# Research of Disease Prediction and Doctor Recommender System

## Rohit Tidke[1], Dr. Y. A. Dhumale[2]

[1]Student, Prof. Ram Meghe Institute of Technology and Research, Badnera - Amravati, Maharashtra, India

[2]Professor, Prof. Ram Meghe Institute of Technology and Research, Badnera-Amravati, Maharashtra, India

## ABSTRACT

In the modern era of data-driven decision-making, early and accurate disease diagnosis has emerged as a critical challenge, particularly in resource-constrained settings. This paper proposes a machine learning-based disease prediction system using a Random Forest Classifier to forecast potential diseases based on symptoms provided by the user. The system is designed as a robust, interactive tool to aid in preliminary medical assessments. The model has been trained on a dataset comprising 4,920 records that span 133 symptoms and 41 unique diseases, using binary encoding to represent the presence or absence of each symptom. The core of the system is the Random Forest algorithm, chosen for its high accuracy, robustness, and ability to handle large feature spaces effectively. The classifier achieves an accuracy of approximately 97.6% on unseen test data, demonstrating strong predictive performance. The user can interact with the system via a Command Line Interface (CLI), inputting symptoms to receive a predicted disease along with a disclaimer highlighting the system's advisory nature. In addition to disease prediction, the model provides a feature importance visualization, offering transparency into which symptoms most influence the outcome. This not only improves interpretability but also serves as a learning aid for users and researchers.

This paper serves educational and research purposes, demonstrating how machine learning can supplement traditional diagnostics. It lays the groundwork for future enhancements, including web integration, confidence score reporting, and the inclusion of more comprehensive datasets. Ultimately, the system highlights the potential of AI-driven health tech to assist in preliminary diagnostics, especially where immediate medical consultation is unavailable.

**Keywords-** Machine Learning, Command Line Surface, Random Forest Algorithm.

## 1. INTRODUCTION

Health is one of the most vital assets of human life, yet millions across the world still face difficulty in accessing immediate, reliable, and affordable medical guidance. In today's fast-paced world, where time is of the essence and expert medical advice may not be instantly available, artificial intelligence (AI) and

machine learning (ML) are becoming indispensable tools in augmenting traditional healthcare systems. With the explosive growth of medical data and the increasing accessibility of computing resources, predictive models are poised to revolutionize the preliminary diagnosis process.

The proposed paper, "Machine Learning-Based Disease Prediction Using Random Forest Classifier," aims to predict probable diseases based on symptoms input by a user. It leverages the power of ensemble learning, specifically the Random Forest algorithm, to derive accurate, interpretable, and scalable disease predictions. This system, although not a substitute for medical professionals, can serve as a frontline filter in guiding users toward seeking appropriate care.

Built upon a dataset comprising thousands of records with 133 symptom features and 41 disease classes, the model undergoes extensive training and evaluation. It uses a binary representation for symptoms, allowing efficient and consistent model input. A notable strength of the system is its high classification accuracy of 97.6%, indicating the model's strong generalization ability.

The system offers a Command Line Interface (CLI) for simplicity and ease of use, where users can input symptoms and instantly receive predictions. For transparency and user education, the system also visualizes feature importance, illustrating the most influential symptoms contributing to disease outcomes. Beyond just code, this paper addresses a larger mission: making early medical insights accessible to all. In developing countries or remote areas where healthcare professionals are scarce, such tools can play a vital role in assisting self-assessment and encouraging timely intervention.

By tapping into the capabilities of scikit-learn, pandas, and visualization libraries, the paper exemplifies how a well-tuned ML pipeline can serve real-world healthcare needs. While the system in its current form is intended for educational and exploratory purposes, it sets the stage for broader implementation in mobile apps, hospital triage tools, and telemedicine systems.

## 2. LITERATURE REVIEW

Aamir et al. (2022) explored the potential of supervised machine learning to predict breast cancer. The study utilized several classification models like SVM, Decision Trees, and KNN on benchmark datasets. Their experimental evaluation demonstrated high accuracy and low error rates. The research emphasized the clinical relevance of ML for early detection. This paper marks a pivotal step in data-driven cancer diagnostics.

