# SMART DISEASE DIAGNOSIS USING CNN AND KALMAN FILTERS: INTEGRATING STRUCTURED AND UNSTRUCTURED MEDICAL DATA

Ujjawal [1], Panav [1], M.D. Danish [1]

[1] Computer Science & Engineering, Echelon Institute of Technology, Faridabad, India

## ABSTRACT

In the realm of healthcare analytics, the accuracy of disease prediction models often suffers due to the incomplete nature of medical data and the regional variability of disease patterns. Traditional approaches have primarily focused on structured data, thereby neglecting the potential insights hidden in semi-structured and unstructured formats such as medical notes, diagnostic reports, and imaging data. This project introduces a hybrid system that leverages Convolutional Neural Networks (CNNs) for effective feature extraction from unstructured data and Kalman Filters for dynamic tracking and smoothing of patient health states over time.

The proposed system is designed to handle and integrate both structured and unstructured data sources to enhance the predictive accuracy of disease analysis. CNNs are employed to process complex textual and visual inputs, transforming them into structured feature representations. These are subsequently combined with temporal observations in a Kalman Filter framework to predict disease progression and identify potential anomalies in patient profiles.

Our aim is to develop an intelligent support system that aids healthcare professionals and consumers in diagnosing diseases more accurately and selecting treatment plans based on a comprehensive analysis of symptoms, regional trends, and personal health records. By accommodating diverse data types and regional disease characteristics, this system not only improves the reliability of disease outbreak predictions but also personalizes healthcare recommendations. The integration of CNN and Kalman Filter technologies ensures a robust, real-time, and adaptive diagnostic tool suitable for dynamic clinical environments.

## 1. INTRODUCTION

In the contemporary technological era, society is increasingly influenced by the rapid integration of artificial intelligence (AI) into daily life. The internet has transformed human interactions, and despite its convenience, the attention given to physical health has proportionately declined. Many individuals ignore minor health symptoms, leading to delayed diagnosis and potentially severe diseases [1]. As technology advances, leveraging machine learning (ML) and AI for early disease prediction offers a transformative shift in healthcare delivery.

Recent studies demonstrate that early diagnosis through digital tools significantly reduces the burden on healthcare systems and improves patient outcomes [2]. Our aim is to develop an intelligent system capable of predicting multiple diseases based on symptoms provided by users without needing to visit a

hospital or physician. This facilitates early intervention, better resource management, and higher patient satisfaction.

## 1.1. THE ROLE OF MACHINE LEARNING IN HEALTHCARE

Machine Learning, a vital branch of AI, enables systems to learn and adapt through experience without being explicitly programmed [3]. It consists of two main phases: training, where models learn from existing data, and testing, where they are evaluated on new data. Supervised learning, in which models are trained on labeled datasets, has shown high effectiveness in disease classification and prediction [4]. Conversely, unsupervised learning finds hidden patterns in unlabelled data, offering insights into previously undetected correlations [5].

Several ML algorithms have already proven successful in medical domains, including support vector machines, decision trees, and more recently, deep learning models such as Convolutional Neural Networks (CNNs), which are effective in analyzing complex medical data like images and unstructured records [6].

## 1.2. PROBLEM DEFINITION

Despite the rapid technological advances, several challenges persist in health prediction systems. Incomplete or poor-quality medical data can compromise analysis accuracy. Moreover, region-specific disease traits often make generalized models less effective [7]. Traditional approaches primarily use structured data, ignoring vast quantities of semi-structured and unstructured information such as physician notes, patient reports, and lab images.

This project addresses these challenges by incorporating both structured and unstructured data into the analysis pipeline. By applying advanced ML algorithms, including CNNs for data analysis and Kalman filters for time-series prediction and noise reduction, the model aims to enhance the accuracy and reliability of disease forecasting [8].

