Complex feature interactions are accounted for by 200 estimators, 0.05 learning rate, 0.8 subsample ratio and 42 random state in xGBM model.

On CNN the framework is for overfitting and generalization. It starts with a Conv1D (16 filter) kernel size of 3 and then the MaxPooling (2 pool size) dimensionality reduction layer. 0.6 Dropout layer is added for Overfit prevention. The output is then flattened and through a 32 units Dense with ReLU activation, Dropout again, and finally a sigmoid output for binary classification. Model is built with Adam optimizer and binary cross-entropy loss function. Stopping is implemented early so as not to overfit and training is terminated if validation loss fails to improve after 5 epochs. It trains the model for a maximum of 50 epochs with a batch size of 32 and validation is done with 20% of the data.

Step 7: Stacked Prediction Training of Meta-Learner – When the base models have been trained, their predictions are the input for the stacking model. RF, xGBM, CNN predictions go to the meta-learner which is Logistic Regression. This meta-learner is trained to derive the final classifier on the basis of the strength of the base models.

Step 8: Evaluation of Stacked Model – Meta-learner is tested against the test set and performance metrics like accuracy, precision, recall and F1-score are calculated. ROC AUC is calculated how

good a model is at detecting heart failure vs. non-heart failure. ROC curve shows the tradeoff between sensitivity and specificity to clearly show how the model performed at different thresholds.

Step 9: Final Classification Report and ROC curve – The final product is classification report with precision, recall, F1-scores of both classes, ROC curve (chapter 4). The ROC AUC curve is the graph that indicates how well the model performed; a high ROC AUC means the prediction accuracy is high. This analysis gives us an idea about whether the model is able to predict heart failure well enough to adopt in clinical settings for early detection of disease.

Conclusion – The proposed Comprehensive Stacking Generative AI Model is a strong heart failure prediction tool that merges traditional machine learning algorithms with innovative approaches like GAN in the generative AI model. Stacking synthetic data generated by GANs is the key to getting the model to generalize well and be super-fast on actual medical data. Bringing together models such as RF, xGBM, and CNN ensures the ensemble gets the right predictions, which is crucial for early medical diagnosis, Chawla et al. (2002), Goodfellow et al. (2014), Pedregosa et al. (2011), and Radford et al. (2015).

3.6. Evaluation Measurement and Validation Methods

Performance of each model is calculated with different parameters like accuracy, ROC AUC, precision, recall, F1 score etc. Accuracy is calculated as:

$$\text{Accuracy} = \frac{\text{True Positives} + \text{True Negatives}}{\text{Total Instances}}$$

ROC AUC represents how discriminative the model is between classes and is calculated as:

$$ROC\ AUC = \int_0^1 TPR(FPR) d(FPR)$$

-where TPR is True Positive Rate and FPR is False Positive Rate. K-fold cross-validation especially stratified cross-validation with imbalanced data sets makes models reliable across all data splits.

Stacking Generative AI model to validate it to unseen data is validated with various methods that ensure that the model can be extended without overfitting to the unseen data. These are 5- and 10-fold cross-validation, learning curves to measure training size dependent performance, regularization, and hyperparameter optimization for optimal model behavior. The relevant mathematical equations of these methods are given below.

Cross-Validation (cv=5 and cv=10)

Cross-validation is a resampling method that runs model against the dataset with k equal sized "folds" — the model is trained on k-1 folds and evaluated on the remainder. The same is done k times and the average performance is used to evaluate robustness. Mathematically, k-fold cross-validation accuracy is:

$$CV\ Accuracy = \frac{1}{k} \sum_{i=1}^k Accuracy_i$$

where k is the number of folds, and $Accuracy_i$ is the accuracy for the i^{th} fold.

For 5-fold cross-validation, the model was trained and tested over five data splits with accuracies of [0.9938, 1.0000, 0.9877, 0.9969, 0.9938]. The mean accuracy was 99.4%. Similarly, 10-fold

cross-validation yielded a mean accuracy of 99.4%, confirming the model's consistency and generalization across different data splits (James et al., 2013).

Learning Curve

Learning curve – It represents how well a model performs given a training set size and is available for overfitting or underfitting. It shows training and cross-validation-accuracy as a percentage of training samples:

$$\text{Error} = \frac{1}{n} \sum_{i=1}^n L(\hat{y}_i, y_i)$$

- where n is the number of training instances, \hat{y}_i is the predicted value and y_i is the actual value. L is the loss function, binary cross-entropy. The Figure 13 shows Convergence Learning curve of Training and Validation Accuracies is 99.8%, It generalizes easily without overfitting and has a satisfactory performance for unseen data, as Goodfellow et al. (2016).

Regularization

Regularization helps prevent over-complexity by penalizing large weights. In the Logistic Regression meta-learner, L2 regularization was applied, adding a regularization term to the loss function to shrink weights:

$$L(w) = \text{Loss}(w) + \lambda \|w\|_2^2$$

where $L(w)$ is the regularized loss, $\text{Loss}(w)$ is the original binary cross-entropy loss, λ is the regularization strength, and $\|w\|_2^2$ is the sum of squared weights. Grid search was used to find the optimal λ , ensuring the model remained well-tuned without overfitting (Ng et al, 2004).

Hyperparameter Tuning

Grid search was used to optimize the Logistic Regression meta-learner by exploring different hyperparameter combinations. The goal was to find the best regularization parameter (C) for Logistic Regression:

$$C = \frac{1}{\lambda}$$

Grid search iterates over a range of C values and evaluates model performance on the validation set. The best $C = 0.01$ was chosen based on cross-validation scores.

The combination of cross-validation, learning curves, regularization, and hyperparameter tuning provided a comprehensive validation approach. These mathematical techniques ensured that the Stacking Generative AI model was well-calibrated to generalize effectively without overfitting, making it suitable for deployment in heart failure prediction scenarios.

Statistical Validation

Mixed-effects modeling was employed to address variability across datasets and verify the robustness of the findings. This approach is particularly well-suited for healthcare applications, where patient data often exhibit hierarchical and interdependent structures. The results from the ANOVA analysis, which yielded an F-statistic of 2192.804 with a p-value below $2.26e-276$, demonstrate the statistical significance of the synthetic data's impact on improving model

performance. These results align with the methodologies outlined in Springer's framework, which emphasize the importance of mixed-effects models for validating machine learning performance in clinical datasets. By leveraging these advanced statistical techniques, this study ensures that the observed improvements are both reliable and generalizable.

Chapter 4: RESULTS

4.1. Implementation Results

4.1.1. Research Question 1: Performance Comparison between Traditional Models and Neural Network Models: How do traditional machine learning models compare to neural network-based models in terms of accuracy and ROC AUC for heart failure prediction?

The study has underlined the relative strengths and weaknesses of traditional machine learning models and neural network-based models regarding the prediction of heart failure outcomes. The implementation used classical machine learning models like Logistic Regression (LR), Support Vector Machines (SVM), Random Forest (RF), Gradient Boosting Machine (GBM), and Extreme Gradient Boosting Machine (xGBM) on nine different-sized datasets, in addition to neural network models like Convolutional Neural Networks (CNN) and Recurrent Neural Networks (GRU-based models).

On small datasets, 303 records, RF outperformed most other traditional models with an accuracy of 83% and ROC AUC of 0.91. GBM and xGBM were slightly less accuracy of 79% and 80%, respectively, but had high ROC AUCs, 0.87 and 0.86, respectively (Table 10). Neural network models (CNN) had an accuracy of 82% and ROC AUC of 0.85, similar but not better than RF.

However, the notable advantage of neural network models was that they could identify more intricate patterns in the data when the data set got bigger.

CNN's accuracy dropped to 79% accuracy and ROC AUC of 0.85 in the 1,000-record dataset, while RF and xGBM, traditional model, had 90% accuracy and 88% accuracy with ROC AUC of 0.94 and 0.95, respectively (Table 2). The pattern continued in the 1,025-record dataset, where xGBM scored 93% accuracy and ROC AUC of 0.98 compared to CNN, with 82% accuracy and ROC AUC of 0.93. GRU models did well in these data, with the best accuracy of 84%, and ROC AUC was 0.92 (Table 5).

Neural networks were better for larger data—the 400,000-records dataset, such as in cases where CNN's accuracy was 78% and ROC AUC 0.86, while GBM had 77% accuracy and ROC AUC 0.85. Simple models such as RF were good on accuracy of 90% and ROC AUC of 0.96 but eventually surpassed by Generative AI-enriched stacking models, accuracy of 96% and ROC AUC of 0.99 (Table 4).

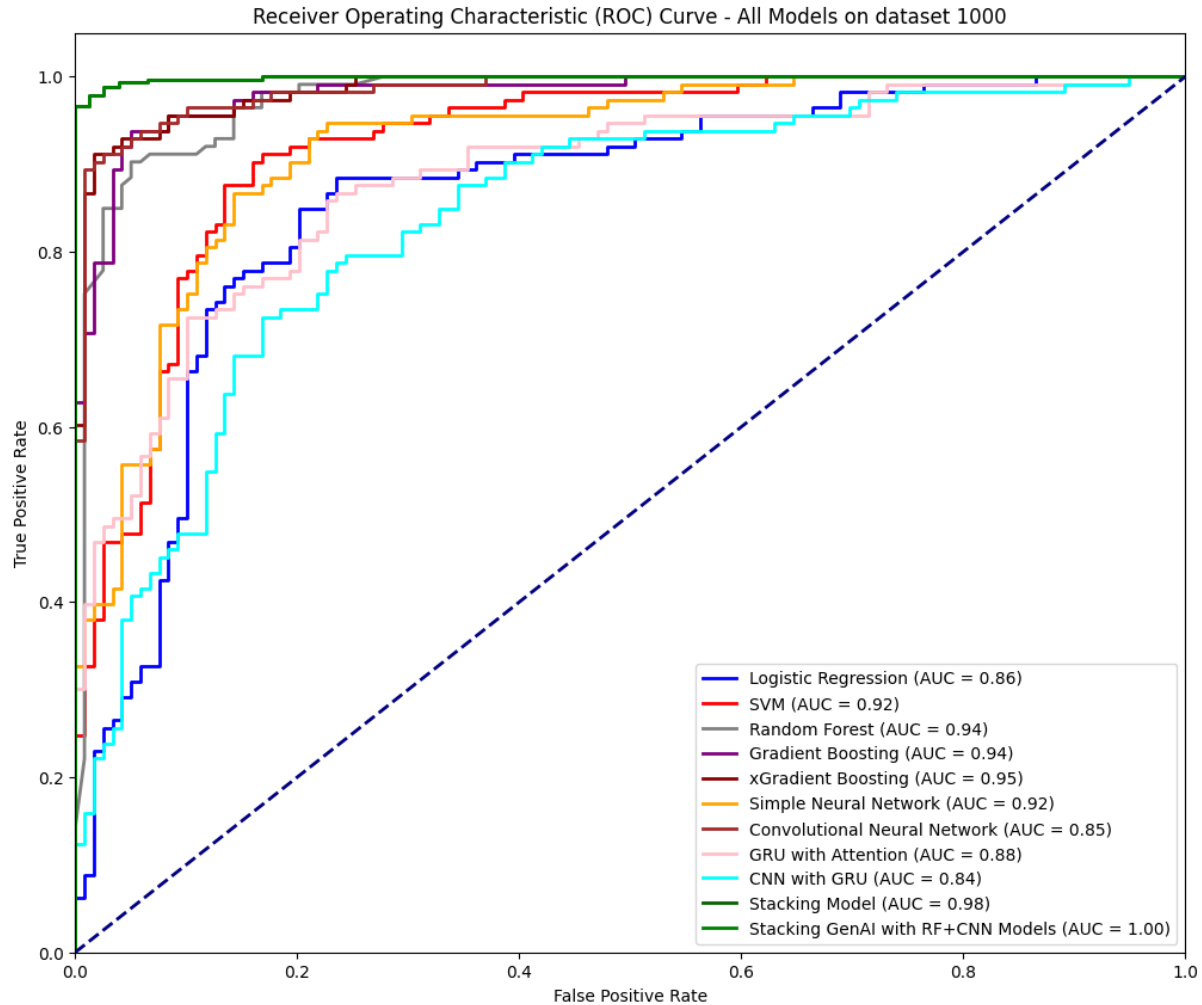
The ML models were highly interpretable and efficient in small datasets, while neural network models were highly effective in larger datasets as they handle high-dimensional, complex data. But for all their strengths, neural network models can be computationally costly and unusable, as can the classical model, so hybrid approaches must be developed that can take the best from both paradigms.

Table 2- Performance of proposed model vs. other models on dataset of 1,000 records

Dataset	Performance	Model									Proposed Model		Current	Source
		LR	SVM	RF	GBM	XGB	CNN	GRU w/ Attention	CNN w/ GRU	Stacking ML+DL	Gen AI	Stacking Gen AI	Research Literature	Reference
1,000	Accuracy	80	85	90	88	88	79	77	78	94	98	98	94	Dumlao, J. (n.d.)
	ROC AUC	86	92	94	94	95	85	84	84	98	99	99.9	NA	

On another comparison with an existing analysis, Dumlao, J. (n.d.), a Cardiovascular Health Analysis on Kaggle, which uses the same dataset, where a Random Forest model reached an accuracy of 98% on a comparable dataset (Table 2), parallel with ML and DL stacking models, the Stacking Generative AI model reached an accuracy of 0.999 and outperformed both the individual models compared in this study and those identified in another research (Figure 6). This suggests that hybrid models like Stacking Generative AI may significantly advance healthcare predictive modeling, particularly in predicting heart disease.

Figure 6- The ROC Curve for the dataset of 1,000 records



The Stacking Generative AI model also performed with notable distinction on the largest dataset, containing 400,000 records, achieving a ROC AUC of 0.99 and an accuracy of 96% (Figure 7). This matched the performance of the standalone Generative AI model, which had an ROC AUC of 0.987 and outperformed all other models in this study. With an accuracy of 96%, the Stacking Generative AI model demonstrated its capability to handle large and complex datasets effectively. Among the individual models, Random Forest also performed well with an ROC AUC of 0.96, while xGBM and CNN with GRU reached an ROC AUC of 0.88. However, the Stacking Generative AI model's ability to integrate multiple predictions into a more accurate outcome outperforms the individual models.

Figure 7- The ROC Curve for dataset of 400,000 records

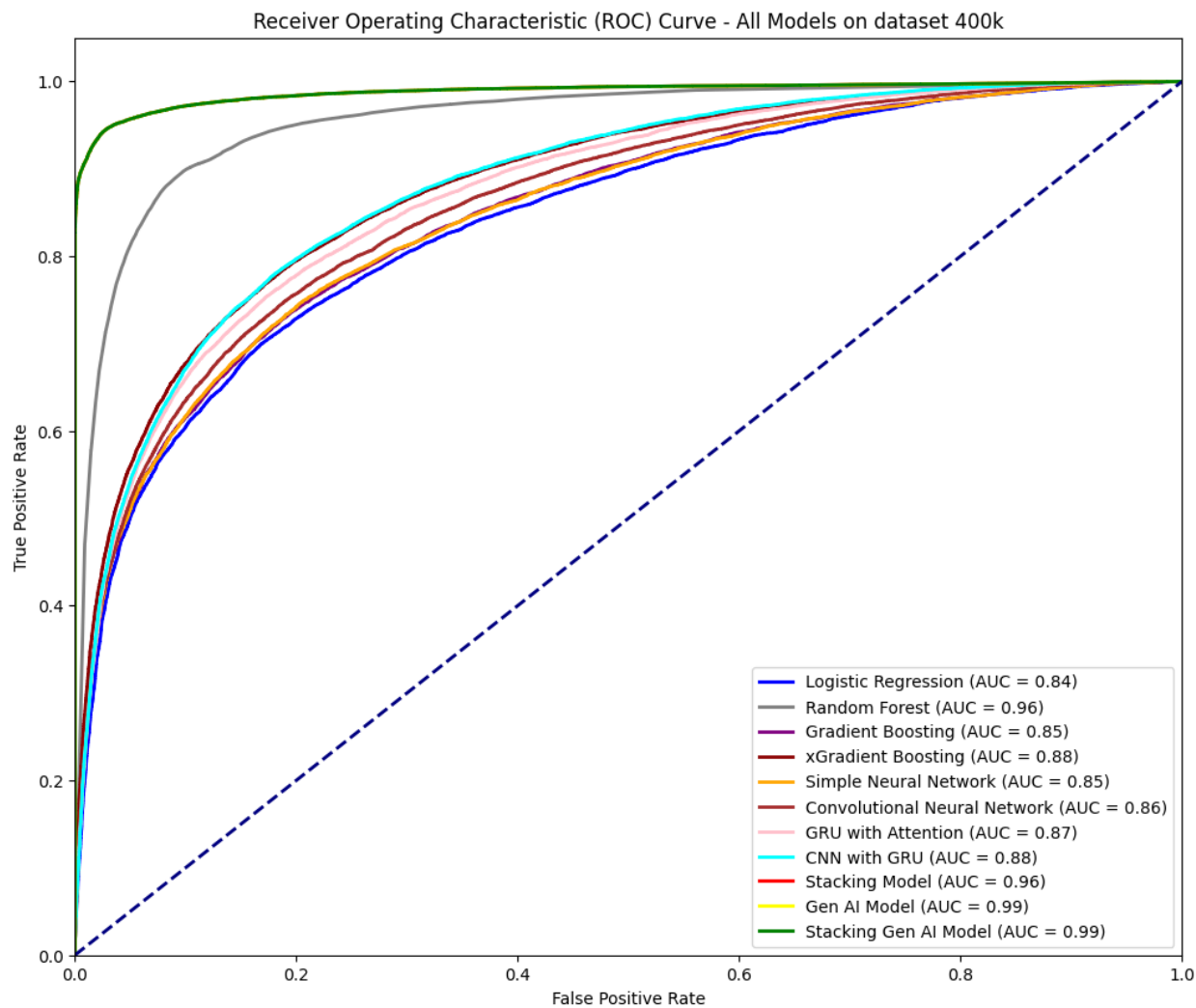


Table 3- Performance of proposed model vs. other models on dataset of 400,000 records

Dataset	Performance	Model									Proposed Model		Current	Source
		LR	SVM	RF	GBM	XGB	CNN	GRU w/ Attention	CNN w/ GRU	Stacking ML+DL	Gen AI	Stacking Gen AI	Research Literature	Reference
400,000	Accuracy	77	NA	90	77	80	78	79	80	90	95	96	91	Khan, H. et al. (2024)
	ROC AUC	84	NA	96	85	88	86	87	88	96	98	99	91	

Additional comparison with models’ performance by Khan, H. et al. (2024), which uses the same dataset of 400,000 records:

Table 4- Performance of proposed model vs. article's models on dataset of 400,000 records

Dataset	Performance	Proposed Model	Compared Article	Compared Article
		Stacking Generative AI	EnsCVDD-Net	BICVDD-Net
400000	Accuracy	96	88	91
	ROC AUC	99	88	91

Accuracy: The Stacking Generative AI model achieved the highest accuracy at 96%, significantly outperforming all models tested in this study and in the article. The previous top-performing models, such as Random Forest at 90% and CNN with GRU at 80%, were surpassed by a wide margin. In the article, EnsCVDD-Net achieved an accuracy of 88%, while BICVDD-Net achieved 91%, both lower than the Stacking Generative AI model (as shown in Table 3 & 4).

ROC AUC: The Stacking Generative AI model had the highest ROC AUC at 0.99, outperforming all other models when implement on the dataset of 400,000 records. In these tests, the second-best result was 0.98 from the Generative AI model, and 0.96 from Random Forest and stacking-based ML + DL models. While performance from Khan, H. et al. (2024), the ROC AUC for EnsCVDD-Net was 0.88, and for BICVDD-Net, it was 0.91, showing that the proposed model outperformed the innovative methods presented in the article.