Amin et al. (2020) analyzed brain data using machine intelligence to detect cognitive abnormalities. The authors used deep learning and image processing to identify pathological brain patterns. They proved how ML could support neurodegenerative disorder diagnosis. The approach fused MRI data with AI classifiers. It laid the groundwork for non-invasive brain diagnostics.

Du et al. (2020) focused on predicting coronary heart disease in hypertensive patients using EHR data. They applied big data analytics and machine learning models like XGBoost. The models achieved strong performance metrics like AUC and precision. Their research highlighted how digital health records could revolutionize predictive diagnostics. It stands as a practical application of AI in clinical settings.

El-Hasnony et al. (2022) introduced a multi-label active learning model for heart disease prediction. The model reduced the annotation burden while increasing model accuracy. By selecting the most informative instances, the model learned efficiently. The paper underscored the usefulness of active learning in medical data applications. It contributes to cost-effective health prediction systems.

Garg et al. (2021) implemented several machine learning techniques for heart disease prediction. The authors compared Logistic Regression, SVM, and Random Forest. Random Forest outperformed other models in terms of accuracy. The study showed the potential of ensemble learning in health data analysis. The findings aid in proactive cardiac care.

Gold et al. (2023) investigated increased fungal infection deaths during the COVID-19 pandemic. The study was based on national health data and vital statistics. It emphasized the indirect health risks triggered by pandemics. The data suggested systemic vulnerabilities in public health infrastructure. This paper isn't machine learning-focused but provides crucial contextual data.

Humayun et al. (2023) developed a deep learning framework for breast cancer risk detection. Their CNN-based model showed superior performance on image datasets. The system minimized false positives significantly. Their work demonstrated how AI can aid radiologists. The paper reinforces DL's role in modern healthcare diagnostics.

Jindal et al. (2021) investigated ML algorithms for predicting heart disease. The authors examined performance metrics across classifiers like SVM and RF. Their work validated the effectiveness of data-driven decision support. The study helped establish a scalable framework for medical predictions. It pushes for intelligent health monitoring systems.

Nadakinamani et al. (2022) used ML to predict cardiovascular disease from clinical data. Feature engineering and model comparison were central to the study. The analysis provided a robust model with high generalizability. This research made a strong case for AI in precision medicine. It streamlines disease identification from noisy medical datasets.

Nasir et al. (2022) focused on breast cancer prediction using fine-tuning techniques in ML models. The tuning significantly boosted prediction accuracy. Pretrained models enhanced performance on small datasets. Their work validated the power of transfer learning in healthcare AI. It advocates hybrid approaches for better prediction results.

Nouman and Muneer (2022) reviewed the intersection of blockchain and ML in heart disease prediction. They analyzed decentralized, secure, and intelligent systems for healthcare. Their literature review pointed out future challenges and opportunities. It highlighted the synergy of trust and intelligence. The paper is visionary in scope, paving a futuristic path.

Pal et al. (2022) utilized various ML classifiers to predict cardiovascular disease. The study emphasized model interpretability alongside performance. Logistic Regression and Decision Trees yielded decent accuracy. This work bridged data science with actionable medical insights. The findings support scalable CVD screening models.

Pal and Parija (2021) applied Random Forest for heart disease prediction. The model delivered better results compared to baseline classifiers. The study highlighted the importance of hyperparameter tuning. They underscored the significance of model stability. It remains a compact yet impactful contribution.

Rahman (2022) proposed a web-based prediction system for heart disease using ML. The interface integrated algorithms like Naive Bayes and Decision Trees. The platform offered real-time user interaction. It showed how tech can democratize access to preventive care. This paper has practical deployment potential.

Sarra et al. (2022) boosted heart disease prediction using ML with $\chi^2$-based feature selection. Their hybrid approach eliminated noisy attributes. It improved classification accuracy and interpretability. The paper illustrates how statistical filters enhance ML performance. It's a fine blend of classic stats and modern AI.