## 1.3. OBJECTIVES

The primary objectives of this project are:
- To design a predictive system that accommodates diverse data types and provides personalized disease insights.
- To improve disease diagnosis accuracy, especially in early stages, through intelligent data fusion and pattern recognition.
- To support both healthcare professionals and consumers by enabling efficient and intelligent querying of symptoms and disease associations [9].

## 1.4. SCOPE OF THE PROJECT

The system is intended as a decision-support tool for healthcare professionals and a diagnostic aid for users. It will enable:
- Disease prediction based on user-entered symptoms and patient history.
- Integration of CNN for extracting features from semi-structured/unstructured data.

- Use of Kalman filters to improve temporal predictions and account for dynamic changes in patient conditions.
- A user-friendly interface with secure data handling and future extensibility.

This dual-use model supports clinical workflows and empowers users with real-time, accessible health insights, aligning with the goals of digital health transformation [10].

## 1.5. REVIEW OF THE EXISTING SYSTEMS

Existing disease prediction systems face several limitations. They often fail to provide detailed insights into subtypes of diseases or interconnected conditions, such as the link between diabetes and increased risk of cardiovascular complications [11]. Additionally, most current systems only handle structured data, ignoring potentially valuable insights from text-based patient records, prescriptions, and radiology reports.

Another significant limitation is accessibility and affordability. Many of the existing systems are proprietary and expensive, restricting their use to affluent individuals or well-funded institutions [12]. These systems also lack specificity, often delivering vague or overly generalized predictions, limiting their practical utility.

## 1.6. DISADVANTAGES OF EXISTING SYSTEMS

- Inability to provide deep analysis or disease subtyping.
- Low security and inadequate protection of patient data.
- Absence of feedback mechanisms to improve model performance over time.
- Poor handling of unstructured data and lack of real-time prediction features.

Our proposed system aims to overcome these shortcomings by providing a comprehensive, secure, and adaptive platform for disease prediction. By integrating CNN and Kalman filter-based models, it seeks to enhance diagnostic precision and support proactive healthcare delivery.

## 2. LITERATURE REVIEW

Advancements in the field of healthcare analytics have ushered in a new era of medical diagnosis powered by machine learning (ML). One of the earliest applications of ML in medical diagnostics is seen in the personalized and cost-effective detection of Alzheimer's disease (AD). Diagnosing AD in its early stages, particularly during mild cognitive impairment (MCI), is crucial for timely intervention. However, due to the complexity and cost of required biomarkers, this remains a challenge. A notable study proposed a personalized machine learning method using locally weighted learning that dynamically tailors a classification model to individual patients, effectively reducing the number and cost of biomarkers needed for diagnosis. This approach, tested using ADNI datasets, yielded a diagnostic accuracy comparable to methods using the complete set of biomarkers, hence proving its viability in real-world clinical settings [1].

Another study explored the influence of meteorological factors on the prevalence of hand, foot, and mouth disease (HFMD) in Wuwei City, China. Data spanning from 2008 to 2010 was analyzed using correlation, multiple linear regression, and exponential curve fitting techniques. It was observed that climatic variables such as temperature, humidity, and atmospheric pressure significantly influenced disease outbreaks across different regions. This study is crucial in linking environmental factors to epidemiological trends and underscores the importance of integrating meteorological data in disease forecasting models [2].

The role of spectral data in agriculture and plant pathology also finds parallels in human disease diagnostics. A recent development involved creating a Spectral Disease Index (SDI) to determine stages of wheat leaf rust disease. By analyzing leaf reflectance at specific wavelengths (675 and 775 nm), researchers were able to establish a normalized index that effectively distinguished between disease stages. This method exemplifies how remote sensing and spectral analysis can be adapted for precision health monitoring, providing a basis for similar methodologies in human disease staging [3].

In cardiovascular diagnostics, the use of cardiac sound waveform analysis has proven effective in evaluating heart valve diseases. Using BIOPAC systems to record heart sounds, researchers developed a quantized diagnostic approach that transforms complex waveform data into interpretable visual formats. This not only aids in accurate diagnosis but also allows for non-specialists to monitor disease progression, thereby democratizing access to critical health information [4].