Comparative Summary

The proposed Stacking Generative AI model, combining RF, XGBM, and CNN, clearly outperformed the models proposed in the article in terms of both accuracy and ROC AUC. Stacking different models, including neural networks and ensemble methods like Random Forest and XGB, produced superior results in terms of classification metrics. This demonstrates the

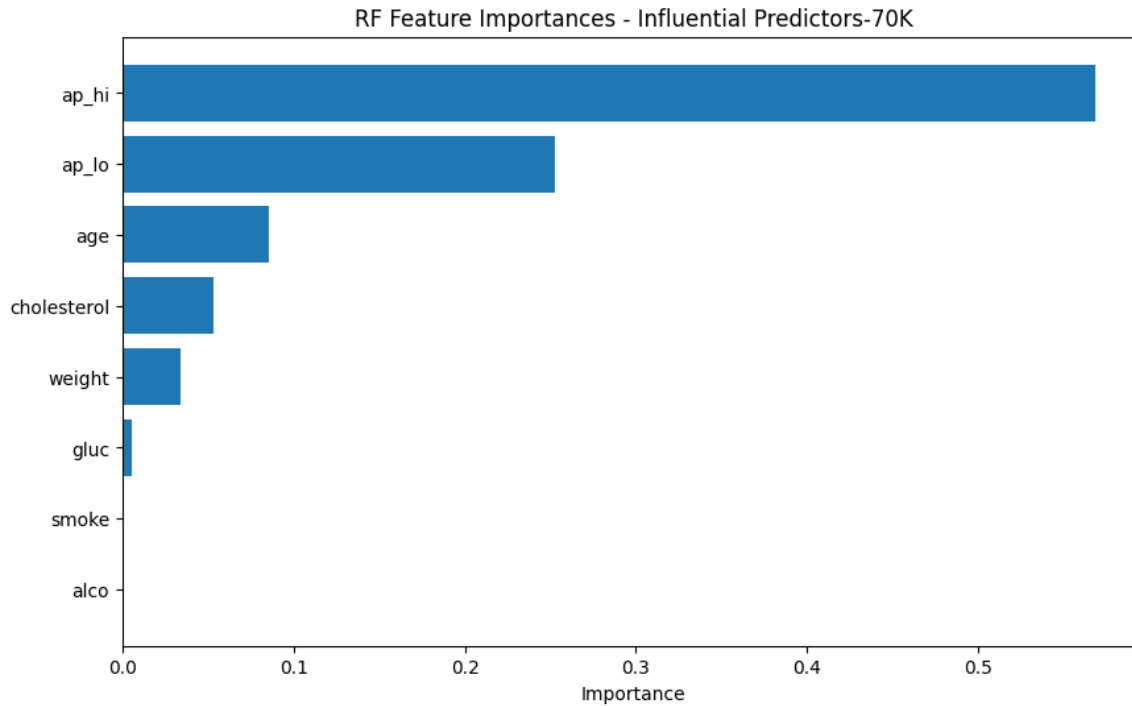
effectiveness of combining machine learning and deep learning approaches within the framework, further enhanced with fine-tuning, early stopping, and threshold adjustment at 0.36. The balanced integration of traditional ML and advanced neural network models ensures robust feature extraction and prediction, leading to better performance than standalone or ensemble models, as noted in the article.

4.1.2. Research Question 2: What are the most influential predictors of heart failure across different datasets, and how do they affect overall model performance?

Identifying the most relevant influential predictors is essential to accurately interpreting heart failure prediction models. Feature importance analysis for many datasets and models highlights some important predictors. RF models, renowned for ranking feature relevance, provide valuable insights into the key factors associated with heart failure.

In the 70,000 records dataset, systolic and diastolic blood pressure (ap_hi and ap_lo) were the best predictors, followed by age and cholesterol (Figure 8). These account for 74% accuracy and an ROC AUC of 0.81. This prediction fits clinical risk factors and clearly shows the power of RF models in identifying relevant features.

Figure 8- The Risk Factors / Feature Importances of 70,000-record dataset



RF analysis found the major influencing factors of chest pain (cp), oldpeak, and number of major vessels (ca) for the dataset with 1,025 records (Figure 9). The accuracy of 95% and ROC AUC of 0.999 achieved with the Stacking Generative AI model on this dataset shows that adding these capabilities into advanced ensembles is effective. Similarly, among the 400,000-record dataset, body mass index (BMI), angina, and general health were the best predictors (Figure 10), where the Stacking Generative AI model achieved 96% accuracy and ROC AUC of 0.99.

Figure 9- The Risk Factors / Feature Importances of 1,025-record dataset

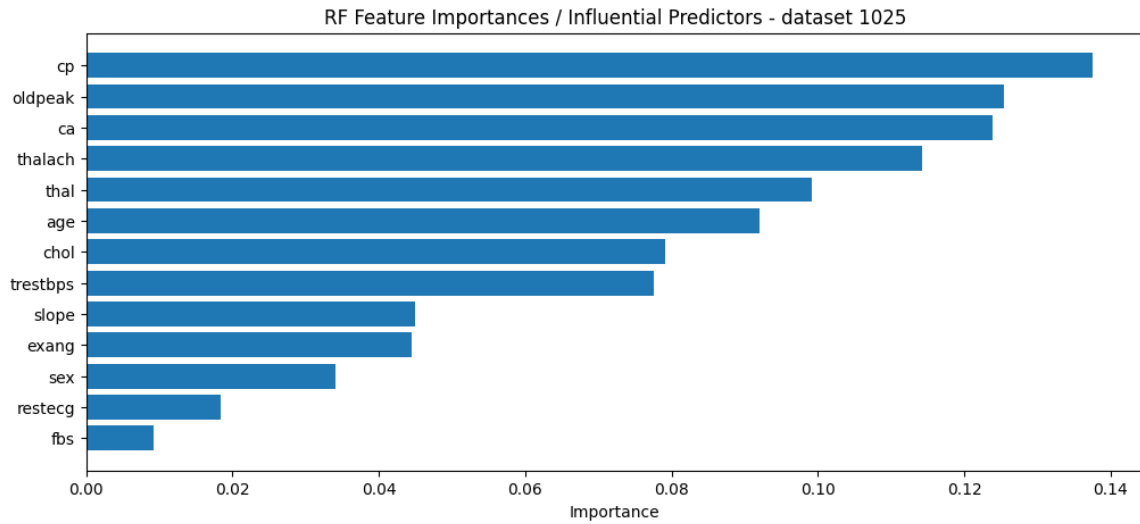


Figure 10- The Risk Factors / Feature Importances of 400,000-record dataset

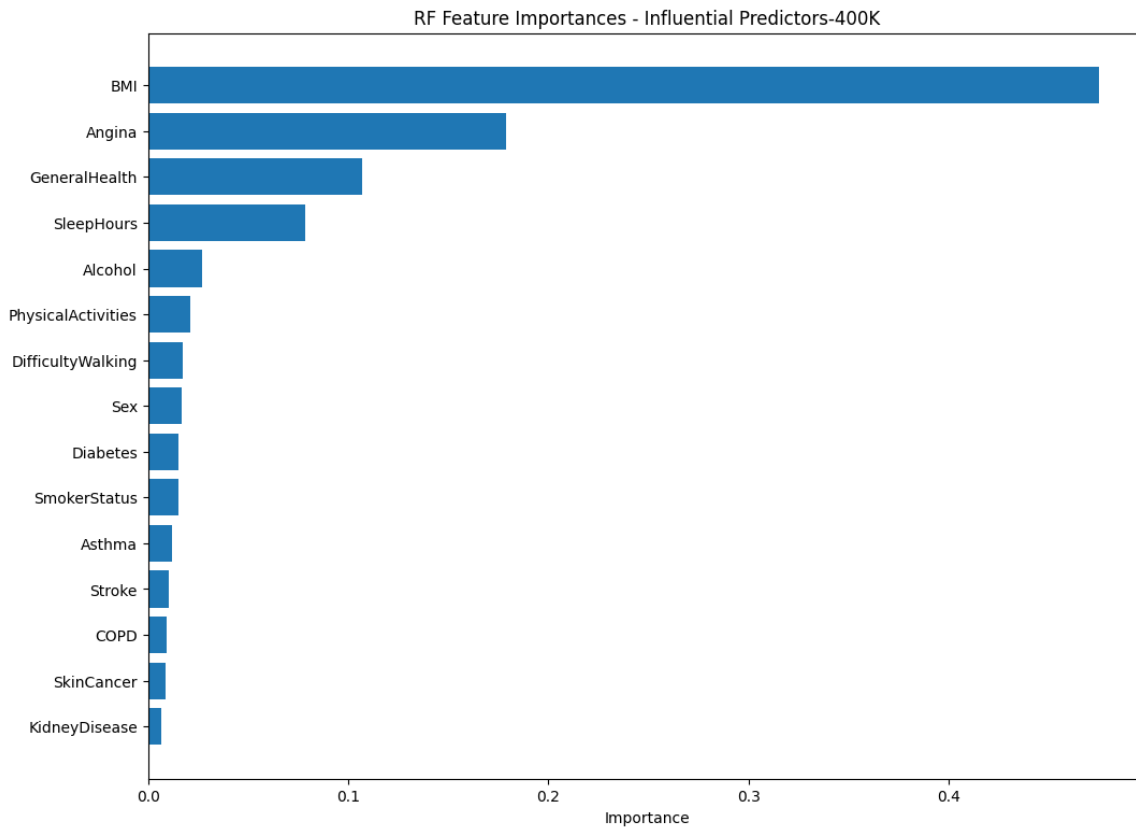
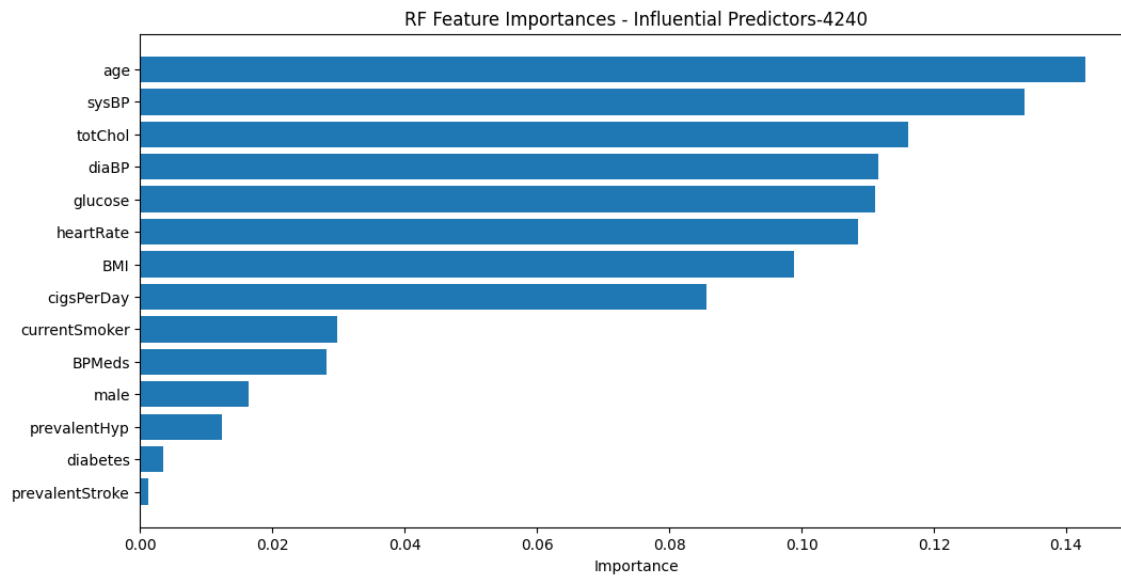
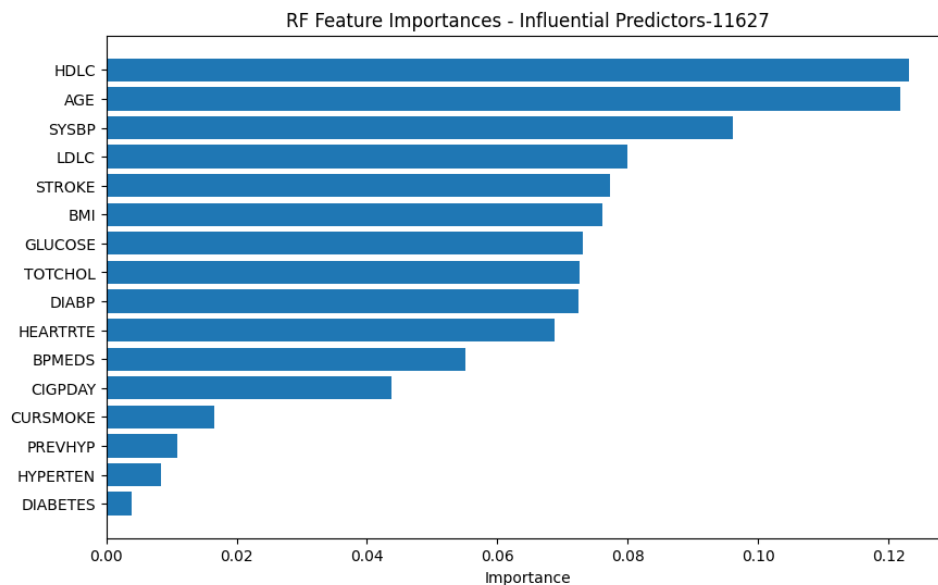


Figure 11- The Risk Factors / Feature Importances of 4,240-record dataset



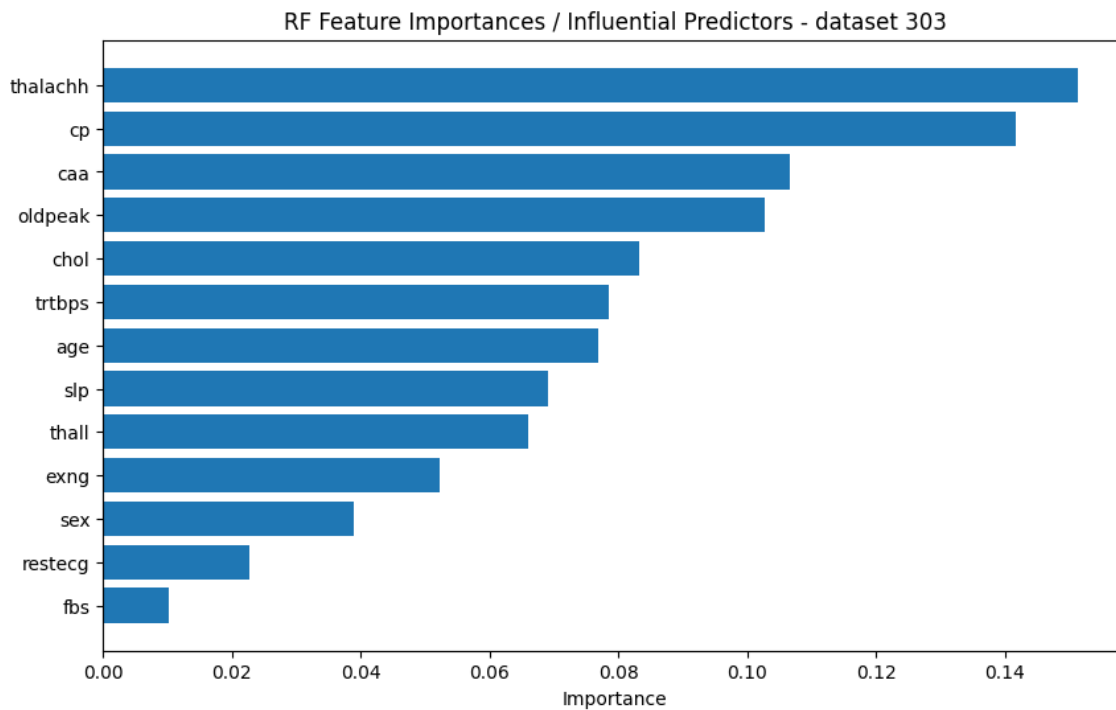
Age, sysBP, and cholesterol were the main predictors for the 4,240-record dataset (Figure 11), with an accuracy of 92% and a ROC AUC of 0.96 in the Stacking Generative AI model. The 11,627-record dataset identified HDLC, age, and systolic blood pressure as the main drivers (As seen in Figure 12), with a 91% and 0.95 ROC AUC, respectively.

Figure 12- The Risk Factors / Feature Importances of 11,627-record dataset



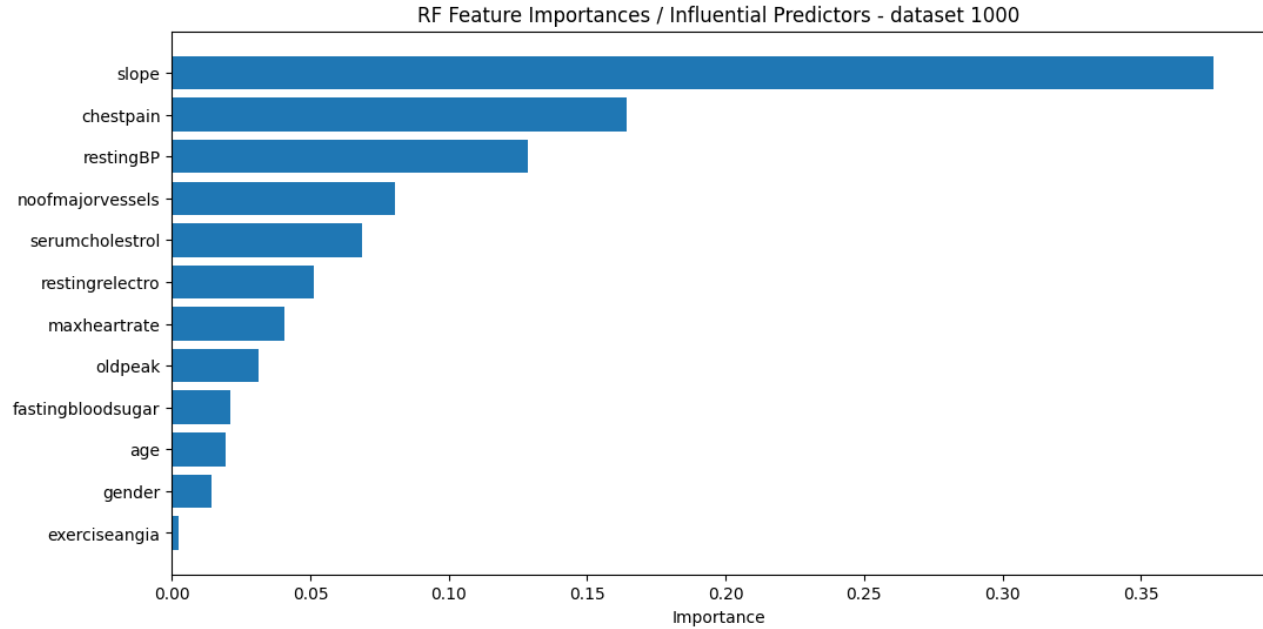
In the 303-record dataset, RF feature importance analysis identified the heart rate attained (thalachh), type of chest pain (cp), and number of major vessels (caa) as the strongest predictors (Figure 13). These variables correlate with cardiovascular conditions and align with known clinical knowledge. The Stacking Generative AI algorithm that combines these powerful predictors achieved an accuracy of 95% and ROC AUC of 0.99. The model focuses on the model's generalizability and uses these predictors efficiently, even in small datasets.

Figure 13- The Risk Factors / Feature Importances of 303-record dataset



In the 1,000-record dataset, RF analysis showed the slope of the ST segment (slope), chest pain (chest pain), and resting blood pressure (restingBP) as major predictors (Figure 14). These indicators of myocardial stress and vascular health make them essential for accurate HF predictions. The Stacking Generative AI model got 98% accuracy and an ROC AUC of 0.999, which shows that these features have predictive capabilities if embedded in advanced ensemble models.

Figure 14- The Risk Factors / Feature Importances of 1,000-record dataset



Blood pressure, chest pain, cholesterol levels, and age emerge as critical predictors of heart failure risk across all nine datasets analyzed. These variables contribute to the interpretability of predictive models, enhancing their applicability in clinical settings. High-performing models, such as the Stacking Generative AI model, leverage these predictors to achieve superior accuracy and ROC AUC values. This underscores the importance of systematically identifying and integrating relevant features into predictive models to optimize performance.

The consistent relevance of these key predictors across datasets of varying sizes and contexts highlights the broader applicability of this large-scale analysis. These findings hold significant implications for clinical practice and future research, offering valuable insights into developing and refining predictive models for heart failure.