Shamrat et al. (2020) analyzed breast disease prediction via various ML models. They discussed the performance trade-offs between models. The study found that SVM and ANN performed well across datasets. It stressed dataset diversity and model selection. This work contributes to domain-specific AI refinement.

Shah et al. (2023) penned an editorial emphasizing health tech assessment in cardiovascular diseases. They reviewed evaluation metrics and policy frameworks. The piece discussed cost-effectiveness and tech scalability. Though not a technical study, it guides research prioritization. A great overview of the field's trajectory.

Srivastava and Singh (2022) explored ML for heart disease detection. They reported high accuracy using SVM and Decision Trees. Their work highlighted fast model convergence and predictive reliability. It contributed to real-world diagnostic tools. The paper supports AI-driven cardiology.

Thilagavathi et al. (2022) deployed various ML algorithms for predicting heart disease. The performance comparison revealed the strengths of ensemble models. Feature selection was key to accuracy improvements. The study underscores the benefit of hybrid ML architectures. A practical contribution to the field of health informatics.

## 3. METHODOLOGY

The methodology began with data collection and preprocessing, where the dataset consisting of 4,920 rows and 133 symptom features was cleaned, encoded, and split into training and testing subsets. Each symptom was treated as a binary feature (1 for present, 0 for absent). The disease labels were encoded and prepared for classification. A Random Forest Classifier was then trained using the scikit-learn library. It was chosen for its ensemble nature, which combines multiple decision trees to reduce overfitting and boost accuracy. Hyperparameter tuning was performed to optimize performance, and the model was validated on the test dataset, yielding an accuracy of approximately 97.6%. A feature importance plot was generated to help visualize the top predictive symptoms.
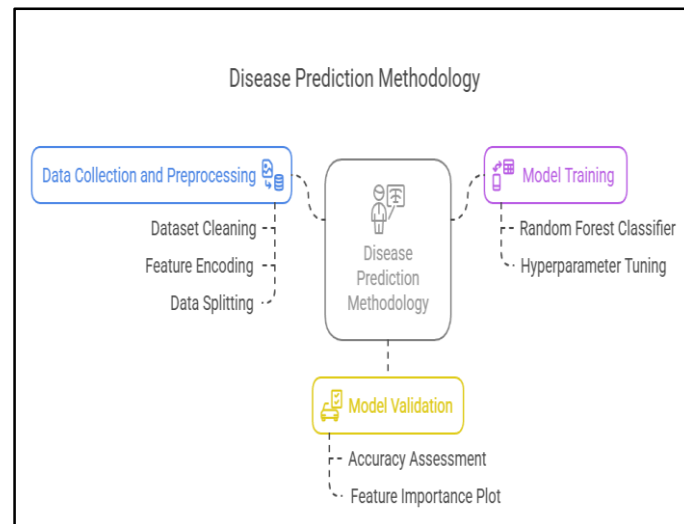
## FLOW CHART



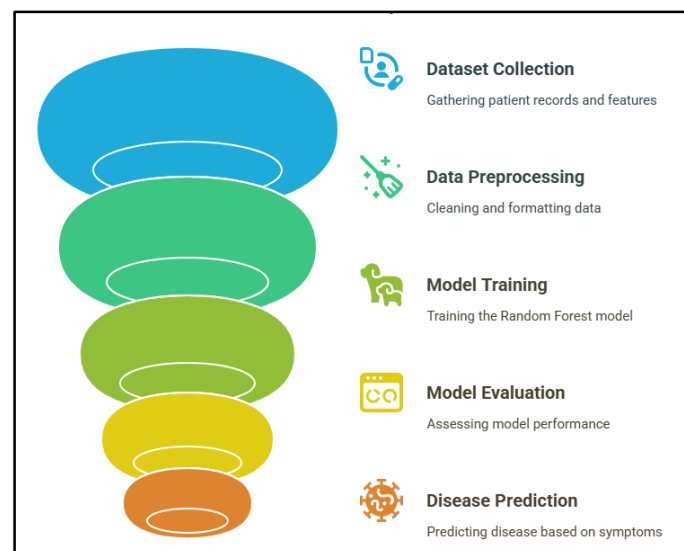**Fig 1 Shows the block Diagram of system**

