

Explainable Artificial Intelligence in Healthcare: Methods, Challenges, and a Conceptual Framework

Poonam Sahibani , Anilkumar Munani

IT Engineering Department, Sardar Patel College of Engineering and Technology, Anand, India

Email: poonamsahibani10@gmail.com

Computer Engineering Department, Bhailalbhai & Bhikhabhai Institute of Technology, Anand, India

Email: almunani@bbit.ac.in

Abstract:

Explainable Artificial Intelligence (XAI) is vital for enabling trust, transparency, and accountability in medical decision support systems. This paper provides a conceptual survey of core XAI techniques—such as LIME, SHAP, Integrated Gradients, Grad-CAM, counterfactual explanations, and concept activation—and analyzes their theoretical foundations, strengths, and limitations in healthcare settings. We further discuss major challenges unique to clinical deployment: explanation validation, human-AI interaction, trust calibration, and regulatory compliance. To guide researchers and practitioners, we propose a conceptual framework mapping XAI methods by clinical task (diagnosis, prognosis, treatment decision), data modality (tabular, imaging, multimodal), and explanation scope (local vs global). We conclude with recommendations for future work: benchmark datasets for explanation evaluation, multimodal interpretation, human-centered studies, and integration with causal reasoning.

Keywords — Explainable AI, Interpretability, LIME, SHAP, Grad-CAM, Counterfactual Explanations, Healthcare, Trust, Clinical Decision Support.

I. INTRODUCTION

Artificial Intelligence (AI) holds great promise in healthcare—ranging from medical image analysis (e.g. tumor detection, segmentation) to predictive models on electronic health records (EHRs) for disease risk or patient outcome forecasting. However, the “black-box” nature of many high-performance models (especially deep learning) poses a significant barrier to adoption in clinical settings, where decisions must be interpretable, auditable, and justifiable to medical experts and regulators. Explainable AI (XAI) aims to reconcile model complexity with human understanding by providing insights into why a model makes a particular decision or prediction.

In healthcare, the stakes are high: erroneous or uninterpretable predictions can risk patient safety,

raise legal liability, and erode clinician trust. Hence, explainability is not a luxury but a necessity. The demand for transparent models is reinforced by regulatory trends (e.g. GDPR’s “right to explanation,” FDA oversight of AI/ML in medical devices) and ethical imperatives of fairness and accountability.

This paper seeks to provide a conceptual and technical overview of XAI methods applied to healthcare use-cases, highlighting their theoretical principles, suitability for different data modalities, and the unique challenges of clinical integration. We also propose a framework to guide selection and combination of XAI methods in healthcare settings.

The rest of this paper is organized as follows: Section II reviews key XAI methods and their theoretical foundations. Section III discusses major challenges in medical deployment. Section IV

presents the proposed framework mapping XAI techniques by task, data type, and interpretability scope. Section V outlines comparative analysis and best practices. Section VI suggests future directions. Section VII concludes.

METHODS OF EXPLAINABLE AI IN HEALTHCARE

All paragraphs must be indented. All paragraphs must be justified, i.e. both left-justified and right-justified.

There is a broad taxonomy of XAI methods; here we cover major classes that are relevant to medical applications: perturbation/surrogate-based, Shapley-based, gradient / saliency-based, counterfactual / example-based, and concept-level / global methods.

A. Perturbation / Surrogate-based Methods (e.g., LIME)

LIME (Local Interpretable Model-Agnostic Explanations) is a popular approach that approximates the local decision boundary of a black-box model by sampling perturbed inputs and fitting an interpretable surrogate (e.g., linear model) weighted by proximity [4]. The coefficients of the surrogate model indicate feature importance locally. LIME is model-agnostic and works for tabular, text, or imaging features (after embedding). Its advantages include flexibility and ease of implementation. However, its fidelity is not guaranteed when the local decision boundary is highly nonlinear, and explanations may be unstable (different perturbations can yield different explanations). Also, it offers local explanations only, and lacks a global view of the model.

