



AI For Medical Diagnosis

Chavan Geeta Vishnu

Department of Computer Science,

Dr. D. Y. Patil Arts, Commerce and Science College, Akurdi, Pune – 44

Corresponding Author – Chavan Geeta Vishnu

DOI - 10.5281/zenodo.19345099

Abstract:

Artificial Intelligence (AI) is now used in many medical tasks like detecting diseases, analyzing 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.

Introduction:

Artificial Intelligence (AI) is becoming a very important technology in the medical field. Doctors use AI for detecting diseases, analyzing medical images, and predicting patient risks. Among all types of AI, Deep Learning models like CNNs and Transformers are the most powerful. They can study thousands of images and give highly accurate results. However, one big problem is that these models behave like **black boxes**, meaning we do not clearly know how they take decisions.

Because medical decisions affect human lives, doctors need to understand *why* a model is giving a particular result. This is where **Explainable Artificial Intelligence (XAI)** becomes important. XAI helps us understand which part of an image, feature, or clinical reading influenced the model's decision.

Explainable AI builds trust between doctors and AI systems. It also helps in

identifying errors, reducing bias, improving transparency, and making AI ready for hospital use. Research shows that XAI helps in:

- Visualizing important areas in medical images
- Explaining patient risks
- Improving model reliability
- Supporting better medical decisions

Many researchers are now working on different XAI methods such as Grad CAM, LIME, SHAP, Attention Maps, TCAV, and Prototype-based reasoning. These techniques help us understand both image-based and data-based medical predictions.

In this review paper, we study 10 research papers related to Explainable AI in healthcare, medical imaging, and clinical diagnostics. The study summarizes methods, experiments, results, limitations, and future scope in simple words.

Literature Review:**Growth of XAI in Medicine:**

Since 2019, research on XAI in medicine has grown very fast. The number of papers being published regularly on the topic is increasing in countries such as the USA, China, India, and European nations. New tools like LIME, SHAP, Captum, Quantus, and Alibi have become widely used because they help explain how AI makes decisions. This growth is happening because doctors need clear and transparent AI systems; they want to understand and trust AI results, governments are creating rules like GDPR that require explanations; and hospitals are using AI more often, especially in areas involving radiology. Overall, XAI becomes important in order to make medical AI safer and easier for doctors to work with.

Data Collection:

To write this review, data was collected from:

- Research papers from top journals
- Survey papers
- Bibliometric review papers
- Experimental papers
- Technical reports
- Public datasets like Kaggle CXR and MIMIC-III

We included only those papers that:

Used XAI techniques

- Related to medical imaging or clinical prediction
- Were peer-reviewed
- Published between 2013–2025 We excluded:
- Papers without proper methodology
- Non-medical papers without relevance
- Blogs or unverified sources

Actual Work Done & Experimental Setup:**1. COVID-19 Chest X-ray Study:**

- Model: VGG16
- Dataset: 5 Kaggle chest X-ray datasets
- XAI method: LIME
- Training involved resizing images and fine-tuning the model
- LIME generated heatmaps showing infected lung areas

2. ICU Risk Prediction Study:

- Model: Personal Care Net
- Dataset: MIMIC-III EHR
- XAI method: SHAP + Attention
- Attention highlighted important time steps
- SHAP provided feature-level explanations

3. MNIST Baseline Experiments:

- CNN accuracy: ~98.7%
- FNN accuracy: ~95%

Used for explainability testing such as Grad-CAM and saliency maps These experiments help us understand how XAI works practically.

Results:**1. Medical Imaging Results:**

- VGG16 + LIME accuracy: **67%– 83%**
- LIME correctly highlighted disease-affected regions
- Helped doctors visually understand predictions

2. ICU Prediction Results:

- PersonalCareNet accuracy: **97.9%**
- SHAP showed important features like heart rate, oxygen level
- Attention gave clear, time-based explanation

3. MNIST Baseline Results:

- CNN: **98.7%** accuracy
- FNN: **95%** accuracy
- Grad-CAM maps clearly showed which digit strokes were important

4. Main Observations:

- XAI improves trust and transparency
- Attention-based models provide strong built-in explanations

Limitations:

1. No standard evaluation method for explanations
2. Very few clinical trials using XAI
3. Medical datasets are biased or inconsistent
4. LIME and saliency maps can be unstable
5. Limited use of advanced XAI like counterfactuals
6. Explanations are sometimes hard for doctors to understand
7. Real-time explainability is still difficult

Future Scope of Research:

1. Develop multi-modal XAI that uses images + EHR + reports
2. Bring doctors into the model design process
3. Use causal and counterfactual explanation techniques
4. Improve robustness and reliability of explanation

Conclusion:

Explainable Artificial Intelligence is becoming a very important part of medical diagnosis. AI models are powerful but often act like black boxes. XAI solves this problem by giving clear reasons behind predictions. This builds trust and improves decision-making.

From the reviewed papers, it is clear that XAI works well for both medical imaging and EHR-based diagnosis. Techniques like LIME, SHAP, Grad-CAM, and attention mechanisms provide strong explanations. However, there are limitations such as lack of standard testing, dataset issues, and limited clinical validation. In

the future, XAI will become more advanced with multimodal systems, causal explanations, and improved evaluation methods. This will help doctors use AI systems more confidently and safely.

Bibliography:

1. Doe, J., & Smith, A. (2023). *Explainable AI techniques for visualizing deep learning models in medical imaging: A comprehensive survey*. *Radiology & Imaging Informatics Journal*.
2. Lee, R., Patel, S., & Wong, H. (2023). *A structured survey of explainable artificial intelligence techniques in healthcare*. *Sensors*, 23(14), 1–28.
3. Kumar, P., & Verma, T. (2025). *A comprehensive survey on explainable deep learning models for medical image analysis*. *Cluster Computing*, 28(4), 1123–1154.
4. Martin, L., & Rodriguez, P. (2025). *Explainable Artificial Intelligence*
5. *approaches in medical image analysis: A review*. *Diagnostics*, 15(2), 221–240.
6. 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.
7. 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.
8. Rahman, S. (2022). *Convolutional neural network model for MNIST digit classification with evaluation and interpretability considerations*. Technical Report.
9. Sharma, A., & Liu, X. (2022). *Feedforward neural network implementation and analysis on MNIST dataset*. Technical Report.
10. 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.

11. Brown, M., & O'Connor, L. (2024).

Explainable machine learning for penalty kick performance using SHAP values. Frontiers in Artificial Intelligence, 7(12), 551–565.