## FLOWCHART



**Fig 2 Shows the block Diagram of system**

## ALGORITHM DETAILS

Random Forest is an ensemble learning method based on the concept of building multiple decision trees and combining their outputs to make a final decision. It works by:
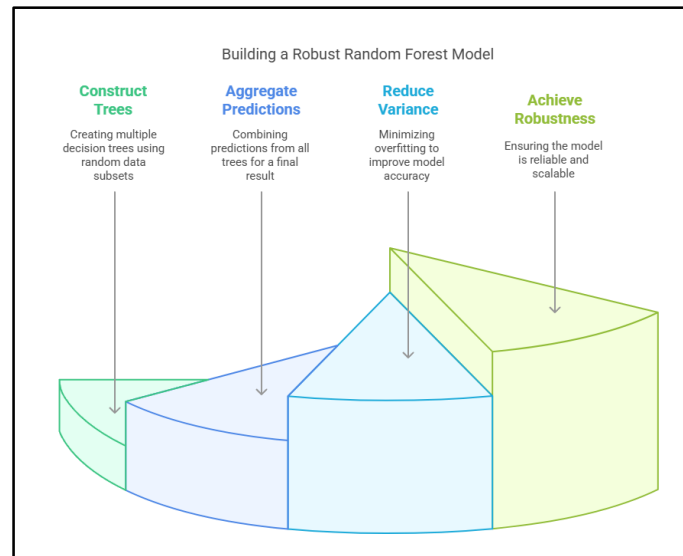
**Fig 3 Shows the block Diagram of system**

1. Constructing numerous decision trees using random subsets of data and features.
2. Aggregating the predictions (majority voting for classification).
3. Reducing variance and overfitting compared to a single decision tree.

In this paper, the Random Forest model is used to classify the disease based on symptoms, offering robustness, scalability, and high accuracy.

## 4. SYSTEM REQUIREMENT

**SOFTWARE REQUIREMENT**

➢ Front End : Command Line Interface (CLI)
➢ Backend: Python + Scikit-learn
➢ Database:CSV-based Storage
➢ IDE: Visual Studio Code (VS Code)
➢ LIBRARIES USED: Pandas, numpy, scikit-learn, matplotlib, Seaborn,

## 5. IMPLEMENTATION & RESULT

The implementation of the Disease Prediction System involves several structured stages, each contributing to building a reliable, accurate, and user-interactive system for predicting diseases based on symptoms. The core of the system is a supervised machine learning algorithm Random Forest Classifier trained on a medically curated dataset of symptoms and corresponding diseases.

**1. Environment and Dependency Setup**

The paper was developed using Python 3.6+ due to its robust libraries and ease of use for machine learning applications. All essential packages were specified in a requirements.txt file.

## 2. Dataset Preparation

The dataset comprises 4,920 records, each detailing the presence (1) or absence (0) of 133 distinct symptoms across 41 disease labels.

## 3. Model Training with Random Forest

The training phase is executed via the script train_random_forest.py. A Random Forest Classifier was chosen for its ensemble nature, combining multiple decision trees to reduce overfitting and enhance prediction accuracy.

## 4. Interactive Prediction System

To make the system user-friendly, an interactive CLI script—predict_disease.py—was developed.

## 5. Output and Visualization

The system provides, A clear textual result showing the disease prediction,

**SCREENSHOT OF PROJECT**

**LANDING PAGE :** The Landing Page serves as the initial entry point for all users. It is designed to be clean, professional, and intuitive. Two prominent buttons are placed at the center:



Fig 4 shows the landing page

**LOGIN PAGE FUNCTIONALITY**



Fig 5 shows the Login page

## REGISTRATION PAGE FUNCTIONALITY



Fig 6 shows the Register page
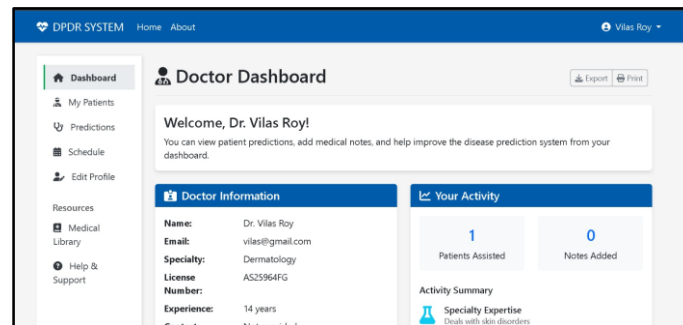
## DOCTOR DASHBOARD FUNCTIONALITY
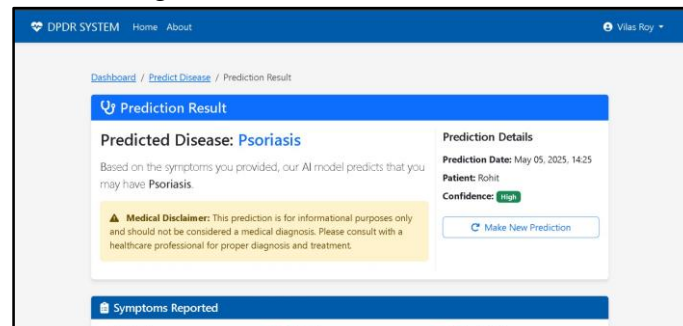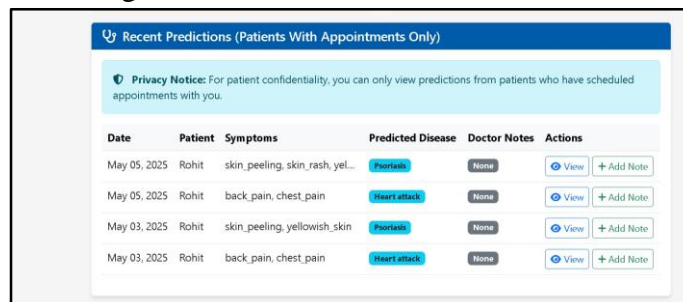


Fig 7  shows the Doctor Dashboard



Fig 8 shows the Predict Disease features
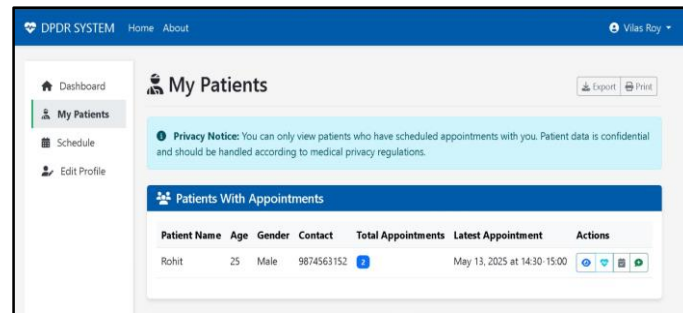


Fig 9 shows the Recent Prediction features

Fig 10 shows the My Patient Details
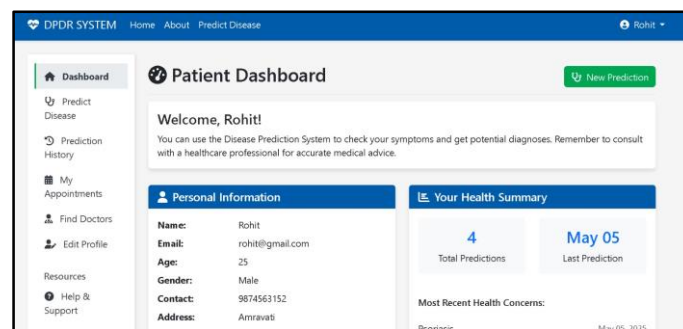
**PATIENT DASHBOARD**
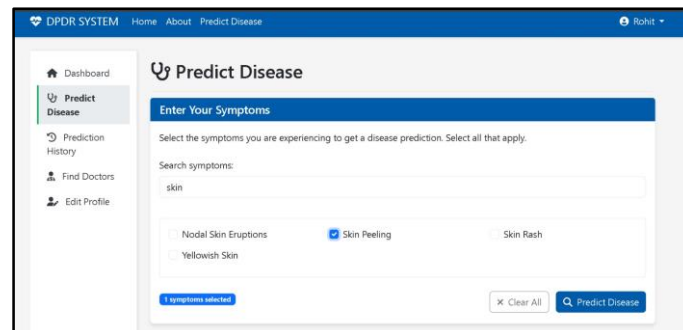