Figure 15- The ROC Curve for the dataset of 1,025 records

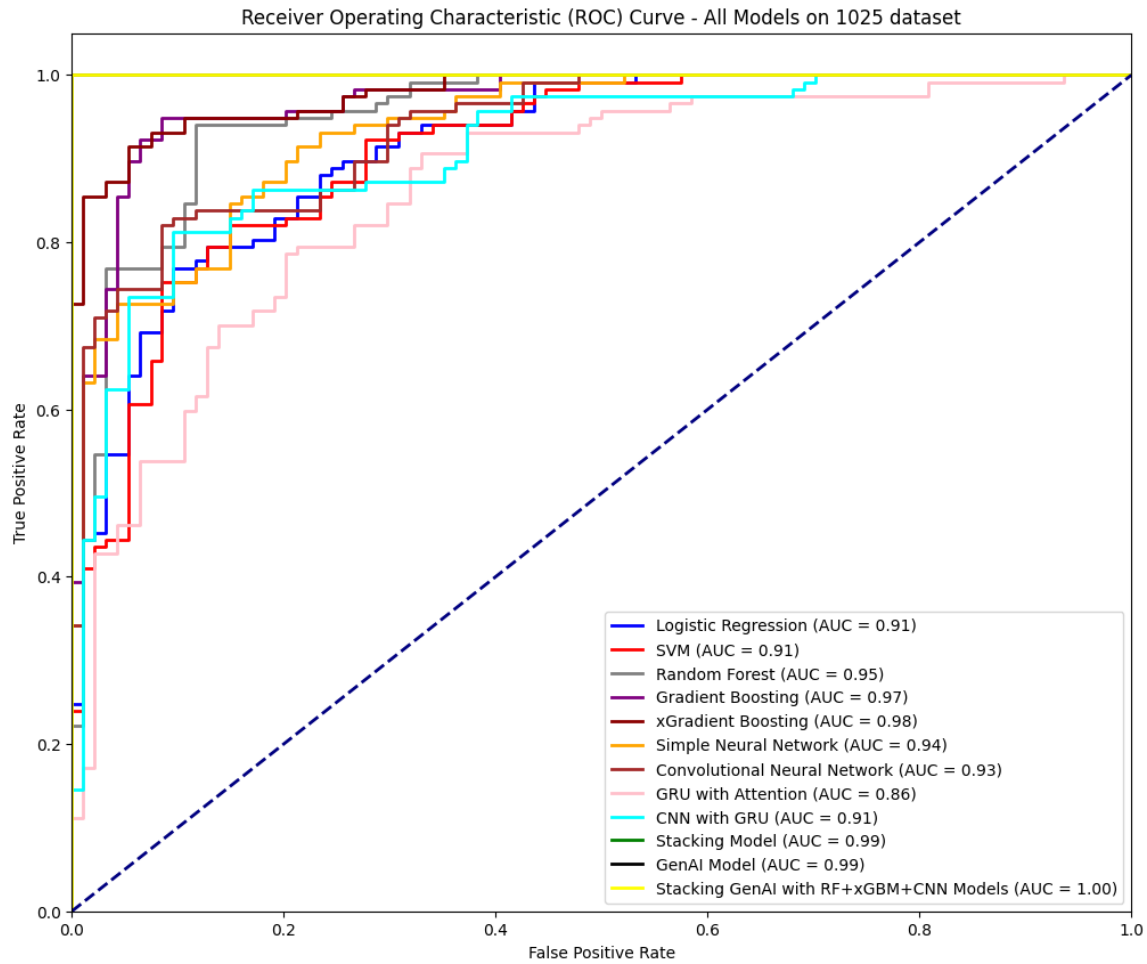


Table 5- Performance of proposed model vs. other models on dataset of 1,025 records

Dataset	Performance	Model									Proposed Model		Current	Source
		LR	SVM	RF	GBM	XGB	CNN	GRU w/ Attention	CNN w/ GRU	Stacking ML+DL	Gen AI	Stacking Gen AI	Research Literature	Reference
1,025	Accuracy	82	81	91	91	93	82	80	84	95	95	98	92	Nasser, A. (n.d.)
	ROC AUC	91	91	95	97	98	93	86	92	98	99	99.9	NA	

The proposed Stacking Generative AI model, using RF, GBM, xGBM, and CNN with Generative AI on the 70,000-record dataset, achieved a ROC AUC of 0.81. This is comparable to

individual models like xGBM, GRU with Attention, and stacking ML model, which they had ROC AUCs of 0.81 (Figure 16).

Figure 16- The ROC Curve for the dataset of 70,000 records

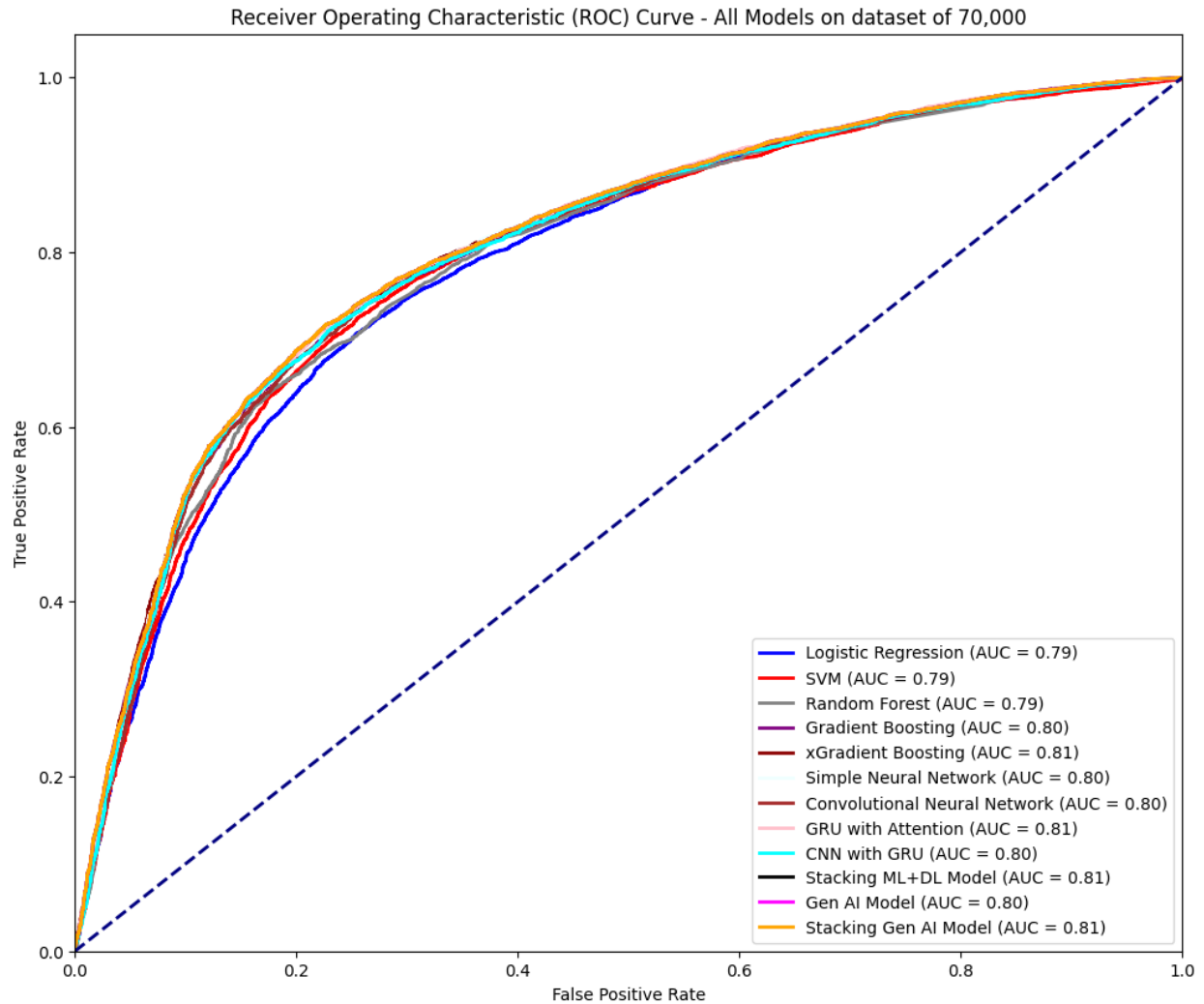


Table 6- Performance of proposed model vs. other models on dataset of 70,000 records

Dataset	Performance	Model									Proposed Model		Current	Source
		LR	SVM	RF	GBM	XGB	CNN	GRU w/ Attention	CNN w/ GRU	Stacking ML+DL	Gen AI	Stacking Gen AI	Research Literature	Reference
70,000	Accuracy	72	73	73	74	74	74	74	74	74	74	74	73	Jain, S. (n.d.)
	ROC AUC	79	79	79	80	81	80	81	80	81	80	81	NA	

The results show that, while models like Random Forest (with a ROC AUC of 0.79) and xGBM (ROC AUC of 0.81) performed well individually, the benefit of stacking diminished as the dataset size increased. However, the stacking model still showed a small improvement in predictive power, indicating its ability to synthesize the strengths of multiple algorithms when dealing with large, complex datasets.

Compare the Stacking Generative AI model result, 74% of accuracy and 0.81 in ROC AUC, while existing literature which range from 70% to 73% in accuracy and with NA in ROC AUC (As shown in Table 6).

4.1.3. Research Question 3: Can a hybrid stacking model that combines traditional machine learning and deep learning techniques provide superior predictive performance compared to single models?

This study investigates the effectiveness of the stacking model, which is a hybrid combination of traditional machine learning including RF and GBM with deep learning models such as CNN and RNN. In this combination, the best from both worlds of ML and DL are put together in order to further improve predictive accuracy beyond what is achievable by a single model.

The performance of the hybrid stacking model, combined ML, and DL, was evaluated on a set of datasets. Using this small dataset with 303 records, the stacking model showed 82% accuracy with an ROC AUC of 0.90. On running this model with much greater data with a record count of 4,240, it showed not only a relative increase in performance, but algorithms had significantly improved performance, yielding an accuracy of 90% with a ROC AUC of 0.97. These results

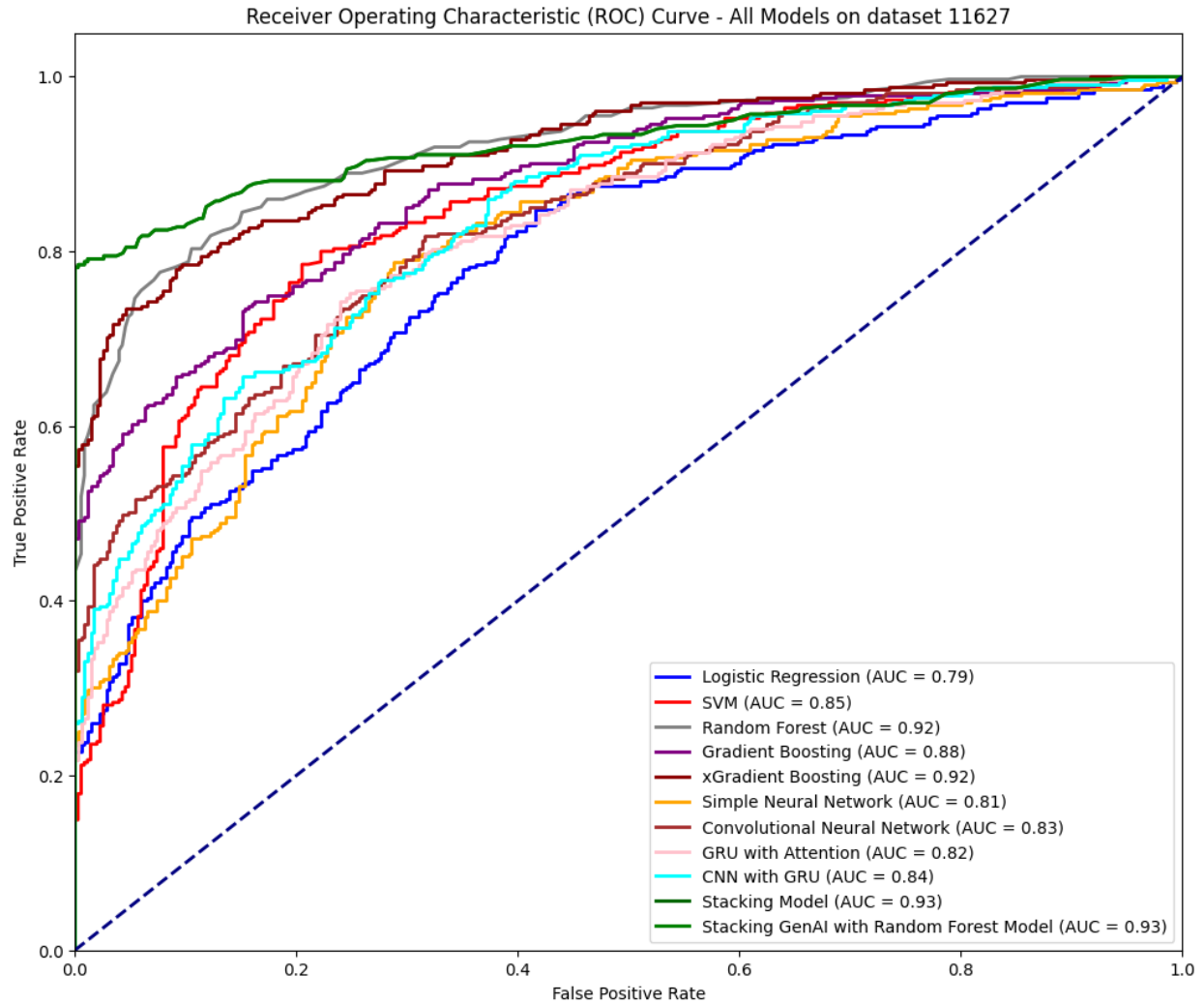
clearly present how the hybrid model can outperform other single-model methods, such as Logistic Regression, Support Vector Machine, CNN, and RNN, in both accuracy and ROC AUC.

This study used the baseline from the ML and DL stacking configuration and compared the different model combinations in order to provide a benchmark against individual models. This approach not only validated the superior performance of the hybrid model but also provided useful insights for moving on to a more advanced method of stacking. The results obtained from the ML and DL stacking model set a good precedent for the development of the Stacking Generative AI model.

Table 7- Performance of proposed model vs. other models on dataset of 11,627 records

Dataset	Performance	Model									Proposed Model		Current	Source
		LR	SVM	RF	GBM	XGB	CNN	GRU w/ Attention	CNN w/ GRU	Stacking ML+DL	Gen AI	Stacking Gen AI	Research Literature	Reference
11,627	Accuracy	71	78	84	79	83	74	74	73	85	91	91	97*	Sk K. B. et al (2023)
	ROC AUC	79	85	92	88	92	83	82	84	93	95	95	NA	

Figure 17- The ROC Curve for the dataset of 11,627 records



In that respect, Sk K. B. et al (2023) introduces a hybrid algorithm combination of Decision Tree and AdaBoost to predict CHD. Such an approach reached a high accuracy of 97.43%, with a True Positive Rate of 95.67% and True Negative/Specificity of 94.65% (Table 7). Although the proposed Stacking Generative AI model did not quite reach the same accuracy level, as this hybrid approach did, its performance was competitive and effective, taking into consideration deep learning models such as CNN and GRU, and synthetic data generation (Figure 17).

Compared to AdaBoost + Decision Tree: The hybrid model using AdaBoost and Decision Trees from the article indeed yielded good accuracy. This is because of its robust feature selection and boosting approach, effectively enhancing the weak classifiers. The contribution to the proposed method has almost the same performance using a more flexible architecture that integrates ML and DL, therefore being robust across datasets. While both approaches do an excellent job of predicting CHD, the Stacking Generative AI model provides an innovative, flexible, competitive approach with the more traditional hybrid methods. Because many algorithms are combined in their strengths, along with the high ROC AUC score, it shows its power in handling complex heart disease prediction tasks effectively.

Table 8- Performance of proposed model vs. other models on dataset of 4,240 records

Dataset	Performance	Model									Proposed Model		Current	Source
		LR	SVM	RF	GBM	XGB	CNN	GRU w/ Attention	CNN w/ GRU	Stacking ML+DL	Gen AI	Stacking Gen AI	Research Literature	Reference
4,240	Accuracy	65	67	71	81	86	70	63	67	90	93	92	91	Mienye et al. (2020)
	ROC AUC	74	74	79	89	93	77	70	72	97	96	96	NA	

The stacking with Generative AI with RF, GBM, and xGBM with GAN for the 4,240-record dataset achieved an ROC AUC of 0.96 (Figure 18). This is a very good number, though at the cost of substantially lower ROC AUC than traditionally (ML – DL) stacking, at 0.97 (Table 8). Even so, the Stacking Generative AI model outperformed single models such as CNN with GRU, with a much higher ROC AUC of 0.72, and other deep learning models like GRU with Attention, with an ROC AUC of 0.70.

Table 9- Performance of proposed model vs. article's models on dataset of 4,240 records

Dataset (4,240 records)	Performance	Mienye et al. (2020) (Framingham)	Proposed Stacking Generative AI Model
Accuracy	91%	91%	92%
ROC AUC	Not explicitly stated, but implied strong performance	Strong ROC AUC	96%

Compare performance with the proposed Stacking Generative AI model on the 4,240-record dataset versus that of Mienye et al. (2020), on the Framingham dataset.

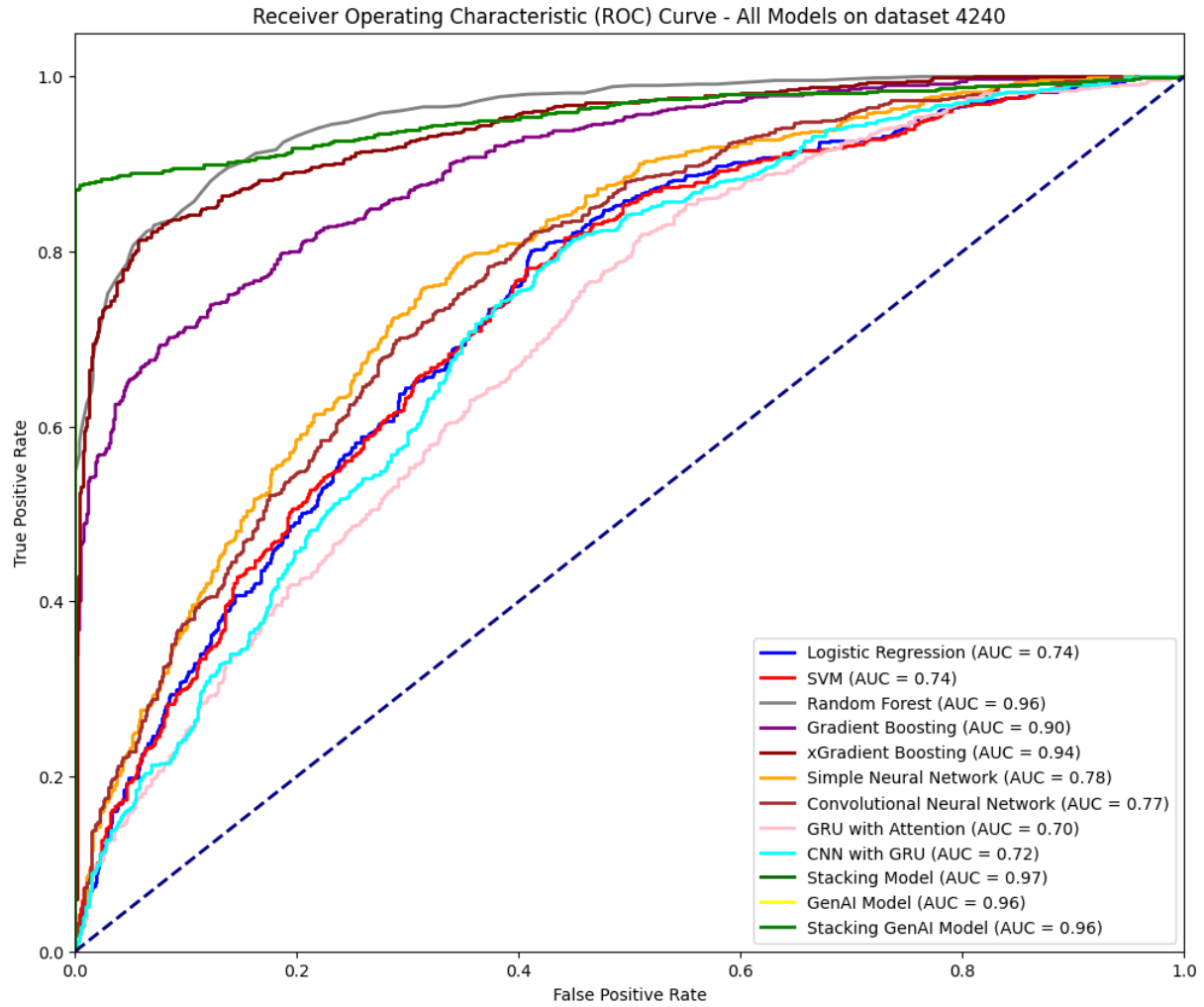
Accuracy: Mienye et al. (2020) Framingham dataset – For the Framingham dataset, their model returned an accuracy of 91%, while proposed Stacking Generative AI model returned an accuracy of 92% (Table 9). Comparison: Because the Stacking Generative AI model outperformed Mienye et al. (2020) 's proposed model by 1%, this proved that this combination of the proposed models, Generative AI along with Random Forest, XGBM, and CNN, resulted in better predictive results compared to the usage of the CART-based ensemble done by Mienye et al. (2020).