B. Shapley-based Methods (SHAP)

SHAP (SHapley Additive exPlanations) is based on cooperative game theory: every feature is assigned a Shapley value representing its average marginal contribution across all possible subsets [5]. SHAP defines an additive explanation model that decomposes the prediction into feature contributions summing to the output (plus a baseline). It offers strong axiomatic guarantees (local accuracy, consistency). Versions such as TreeSHAP, DeepSHAP optimize computations for tree ensembles or neural networks. In healthcare, SHAP is widely used for both local and aggregated global explanations, e.g. highlighting which biomarkers generally contribute most to risk predictions [6]. The main drawbacks are

computational cost (exact Shapley is exponential) and the sensitivity to the choice of baseline or reference distribution.

C. Gradient / Saliency-based Methods

For differentiable models (especially neural networks), gradient-based explanations compute partial derivatives of output with respect to inputs. Some variants:

Integrated Gradients (IG): Compute the path integral of gradients from a baseline input to the actual input [7]. This yields attributions satisfying completeness (sum equals model output difference).

Guided Backpropagation / Saliency maps: Basic saliency uses the gradient magnitude; guided backprop suppresses negative flows to clean up maps.

Grad-CAM / CAM: Originally for CNNs, Grad-CAM computes gradients at last convolutional layers to generate class activation maps, highlighting spatial regions influencing prediction [8]. Grad-CAM is extensively used in medical imaging (MRI, CT, pathology) to show heatmaps of relevant structures [9].

These methods are efficient (one backward pass) and produce intuitive visual explanations. However, they face issues like noise, gradient saturation, sensitivity to baseline, and inability to explain non-differentiable models. Also, explanations are local.

D. Counterfactual and Example-Based Explanations

Counterfactual explanations answer “what minimal perturbation would flip the model’s decision?” (e.g. “If glucose had been 10 mg/dL lower, the model would predict no diabetes”) [10]. These are intuitive and actionable. Generating valid counterfactuals under clinical constraints (e.g. realistic feature changes) is nontrivial. Prototype or example-based explanations instead point to similar instances (nearest neighbors) as reference points. Anchors extend LIME by building local rules that “anchor” the prediction robustly [11]. These methods are local, but conceptually compelling for clinicians.

E. Concept-Level and Global Explanation Methods

These methods provide higher-level explanations or global model insight. For example:

TCAV (Testing with Concept Activation Vectors): Measures sensitivity of model class

output to user-defined human concepts (e.g. “texture of lesion”) [12].

Surrogate models / Rule extraction: Train a simpler interpretable model (e.g. decision tree, rule lists) approximating the black box; the surrogate gives a global interpretability perspective.

Global SHAP summaries: Aggregate local SHAP values (mean absolute) across many instances to rank feature importance globally.

These methods offer global interpretability, helpful for auditing, regulatory reporting, and understanding model trends.

II. CHALLENGES OF XAI IN HEALTHCARE

Deploying XAI in clinical settings faces several domain-specific challenges:

A. Trust and Clinical Validation

Clinicians must trust that the provided explanation is faithful to the model logic and medically plausible. However, some explanations (e.g. saliency maps highlighting irrelevant regions) may mislead [13]. Without empirical evaluations (human-in-the-loop studies), trust remains weak.

B. Evaluation of Explanations

There is no consensus on metrics for explanation quality. Proposed metrics include faithfulness (does removing important features change predictions), infidelity, sensitivity, and stability. In medical domains, comparing explanations to known ground truth (e.g. annotated lesions) is possible only in certain imaging datasets [14]. Human evaluation (clinician scoring) is necessary but expensive. Some works examine using Delete-and-Retrain or Remove-and-Retrain to test attribution fidelity [15].

C. Clinical Usability and Human-AI Interaction

Explanations must be concise, intelligible, and integrated into workflow. Overly detailed explanations can overwhelm clinicians. The “human in the loop” paradigm is often necessary: experts can override or adjust AI suggestions [16]. Design of explanation UI (heatmaps, textual summaries) matters significantly.

D. Data Complexity, Multimodality, and Model Scale

Clinical data is often multimodal (EHR + imaging + genomics). Combining explanations across modalities is underexplored. Large models (e.g. multi-modal transformers) complicate explanation generation and make real-time explanation difficult.

E. Ethical, Privacy, and Regulatory Constraints

Explanations must not leak private data or violate patient confidentiality. Regulatory bodies like FDA and GDPR may require “meaningful explanations” for automated decisions. Explaining AI in a way that meets regulatory standards is nontrivial [3]. Also, incorrect explanations could lead to liability issues.