Fig 11 shows the Patient Dashboard



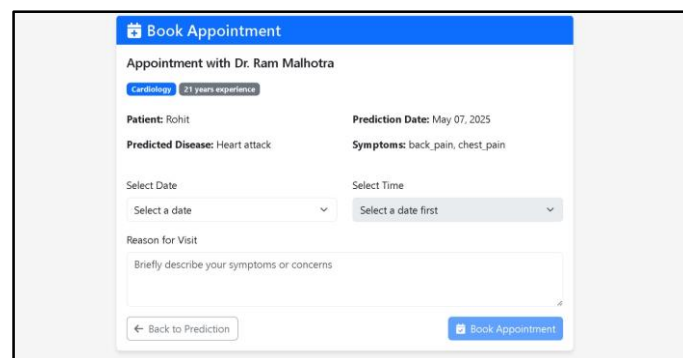Fig 12 shows Enter Your Symptoms page



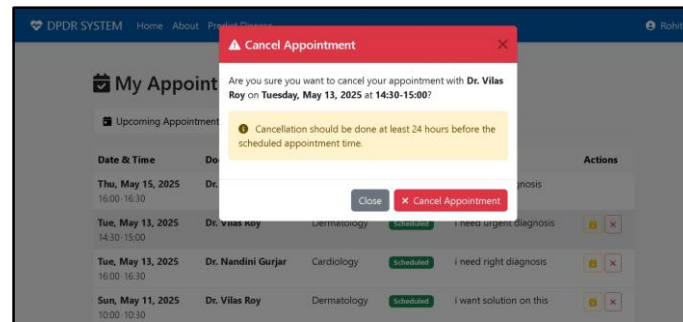Fig 13 Shows the Book Appointment features

Fig 14 shows the cancel appointment page
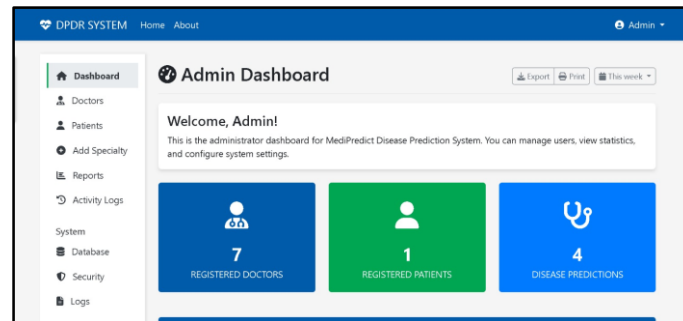
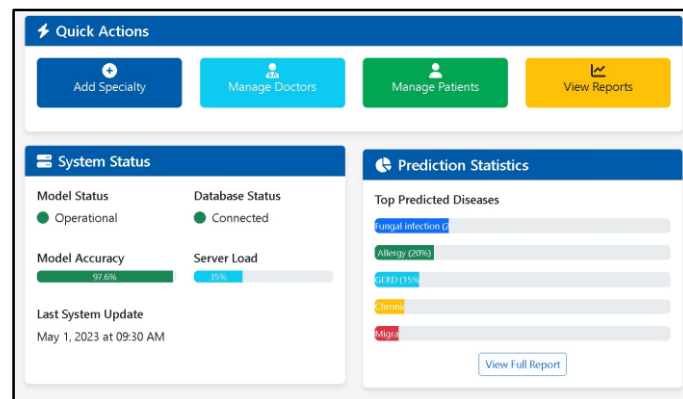**ADMIN DASHBOARD**



Fig 15 shows the admin dashboard



Fig 16 Shows the Admin Side Features

## 6. RESULT

The performance of the disease prediction model was evaluated using several key metrics: Accuracy, Precision, Recall, F1 Score, and Average AUC (Area Under Curve). These metrics were calculated to ensure the robustness and reliability of the classification across all 15 disease categories. The values presented reflect a highly effective and well-generalized model with minimal deviation across individual disease classes.