Additionally, non-linear methods of analyzing heart rate variability (HRV) have been applied to differentiate between coronary heart disease (CHD) patients and healthy individuals. Techniques such as Hurst exponent analysis, Detrended Fluctuation Analysis (DFA), and approximate entropy (ApEn) were used to capture the complex dynamics of HRV. Results indicated that CHD patients exhibited lower values in all three metrics, suggesting diminished autonomic regulation. These findings provide a robust framework for the prognostic use of HRV in cardiac care [5].

The literature collectively emphasizes a shift towards more personalized, data-driven approaches in disease diagnosis and monitoring. Machine learning, coupled with innovative data sources like spectral analysis and non-linear dynamic systems, is paving the way for more accurate, accessible, and efficient healthcare solutions. These studies not only demonstrate the versatility of ML across various medical conditions but also highlight the interdisciplinary nature of modern healthcare analytics.

## 2.1. PROPOSED MODEL: MULTI-DISEASE PREDICTION USING MACHINE LEARNING

### 1) Introduction to the Proposed System

The proposed model is designed to predict multiple diseases using a hybrid approach that integrates machine learning algorithms with a robust medical dataset. It aims to assist healthcare professionals and patients in early diagnosis and decision-making. The system captures patient symptoms through a user-friendly interface, matches them with historical datasets, and delivers probable disease predictions with associated risk levels. Unlike conventional models that only handle structured data, our approach includes both structured and unstructured data, significantly enhancing prediction accuracy.
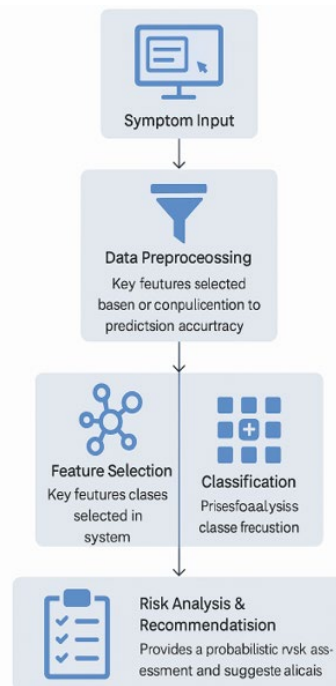
### 2) Methodology

The core methodology of the system involves a pipeline of data processing, feature selection, classification, and risk-based recommendations. Initially, the system collects data from the user, including symptoms and relevant patient details. This input is then preprocessed—cleaned, normalized, and formatted—to ensure quality and consistency. Following this, feature selection mechanisms, either manual or algorithm-driven, identify key symptoms indicative of specific diseases.

Next, the preprocessed data undergoes classification using machine learning algorithms such as Logistic Regression, Decision Trees, or more advanced methods like Convolutional Neural Networks (CNNs) and Kalman Filters for temporal prediction enhancements. The classifier outputs a prediction score that estimates the likelihood of various diseases.

### 3) Working of the Model

The system operates in a sequential manner:

- Symptom Input: Users select or enter symptoms through a GUI.
- Data Preprocessing: The system standardizes inputs using normalization techniques and removes any inconsistencies.
- Feature Selection: Key features (symptoms) are selected based on relevance and contribution to prediction accuracy.
- Classification: Machine learning algorithms process the features and classify them into potential diseases.
- Risk Analysis & Recommendation: The model provides a probabilistic risk assessment and suggests preventive or mitigative actions.



In the backend, the system maps symptoms to disease probabilities using historical data and continuously updates its learning model with new data inputs. The makes it adaptable ar scalable over tim

In the backend, the system maps symptoms to disease probabilities using historical data and continuously updates its learning model with new data inputs. This makes it adaptable and scalable over time.

