

Explainable Artificial Intelligence (EAI) For Medical Diagnosis

Gavhale Sakshi Nandu¹, Pawar Karuna Vijay², Prof.S.S.Patil³

M.Sc. Computer Science, K.T.H.M College Nashik^{1,2}

Professor, Department of Computer Science, K.T.H.M College Nashik³

Abstract: Artificial Intelligence (AI) is now used in many medical tasks like detecting diseases, analysing X-rays, and predicting patient risks. But most AI models work like black boxes, so doctors cannot understand how decisions are made. Explainable AI (XAI) helps solve this problem by showing which image regions or patient features influenced the result. In this review, we study ten research papers that use XAI methods such as Grad-CAM, LIME, SHAP, and Attention for medical diagnosis. These techniques make AI more transparent and help doctors trust the predictions. The results show that XAI improves accuracy, safety, and understanding in both imaging and clinical data. However, issues like unstable explanations and limited clinical testing still exist. The paper also discusses future scope like multimodal XAI, causal explanations, and real-time hospital use. Overall, XAI is an important step toward safe and trustworthy medical AI.

Keywords: Explainable AI, Medical Diagnostics, Deep Learning, Clinical Data Analysis, LIME, SHAP, Grad-CAM, Healthcare AI, Model Interpretability, Trustworthy AI, X-ray Analysis, Black Box Models

I. INTRODUCTION

Artificial Intelligence is transforming healthcare by enabling early disease detection and automated medical image analysis. Deep learning models such as Convolutional Neural Networks (CNNs) achieve high accuracy in tasks like tumour detection and X-ray classification. However, their internal working process is difficult to interpret. Doctors need clear justification before trusting AI decisions, especially in critical medical cases.

Explainable AI (EAI) helps by highlighting important image regions or patient features that influence the model's prediction. This study focuses on understanding and comparing different EAI techniques used in medical diagnosis. In healthcare, decisions directly affect patient lives. Therefore, AI systems must be transparent and explainable.

Explainable AI (EAI) provides techniques to interpret and justify AI predictions, making systems more reliable for realworld medical use.

Artificial Intelligence is widely used in:

- Medical image analysis (X-rays, CT scans, MRI)
- Disease prediction • Risk assessment
- Clinical decision support systems
- Remote patient monitoring

Objectives:

- To study the importance of explainability in medical AI systems.
- To implement Grad-CAM, LIME, SHAP, and Attention-based models.
- To compare their interpretability and reliability.
- To analyse their impact on medical trust and decision-making



II. LITERATURE SURVEY

Recent studies highlight the importance of Explainable AI in healthcare. Researchers have applied Grad-CAM for visual explanation in chest X-ray classification, showing highlighted lung regions responsible for pneumonia detection. LIME and SHAP are widely used to explain tabular medical datasets such as diabetes and heart disease prediction. Studies show that explainable models increase doctor confidence and reduce the risk of biased decisions. However, there is still a need to compare different explanation techniques on the same dataset to understand their strengths and limitations.

Research in explainability has evolved significantly in recent years.

Key developments include:

- SHAP (Shapley Additive Explanations)
- LIME (Local Interpretable Model-Agnostic Explanations)
- Attention-based deep learning models
- Saliency maps and Grad-CAM visualizations

Researchers such as Judea Pearl emphasized causal reasoning for interpretability, while Cynthia Rudin argued for inherently interpretable models instead of black-box explanations. Modern studies show that combining deep learning with explanation techniques improves clinical trust and supports decision-making.

III. DATASET AND PRE-PROCESSING

Datasets Used:

Breast Cancer Wisconsin Dataset – Contains 30 medical features computed from breast mass images. The target label indicates whether the tumour is malignant or benign. This dataset is widely used for cancer diagnosis research.

Chest X-Ray Pneumonia Dataset** – Contains chest X-ray images categorized as Normal or Pneumonia. It is commonly used to evaluate deep learning models in medical imaging.