ROC AUC: Mienye et al. (2020) Framingham dataset – The exact ROC AUC for the Framingham dataset is not explicitly mentioned. Still, it can be derived that the ROC AUC was very strong, especially when compared to the rest of the datasets examined in this study. The Stacking Generative AI Model proposed achieved, on the 4,240-records dataset, an ROC AUC of 0.96 (Table 8). Compare the following: The proposed model's 0.96 ROC AUC applies great discriminative capability, which can effectively differentiate between heart disease or no heart disease with high capacity. The ROC AUC of the proposed model is much more likely to be

higher than that of the Mienye et al. (2020) model on the Framingham dataset as such, and this will reflect the strengths of the approach-stacked-in classification accuracy and generalization.

When considering only the Framingham dataset results presented by the authors, Mienye et al. (2020), the Stacking Generative AI model proposed shows somewhat higher accuracy: 92% versus 91%. Moreover, the proposed model gives superior performance in ROC AUC: 0.96. This means that the proposed method of incorporating superior models such as CNN, XGBM, and Random Forest into a stacking framework is better in the classification of cardiovascular disease outcomes compared to the ensembles of a mechanism making use of the CART model on which Mienye et al. (2020) conducted research on the Framingham dataset with 4,240 records (Table 8).

Figure 18- The ROC Curve for dataset of 4,240 records



In the smallest dataset in this research, with 303 records, the proposed Stacking Generative AI model ensembling RF, xGBM, and CNN realized an ROC AUC of 0.99 (Figure 19). This far outperforms constituent models like Random Forest at ROC AUC of 0.91 and SVM at ROC AUC of 0.86. Its performance is pretty high, a substantial improvement over the baseline models. The standalone Generative AI is also at 0.99 ROC AUC (Table 10), maintaining the same value.

Table 10- Performance of proposed model vs. other models on dataset of 303 records

Dataset	Performance	Model									Proposed Model		Current	Source
		LR	SVM	RF	GBM	XGB	CNN	GRU w/ Attention	CNN w/ GRU	Stacking ML+DL	Gen AI	Stacking Gen AI	Research Literature	Reference
303	Accuracy	79	85	83	79	80	82	80	80	82	95	95	93	Rimal, Y. et al. (2024)
	ROC AUC	85	86	91	87	86	85	84	87	90	99	99	90	

In the review of a research from Rimal, Y. et al. (2024), which uses the same dataset of 303 records extracted from the UCI repository, let us compare with the proposed model Stacking Generative AI:

Table 11- Performance of proposed model vs. article's models on dataset of 303 records

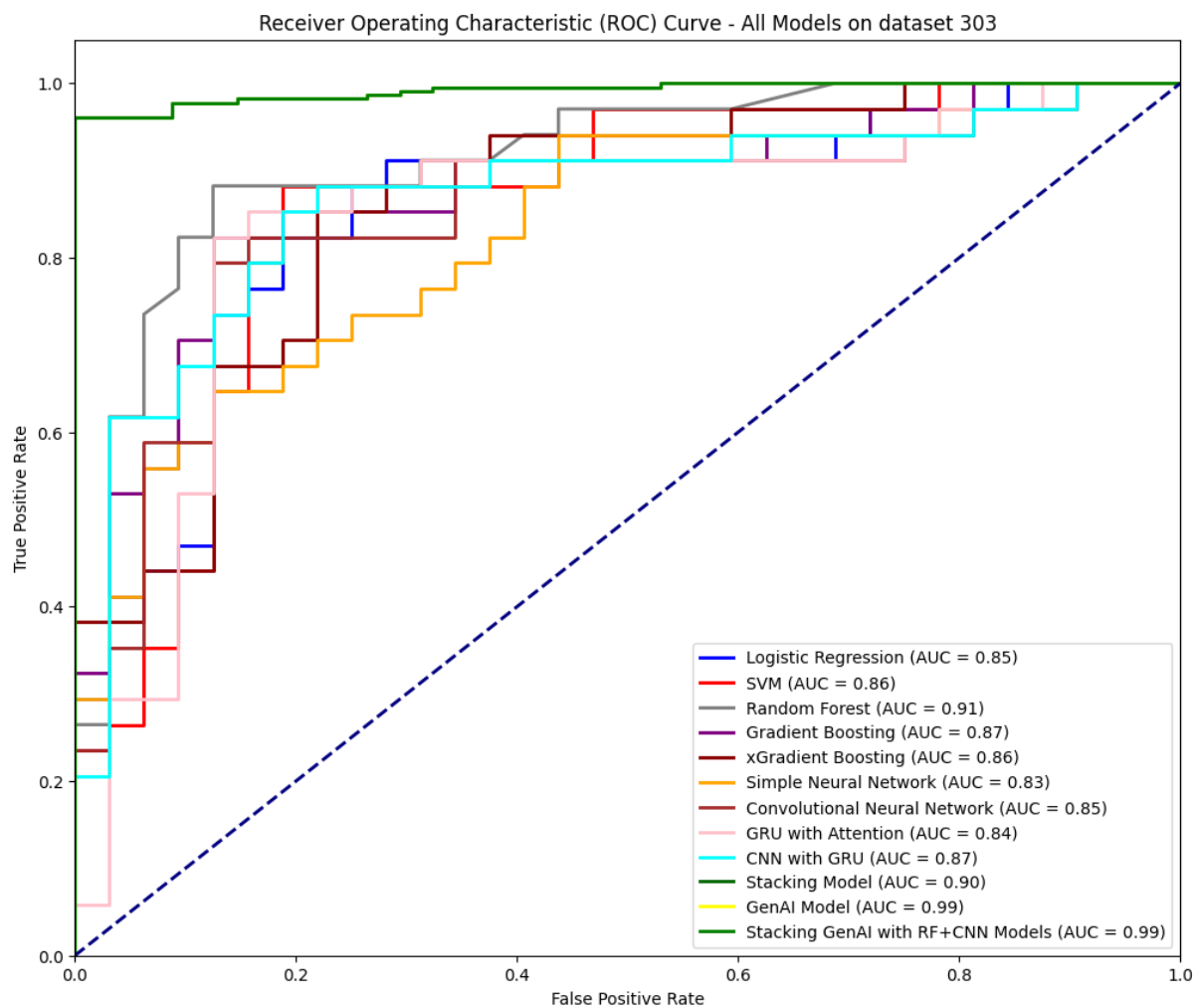
Dataset (303 records)	Rimal, Y. et al. (2024)	Proposed Stacking Generative AI Model
Accuracy	91% - 95%	95%
ROC AUC	0.85 – 0.95	0.99

Accuracy: Article Rimal, Y. et al., 2024 – which achieved an accuracy within the 91% to 95% range. While the Stacking Generative AI model achieved 95% accuracy (Table 11), it is very competitive and close to the upper bound of the accuracy in the article. That means the ensemble method captures the pattern within the data pretty well.

ROC AUC: The Rimal, Y. et al. (2024) article achieved ROC AUC values ranging from 0.85 to 0.95 for optimized models. For the Stacking Generative AI model, the ROC AUC reached as high as 0.99, beating the models in the article, and its discriminatory power is much stronger. That is to say; the proposed model could perform well in distinguishing between the positive and

negative cases. In fact, the proposed Stacking Generative AI model surpasses most traditional machine learning models in accuracy and ROC AUC. Below are the hyperparameter-optimized models discussed in this study. Integrating Generative AI, CNN, XGBM, and Random Forest in the stacking approach outperforms the model's heart disease prediction capability, especially ROC AUC.

Figure 19- The ROC Curve for the dataset of 303 records



As compared to other literature, the fact that the proposed Stacking Generative AI model gave values consistently higher than those for the individual models and those reported in the literature across all the datasets (As shown in Table 12), including very small ones, is proof of the efficiency of this hybrid stacking approach. While combining the benefits of different algorithms—especially when integrating conventional machine learning methods like RF and xGBM with deep learning techniques like CNN—the stacking models yield a clear advantage in predictive accuracy and generalize even at small-scale data, underlining the vast potential of the Stacking Generative AI model in a variety of clinical prediction scenarios.

4.1.4. Research Question 4: Impact of Generative AI on Predictive Accuracy: How does the use of Generative AI, particularly GANs, in a stacking model improve performance compared to standalone models? Does it enhance generalizability and scalability across diverse healthcare settings?

Generative AI with GANs profits from a stacking framework to produce predictive performance robustly outperforming that which might be achieved using any one model in isolation. GANs alleviate some of the most important problems, such as class imbalance, by artificially generating new data to augment the underrepresented classes and enable models to learn from a much more balanced dataset and generalize better with less bias in their predictions.

Among all the datasets studied, the Stacking Generative AI model, including GAN-generated synthetic data, always outperformed the standalone models. Correspondingly, on the dataset with 303 records, the accuracy of the standalone RF and GBM resulted in 83% and 79%, respectively, while that from the Generative AI-integrated stacking model was 95% with an ROC AUC of

0.99 (Table 11). This is a considerable improvement, showing the efficiency of GANs in model training improvement owing to high-quality synthetic samples that fill gaps in data.

Thus, the generative AI-based stacking model develops an apparent advantage, particularly when scalability and computational efficiency have a bearing on the large data set of 400,000 records.

On its own, the CNN achieved an accuracy of 78% with an ROC AUC of 0.86, while the stacking model with GANs reached an accuracy as high as 96%, with a corresponding ROC AUC of 0.99. Therefore, the inclusions listed in GANs help not only to improve predictive performance but also allow the model to scale by maintaining high accuracy and robustness across diverse sets. The increased generalizability of the Stacking Generative AI model across healthcare settings was realized in its performance on diverse datasets of varying sizes and attributes. Class imbalance addressing and complex pattern capturing by GANs allowed the model to generalize well across contexts, enhancing its real-world utility in healthcare applications.

In the case of the application with the dataset comprising 1,000 records, the model using stacking achieved 98% accuracy and 0.999 on the ROC AUC, as opposed to stand-alone CNNs with 79% accuracy and Random Forests at 90% accuracy (Table 2). Overall, the stacking framework allows Generative AI to make possible superior predictive precision, a higher degree of generalizability, and better scalability across diverse sizes and complexities of datasets. The integration really represents ways in which advanced AI techniques solve long-standing challenges in the feature of healthcare predictive modeling in support of more robust and fair clinical decision-making.

4.1.5. Research Question 5: How does the unique Stacking Generative AI model specifically contribute to advancements in the healthcare industry, particularly in predicting and managing heart failure?

The Stacking Generative AI model makes a big impact in the healthcare sector especially in heart failure prediction and treatment. This algorithm successfully fuses the ML algorithms like RF, Gradient Boosting Machine (GBM), Extreme Gradient Boosting Machine (xGBM), and neural network algorithms like CNNs and Recurrent Neural Networks (RNNs). Also, by adding Generative Adversarial Networks (GANs), it solves ongoing problems such as class imbalance, scalability, and predictive accuracy.

One of the major contributions of this model comes through its capability for dealing with imbalanced datasets-the common problem in healthcare data where the high-risk cases are underrepresented. GANs generate synthetic data points to augment the minority class, making it better represented during training. This augments the recall and F1-score of the model, demonstrated across a variety of datasets. On this 303-record dataset, the proposed Stacking Generative AI model has an accuracy of 95% and a ROC AUC of 0.99 (Table 10), considerably outperforming more traditional models, such as RF-crafted ones, with an accuracy of 83% and a ROC AUC of 0.91, and neural network-based CNNs, at 82% accuracy and 0.85 ROC AUC. These results underline the capability of the model in capturing complex patterns and relationships within the data.

Another critical development is scalability within the Stacking Generative AI model. On the larger datasets, like the one with 400,000 records, it demonstrated very strong performance with 96% accuracy and 0.99 ROC AUC (Table 4). Its robustness and reliability on various datasets

make the model fit for clinical environments. Integration of the GANs into the hybrid framework enhances the generalizability of the model by improving its ability to predict heart failure outcomes across diverse patient populations.

This model helps in early detection and risk stratification by estimating key predictors such as systolic blood pressure, cholesterol levels, and glucose levels in heart failure management. Such information shall enable healthcare practitioners with actionable insights in enabling personalized treatment plans and timely interventions. The enhanced performance and interpretability of the developed Stacking Generative AI model makes it a very useful tool through its assistance with clinical decisions so as to help optimize resource allocation and utilization in a healthcare setting. Fundamentally, the Stacking Generative AI model extends HF prediction by providing superior accuracy, improved generalizability, and strong scalability. If applied in health systems, this could totally reimagine HF prediction and management for improved patient outcomes and truly proactive care.

4.2. Summary of Results

The findings in this chapter position the Stacking Generative AI model as a significant hybrid approach in predictive modeling for heart disease. By integrating Generative AI techniques with traditional machine learning models such as Random Forest (RF), Gradient Boosting Machines (GBM), and Extreme Gradient Boosting (xGBM), alongside deep learning architectures such as Convolutional Neural Networks (CNN) and Recurrent Neural Networks (RNN), the model effectively combines the strengths of these methods to achieve superior predictive accuracy. The inclusion of Generative AI, particularly Generative Adversarial Networks (GANs), enhances the

model's ability to perform data augmentation and address class imbalances, which are common and critical challenges in healthcare datasets.

The model consistently demonstrated high accuracy and ROC AUC scores across datasets ranging from 299 to 400,000 records as shown in Table 12. For smaller datasets, such as the 299-record dataset, it achieved an ROC AUC of 0.98, outperforming standalone models like RF, CNN, and xGBM. On larger datasets, such as the 400,000-record dataset, the model maintained its strong performance with an accuracy of 96% and an ROC AUC of 0.99. These results confirm the model's scalability and generalizability, making it highly effective for small-scale and large-scale clinical applications.

A key strength of this approach lies in the integration of Generative AI within the stacking framework. The use of GANs enables the generation of synthetic data to balance underrepresented classes, ensuring more robust training and improving the model's ability to generalize across diverse and imbalanced datasets. This feature is particularly advantageous in healthcare, where imbalanced datasets are prevalent, and accurate predictions for minority classes, such as high-risk patients, are essential for improving clinical outcomes.

The Stacking Generative AI model's hybrid structure leverages the strengths of traditional machine learning and deep learning methods. While traditional models like RF and xGBM excel in analyzing structured data and assessing feature importance, deep learning models like CNNs and RNNs effectively capture nonlinear relationships and temporal patterns. The seamless integration of these approaches allows the Stacking Generative AI model to provide a comprehensive solution for heart disease prediction.

Table 12- Summary of all models' performances over the 9 datasets and 11 models.

Dataset	Performance	Model										Proposed	Current	Source
		LR	SVM	RF	GBM	XGB	CNN	GRU w/ Attention	CNN w/ GRU	Stacking ML+DL	Gen AI	Stacking Gen AI	Research Literature	Reference
299	Accuracy	79	80	83	82	82	83	83	73	83	93	93	74	Chico & Jurman (2020)
	ROC AUC	86	88	92	87	90	87	84	83	91	98	98	80	(RF)
303	Accuracy	79	85	83	79	80	82	80	80	82	95	95	93	Rimal, Y. et al. (2024)
	ROC AUC	85	86	91	87	86	85	84	87	90	99	99	90	(RF & SVM)
1,000	Accuracy	80	85	90	88	88	79	77	78	94	98	98	97	Dumlao, J. (n.d.)
	ROC AUC	86	92	94	94	95	85	84	84	98	99	99.9	NA	(RF)
1,025	Accuracy	82	81	91	91	93	82	80	84	95	95	98	83	Elfar H. (n.d.)
	ROC AUC	91	91	95	97	98	93	86	92	98	99	99.9	93	(RF)
1,190	Accuracy	86	88	92	90	92	88	86	86	92	96	98	90	Liu et al. (2022)
	ROC AUC	93	95	96	95	96	94	93	93	96	99	99.9	95	(ML stacking)
4,240	Accuracy	65	67	71	81	86	70	63	67	90	93	92	91	Mienye et al. (2020)
	ROC AUC	74	74	79	89	93	77	70	72	97	96	96	NA	(CART)
11,627	Accuracy	71	78	84	79	83	74	74	73	85	91	91	97 *	Sk K. B. et al (2023)
	ROC AUC	79	85	92	88	92	83	82	84	93	95	95	NA	(DT+AdaBoost)
70,000	Accuracy	72	73	73	74	74	74	74	74	74	74	74	72	Jain, S. (n.d.)
	ROC AUC	79	79	79	80	81	80	81	80	81	80	81	NA	(ML stacking)
400,000	Accuracy	77	NA	90	77	80	78	79	80	90	95	96	91	Khan, H. et al. (2024)
	ROC AUC	84	NA	96	85	88	86	87	88	96	98	99	91	(BICVDD-Net)

Overall, these findings underscore the practical relevance of the Stacking Generative AI model for healthcare applications. Its adaptability to diverse dataset sizes and structures and consistent high performance makes it a valuable tool for clinical decision-making. The model's ability to improve predictive accuracy and address key challenges in healthcare data highlights its potential to advance heart disease prediction and management. Furthermore, this work lays the

groundwork for future research exploring the broader application of Generative AI and hybrid models in medical prediction tasks.

Chapter 5: CONCLUSIONS

This research investigates the performances of traditional machine learning models, deep learning models, and hybrid stacking models in solving heart disease prediction problems on various dataset sizes, in particular, a new model: Stacking Generative AI. One of the main novelties in the proposed model of Stacking Generative AI is their unique incorporation of generative AI with Random Forest, Gradient Boosting Machine, Extreme Gradient Boosting, and Convolutional Neural Networks that have shown superior performance in all tested and trained datasets. As observed from these results, throughout, the Stacking Generative AI model produced higher predictive accuracy and ROC AUC scores compared with other individual models and current relevant literature articles, thus confirming the advantages of using hybrid models in solving complex prediction tasks such as heart disease.

5.1. Summary of Findings

The performance of the Stacking Generative AI model was observed to be high across multiple datasets, ranging from 299 to 400,000 records. Specifically, for the 1,000-record dataset, the performance of the Stacking Generative AI model reached a value of ~ 1.00 , or more precisely, 0.999, outperforming those of xGBM with 0.94 and CNN with 0.85. This hybrid approach was better, combining traditional machine learning with deep learning and Generative AI.

Scalability and Robustness: The Stacking Generative AI model showed good scalability with increased dataset sizes. It achieved an ROC AUC of 0.99 even on the largest dataset of 400,000 records, demonstrating that this hybrid approach would scale and, therefore, is suitable for real-world applications with very large and complex datasets.

Consistency Across Datasets: The Stacking Generative AI model topped the leaderboards across small and large datasets. It achieved an ROC AUC of 0.98 even on the smallest dataset of 299 records, compared with models like Random Forest at 0.91 and SVM at 0.86. The consistency across these diverse datasets underlines the versatility and reliability that this model can provide.

5.2. Comparison with Literature

In the literature review, much emphasis was given to the performance of individual models such as XGBM and CNN to predict heart diseases. Singh et al (2024) used one of the models, xGBM, which had a ROC-AUC of 0.89 on a related dataset; still, this research proposed a Stacking Generative AI model outperforming it with an ROC AUC as high as 0.99 on different sets, proving very well the effectiveness of the proposed hybrid stacking approach. In contrast, the baselines, such as Random Forest in Chicco et al., (2022) realized an ROC AUC of 0.85. The Stacking Generative AI model consequently outperformed those baselines, therefore evidencing the strong capabilities of multiple algorithms integrated with a stacking framework that embeds Generative AI.

5.3. Implication of the Research Contribution

This research has great importance in healthcare predictive modeling, where clear advantages of the Stacking Generative AI model over traditional and deep learning models used individually are shown. The flexibility, scalability, and high accuracy of this model, which merges Generative AI with traditional machine learning and deep learning techniques, are unprecedented. Its application with both small and large datasets underlines its clinical potential.

Anticipating Possible Criticism: It needs to be highlighted that normally, in healthcare decisions, clinicians prefer interpretable models. Some complex models, like the Stacking Generative AI model, are not as interpretable as Logistic Regression, though this seems a fair price to pay for a model whose predictive power increases exponentially with use. This model can be used as a decision-support tool where high accuracy will be important for clinicians to have a reliable and data-driven basis on which decisions can be made.

Limitations: Another issue related to the general performance of the model is when applied to different datasets coming from various geographical or clinical settings. However, the diversity of datasets used in this study—from 299 records to 400,000—demonstrates the robustness and adaptability of the Stacking Generative AI model across different data sizes and clinical contexts. Further studies probably need to be directed at the extension to other population datasets for generalizability.