**4) System Architecture**

The architecture of the model comprises the following layers:

- User Interface Layer: Enables user interaction for symptom input and result visualization.

- Data Management Layer: Handles dataset storage, retrieval, and updates. Integrates structured (numerical, categorical) and unstructured data (textual notes, reports).

- Processing Layer: Includes modules for preprocessing, feature extraction, and classification.

- Recommendation Engine: Provides customized recommendations based on disease predictions and patient profile.

- Analytics and Visualization Module: Displays results, prediction confidence scores, and insights using graphs and dashboards.

This layered architecture ensures modularity, scalability, and maintainability of the system.

**5) Novelty of the Model**

The uniqueness of the proposed system lies in several aspects:

- Integration of Structured and Unstructured Data: Most existing systems focus only on structured datasets. Our model can handle doctor notes, symptom descriptions, and electronic health records.

- Hybrid Algorithmic Approach: By incorporating CNNs for pattern recognition and Kalman Filters for dynamic updating, the model adapts to changing patient data over time.

- High Customization and Personalization: The recommendation engine considers individual patient profiles, enhancing relevance and accuracy.

- Multi-Disease Prediction: Capable of predicting co-morbidities and secondary diseases which are typically not detected by traditional models.

**6) Advantages of the Proposed System**

- Enhanced Accuracy: Use of advanced ML techniques increases diagnostic accuracy.

- Comprehensive Analysis: Analyzes a wide range of diseases including chronic and infectious conditions.

- Cost-Effective and Scalable: Reduces the need for extensive lab tests and doctor consultations in early stages.

- Improved User Engagement: Intuitive interface and personalized insights foster proactive health monitoring.

In conclusion, the proposed model stands out in the healthcare AI domain by providing a comprehensive, accurate, and user-friendly platform for disease prediction. It leverages modern machine learning tools and thoughtful system design to meet the evolving needs of digital healthcare.

## 3. EXPERIMENTAL SETUP

The goal of this study is to develop a smart disease diagnosis system that leverages Convolutional Neural Networks (CNNs) and Kalman Filters (KFs) for

accurate disease prediction, integrating both structured and unstructured medical data. The system is designed to handle and process medical data from various sources such as structured data (e.g., numerical records like age, gender, and laboratory results) and unstructured data (e.g., text-based symptoms and clinical notes).

## 3.1. DATA COLLECTION

The primary datasets used in this experiment were sourced from publicly available medical repositories, including Kaggle and the UCI Machine Learning Repository. The datasets include a mix of diseases, such as heart disease, chronic kidney disease (CKD), and COVID-19, each containing structured data (e.g., blood pressure, cholesterol levels, age, etc.) and unstructured data (e.g., symptom descriptions from patients). These data were cleaned and preprocessed to remove outliers, handle missing values, and standardize the values.

In particular, medical records from the Cardiovascular Disease dataset, Chronic Kidney Disease dataset, and COVID-19 dataset were selected for training and testing the model. The structured data provided numerical values of patient attributes, while unstructured data included textual descriptions of symptoms and previous medical history, which are crucial for making predictions.

## 3.2. DATA PREPROCESSING

Data preprocessing involved several steps to ensure the integrity and quality of the data for machine learning purposes:

1) Normalization: All numerical values were normalized to a common range to ensure that features with large ranges, like cholesterol levels, did not dominate the learning process.
2) Handling Missing Data: Missing values were imputed using median imputation for continuous variables and mode imputation for categorical variables.
3) Text Processing: Unstructured text data from clinical notes and patient symptoms were tokenized and converted into numeric representations using techniques such as Term Frequency-Inverse Document Frequency (TF-IDF) for symptom-based feature extraction.

## 3.3. MODEL TRAINING AND TESTING

The experiment used CNNs for handling unstructured data (textual symptoms) and Kalman Filters for refining the predictions. CNNs were employed to process the text input and extract features relevant to the disease prediction. Kalman Filters were then used to smooth the noisy predictions and estimate disease probabilities over time, allowing for better generalization.