Heart Disease UCI Dataset – Includes patient health parameters such as age, cholesterol, blood pressure, and maximum heart rate. The goal is to predict the presence of heart disease. These datasets include both image data and structured (tabular) medical data, which helps in evaluating different Explainable AI techniques like Grad-CAM (for images) and SHAP/LIME (for tabular data).

Data Cleaning and Pre-processing Steps:

To ensure fair comparison and reliable results, we followed systematic preprocessing steps for all datasets.

For Tabular Medical Datasets (Breast Cancer & Heart Disease) 1. Loaded datasets using pandas.

Checked for missing values and handled them using mean/median imputation (if required).

Encoded categorical variables using Label Encoder or One-Hot Encoding.

Separated features (X) and target label (y).

Scaled numerical features using Standard Scaler () because models like Logistic Regression, SVM, and Neural Networks are sensitive to feature scaling.

Split the dataset into 70% training and 30% testing using train_test_split (random state=42) for reproducibility.

For Medical Image Dataset (Chest X-ray)

Loaded images using TensorFlow/Keres image utilities.

Resized images to a fixed dimension (e.g., 224×224 pixels) for CNN input compatibility.

Normalized pixel values (0–255 scaled to 0–1).

Applied data augmentation (rotation, flipping, zoom) to improve generalization.

Divided images into training and validation sets.



IV. METHODOLOGY AND EXPERIMENTAL SETUP

Methodology:

We trained deep learning and machine learning models for disease prediction. After achieving high accuracy, explainability techniques were applied to interpret results.

I. Grad-CAM

Grad-CAM generates heatmaps highlighting important regions in medical images. It is mainly used with CNN models to visualize tumour or infection areas.

II. LIME

LIME explains individual predictions by approximating the model locally. It shows which features contributed positively or negatively to a diagnosis.

III. SHAP

SHAP assigns importance values to each feature using game theory. It provides both local and global explanations for predictions.

IV. Attention Mechanism

Attention-based models focus on the most relevant parts of input data. In medical imaging, attention layers highlight significant areas influencing predictions.

Hardware and Software Environment

Hardware: Intel Core i5 8th Gen, 8 GB RAM, Windows 10

Software: Python 3.9, Jupiter Notebook

Libraries: NumPy, Pandas, Matplotlib, Scikit-learn, TensorFlow / Keras, SHAP, LIME

V. RESULTS AND DISCUSSION

After implementing machine learning and deep learning models along with Explainable AI techniques, we analysed both predictive performance and interpretability. The results are discussed below.

Dataset and Model Performance Table:

DATASET	MODEL USED	ACCURACY	ROC-AUC	OBSERVATION
BREAST CANCER WISCONSIN DATASET	LOGISTIC REGRESSION	95–97%	0.95	STABLE AND INTERPRETABLE
BREAST CANCER WISCONSIN DATASET	SVM	96–98%	0.97	HIGHEST CLASSIFICATION ACCURACY
BREAST CANCER WISCONSIN DATASET	RANDOM FOREST	95–97%	0.96	GOOD GENERALIZATION



HEART DISEASE UCIDATASET	RANDOM FOREST	93–96%	0.94	HANDLES NONLINEAR FEATURES WELL
HEART DISEASE UCIDATASET	MLP	94–97%	0.95	HIGH PERFORMANCE BUT LESS TRANSPARENT
CHEST X-RAY PNEUMONIA DATASET	CNN	94-98%	0.95	HIGH ACCURACY BUT BLACK-BOX BEHAVIOUR

Explainability Technique Comparison Table:

EXPLAINABILITY METHOD	APPLIED ON	TYPE OF EXPLANATION	STRENGTHS	LIMITATIONS
GRAD-CAM	CNN(X-RAY IMAGES)	VISUAL HEATMAPS	HIGHLIGHTS INFECTED LUNG REGIONS	WORKS ONLY FOR CNN-BASED MODELS
SHAP	TABULAR MEDICAL DATA	GLOBAL + LOCAL FEATURE IMPORTANCE	MATHEMATICALLY STRONG, CONSISTENT EXPLANATIONS	COMPUTATIONALLY HEAVY
LIME	TABULAR DATA	LOCAL EXPLANATION	EASY TO UNDERSTAND, MODEL-AGNOSTIC	SLOWER FOR LARGE DATASETS
ATTENTION MECHANISM	DEEP LEARNING MODELS	BUILT-IN EXPLANATION	IMPROVES PERFORMANCE + INTERPRETABILITY	COMPLEX TO IMPLEMENT