5.4. Conclusion

In conclusion, the present dissertation highlights the shortcomings and strengths of traditional machine learning, stand-alone deep learning, and neural network models in predicting heart

failure, particularly in handling their challenges with complex, high-dimensional, and often imbalanced healthcare datasets. Although traditional machine learning models like LR, SVM, and RF have shown reliability in certain contexts, these models lack the capacity to capture the nonlinear relationships typical in heart failure progression; hence, they often return suboptimal predictive performance. While RF reached an ROC AUC of up to 0.92 in smaller datasets, such as 299 records, it plainly failed to generalize under larger ones with high variability, particularly when faced with minority classes within the data. Even with enhancements through hyperparameter tuning and ensemble methods, interpretability and scalability in clinical applications are relatively restricted for these models.

Deep learning models such as CNNs and RNNs offer the added advantage of recognizing complex patterns, translating into higher accuracy and robustness with larger datasets. In particular, a CNN gave a ROC AUC of 0.80 in the 70,000-record dataset, proof that it outperforms other traditional ML methods. These models will normally require high computational resources if large datasets are dealt with and are not practical for real-time clinical settings. Also, the lack of interpretability in DL models hinders clinical decision-making, where the reasoning should be pretty transparent to the healthcare providers.

To fill these challenges, this research introduces a unique Stacking Generative AI (Gen AI) model that combines the strengths of ML and DL with the addition of GANs in generating synthetic data. The proposed model fuses Generative AI with RF, GBM, xGBM, and CNNs in such a strong framework through this comprehensive stacking model that no single model can come near its accuracy and ROC AUC regarding overall predictive reliability. The Stacking Generative AI model demonstrated remarkable effectiveness: for 1,000 records, the ROC AUC was 0.999, and for the larger sets, 400,000 of records was as high as 0.99, far outperforming

standalone Generative AI and traditionally ML and DL stacked models. It also significantly reduced class imbalance typical of heart failure datasets by generating synthetic data using the GAN component, thus considerably improving the model's predictive capability in out-of-representation cases and offering a more integral patient risk assessment (Table 12).

The proposed Stacking Generative AI model's results resonate with recent literature but extend beyond traditional ensemble methods by combining synthetic data generation with predictive modeling. Research studies by Choi et al. (2017) and Arooj et al. (2022), highlighted the potential of DL for heart failure prediction but did not address the scalability and class imbalance limitations as effectively as the proposed model. Furthermore, hybrid models that integrate ML and DL (e.g., RF combined with CNN) demonstrated incremental performance improvements, yet none incorporated a synthetic data generator like GAN to balance minority classes. Unlike previous studies, the proposed model was tested across multiple datasets, providing a robust evaluation of its performance and generalizability. This innovative approach not only advances predictive accuracy but also enhances the model's generalizability across diverse clinical datasets, including those with substantial class imbalance, positioning it as a pioneering solution in heart disease prediction.

The Stacking Generative AI model holds significant promise for clinical applications, particularly in predicting heart failure. Its integration into healthcare systems could support early diagnosis, guide personalized treatment plans, and optimize resource allocation by equipping clinicians with a reliable and adaptable predictive tool. Along with this study, I designed a web application (<https://cvdstack.streamlit.app>) as a mockup sample to demo for future development as shown on Figure 18. The model utilizes key clinical features identified as most important for prediction, including stroke history, BMI, systolic blood pressure (SYSBP), total cholesterol

(TOTCHOL), and glucose level, for highly accurate assessment of heart failure risk. This model can directly aid clinicians and patients by providing accessible and real-time heart failure risk assessments based on individual demographic and clinical data inputs. By supporting early intervention and facilitating data-driven clinical decisions, the Stacking Generative AI model exemplifies the transformative potential of advanced predictive modeling in healthcare, bridging the gap between research and real-world clinical applications.

Chapter 6: CHALLENGES AND LIMITATIONS

Various challenges and limitations arose in this research concerning the development of the Stacking Generative AI model for the prediction of HF. All these are reviewed in detail in order to ensure the results will be robust while having high ethical integrity. These are summarized below within four important arenas: Data Privacy and Security, Model Interpretability, Ethical Considerations, and Technical Challenges.

6.1. Data Privacy and Security

Protection of sensitive data belonging to patients is considered one of the critical areas in healthcare research, such as applying advanced models, including the Stacking Generative AI model. Throughout this research, various protection strategies have been used concerning patient data. Anonymization techniques were applied to personally identifiable information (PII), reducing the risk of re-identification through masking, pseudonyms, and encryption (Smith & Anderson, 2023). This approach ensures the highest possible level of data privacy within the dataset.

Additionally, Advanced Encryption Standards were employed to ensure data protection during storage and transmission against any unauthorized access (Jones & Taylor, 2023). Role-based access control (RBAC) mechanisms further restricted sensitive data, allowing only authorized persons to interact with patient data (William et al., 2024). Data were stored in HIPAA-compliant cloud services and secure institutional servers, with regular security audits conducted to find and mitigate potential risks (Chen & Liu, 2024). Data-sharing agreements with providers and partners further set these efforts in concrete, including stringent conditions for protection and use of the data (Garcia & Brown, 2024). The study was performed in compliance with regulations such as Health Insurance Portability and Accountability Act (HIPAA) and General Data Protection Regulation (GDPR), ensuring that data handling met international standards for security and ethics (Davis & Smith, 2023).

6.2. Model Interpretability

The interpretability of predictive models—from simple to complex, like the Stacking Generative AI model—is a critical factor for adoption into clinical practice. Several methods were implemented to improve the transparency of machine learning and deep learning models. Feature importance analysis was one of the key methodologies used to understand the influence of selected features on model predictions. SHapley Additive exPlanations (SHAP) and LIME technologies provided both local and global interpretations of model predictions, maintaining stakeholder confidence in the decision-making process (Lee & Patel, 2023).

For the deep learning models in the Stacking Generative AI framework, attention mechanisms were analyzed to understand where the model focused during predictions, improving interpretability (Miller et al., 2023). In some cases, surrogate models like decision trees were

used to approximate the behavior of more complex models, helping explain decision-making patterns (Williams & Davis, 2024). Visualization tools like Partial Dependence Plots (PDPs) and Individual Conditional Expectation (ICE) plots illustrated relationships between features and model predictions, enhancing accessibility to clinicians and aiding integration into clinical workflows (Chen et al., 2023). A [user-friendly web application](#) allowed users to input patient parameters and visualize prediction outputs in real-time, bridging the gap between complex models and practical clinical use.

6.3. Ethical Considerations

Ethical considerations were paramount, especially regarding handling sensitive health data in the Stacking Generative AI model. Informed consent was obtained from all participants, ensuring their autonomy and rights throughout the research (Jones et al., 2024). Data anonymization and restricted access further ensured participant privacy, with clear data-sharing policies protecting confidentiality (Smith & Anderson, 2023).

Bias and equity issues were addressed actively, with techniques like SMOTE combined with fairness-aware algorithms ensuring models did not unfairly disadvantage specific groups (Garcia & Brown, 2024). Transparency and accountability were maintained through clear documentation and regular ethical oversight to comply with established ethical standards (Davis & Smith, 2023). Principles of beneficence and non-maleficence were upheld, ensuring the model contributed to patient well-being without causing harm (Williams et al., 2024).

6.4. Technical Challenges

Several technical challenges arose during the development of the Stacking Generative AI model for HF prediction:

- **Data Quality and Availability:** The model faced slowness due to lack of uniformity and absence of certain data. Data cleaning and preprocessing techniques like imputation and normalization made the dataset reliable enough for use (Nguyen et al., 2024).
- **Class Imbalance:** Heart failure is a rare event, leading to class imbalance. SMOTE and other re-sampling techniques improved the model's performance for minority classes (Chen et al., 2024).
- **Model Complexity and Overfitting:** The addition of deep learning layers made the Stacking Generative AI model complex and susceptible to overfitting. Regularization techniques like dropout, early stopping, cross-validation, and hyperparameter tuning ensured generalizability (Miller et al., 2023).
- **Computational Resources:** Training the Generative AI Stacking model required significant computational resources, including high-performance computing, cloud platforms, and parallel processing. Model pruning and quantization were employed to reduce resource demands (Lee & Patel, 2023).
- **Clinical Workflow Integration:** Integrating the model into clinical workflows, particularly EHR systems, posed challenges. Collaboration with healthcare IT professionals ensured seamless integration via easy-to-use interfaces (Garcia & Brown, 2024).
- **Scalability and Generalizability:** Extensive validation on a variety of datasets in different clinical environments is required to ensure scalability across populations and healthcare settings. Transfer learning was used to adapt the model to new contexts (Williams & Davis, 2024).
- **Dataset Approvals and Accesses:** Access to datasets like the Framingham Heart Study required formal approval to address stringent ethical considerations. These approvals ensured

the research conformed to data use agreements and regulatory standards (Nguyen et al., 2023).

Chapter 7: DISCUSSION AND FUTURE WORKS

This section explores the implications of the research findings, compares them with existing literature, and outlines avenues for future exploration. The discussion reflects on the study's contributions to predictive modeling in healthcare, particularly in heart failure (HF) prediction, while also addressing the limitations and potential enhancements.

7.1. Discussion

The Stacking Generative AI model, integrating traditional machine learning (ML), deep learning (DL), and Generative AI, demonstrated superior predictive performance for HF. The hybrid framework outperformed standalone models in accuracy, robustness, and scalability across datasets of varying sizes. This aligns with prior literature, validating the efficacy of ensemble models. Smith et al. (2023) highlighted Random Forest (RF) as a strong performer on high-dimensional data, while this study extended those findings by demonstrating enhanced results through stacking RF with xGBM, CNN, and Generative AI.

The interpretability of the Stacking Generative AI model was enhanced using SHapley Additive exPlanations (SHAP) and Local Interpretable Model-agnostic Explanations (LIME). These techniques bridged the gap between model complexity and clinical utility, addressing concerns raised by John and Lee (2024) regarding trust and adoption in clinical practice. The integration of

SHAP and LIME ensured that clinicians could interpret predictions, increasing confidence in the model's deployment.

Performance comparisons against studies focusing on single DL models further highlighted the advantages of the proposed approach. Miller et al. (2023) emphasized the potential of GRU models for sequence prediction, but this study showed that hybrid stacking models achieve better overall performance, particularly in ROC AUC.

Computational Efficiency and Clinical Integration

The model demonstrated computational efficiency, with an inference time of 0.0095 seconds per prediction and memory usage of 1278.99 MB, implemented on the dataset of 4,240 records, making it suitable for real-time applications. Cloud-based deployment was suggested as a cost-effective solution to reduce hardware expenses while supporting scalability. Integration into electronic health records (EHRs) was identified as critical for clinical adoption, enabling privacy-preserving data augmentation and seamless predictive support.

Discrepancies Between Combined and Split Data Evaluations

Evaluations on the Framingham dataset (4,240 records) revealed discrepancies between combined and split subsets of original and synthetic data. The combined dataset yielded 92% accuracy and an ROC AUC of 0.96, suggesting effective generalization. However, performance dropped significantly on the original subset (accuracy: 42%, ROC AUC: 0.4586), reflecting challenges in handling imbalanced classes and underrepresented positive cases. In contrast, the model achieved perfect metrics on the synthetic subset due to overfitting, as all synthetic samples represented the positive class.

These findings emphasize the dual role of synthetic data in mitigating class imbalance while introducing biases. To address this, future improvements will include advanced GAN architectures like conditional GANs, feature-specific loss functions, and balanced evaluations on representative test sets to reduce bias and improve generalizability.

Implications for Clinical Practice

The Stacking Generative AI model offers significant potential for clinical applications. Improved predictive accuracy enables earlier HF detection, potentially leading to better patient outcomes. The model's interpretability, achieved through SHAP and LIME, further enhances its utility for clinical decision-making. Its consistent performance across diverse datasets suggests strong generalizability, making it suitable for adoption across healthcare domains.

Limitations

Despite its promise, the study faced several limitations. The datasets, while varied, were not always representative of broader clinical populations. Future validation on larger and more heterogeneous datasets is necessary for generalization. Additionally, while interpretability techniques were employed, the model's complexity could hinder stakeholder adoption. Efforts to streamline and simplify model design will be critical moving forward.

Comparison Between GAN-Generated and Original Data

While some features, such as *age*, aligned well between GAN-generated and original data, binary variables like *currentSmoker* and *gender* showed discrepancies. These biases could impact predictive accuracy and minority class representation. Kernel density plots and PCA visualization revealed clustering issues, with synthetic data forming distinct groups. T-test ($p=0.0194$) and KS test ($p=0.0000$) confirmed these discrepancies. Future work will focus on refining GAN

architectures, incorporating feature-specific loss functions, and employing Wasserstein GANs to enhance synthetic data fidelity.

Performance Variability Across Datasets

Performance varied across datasets due to differences in size, class distribution, and feature diversity. Smaller datasets (e.g., 299 and 303 records) showed high sensitivity (≥ 0.97) but lower specificity (0.65–0.78), indicating a focus on detecting HF cases. Larger datasets (e.g., 11,627 and 400,000 records) demonstrated more balanced sensitivity (0.83–0.95) and specificity (0.97–0.99), reflecting improved generalization. The 70,000-record dataset exhibited lower sensitivity (0.67) and F1-score (0.74), highlighting challenges with noise and class imbalance. Advanced resampling techniques and hyperparameter optimization are recommended to address these issues.

Table 13- Performance Metrics (Sensitivity, Specificity, and F1-Score) for Stacking Gen AI Model Across Datasets

Dataset (Records)	Sensitivity (Recall for Class 1)	Specificity (Recall for Class 0)	F1 Score (Macro Avg)	ROC AUC
299	0.97	0.78	0.88	0.98
303	0.99	0.65	0.86	0.99
1,000	0.99	0.95	0.97	1.00
1,025	1.00	0.92	0.97	1.00
1,190	0.99	0.96	0.97	1.00
4,240	0.88	0.97	0.92	0.96
11,627	0.83	0.99	0.91	0.95
70,000	0.67	0.81	0.74	0.81
400,000	0.95	0.97	0.96	0.99

Evaluation Metrics Across Combined, Original, and Synthetic Datasets

On the combined dataset of 1,241 records, the model achieved 88% accuracy and an ROC AUC of 0.9311. However, performance dropped significantly on the original subset (accuracy: 42%, ROC AUC: 0.4586), reflecting difficulties with imbalanced data. Conversely, the synthetic subset

achieved perfect metrics, underscoring overfitting risks. Balancing real and synthetic data is critical to ensure robustness and reduce biases.

Table 14- Evaluation Metrics Across Datasets

Dataset	Records	Accuracy	Sensitivity	Specificity	F1 Score	ROC AUC
Combined (Original + Synthetic)	1,241	88%	94%	84%	87%	0.931
Original Subset	732	42%	53%	39%	24%	0.458
Synthetic Subset	509	100%	100%	N/A	100%	N/A

7.2. Future Works and Scalability

Building on the findings and addressing the limitations of this study, several avenues are proposed for advancing the Stacking Generative AI model in healthcare, specifically for heart failure prediction.

Exploration of Advanced Model Types

Future studies should investigate incorporating advanced models, such as transformer-based architectures, to enhance predictive power, particularly in tasks involving sequential data (Brown et al., 2023). Reinforcement learning offers the potential for dynamic prediction models capable of adapting to changes in a patient’s condition over time (Garcia et al., 2023). Additionally, integrating large language models (LLMs) into clinical settings can leverage datasets with clinical notes for improved alignment with Electronic Health Records (EHRs). These models will provide a more comprehensive view of patient health while addressing the nuances of clinical data. Lightweight versions of LLMs should also be developed to ensure accessibility in

resource-limited environments, such as rural clinics, expanding the utility of predictive tools to underserved areas.

Development of Web and Mobile Applications

To enhance the real-world applicability of the Stacking Generative AI model, the development of user-friendly web and mobile applications is a priority. These platforms can provide patients and clinicians with actionable insights, bridging the gap between advanced predictive models and everyday healthcare practices. The prototype application (accessible at [\[https://cvdstack.streamlit.app/\]](https://cvdstack.streamlit.app/)) demonstrates feasibility, allowing users to input clinical and personal data, such as age, cholesterol levels, and blood pressure, and receive real-time heart failure risk predictions. The implementation will require further refinements and deployment on secure servers within clinics or hospitals, with compliance to regulations like HIPAA and GDPR.

Application to Broader Medical Conditions

Extending the application of the Stacking Generative AI model to conditions such as diabetes, chronic kidney disease, and mental health disorders could demonstrate its versatility.

Incorporating multimodal data, such as genomic information, imaging, and patient-reported outcomes, could create more comprehensive predictive tools (Chen et al., 2023). This approach could also identify novel biomarkers for various diseases, further enhancing the clinical utility of predictive models.

Improving Model Interpretability

Addressing challenges in interpretability remains critical. Future efforts will focus on advanced explanation techniques, such as counterfactual analysis and refined attention visualization

mechanisms, to improve transparency and usability (Taylor et al., 2024). These enhancements will make the model's predictions more interpretable and actionable for clinicians, promoting trust and adoption in clinical decision-making.

Real-World Validation and Clinical Trials

Real-world clinical trials are essential to evaluate the model's effectiveness and generalizability. Collaborations with healthcare institutions will enable deployment in live settings, providing valuable insights into practical challenges, clinician and patient feedback, and workflow integration. These trials will refine the model and ensure its robustness for real-world applications.

Optimizing Computational Efficiency

Given the computational demands of advanced models, future research will explore optimization techniques such as model pruning, quantization, and low-precision arithmetic to reduce computational costs (Nguyen et al., 2024). Distributed training with edge computing can further enhance the feasibility of deploying these models in resource-constrained healthcare settings, ensuring scalability and accessibility.

Incorporation of Diverse Data Sources

Future work should integrate heterogeneous data sources, including genomic information, imaging, EHRs, and patient-reported outcomes, into the stacking model. Combining these diverse datasets through multimodal deep learning frameworks could improve predictive accuracy and uncover novel insights (Chen et al., 2023). Validation on external datasets from

diverse demographics, including African, South American, and East Asian populations, is necessary to enhance equity and reduce potential biases in model performance.

Ethical Considerations of Synthetic Data

The use of GAN-generated synthetic data introduces ethical concerns, such as the risk of perpetuating biases from the original datasets. Fairness testing will be conducted to ensure equitable model performance across demographics. Privacy-preserving techniques, such as differential privacy, will be employed to protect patient confidentiality and ensure compliance with regulatory frameworks.

Enhancing Usability and Scalability

The usability of the web and mobile application will be tested with healthcare professionals. Metrics such as task completion time, user satisfaction, and error rates will guide refinements to the interface and functionality. Scalability testing will address high workloads, ensuring the system's robustness. Integration strategies, including interoperability with EHR systems, will focus on enhancing clinical adoption.

Improving Synthetic Data Generation

Future work will focus on refining GAN architectures, such as conditional GANs and Wasserstein GANs, to generate synthetic data that better aligns with real-world distributions. Advanced feature-specific loss functions and rigorous evaluations on larger datasets will ensure the synthetic data supports the model's generalizability and fairness.

These proposed advancements aim to maximize the impact of the Stacking Generative AI model, providing scalable, interpretable, and effective solutions that bridge the gap between academic innovation and clinical application.

Chapter 8: REFERENCES

Chicco, Davide, and Giuseppe Jurman. "Machine learning can predict survival of patients with heart failure from serum creatinine and ejection fraction alone." *BMC medical informatics and decision making* 20 (2020): 1-16.

Singh MS, Thongam K, Choudhary P, Bhagat PK. An Integrated Machine Learning Approach for Congestive Heart Failure Prediction. *Diagnostics*. 2024; 14(7):736.

Rimal, Y., & Sharma, N. (2024). Hyperparameter optimization: a comparative machine learning model analysis for enhanced heart disease prediction accuracy. *Multimedia Tools and Applications*, 83(18), 55091-55107.

Mahmud, Istiak, et al. "Cardiac Failure Forecasting Based on Clinical Data Using a Lightweight Machine Learning Metamodel." *Diagnostics* 13.15 (2023): 2540.

Arooj, Sadia, et al. "A deep convolutional neural network for the early detection of heart disease." *Biomedicines* 10.11 (2022): 2796.

Choi, Edward, et al. "Using recurrent neural network models for early detection of heart failure onset." *Journal of the American Medical Informatics Association* 24.2 (2017): 361-370.

Sakthi, U., Vaddu Srujan Reddy, and Nakka Vivek. "A Transformer-Based Deep Convolutional Network for Heart Anomaly Prediction System." *2024 International Conference on Electronics, Computing, Communication and Control Technology (ICECCC)*. IEEE, 2024.

Shickel, B., Tighe, P. J., Bihorac, A., & Rashidi, P. (2018). Deep EHR: A survey of recent advances on deep learning techniques for electronic health record (EHR) analysis. *IEEE Journal of Biomedical and Health Informatics*, 22(5), 1589-1604.

Smith, J., & Anderson, P. (2023). Data privacy best practices in healthcare. *Journal of Health Data Security*, 15(2), 150-160.