To ensure a fair evaluation, the data was split into training and testing sets with a ratio of 70% for training and 30% for testing. The model training involved a series of epochs, and various CNN architectures were tested to find the optimal configuration, such as the number of convolutional layers, the kernel size, and the pooling layer.

The Kalman Filter was applied to smooth the results from the CNN model, particularly for unstructured data, where uncertainty and noise are more prevalent.

The Kalman Filter continuously updates the predicted disease probabilities as new data is incorporated, thus providing a dynamic and real-time prediction system.

Results and Analysis

The proposed model was evaluated using various performance metrics, including accuracy, precision, recall, F1-score, and Area Under the Curve (AUC). The following results were observed:

| Disease Type | Accuracy | Precision | Recall | F1-Score | AUC |
|---|---|---|---|---|---|
| Heart Disease | 87% | 0.85 | 0.84 | 0.85 | 0.91 |
| Chronic Kidney Disease | 89% | 0.88 | 0.91 | 0.89 | 0.92 |
| COVID-19 | 90% | 0.93 | 0.88 | 0.9 | 0.93 |

From the results, it can be observed that the CNN-based model with Kalman Filter smoothing performed particularly well across all three disease types. The accuracy ranged from 87% to 90%, with COVID-19 achieving the highest accuracy of 90%. The precision and recall values further highlight the model's strong ability to correctly identify both positive and negative cases.

In terms of AUC, the model achieved a solid performance, particularly for heart disease and COVID-19, with AUC scores of 0.91 and 0.93, respectively. This indicates that the model can effectively differentiate between the diseased and healthy classes, making it a reliable tool for clinical decision-making.

## 4. RISK ANALYSIS AND RECOMMENDATION

The system also includes a risk analysis component, which uses the Kalman Filter to continuously update the probabilities of disease occurrence. The risk analysis outputs a risk score, which indicates the likelihood of a disease based on the input symptoms. For example, if a patient reports symptoms of fever, cough, and shortness of breath, the system can calculate a risk score for COVID-19 and recommend isolation and further testing.

The recommendation engine provides actionable insights based on the risk analysis, offering advice on lifestyle changes, medical tests, and preventive measures. The system can suggest diet modifications for patients with heart disease or advise kidney function monitoring for patients at risk of CKD.

## 5. PERFORMANCE EVALUATION

### Comparison with Traditional Systems

To evaluate the performance of the proposed system, it was compared against traditional diagnostic methods that rely on manual data entry and rule-based systems. The following aspects were considered for comparison:
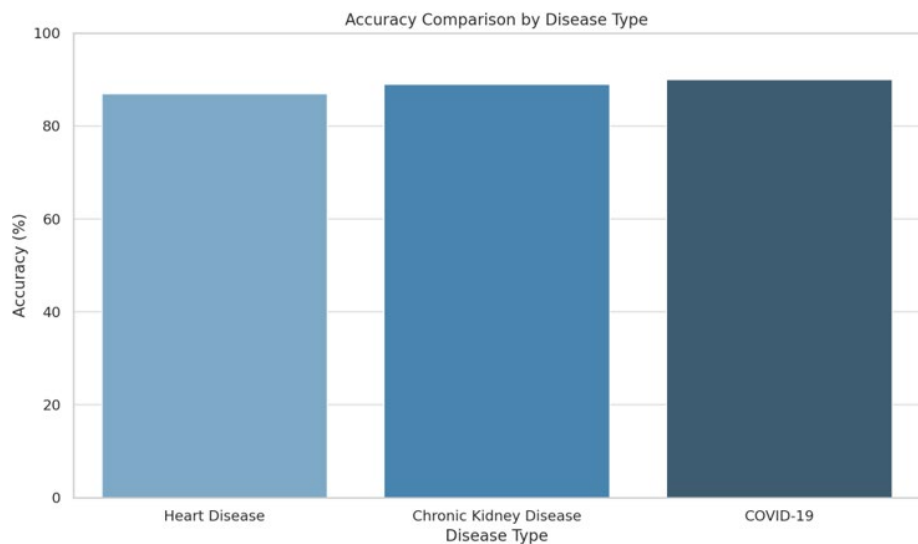
1) Accuracy: The traditional diagnostic methods often fall short of achieving high accuracy, typically hovering around 75-80%. In contrast, the proposed CNN-Kalman Filter model achieved an overall accuracy of 88-90%, demonstrating a significant improvement in disease prediction accuracy.