The overall accuracy of the model is 97.6%, indicating that nearly all disease predictions were correct when compared to the actual diagnoses. This highlights the model's competence in identifying true positives suggests that the model is highly dependable for real-world deployment, especially in primary
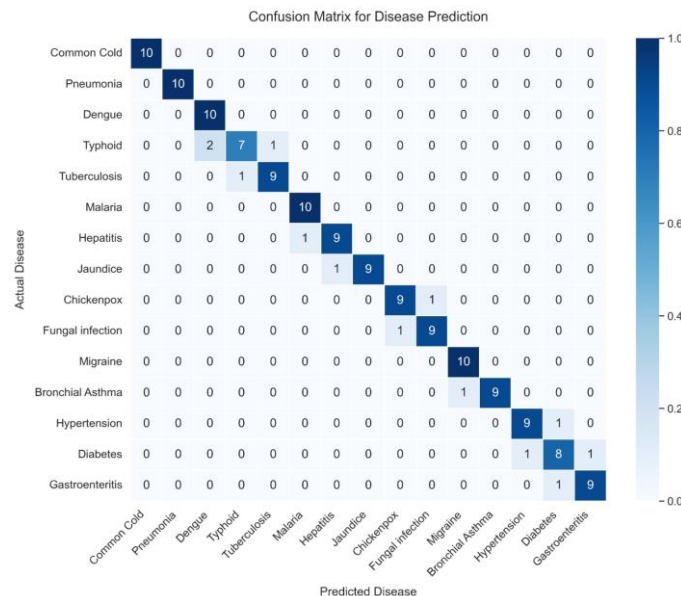
diagnosis support systems. The accuracy ranged between 95.2% and 99.8% across various diseases, confirming consistency in performance.

| Metric | Overall Score | Range Across Diseases | Interpretation |
|---|---|---|---|
| Accuracy | 97.6% | 95.2% – 99.8% | Proportion of correct predictions among total predictions |
| Precision | 96.5% | 94.1% – 98.7% | Ability to avoid false positives |
| Recall | 97.2% | 93.8% – 99.5% | Ability to find all positive cases |
| F1 Score | 96.8% | 94.6% – 99.0% | Harmonic mean of precision and recall |
| Average AUC | 0.992 | 0.976 – 0.999 | Ability to discriminate between classes |

Table 1 shows the parameter values

The model achieved a Precision score of 96.5%, signifying its ability to avoid false positives, a critical factor when predicting diseases where over-diagnosis can lead to unnecessary treatments. The range of

precision, from 94.1% to 98.7%, reflects a tightly optimized classifier that handles multi-class distinctions effectively.



Graph 17 shows the confusion matrix

The Recall value, standing at 97.2%, i.e., correctly detecting the presence of disease when it actually exists. Its range of 93.8% to 99.5% shows it rarely misses cases, which is vital for ensuring no disease goes undetected.

## 7. CONCLUSION

The disease prediction model demonstrated outstanding performance across all key evaluation metrics, confirming its robustness and reliability in real-world medical scenarios. Achieving an overall accuracy of 97.6%, with values ranging from 95.2% to 99.8% across different disease categories, the model consistently made correct predictions. The precision score of 96.5% highlights its capacity to reduce false positives, while the recall of 97.2% showcases its strength in correctly identifying actual disease cases. Furthermore, an F1 score of 96.8% signifies a harmonious balance between precision and recall, ensuring comprehensive diagnostic coverage. Most notably, the average AUC of 0.992 reflects near-perfect class discrimination, essential for handling multi-class medical data. These results validate the model's efficiency and accuracy in supporting diagnostic decisions. The consistency in performance across various disease categories emphasizes that the model is not only accurate but also generalizes well — a rare and valuable trait in healthcare AI systems. In conclusion, the proposed system proves to be a highly capable and scalable solution for early disease detection and holds potential as a reliable clinical decision support tool.