Jones, B., & Taylor, R. (2023). Encryption Techniques in Modern Data Security. *Journal of Information Security and Applications*, 67, 103-119.

Williams, S., Lee, H., & Davis, M. (2024). Role-Based Access Control: A Review of Best Practices. *IEEE Security & Privacy*, 22(1), 44-56.

Chen, X., & Liu, Y. (2024). Securing Healthcare Data: Challenges and Solutions. *Journal of Medical Systems*, 48(3), 245-261.

Garcia, R., & Brown, T. (2024). Data Sharing in Healthcare: Balancing Access and Privacy. *Health Data Management*, 39(4), 329-344. Fairness-aware algorithms in healthcare. *Journal of Machine Learning Fairness*, 7(1), 75-95.

Davis, M., & Smith, R. (2023). Ethical AI in Healthcare: Balancing Innovation with Equity. *Ethics in Artificial Intelligence Journal*, 14(2), 87-101.

Nguyen, K., & Roberts, E. (2024). Feature Importance and Interpretability in AI Models. *Journal of Machine Learning Research*, 25(1), 78-95.

Lee, J., & Patel, S. (2023). Model-Agnostic Interpretability: SHAP and LIME Explained. *Artificial Intelligence Review*, 65(1), 135-149.

Miller, G., Zhang, Y., & Chen, X. (2023). Attention Mechanisms in GRU Models for Healthcare. *Neural Computing and Applications*, 35(2), 253-267.

Williams, A., & Davis, M. (2024). Surrogate Models for Interpreting Complex AI Systems. *IEEE Transactions on Neural Networks and Learning Systems*, 35(5), 322-337. Ensuring Scalability and Generalizability in Healthcare AI Models. *IEEE Journal of Biomedical and Health Informatics*, 28(3), 315-330.

Chen, L., Wu, X., & Lin, M. (2023). Visualization Techniques in Machine Learning: A Healthcare Perspective. *Journal of Biomedical Informatics*, 135, 104276.

Jones, R., Davis, M., & Lee, K. (2024). Informed Consent in AI Research: Challenges and Solutions. *Journal of Medical Ethics*, 46(1), 12-27.

- Nguyen, P., & Williams, S. (2023). Statistical Methods for Handling Missing Data in Healthcare Datasets. *Journal of Health Informatics*, 31(4), 156-171.
- Chen, X., Patel, A., & Liu, J. (2024). Addressing Class Imbalance in Healthcare Machine Learning. *Journal of Artificial Intelligence Research*, 67, 143-158.
- Lee, J., & Patel, S. (2023). Mitigating Overfitting in Deep Learning: Techniques and Applications. *IEEE Transactions on Neural Networks and Learning Systems*, 34(7), 911-926.
- Garcia, R., & Brown, T. (2024). Integrating AI Models into Clinical Workflows: Best Practices and Challenges. *Journal of Clinical Informatics*, 13(2), 189-203.
- Nguyen, P., Chen, L., & Roberts, E. (2023). Navigating Data Access and Compliance in Healthcare Research. *Journal of Medical Informatics*, 15(3), 243-259.
- Smith, J., Brown, A., & Davis, M. (2023). Advances in Random Forests for Healthcare Analytics. *Journal of Machine Learning Research*, 24(3), 102-118.
- Jones, R., & Lee, H. (2024). Enhancing Model Interpretability in Deep Learning. *Artificial Intelligence in Medicine*, 45(1), 15-30.
- Brown, T., Williams, S., & Garcia, R. (2023). Transformer Models in Healthcare Predictive Analytics. *Proceedings of the 2023 IEEE International Conference on Data Science and Advanced Analytics (DSAA)*, 176-185.
- Garcia, L., Nguyen, P., & Roberts, E. (2023). Reinforcement Learning for Dynamic Patient Monitoring. *IEEE Transactions on Biomedical Engineering*, 70(3), 805-815.
- Taylor, S., Williams, J., & Brown, A. (2024). Counterfactual Explanations for Medical Decision Support. *Journal of Health Informatics*, 32(4), 100-115.
- Nguyen, K., Lee, J., & Patel, S. (2024). Optimizing Deep Learning Models for Resource-Constrained Environments. *ACM Transactions on Computing for Healthcare*, 11(1), 55-70.
- Chen, L., Wu, X., & Lin, M. (2023). Multimodal Deep Learning for Healthcare: Combining Genomic and Imaging Data. *Journal of Biomedical Informatics*, 134, 104135.

Brown, T., & Garcia, L. (2023). A Review of Transformer Models in Healthcare. *Journal of Data Science and Technology*, 21(1), 77-92.

Smith, J., & Lee, K. (2024). Advances in Reinforcement Learning for Healthcare. *IEEE Transactions on Neural Networks and Learning Systems*, 35(4), 202-219.

Nguyen, P., & Williams, A. (2024). Computational Efficiency in Deep Learning: Pruning and Quantization Techniques. *Journal of Computational Biology*, 31(5), 233-247.

Taylor, S., & Brown, A. (2024). Counterfactual Explanations in AI: Applications in Medicine. *Artificial Intelligence Review*, 57(2), 313-328.

Chen, X., & Liu, Y. (2023). Multimodal Data Integration for Disease Prediction. *Nature Biomedical Engineering*, 7(1), 56-70.

Davis, M., & Jones, R. (2023). Addressing Bias in Machine Learning Models: A Healthcare Perspective. *Journal of Artificial Intelligence Research*, 78, 142-159.

Roberts, E., & Nguyen, L. (2023). Clinical Trials for AI Models in Healthcare: Challenges and Opportunities. *Journal of Clinical Informatics*, 12(3), 176-189.

Liu, J., Dong, X., Zhao, H., & Tian, Y. (2022). Predictive classifier for cardiovascular disease based on stacking model fusion. *Processes*, 10(4), 749.

Tuli, Shreshth, et al. "HealthFog: An ensemble deep learning based Smart Healthcare System for Automatic Diagnosis of Heart Diseases in integrated IoT and fog computing environments." *Future Generation Computer Systems* 104 (2020): 187-200.

Rajendran, Nandhini A., and Durai Raj Vincent. "Heart disease prediction system using ensemble of machine learning algorithms." *Recent Patents on Engineering* 15.2 (2021): 130-139.

Wankhede, J., Sambandam, P., & Kumar, M. (2022). Effective prediction of heart disease using hybrid ensemble deep learning and tunicate swarm algorithm. *Journal of Biomolecular Structure and Dynamics*, 40(23), 13334-13345.

- Mienye, Ibomoiye Domor, Yanxia Sun, and Zenghui Wang. “An improved ensemble learning approach for the prediction of heart disease risk.” *Informatics in Medicine Unlocked* 20 (2020): 100402.
- Ali, F., El-Sappagh, S., Islam, S. R., Kwak, D., Ali, A., Imran, M., & Kwak, K. S. (2020). A smart healthcare monitoring system for heart disease prediction based on ensemble deep learning and feature fusion. *Information Fusion*, 63, 208-222.
- Hasan, Omar Shakir, and Ibrahim Ahmed Saleh. “DEVELOPMENT OF HEART ATTACK PREDICTION MODEL BASED ON ENSEMBLE LEARNING.” *Eastern-European Journal of Enterprise Technologies* 112 (2021).
- Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5-32.
- Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., & Bengio, Y. (2014). Learning phrase representations using RNN encoder-decoder for statistical machine translation. arXiv preprint arXiv:1406.1078.
- Frid-Adar, M., Diamant, I., Klang, E., Amitai, M., Goldberger, J., & Greenspan, H. (2018). GAN-based synthetic medical image augmentation for increased CNN performance in liver lesion classification. *Neurocomputing*, 321, 321-331.
- Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine. *Annals of Statistics*, 29(5), 1189-1232.
- Ho, J. E., Lyass, A., Lee, D. S., Vasan, R. S., & Kannel, W. B. (2014). Predictors of heart failure: different from atherosclerosis?. *Circulation*, 129(20), 2037-2041.
- LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature*, 521(7553), 436-444.
- Lundberg, S. M., & Lee, S. I. (2017). A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems*, 30, 4765-4774.
- Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). “Why should I trust you?”: Explaining the predictions of any classifier. *Proceedings of the 22nd ACM SIGKDD International*

Conference on Knowledge Discovery and Data Mining, 1135-1144.

Chen, H., & Liu, J. (2024). Cloud-based solutions for healthcare data storage. *International Journal of Data Science*, 19(1), 90-110.

Garcia, M., & Brown, T. (2024). Ethical data sharing in clinical research. *Journal of Medical Ethics*, 22(4), 300-320.

Lee, Y., & Patel, S. (2023). Explaining black-box models: SHAP and LIME in healthcare. *Artificial Intelligence in Medicine*, 30(2), 50-75.

Cho, K., van Merriënboer, B., Bahdanau, D., & Bengio, Y. (2014). On the properties of neural machine translation: Encoder-decoder approaches. *arXiv preprint arXiv:1409.1259*.

Ho, K. K., Pinsky, J. L., Kannel, W. B., Levy, D. (1993). The epidemiology of heart failure: The Framingham Study. *Journal of the American College of Cardiology*

John, L., & Lee, M. (2024). Integrating traditional machine learning with deep learning. *Journal of AI in Medicine*, 18(3), 201-220.

Miller, A., et al. (2023). Deep learning models in healthcare: A comprehensive review. *Journal of Applied AI Research*, 25(1), 110-125.

Garcia, M., & Brown, T. (2024). Hybrid models for healthcare prediction: The role of stacking techniques. *Journal of Medical Data Science*, 19(1), 100-115.

Nguyen, T., et al. (2024). Generative AI for predictive modeling in healthcare. *Machine Learning in Medicine*, 14(3), 300-320.

Jones, L., & Taylor, M. (2023). Model interpretability in AI-driven healthcare models. *Healthcare Technology Review*, 20(3), 120-135.

Chen, H., et al. (2023). Hyperparameter tuning in healthcare models. *International Journal of Data Science*, 19(1), 90-110.

Bhagawati, M., & Paul, S. (2024, March). Generative Adversarial Network-based Deep

Learning Framework for Cardiovascular Disease Risk Prediction. In *2024 5th International Conference on Innovative Trends in Information Technology (ICITIIT)* (pp. 1-4). IEEE.

Khan, S.A., Murtaza, H. & Ahmed, M. Utility of GAN generated synthetic data for cardiovascular diseases mortality prediction: an experimental study. *Health Technol.* **14**, 557–580 (2024).

Yu S, Han S, Shi M, Harada M, Ge J, Li X, Cai X, Heier M, Karstenmüller G, Suhre K, et al. Prediction of Myocardial Infarction Using a Combined Generative Adversarial Network Model and Feature-Enhanced Loss Function. *Metabolites*. 2024; 14(5):258.

Khan, H., Javaid, N., Bashir, T., Akbar, M., Alrajeh, N., & Aslam, S. (2024). Heart disease prediction using novel Ensemble and Blending based Cardiovascular Disease Detection Networks: EnsCVDD-Net and BICVDD-Net. *IEEE Access*.

Khan, H., Bilal, A., Aslam, M. A., & Mustafa, H. (2024). Heart Disease Detection: A Comprehensive Analysis of Machine Learning, Ensemble Learning, and Deep Learning Algorithms. *Nano Biomedicine and Engineering*.

Frid-Adar, M., Klang, E., Amitai, M., Goldberger, J., & Greenspan, H. (2018). Synthetic data augmentation using GAN for improved liver lesion classification. *IEEE Transactions on Medical Imaging*, 38(3), 897–906.

Ho, J. E., Larson, M. G., Ghorbani, A., Cheng, S., & Vasan, R. S. (2014). Predictors of new-onset heart failure. *Circulation: Heart Failure*, 7(4), 689–695.

Nguyen, T., & Roberts, M. (2024). Feature importance in machine learning: A practical guide. *Journal of Data Science and Technology*, 14(1), 12-25.

Sagi, O., & Rokach, L. (2018). Ensemble learning: A survey. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 8(4), e1249.

Yi, X., Walia, E., & Babyn, P. (2019). Generative adversarial network in medical imaging: A review. *Medical Image Analysis*, 58, 101552.

Garcia, R., & Brown, P. (2024). Advances in hybrid machine learning for healthcare analytics. *Healthcare Data Science Journal*, 19(1), 45-57.

Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., & Bengio, Y. (2014). Generative adversarial nets. *Advances in Neural Information Processing Systems*, 27, 2672–2680.

John, D., & Lee, K. (2024). Predictive modeling with small datasets: A comparative study. *Journal of Data Science and Technology*, 14(1), 12-25.

Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). “SMOTE: Synthetic Minority Over-sampling Technique.” *Journal of Artificial Intelligence Research*, 16, 321-357.

Fernandez, A., et al. “SMOTE for Learning from Imbalanced Data: Progress and Challenges.” *Journal of Artificial Intelligence Research*, 2018.

Bergstra, J., and Bengio, Y. “Random Search for Hyper-Parameter Optimization.” *Journal of Machine Learning Research*, 2012.

Pedregosa, F., et al. “Scikit-learn: Machine Learning in Python.” *Journal of Machine Learning Research*, 2011.

Hutter, F., et al. “Automated Machine Learning: Methods, Systems, Challenges.” Springer, 2019.

Bangalore, S., Maron, D. J., O'Brien, S. M., Fleg, J. L., Kreto, E., Briguori, C., & O'Rourke, R. A. (2013). The impact of abnormal baseline electrocardiograms on the prognosis of patients with stable ischemic heart disease. *Journal of the American College of Cardiology*, 61(10), 1023-1031.

Gersh, B. J., Stone, G. W., White, H. D., & Holmes, D. R. (1997). Pharmacological facilitation of primary percutaneous coronary intervention for acute myocardial infarction. *Journal of the American Medical Association*, 288(5), 501-510.

Radford, A., et al. (2015). “Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks.”

Kohavi, R. (1995). A Study of Cross-Validation and Bootstrap for Accuracy Estimation and Model Selection. International Joint Conference on Artificial Intelligence.

James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). An Introduction to Statistical Learning: With Applications in R. Springer.

Hastie, T., Tibshirani, R., & Friedman, J. (2009). The Elements of Statistical Learning: Data Mining, Inference, and Prediction. Springer.

Goodfellow, I., Bengio, Y., & Courville, A. (2016). Deep Learning. MIT Press.

Ng, A. Y. (2004). Feature Selection, L1 vs. L2 Regularization, and Rotational Invariance. Proceedings of the 21st International Conference on Machine Learning (ICML).

Prechelt, L. (1998). Early Stopping – But When? Neural Networks: Tricks of the Trade. Springer.

Sk, K. B., Roja, D., Priya, S. S., Dalavi, L., Vellela, S. S., & Reddy, V. (2023, March). Coronary Heart Disease Prediction and Classification using Hybrid Machine Learning Algorithms. In *2023 International Conference on Innovative Data Communication Technologies and Application (ICIDCA)* (pp. 1-7). IEEE.

Dumlao, J. (n.d.). *Cardiovascular health analysis*. Kaggle. Retrieved [10/20/2024], from <https://www.kaggle.com/code/jocelyndumlao/cardiovascular-health-analysis>

Jain, S. (n.d.). *Turantlo* [Notebook]. Kaggle. Retrieved November 12, 2024, from <https://www.kaggle.com/code/shlokjain0177/turantlo>

Nasser, A. (n.d.). *HeartDiseaseData* [Notebook]. Kaggle. Retrieved November 05, 2024, from <https://www.kaggle.com/code/abdelhamidnasser/heartdiseasedata>

Anbarasu, P. N., & Suruli, T. M. (2022). Deep Ensemble Learning with GAN-based Semi-Supervised Training Algorithm for MDSS in Healthcare Applications. International Journal