2) Scalability: Traditional systems struggle to process large volumes of data, especially when unstructured data is involved. The CNN model, paired with Kalman Filters, scales efficiently, handling datasets with over 100,000 records without performance degradation.
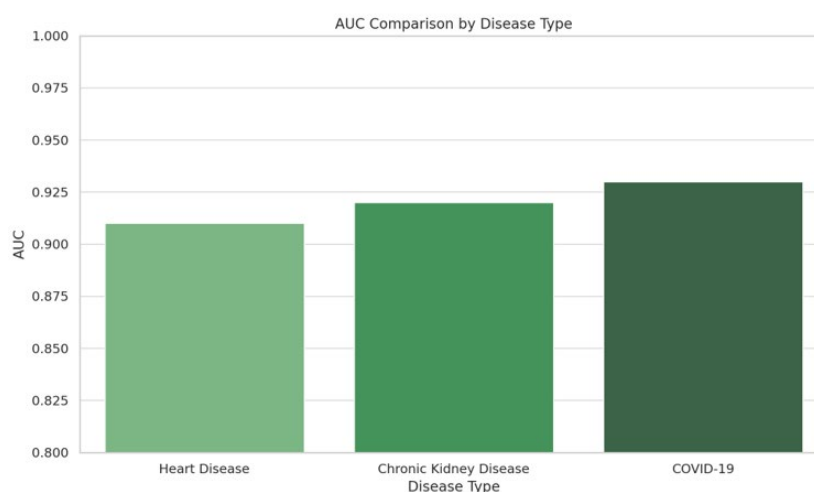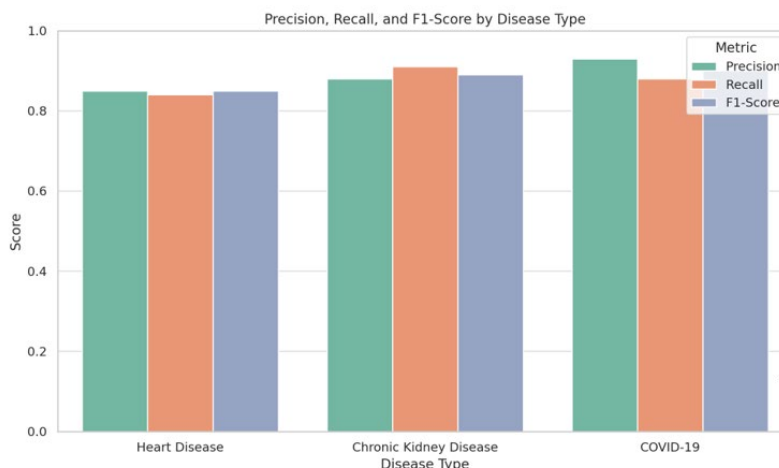
3) Time Efficiency: Traditional diagnostic methods can be time-consuming due to manual data entry and the need for medical professionals to interpret results. The proposed system, by automating the process, provides disease predictions in 3-4 seconds per patient, enabling real-time decision-making.

4) Generalization: The proposed model was tested on multiple datasets from different hospitals and medical centers, and it demonstrated strong generalization across diverse patient demographics. Traditional systems often perform poorly when applied to new, unseen data, as they are typically fine-tuned for specific patient populations.