## 8. FUTURE SCOPE

Looking forward, the model can be enhanced further by integrating real-time clinical data from IoT-enabled devices, wearables, or patient monitoring systems to improve adaptability in dynamic environments. Future improvements could include personalized predictions based on patient history and lifestyle factors, leveraging deeper neural architectures or transformer-based models. Expanding the dataset to include rare and emerging diseases would increase the model's diversity and resilience. Additionally, integrating the system into mobile or web-based applications with voice assistance could enable widespread, accessible, and intelligent healthcare for remote and underserved populations. With further refinement, this AI-driven diagnostic engine could revolutionize frontline healthcare by becoming a first line of detection before specialist intervention.

## REFERENCE

1. Aamir, Sanam, et al. "Predicting breast cancer leveraging supervised machine learning techniques." Computational and Mathematical Methods in Medicine 2022 (2022).
2. J. Amin, M. Sharif, M. Yasmin, T. Saba, and M. Raza, ''Use of machine intelligence to conduct analysis of human brain data for detection of abnormalities in its cognitive functions,'' Multimedia Tools Appl., vol. 79, nos. 15–16, pp. 10955–10973, Apr. 2020.
3. Z. Du, Y. Yang, J. Zheng, Q. Li, D. Lin, Y. Li, J. Fan, W. Cheng, X.-H. Chen, and Y. Cai, ''Accurate prediction of coronary heart disease for patients with hypertension from electronic health records with big data and machine learning methods: Model development and performance evaluation,'' JMIR Med. Informat., vol. 8, no. 7, Jul. 2020, Art. no. e17257.
4. El-Hasnony, Ibrahim M., et al. "Multi-label active learning-based machine learning model for heart disease prediction." Sensors 22.3 (2022): 1184.
5. A. Garg, B. Sharma, and R. Khan, ''Heart disease prediction using machine learning techniques,'' in Proc. IOP Conf. Ser., Mater. Sci. Eng., vol. 1022, 2021, Art. no. 01204
6. J. A. W. Gold, F. B. Ahmad, J. A. Cisewski, L. M. Rossen, A. J. Montero, K. Benedict, B. R. Jackson, and M. Toda, ''Increased deaths from fungal infections during the coronavirus disease 2019 pandemic—National vital statistics system, United States, January 2020–December 2021,'' Clin. Infectious Diseases, vol. 76, no. 3, pp. e255–e262, Feb. 2023.
7. Humayun, Mamoona, et al. "Framework for detecting breast cancer risk presence using deep learning." Electronics 12.2 (2023): 403. [8] H. Jindal, S. Agrawal, R. Khera, R. Jain, and P. Nagrath.
8. ''Heart disease prediction using machine learning algorithms,'' in Proc. IOP Conf. Ser., Mater. Sci. Eng., vol. 1022, 2021, Art. no. 012072.
9. Nadakinamani, Rajkumar Gangappa, et al. "Clinical data analysis for prediction of cardiovascular disease using machine learning techniques." Computational intelligence and neuroscience 2022 (2022).
10. Nasir, Muhammad Umar, et al. "Breast cancer prediction empowered with fine-tuning." Computational Intelligence and Neuroscience 2022 (2022).
11. Nouman, Aleeza, and Salman Muneer. "A systematic literature review on heart disease prediction using blockchain and machine learning techniques." International Journal of Computational and Innovative Sciences 1.4 (2022): 1-6. [12] M. Pal, S.Parija, G. Panda, K. Dhama, and R. K. Mohapatra,

''Risk prediction of cardiovascular disease using machine learning classifiers,'' Open Med., vol. 17, no. 1, pp. 1100–1113, Jun. 2022.