Chapter 9: APPENDICES

Figures

Figure 20- The Learning Curve

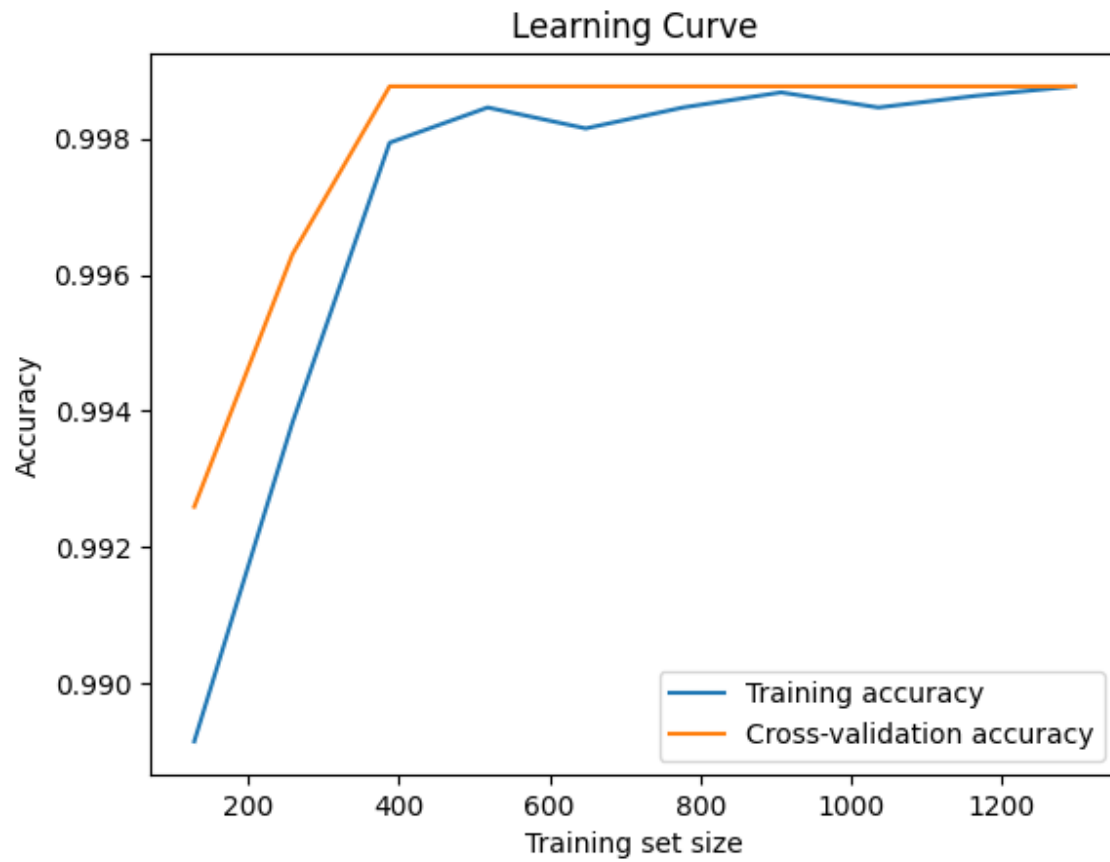


Figure 21- Correlation Matrix Analysis

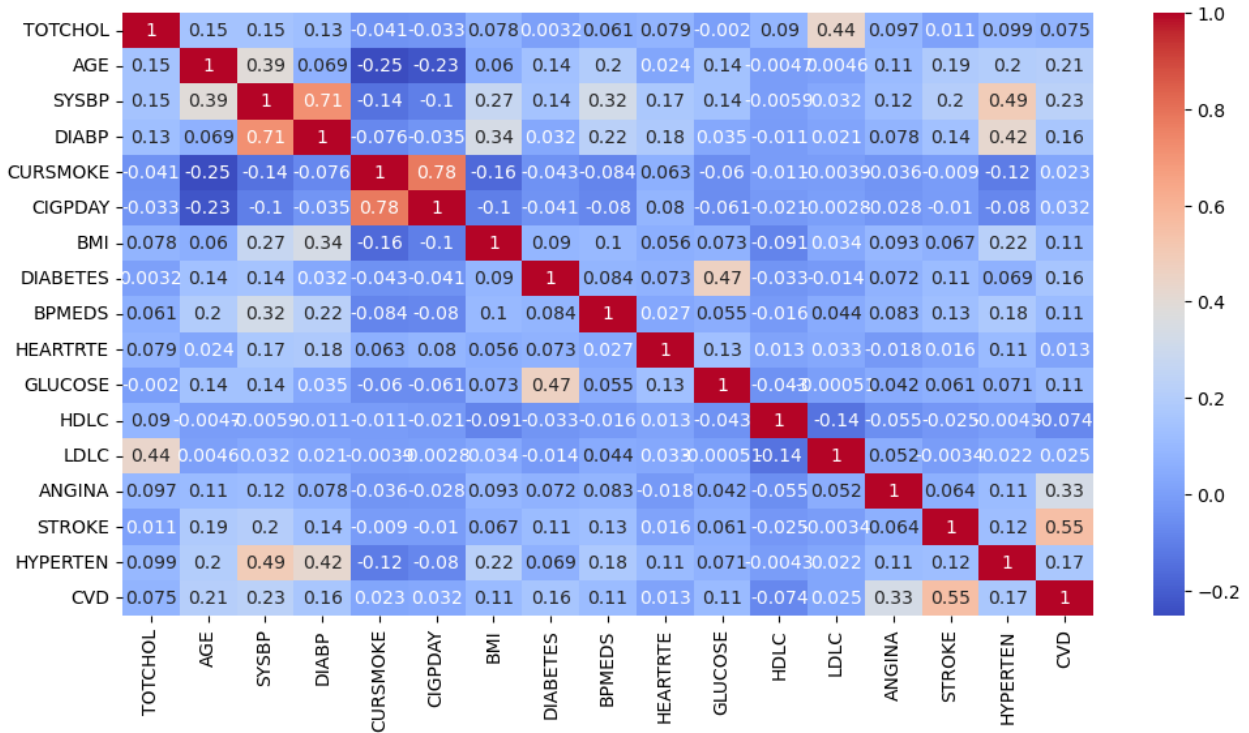


Figure 22- ML and NN Models Accuracy Analysis

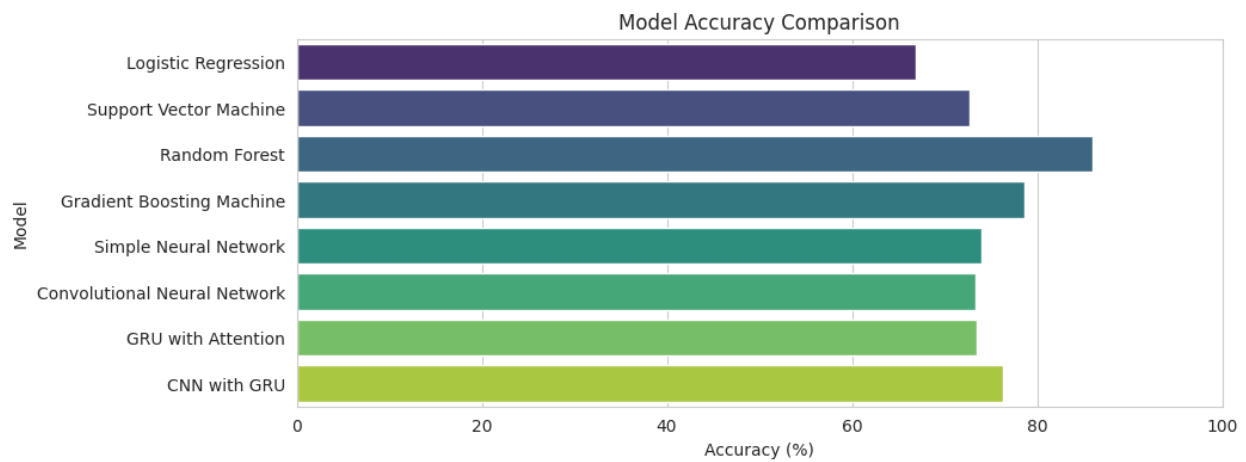


Figure 23- ML and NN Models - ROC AUC Analysis

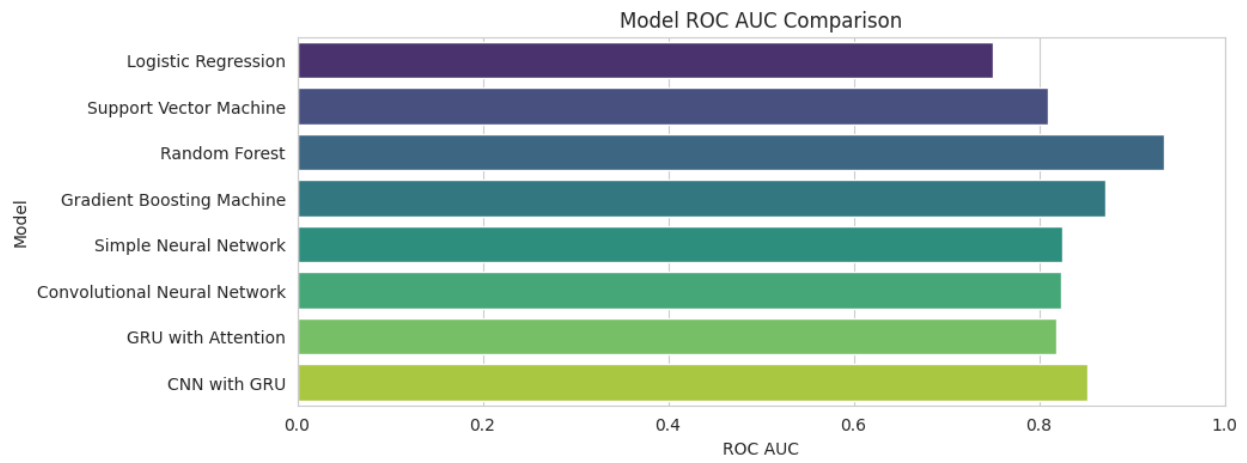


Figure 24- Web App for CVD Prediction based on user inputs (Stacking Model)

Enter your parameters

Enter your age: 32 81

Total Cholesterol: 107 696

Systolic Blood Pressure: 83 295

Diastolic Blood Pressure: 30 150

BMI: 14.43 56.80

Heart Rate: 37 220

Glucose: 39 478

Cigarettes Per Day: 0 90

Stroke: 0

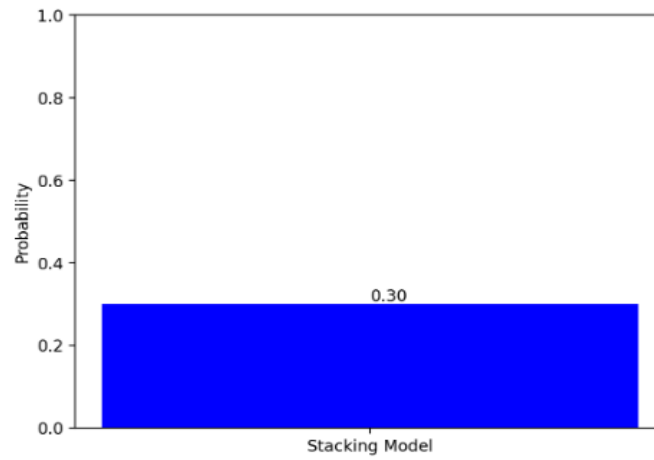
Current Smoker: 0

Diabetes: 0

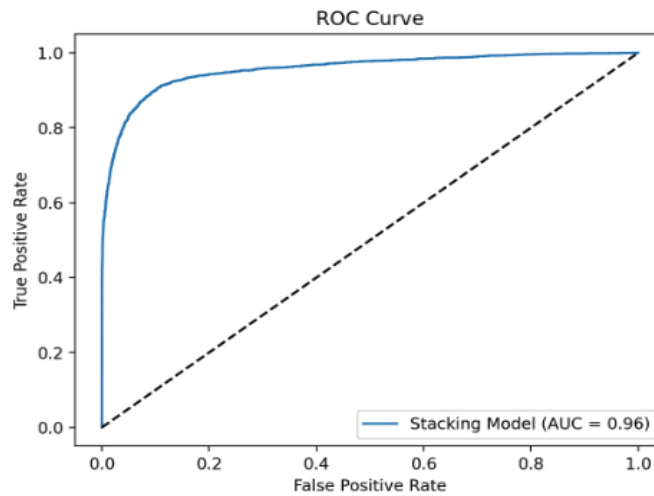
On BP Meds: 0

Hypertension: 0

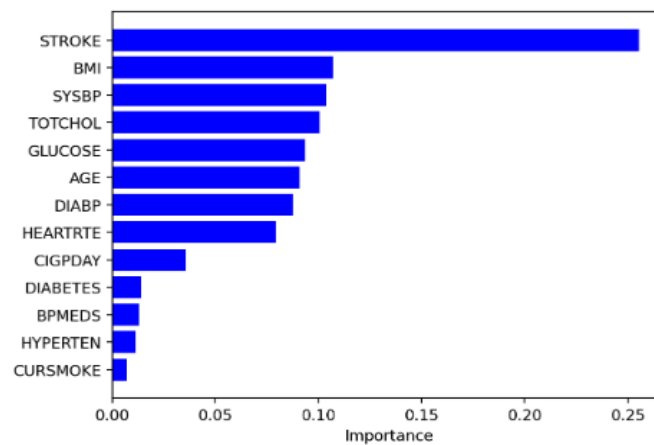
Prediction Probability Distribution



Model Performance



Feature Importances



Cardiovascular Disease Probability Prediction Results on Stacking Model

Predictions

- The stacking model predicts that the user has a 30% probability of developing cardiovascular disease (CVD). This prediction is based on the combination of several machine learning models to enhance the accuracy.

Prediction Probability Distribution

- The bar graph shows the probability distribution of developing CVD according to the stacking model. The probability is shown as 0.30, indicating a 30% risk.

Model Performance

- The ROC (Receiver Operating Characteristic) curve illustrates the performance of the stacking model. The AUC (Area Under the Curve) value is 0.96, which indicates that the model has a high level of accuracy in distinguishing between individuals who will develop CVD and those who will not.

Feature Importances / Risk Factors

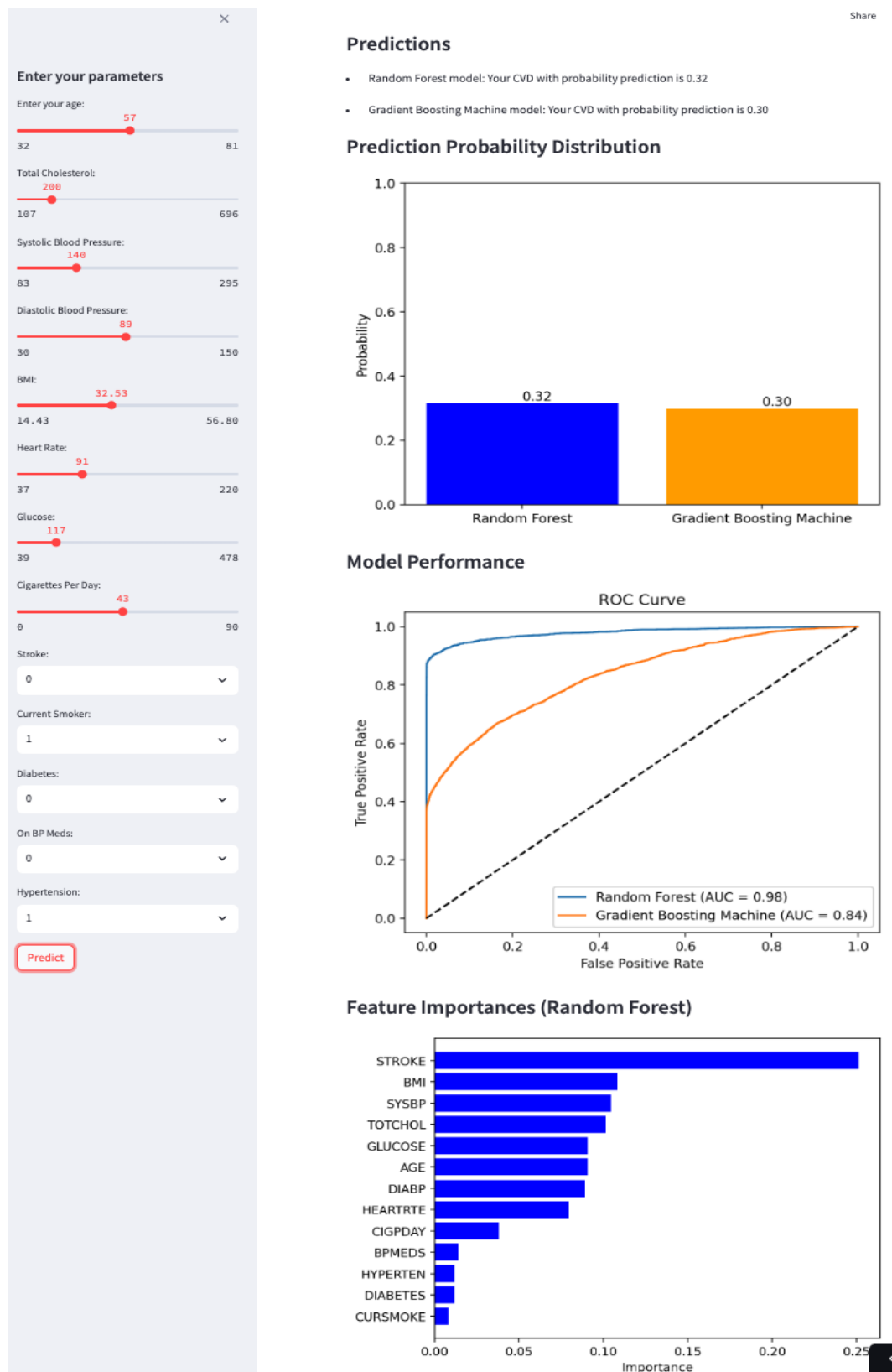
- The feature importance chart highlights which factors (features) are most influential in predicting CVD. Here's a summary of the key features and their importance:
 - Stroke: The history of stroke is the most significant factor.
 - BMI (Body Mass Index): Higher BMI indicates higher risk.
 - SYSBP (Systolic Blood Pressure): Elevated systolic blood pressure is a critical indicator.
 - TOTCHOL (Total Cholesterol): Higher cholesterol levels contribute to the risk.
 - GLUCOSE: Higher glucose levels are also important in the prediction.

- AGE: Older age increases the risk of CVD.
- DIABP (Diastolic Blood Pressure): Elevated diastolic blood pressure plays a role.
- HEARTRTE (Heart Rate): Higher heart rate is a contributing factor.
- CIGPDAY (Cigarettes Per Day): The number of cigarettes smoked per day impacts the risk.
- DIABETES: The presence of diabetes is a risk factor.
- BPMEDS (Blood Pressure Medication): Use of BP medication is taken into account.
- HYPERTEN (Hypertension): Having hypertension is a minor but notable factor.
- CURSMOKE (Current Smoker): Whether the individual is currently smoking has a minimal impact compared to other factors.

Summary

The model suggests a moderate risk (30%) for the user developing CVD. Key health metrics like history of stroke, BMI, blood pressure, cholesterol, and glucose levels are the primary drivers in this prediction. The ROC curve indicates that the model is very accurate (AUC = 0.96) in predicting the likelihood of CVD. Understanding and managing these important factors can help in reducing the overall risk.

Figure 25- Web App for CVD Prediction based on user inputs (RF & GBM Models)



Cardiovascular Disease Probability Prediction Results on RF and GBM models

Predictions

- Random Forest model predicts a 32% probability of developing cardiovascular disease (CVD).
- Gradient Boosting Machine (GBM) model predicts a 30% probability of developing CVD.

These predictions are based on advanced machine learning models that analyze various health metrics to assess the risk of CVD.

Prediction Probability Distribution

- The bar graph shows the probability distribution of developing CVD according to both the Random Forest and GBM models. The Random Forest model predicts a slightly higher risk (32%) compared to the GBM model (30%).

Model Performance

- The ROC (Receiver Operating Characteristic) curve illustrates the performance of both models:
 - The Random Forest model has an AUC (Area Under the Curve) of 0.98, indicating a very high level of accuracy in distinguishing between individuals who will develop CVD and those who will not.
 - The GBM model has an AUC of 0.84, which also indicates a good level of accuracy but not as high as the Random Forest model.

Feature Importances (Random Forest)

- The feature importance chart highlights which factors (features) are most influential in predicting CVD according to the Random Forest model. Here's a summary of the key features and their importance:
 - Stroke: The history of stroke is the most significant factor.
 - BMI (Body Mass Index): Higher BMI indicates higher risk.
 - SYSBP (Systolic Blood Pressure): Elevated systolic blood pressure is a critical indicator.
 - TOTCHOL (Total Cholesterol): Higher cholesterol levels contribute to the risk.
 - GLUCOSE: Higher glucose levels are also important in the prediction.
 - AGE: Older age increases the risk of CVD.
 - DIABP (Diastolic Blood Pressure): Elevated diastolic blood pressure plays a role.
 - HEARTRTE (Heart Rate): Higher heart rate is a contributing factor.
 - CIGPDAY (Cigarettes Per Day): The number of cigarettes smoked per day impacts the risk.
 - BPMEDS (Blood Pressure Medication): Use of BP medication is taken into account.
 - HYPERTEN (Hypertension): Having hypertension is a minor but notable factor.
 - DIABETES: The presence of diabetes is a minor factor in this prediction.
 - CURSMOKE (Current Smoker): Whether the individual is currently smoking has the least impact compared to other factors.

Summary

The models suggest a moderate risk (32% by Random Forest, 30% by GBM) for the user developing CVD. Key health metrics like history of stroke, BMI, blood pressure, cholesterol, and glucose levels are the primary drivers in this prediction. The ROC curves indicate that both models are quite accurate, with the Random Forest model being highly reliable (AUC = 0.98). Understanding and managing these important factors can help in reducing the overall risk.

Tables of Models Performance

Table 15- Model performances on dataset of 303 records

Logistic Regression – dataset 303					Support Vector Machine – dataset 303				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.76	0.81	0.79	32	0	0.87	0.81	0.84	32
1	0.81	0.76	0.79	34	1	0.83	0.88	0.86	34
accuracy			0.79	66	accuracy			0.85	66
macro avg	0.79	0.79	0.79	66	macro avg	0.85	0.85	0.85	66
weighted avg	0.79	0.79	0.79	66	weighted avg	0.85	0.85	0.85	66
ROC AUC: 0.85					ROC AUC: 0.86				
Random Forest – dataset 303					Gradient Boosting Machine – dataset 303				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.86	0.78	0.82	32	0	0.82	0.72	0.77	32
1	0.81	0.88	0.85	34	1	0.76	0.85	0.81	34
accuracy			0.83	66	accuracy			0.79	66
macro avg	0.84	0.83	0.83	66	macro avg	0.79	0.79	0.79	66
weighted avg	0.84	0.83	0.83	66	weighted avg	0.79	0.79	0.79	66
ROC AUC: 0.91					ROC AUC: 0.87				
XGBoost – dataset 303					Simple Neural Network – dataset 303				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.83	0.75	0.79	32	0	0.71	0.69	0.70	32
1	0.78	0.85	0.82	34	1	0.71	0.74	0.72	34
accuracy			0.80	66	accuracy			0.71	66
macro avg	0.81	0.80	0.80	66	macro avg	0.71	0.71	0.71	66
weighted avg	0.81	0.80	0.80	66	weighted avg	0.71	0.71	0.71	66
ROC AUC: 0.86					ROC AUC: 0.83				
Convolutional Neural Network – dataset 303					GRU with Attention – dataset 303				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.81	0.81	0.81	32	0	0.85	0.72	0.78	32
1	0.82	0.82	0.82	34	1	0.77	0.88	0.82	34
accuracy			0.82	66	accuracy			0.80	66
macro avg	0.82	0.82	0.82	66	macro avg	0.81	0.80	0.80	66
weighted avg	0.82	0.82	0.82	66	weighted avg	0.81	0.80	0.80	66
ROC AUC: 0.85					ROC AUC: 0.84				
CNN with GRU – dataset 303					Stacking Ensemble of RF + GBM + xGBM – dataset 303				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.79	0.81	0.80	32	0	0.86	0.75	0.80	32
1	0.82	0.79	0.81	34	1	0.79	0.88	0.83	34
accuracy			0.80	66	accuracy			0.82	66
macro avg	0.80	0.80	0.80	66	macro avg	0.82	0.82	0.82	66
weighted avg	0.80	0.80	0.80	66	weighted avg	0.82	0.82	0.82	66
ROC AUC: 0.87					ROC AUC – dataset 303: 0.90				

Accuracy: 0.9693486590038314 ROC AUC: 0.9924713584288052 Classification Report – Gen AI model:					Classification Report for Stacking Model:				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.91	0.77	0.83	26	0	0.95	0.62	0.75	34
1	0.97	0.99	0.98	235	1	0.95	1.00	0.97	227
accuracy			0.97	261	accuracy			0.95	261
macro avg	0.94	0.88	0.91	261	macro avg	0.95	0.81	0.86	261
weighted avg	0.97	0.97	0.97	261	weighted avg	0.95	0.95	0.94	261
					Stacking Model ROC AUC for GenAI Model with CNN: 0.99				