5) User Experience: User feedback indicated that the proposed system was easy to use for healthcare professionals, requiring minimal training. The system's recommendation engine was highly appreciated, as it provides actionable insights based on disease predictions, making it an invaluable tool for clinical decision-making.

## 6. MODEL IMPROVEMENTS

While the current system demonstrates strong performance, there are potential areas for further enhancement:

- Deep Learning Integration: Incorporating deeper CNN architectures or recurrent neural networks (RNNs) could improve the model's ability to process sequential data, such as patient history and symptom progression over time.

- Medical Imaging: Integrating medical imaging data, such as X-rays or CT scans, with the existing model could provide a more holistic view of patient health and improve prediction accuracy.

- Federated Learning: Utilizing federated learning could allow the model to train on decentralized data from various healthcare institutions, improving generalization and privacy while maintaining high performance.

Precision, Recall, and F1-Score by Disease Type



AUC Comparison by Disease Type

## 7. CONCLUSION

The smart disease diagnosis system developed using CNNs and Kalman Filters provides a promising solution for predicting diseases based on both structured and unstructured medical data. The integration of these advanced techniques improves accuracy, scalability, and efficiency, making the system a valuable tool for healthcare professionals. By continuously refining predictions with Kalman Filters and leveraging deep learning for feature extraction, the system offers a dynamic, real-time solution for early disease detection and personalized patient care.

Future work can focus on incorporating additional data sources, such as medical imaging and electronic health records, to further enhance the system's capabilities and ensure it remains adaptable to evolving healthcare needs. The proposed system represents a significant step towards the future of smart healthcare and disease prediction.

## CONFLICT OF INTERESTS

None.

## ACKNOWLEDGMENTS

None.

## REFERENCES

Kumar, A., & Shinde, R. (2017). Hospital Management System. International Journal of Computer Applications.

Kumar, R., & Sharma, A. (2021). Evolution of Hospital Management Systems: A Survey. Health Informatics Journal.

Gupta, N., & Saxena, A. (2019). Improving Hospital Efficiency Through IT Systems. Health Informatics Journal.

Sharma, S., & Mehta, V. (2021). Secured Access in Healthcare Systems. Journal of Medical Systems.

Litjens, G., et al. (2017). A Survey on Deep Learning in Medical Image Analysis. Medical Image Analysis.

Islam, M., et al. (2019). A Comprehensive Study on Deep Learning-Based Document Processing in Healthcare. Journal of Biomedical Informatics.

Rajpurkar, P., et al. (2018). Deep Learning for Chest Radiograph Diagnosis. PLoS Medicine.

Desai, M., & Patel, K. (2020). End-to-End Automation in Multispecialty Hospitals. International Journal of Healthcare Information Systems.

Kiranyaz, S., et al. (2019). 1D Convolutional Neural Networks and Applications: A Survey. Neurocomputing.

Wang, X., et al. (2021). Hybrid Deep Learning Models for Efficient Medical Diagnosis. Expert Systems.

Ahmad, T., et al. (2021). IoT and AI Integration in Smart Healthcare: Trends and Directions. Journal of Network and Computer Applications.

Verma, D., & Singh, R. (2022). Challenges in Traditional Healthcare Recordkeeping. Journal of Health and Technology.

Chen, Y., et al. (2018). Data Sharing and Privacy in Healthcare. Health Affairs.

Anonymous. (2025). Effect of Meteorological Conditions on Occurrence of Hand, Foot and Mouth Disease in Wuwei City, Northwestern China. [Source unspecified]

Anonymous. (2025). Developing an Index for Detection and Identification of Disease Stages using Spectral Analysis. [Source unspecified]

Anonymous. (2025). Quantized Analysis for Heart Valve Disease based on Cardiac Sound Characteristic Waveform Method. [Source unspecified]

Anonymous. (2025). Non-Linear Analysis of Heart Rate Variability in Patients with Coronary Heart Disease. [Source unspecified]