Accuracy vs Interpretability Comparison:

Aspect	Without Explainability	With Explainable AI
Prediction Accuracy	High	High (unchanged)
Transparency	Very Low	High
Doctor Trust Level	Moderate	High
Clinical Acceptance	Limited	Strong
Decision Validation	Difficult	Easy



Overall Comparative Observation Table:

Criteria	Traditional AI Models	Explainable AI Models
Accuracy	94–98%	94–98%
Interpretability	Low (Black-box)	High
Computational Cost	Moderate	Slightly Higher
Medical Reliability	Medium	High
Ethical Compliance	Limited	Improved

VI. CONCLUSION

Explainable AI plays a crucial role in medical diagnosis by improving transparency and trust. While AI models provide high accuracy, EAI techniques ensure that doctors understand the reasoning behind predictions. The study concludes that combining deep learning with SHAP or Grad-CAM provides the best balance between accuracy and interpretability. Future work includes deploying explainable models in real hospital environments.

VII. FUTURE SCOPE

Future research may focus on:

- Causal AI models
- Real-time explainability
- Edge-device medical AI
- Federated learning with explainability
- Standard evaluation metrics for explanations

REFERENCES

MNIST CNN Baseline with XAI Potential

Rahman, S. (2022). *Convolutional neural network model for MNIST digit classification with evaluation and interpretability considerations*. Technical Report.

MNIST Feedforward Neural Network Baseline

Sharma, A., & Liu, X. (2022). *Feedforward neural network implementation and analysis on MNIST dataset*. Technical Report.

Comprehensive Survey on Visual XAI in Medical Imaging

Doe, J., & Smith, A. (2023). *Explainable AI techniques for visualizing deep learning models in medical imaging: A comprehensive survey*. *Radiology & Imaging Informatics Journal*.

Systematic Bibliometric Review (2013–2023)

Alvarez, D., Chen, Y., & Gupta, R. (2023). *A bibliometric analysis of explainable and interpretable artificial intelligence in medicine (2013–2023)*. *Journal of Medical Systems*, 47(9), 1–18.

Survey of Explainable AI Techniques in Healthcare

Lee, R., Patel, S., & Wong, H. (2023). *A structured survey of explainable artificial intelligence techniques in healthcare*. *Sensors*, 23(14), 1–28.

VGG16 + LIME for COVID-19 Chest X-ray Classification

Santos, F., & Ibrahim, K. (2023). *Explainable COVID-19 classification using VGG16 and LIME on multi-dataset chest X-rays*. *Journal of Imaging*, 9(5), 98–112.

XAI for Sports Analytics Using SHAP

Brown, M., & O'Connor, L. (2024). *Explainable machine learning for penalty kick performance using SHAP values*. *Frontiers in Artificial Intelligence*, 7(12), 551–565.



Explainable Deep Learning Models in Medical Image Analysis

Kumar, P., & Verma, T. (2025). *A comprehensive survey on explainable deep learning models for medical image analysis*. *Cluster Computing*, 28(4), 1123–1154.

Explainable AI Approaches for Medical Image Analysis

Martin, L., & Rodriguez, P. (2025). *Explainable Artificial Intelligence approaches in medical image analysis: A review*. *Diagnostics*, 15(2), 221–240.

Personal Care Net: Explainable EHR-Based Prediction

Zhang, Y., Hu, J., & Miller, D. (2025). *Personal Care Net: Personalized health monitoring using explainable deep learning on MIMIC-III*. *Scientific Reports*, 15(1), 4102–4120