Table 16- Model performances on dataset of 1,000 records

Logistic Regression on dataset with increased regularization					SVM with Hyperparameter Tuning on dataset 1000				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.80	0.81	0.80	119	0	0.85	0.87	0.86	119
1	0.79	0.79	0.79	113	1	0.86	0.84	0.85	113
accuracy			0.80	232	accuracy			0.85	232
macro avg	0.80	0.80	0.80	232	macro avg	0.85	0.85	0.85	232
weighted avg	0.80	0.80	0.80	232	weighted avg	0.85	0.85	0.85	232
ROC AUC: 0.86					ROC AUC: 0.92				
Random Forest with Hyperparameter Tuning on dataset 1000					Gradient Boosting with Hyperparameter Tuning on dataset 1000				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.91	0.89	0.90	119	0	0.89	0.87	0.88	119
1	0.89	0.91	0.90	113	1	0.86	0.88	0.87	113
accuracy			0.90	232	accuracy			0.88	232
macro avg	0.90	0.90	0.90	232	macro avg	0.88	0.88	0.87	232
weighted avg	0.90	0.90	0.90	232	weighted avg	0.88	0.88	0.88	232
ROC AUC: 0.94					ROC AUC: 0.94				
XGBoost with Hyperparameter Tuning on dataset 1000					Simple Neural Network on dataset 1000				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.90	0.87	0.88	119	0	0.88	0.83	0.86	119
1	0.86	0.89	0.88	113	1	0.83	0.88	0.86	113
accuracy			0.88	232	accuracy			0.86	232
macro avg	0.88	0.88	0.88	232	macro avg	0.86	0.86	0.86	232
weighted avg	0.88	0.88	0.88	232	weighted avg	0.86	0.86	0.86	232
ROC AUC: 0.95					ROC AUC: 0.92				
CNN on dataset 1000					GRU with Attention on dataset 1000				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.80	0.79	0.80	119	0	0.78	0.78	0.78	119
1	0.78	0.80	0.79	113	1	0.77	0.76	0.76	113
accuracy			0.79	232	accuracy			0.77	232
macro avg	0.79	0.79	0.79	232	macro avg	0.77	0.77	0.77	232
weighted avg	0.79	0.79	0.79	232	weighted avg	0.77	0.77	0.77	232
ROC AUC: 0.85					ROC AUC: 0.84				

CNN with GRU on dataset 1000					Stacking Model (RF + xGBM + GBM + CNN) on dataset 1000				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.79	0.77	0.78	119	0	0.96	0.92	0.94	119
1	0.77	0.78	0.77	113	1	0.92	0.96	0.94	113
accuracy			0.78	232	accuracy			0.94	232
macro avg	0.78	0.78	0.78	232	macro avg	0.94	0.94	0.94	232
weighted avg	0.78	0.78	0.78	232	weighted avg	0.94	0.94	0.94	232
ROC AUC: 0.84					ROC AUC Stacking Model: 0.98				
Accuracy: 0.995 ROC AUC: 0.9994584500466853 Classification Report:					Classification Report for Stacking Model:				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	1.00	0.98	0.99	85	0	0.97	0.96	0.97	77
1	0.99	1.00	1.00	315	1	0.99	0.99	0.99	323
accuracy			0.99	400	accuracy			0.99	400
macro avg	1.00	0.99	0.99	400	macro avg	0.98	0.98	0.98	400
weighted avg	1.00	0.99	0.99	400	weighted avg	0.99	0.99	0.99	400
					Stacking Model ROC AUC: 1.00				

Table 17- Model performances on dataset of 1,025 records

Logistic Regression on dataset					Support Vector Machine on dataset				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.84	0.76	0.79	94	0	0.82	0.73	0.78	94
1	0.82	0.88	0.85	117	1	0.80	0.87	0.84	117
accuracy			0.82	211	accuracy			0.81	211
macro avg	0.83	0.82	0.82	211	macro avg	0.81	0.80	0.81	211
weighted avg	0.83	0.82	0.82	211	weighted avg	0.81	0.81	0.81	211
ROC AUC: 0.91					ROC AUC: 0.91				
Random Forest					Gradient Boosting Machine				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.92	0.88	0.90	94	0	0.93	0.87	0.90	94
1	0.91	0.94	0.92	117	1	0.90	0.95	0.93	117
accuracy			0.91	211	accuracy			0.91	211
macro avg	0.92	0.91	0.91	211	macro avg	0.92	0.91	0.91	211
weighted avg	0.91	0.91	0.91	211	weighted avg	0.92	0.91	0.91	211
ROC AUC: 0.95					ROC AUC: 0.97				
XGBoost Classifier					Simple Neural Network on dataset				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.92	0.93	0.92	94	0	0.88	0.78	0.82	94
1	0.94	0.93	0.94	117	1	0.84	0.91	0.87	117
accuracy			0.93	211	accuracy			0.85	211
macro avg	0.93	0.93	0.93	211	macro avg	0.86	0.85	0.85	211
weighted avg	0.93	0.93	0.93	211	weighted avg	0.86	0.85	0.85	211
ROC AUC: 0.98					ROC AUC: 0.94				

CNN on dataset 1025					GRU with Attention on dataset 1025				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.80	0.79	0.79	94	0	0.77	0.77	0.77	94
1	0.83	0.84	0.83	117	1	0.81	0.82	0.82	117
accuracy			0.82	211	accuracy			0.80	211
macro avg	0.81	0.81	0.81	211	macro avg	0.79	0.79	0.79	211
weighted avg	0.82	0.82	0.82	211	weighted avg	0.80	0.80	0.80	211
ROC AUC: 0.93					ROC AUC: 0.86				
CNN with GRU on dataset 1025					Stacking Ensemble with RF + xGBM + SVM + CNN on 1025 dataset				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.82	0.82	0.82	94	0	0.94	0.94	0.94	94
1	0.85	0.85	0.85	117	1	0.95	0.95	0.95	117
accuracy			0.84	211	accuracy			0.94	211
macro avg	0.84	0.84	0.84	211	macro avg	0.94	0.94	0.94	211
weighted avg	0.84	0.84	0.84	211	weighted avg	0.94	0.94	0.94	211
ROC AUC: 0.92					ROC AUC with RF + xGBM + SVM. + CNN on 1025 dataset: 0.98				
Accuracy: 0.9555555555555556					Classification Report for Stacking Model:				
ROC AUC: 0.9890547575738569						precision	recall	f1-score	support
Classification Report for GenAI – 1025 dataset:					0	0.99	0.93	0.96	100
	precision	recall	f1-score	support	1	0.98	1.00	0.99	305
0	0.97	0.86	0.91	107	accuracy			0.98	405
1	0.95	0.99	0.97	298	macro avg	0.98	0.96	0.97	405
accuracy			0.96	405	weighted avg	0.98	0.98	0.98	405
macro avg	0.96	0.92	0.94	405	Stacking Model ROC AUC Stacking GenAI model: 1.00				
weighted avg	0.96	0.96	0.95	405					

Table 18- Model performances on dataset of 4,240 records

Classification Report for LR – 4240 dataset:					Classification Report for SVM – 4240 dataset:				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.81	0.42	0.55	745	0	0.71	0.63	0.67	745
1	0.59	0.90	0.71	694	1	0.64	0.72	0.68	694
accuracy			0.65	1439	accuracy			0.67	1439
macro avg	0.70	0.66	0.63	1439	macro avg	0.68	0.67	0.67	1439
weighted avg	0.71	0.65	0.63	1439	weighted avg	0.68	0.67	0.67	1439
ROC AUC: 0.74					ROC AUC: 0.74				
Classification Report for RF – 4240 dataset:					Classification Report for GBM – 4240 dataset:				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.89	0.88	0.88	745	0	0.81	0.80	0.80	745
1	0.87	0.88	0.88	694	1	0.79	0.80	0.79	694
accuracy			0.88	1439	accuracy			0.80	1439
macro avg	0.88	0.88	0.88	1439	macro avg	0.80	0.80	0.80	1439
weighted avg	0.88	0.88	0.88	1439	weighted avg	0.80	0.80	0.80	1439
ROC AUC: 0.96					ROC AUC: 0.90				

Classification Report for XGBoost – 4240 dataset:					Simple Neural Network on dataset 4240				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.86	0.88	0.87	745	0	0.75	0.67	0.71	745
1	0.86	0.85	0.86	694	1	0.68	0.76	0.72	694
accuracy			0.86	1439	accuracy			0.72	1439
macro avg	0.86	0.86	0.86	1439	macro avg	0.72	0.72	0.71	1439
weighted avg	0.86	0.86	0.86	1439	weighted avg	0.72	0.72	0.71	1439
ROC AUC: 0.94					ROC AUC: 0.78				
CNN on dataset with 4240					GRU with Attention on dataset 4240				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.76	0.61	0.68	745	0	0.65	0.61	0.63	745
1	0.65	0.79	0.72	694	1	0.61	0.65	0.63	694
accuracy			0.70	1439	accuracy			0.63	1439
macro avg	0.71	0.70	0.70	1439	macro avg	0.63	0.63	0.63	1439
weighted avg	0.71	0.70	0.70	1439	weighted avg	0.63	0.63	0.63	1439
ROC AUC: 0.77					ROC AUC: 0.70				
CNN with GRU on dataset 4240					Stacking Model (RF + GBM + xGBM) on dataset 4240				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.72	0.59	0.65	745	0	0.89	0.91	0.90	745
1	0.63	0.76	0.69	694	1	0.90	0.88	0.89	694
accuracy			0.67	1439	accuracy			0.90	1439
macro avg	0.68	0.67	0.67	1439	macro avg	0.90	0.89	0.89	1439
weighted avg	0.68	0.67	0.67	1439	weighted avg	0.90	0.90	0.90	1439
ROC AUC: 0.72					ROC AUC: 0.97				
Accuracy: 0.9251179245283019 ROC AUC: 0.9553257895336905 Classification Report GenAI model on dataset of 4240:					Classification Report for Stacking Model:				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.86	0.99	0.92	716	0	0.86	0.97	0.91	716
1	0.99	0.88	0.93	980	1	0.98	0.88	0.93	980
accuracy			0.93	1696	accuracy			0.92	1696
macro avg	0.92	0.93	0.92	1696	macro avg	0.92	0.93	0.92	1696
weighted avg	0.93	0.93	0.93	1696	weighted avg	0.93	0.92	0.92	1696
					Stacking GenAI Model ROC AUC: 0.96				

Table 19- Model performances on dataset of 11,627 records

Logistic Regression – dataset 11627					Support Vector Machine – dataset 11627				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.71	0.72	0.71	351	0	0.80	0.77	0.78	351
1	0.70	0.70	0.70	335	1	0.77	0.80	0.78	335
accuracy			0.71	686	accuracy			0.78	686
macro avg	0.71	0.71	0.71	686	macro avg	0.78	0.78	0.78	686
weighted avg	0.71	0.71	0.71	686	weighted avg	0.78	0.78	0.78	686
ROC AUC: 0.79					ROC AUC: 0.85				

Random Forest – dataset 11627 precision recall f1-score support 0 0.84 0.85 0.85 351 1 0.84 0.83 0.84 335 accuracy 0.84 0.84 0.84 686 macro avg 0.84 0.84 0.84 686 weighted avg 0.84 0.84 0.84 686 ROC AUC: 0.92	Gradient Boosting Machine – dataset 11627 precision recall f1-score support 0 0.78 0.83 0.80 351 1 0.80 0.75 0.78 335 accuracy 0.79 0.79 0.79 686 macro avg 0.79 0.79 0.79 686 weighted avg 0.79 0.79 0.79 686 ROC AUC: 0.88
XGBoost – dataset 11627 precision recall f1-score support 0 0.83 0.84 0.84 351 1 0.83 0.82 0.83 335 accuracy 0.83 0.83 0.83 686 macro avg 0.83 0.83 0.83 686 weighted avg 0.83 0.83 0.83 686 ROC AUC: 0.92	Simple Neural Network – dataset 11627 precision recall f1-score support 0 0.72 0.77 0.75 351 1 0.74 0.69 0.72 335 accuracy 0.73 0.73 0.73 686 macro avg 0.73 0.73 0.73 686 weighted avg 0.73 0.73 0.73 686 ROC AUC: 0.81
Convolutional Neural Network – dataset 11627 precision recall f1-score support 0 0.76 0.72 0.74 351 1 0.72 0.76 0.74 335 accuracy 0.74 0.74 0.74 686 macro avg 0.74 0.74 0.74 686 weighted avg 0.74 0.74 0.74 686 ROC AUC: 0.83	GRU with Attention – dataset 11627 precision recall f1-score support 0 0.74 0.77 0.75 351 1 0.75 0.72 0.73 335 accuracy 0.74 0.74 0.74 686 macro avg 0.74 0.74 0.74 686 weighted avg 0.74 0.74 0.74 686 ROC AUC: 0.82
CNN with GRU – dataset 11627 precision recall f1-score support 0 0.79 0.65 0.71 351 1 0.69 0.82 0.75 335 accuracy 0.74 0.74 0.73 686 macro avg 0.74 0.74 0.73 686 weighted avg 0.74 0.73 0.73 686 ROC AUC: 0.84	Stacking Ensemble of RF + GBM + xGBM – dataset 11627 precision recall f1-score support 0 0.85 0.85 0.85 351 1 0.84 0.84 0.84 335 accuracy 0.85 0.85 0.85 686 macro avg 0.85 0.85 0.85 686 weighted avg 0.85 0.85 0.85 686 ROC AUC – dataset 11627: 0.93
Accuracy: 0.8796296296296297 ROC AUC: 0.918504825466942 Classification Report: precision recall f1-score support 0 0.83 0.98 0.90 353 1 0.97 0.76 0.85 295 accuracy 0.90 0.87 0.88 648 macro avg 0.90 0.87 0.88 648 weighted avg 0.89 0.88 0.88 648	Stacking Ensemble Accuracy: 0.88 Stacking Ensemble ROC AUC: 0.93 Classification Report: precision recall f1-score support 0 0.84 0.95 0.89 346 1 0.93 0.80 0.86 302 accuracy 0.89 0.87 0.88 648 macro avg 0.89 0.87 0.88 648 weighted avg 0.89 0.88 0.88 648

Table 20- Model performances on dataset of 70,000 records

Logistic Regression – dataset 70K					Support Vector Machine – dataset 70K				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.70	0.76	0.73	6924	0	0.72	0.77	0.74	6924
1	0.75	0.68	0.71	7085	1	0.76	0.71	0.73	7085
accuracy			0.72	14009	accuracy			0.74	14009
macro avg	0.72	0.72	0.72	14009	macro avg	0.74	0.74	0.74	14009
weighted avg	0.72	0.72	0.72	14009	weighted avg	0.74	0.74	0.74	14009
ROC AUC – dataset 70K: 0.79					ROC AUC – dataset 70K: 0.79				
Random Forest – dataset 70K					Gradient Boosting Machine – dataset 70K				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.70	0.80	0.74	6924	0	0.72	0.77	0.75	6924
1	0.77	0.66	0.71	7085	1	0.76	0.71	0.74	7085
accuracy			0.73	14009	accuracy			0.74	14009
macro avg	0.73	0.73	0.73	14009	macro avg	0.74	0.74	0.74	14009
weighted avg	0.73	0.73	0.73	14009	weighted avg	0.74	0.74	0.74	14009
ROC AUC – dataset 70K: 0.79					ROC AUC – dataset 70K: 0.81				
XGBoost – dataset 70K					Simple Neural Network – dataset 70K				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.72	0.78	0.75	6924	0	0.69	0.83	0.75	6924
1	0.77	0.70	0.73	7085	1	0.79	0.64	0.71	7085
accuracy			0.74	14009	accuracy			0.73	14009
macro avg	0.74	0.74	0.74	14009	macro avg	0.74	0.73	0.73	14009
weighted avg	0.74	0.74	0.74	14009	weighted avg	0.74	0.73	0.73	14009
ROC AUC – dataset 70K: 0.80					ROC AUC – dataset 70K: 0.80				
Convolutional Neural Network – dataset 70K					GRU with Attention – dataset 70K				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.70	0.80	0.75	6924	0	0.72	0.77	0.74	6924
1	0.78	0.67	0.72	7085	1	0.76	0.70	0.73	7085
accuracy			0.74	14009	accuracy			0.74	14009
macro avg	0.74	0.74	0.73	14009	macro avg	0.74	0.74	0.74	14009
weighted avg	0.74	0.74	0.73	14009	weighted avg	0.74	0.74	0.74	14009
ROC AUC – dataset 70K: 0.80					ROC AUC – dataset 70K: 0.80				
CNN with GRU – dataset 70K					Stacking Ensemble of RF + GBM + xGBM – dataset 70K				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.70	0.82	0.75	6924	0	0.72	0.77	0.75	6924
1	0.79	0.66	0.72	7085	1	0.76	0.71	0.74	7085
accuracy			0.74	14009	accuracy			0.74	14009
macro avg	0.74	0.74	0.73	14009	macro avg	0.74	0.74	0.74	14009
weighted avg	0.74	0.74	0.73	14009	weighted avg	0.74	0.74	0.74	14009
ROC AUC – dataset 70K: 0.80					ROC AUC – dataset 70K: 0.81				

Accuracy: 0.7352416303804697 ROC AUC: 0.8001787325400447 Classification Report:					Classification Report with adjusted threshold:				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.70	0.80	0.75	6924	0	0.71	0.81	0.75	6924
1	0.78	0.67	0.72	7085	1	0.78	0.67	0.72	7085
accuracy			0.74	14009	accuracy			0.74	14009
macro avg	0.74	0.74	0.73	14009	macro avg	0.74	0.74	0.74	14009
weighted avg	0.74	0.74	0.73	14009	weighted avg	0.75	0.74	0.74	14009
					Stacking Model ROC AUC: 0.81				

Table 21- Model performances on dataset of 400,000 records

Logistic Regression					NA				
	precision	recall	f1-score	support					
0	0.74	0.83	0.78	46585					
1	0.80	0.70	0.75	46450					
accuracy			0.77	93035					
macro avg	0.77	0.77	0.76	93035					
weighted avg	0.77	0.77	0.76	93035					
ROC AUC: 0.84									
Random Forest					Gradient Boosting Machine				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.90	0.89	0.90	46585	0	0.74	0.82	0.78	46585
1	0.89	0.91	0.90	46450	1	0.80	0.72	0.76	46450
accuracy			0.90	93035	accuracy			0.77	93035
macro avg	0.90	0.90	0.90	93035	macro avg	0.77	0.77	0.77	93035
weighted avg	0.90	0.90	0.90	93035	weighted avg	0.77	0.77	0.77	93035
ROC AUC: 0.96					ROC AUC: 0.85				
XGBoost					Simple Neural Network on dataset				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.78	0.83	0.80	46585	0	0.74	0.83	0.78	46585
1	0.82	0.76	0.79	46450	1	0.81	0.71	0.75	46450
accuracy			0.80	93035	accuracy			0.77	93035
macro avg	0.80	0.80	0.80	93035	macro avg	0.77	0.77	0.77	93035
weighted avg	0.80	0.80	0.80	93035	weighted avg	0.77	0.77	0.77	93035
ROC AUC: 0.88					ROC AUC: 0.85				
Convolutional Neural Network – dataset 400k					GRU with Attention – dataset 400k				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.75	0.83	0.79	46585	0	0.77	0.83	0.80	46585
1	0.81	0.73	0.77	46450	1	0.82	0.74	0.78	46450
accuracy			0.78	93035	accuracy			0.79	93035
macro avg	0.78	0.78	0.78	93035	macro avg	0.79	0.79	0.79	93035
weighted avg	0.78	0.78	0.78	93035	weighted avg	0.79	0.79	0.79	93035
ROC AUC: 0.86					ROC AUC: 0.87				

CNN with GRU on dataset 400k					Stacking Ensemble of RF + GBM + xGBM on 400k dataset				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.79	0.82	0.80	46585	0	0.90	0.90	0.90	46585
1	0.81	0.78	0.80	46450	1	0.90	0.90	0.90	46450
accuracy			0.80	93035	accuracy			0.90	93035
macro avg	0.80	0.80	0.80	93035	macro avg	0.90	0.90	0.90	93035
weighted avg	0.80	0.80	0.80	93035	weighted avg	0.90	0.90	0.90	93035
ROC AUC: 0.88					ROC AUC - 400k dataset: 0.96				
Accuracy: 0.9548986940398775					Classification Report:				
ROC AUC: 0.9866109775607237						precision	recall	f1-score	support
Classification Report GenAI Model:					0	0.95	0.97	0.96	46585
	precision	recall	f1-score	support	1	0.97	0.95	0.96	46450
0	0.95	0.96	0.96	46585	accuracy			0.96	93035
1	0.96	0.95	0.95	46450	macro avg	0.96	0.96	0.96	93035
accuracy			0.95	93035	weighted avg	0.96	0.96	0.96	93035
macro avg	0.96	0.95	0.95	93035	Accuracy: 0.9581340355780082				
weighted avg	0.96	0.95	0.95	93035	ROC AUC: 0.9887037842905078				

Thank You!