IV. CONCEPTUAL FRAMEWORK FOR HEALTHCARE XAI

We propose a three-axis framework to guide selection and evaluation of XAI methods in clinical contexts:

Clinical Task / Use-Case (Diagnosis, Prognosis, Treatment)

Data Modality (Tabular / Time-series, Imaging, Multimodal)

V. COMPARATIVE INSIGHTS AND DISCUSSION

Model-agnostic vs model-specific: SHAP / LIME are broadly applicable but may trade off fidelity; gradient-based methods exploit model structure with higher efficiency.

Local vs global: Local explanations help individual decisions, while global methods aid model audit and regulatory transparency. Hybrid systems combining both are promising.

Fidelity vs interpretability trade-off: Simpler explanations may not capture model intricacies; more faithful ones may be complex. Always validate explanations (e.g. via perturbation tests).

Concept alignment: Explanations aligned to medical concepts (e.g. “edema,” “contrast

enhancement”) improve interpretability. Methods like TCAV help bridge this gap.

Efficiency: Real-time clinical use demands fast explanation methods; gradient methods are advantageous, but scalable approximations of SHAP are beneficial.

Human-centered design: Explanation UI, modality (text, visuals), and user training are critical but often under-studied in literature.

VI. FUTURE DIRECTIONS

Benchmark datasets with ground-truth explanations (e.g. annotated lesion maps, expert decision rationales) to enable objective evaluation of XAI.

Multimodal explanation methods that unify EHR, imaging, genomics, and text.

Hybrid models of built-in interpretability (self-explainable architectures) that produce explanations inherently.

Longitudinal clinical studies measuring how explanations influence clinician decisions, trust, and patient outcomes.

Regulatory-aligned explanation frameworks: standardization of explanation documentation, certification, and compliance procedures.

Causal and actionable explanations: integrate causal reasoning so that explanations suggest safe clinical interventions.

VII. CONCLUSION

Explainable AI is essential for trustworthy and safe deployment of AI in healthcare. While many methods exist (LIME, SHAP, Grad-CAM, counterfactuals, concept-based), none is universally optimal. Clinically meaningful application requires tailoring explanation techniques to specific use-cases, data modalities, and interpretability needs. Our proposed framework helps structure this alignment. However, challenges remain in validating explanations, integrating them into clinical workflows, scaling to multimodal data, and

meeting regulatory demands. Progress in benchmark datasets, human-centered evaluation, and hybrid interpretable models will be key to realizing XAI's promise in medicine.

REFERENCES

- [1] S. M. Metev and V. P. Veiko, *Laser Assisted Microtechnology*, 2nd ed., R. M. Osgood, Jr., Ed. Berlin, Germany: Springer-Verlag, 1998.
- [2] J. Breckling, Ed., *The Analysis of Directional Time Series: Applications to Wind Speed and Direction*, ser. Lecture Notes in Statistics. Berlin, Germany: Springer, 1989, vol. 61.
- [3] S. Zhang, C. Zhu, J. K. O. Sin, and P. K. T. Mok, “A novel ultrathin elevated channel low-temperature poly-Si TFT,” *IEEE Electron Device Lett.*, vol. 20, pp. 569–571, Nov. 1999.
- [4] M. Wegmuller, J. P. von der Weid, P. Oberson, and N. Gisin, “High resolution fiber distributed measurements with coherent OFDR,” in *Proc. ECOC'00*, 2000, paper 11.3.4, p. 109.
- [5] R. E. Sorace, V. S. Reinhardt, and S. A. Vaughn, “High-speed digital-to-RF converter,” U.S. Patent 5 668 842, Sept. 16, 1997.
- [6] (2002) The IEEE website. [Online]. Available: <http://www.ieee.org/>
- [7] M. Shell. (2002) IEEEtran homepage on CTAN. [Online]. Available: <http://www.ctan.org/tex-archive/macros/latex/contrib/supported/IEEEtran/>
- [8] *FLEXChip Signal Processor (MC68175/D)*, Motorola, 1996.
- [9] “PDCA12-70 data sheet,” Opto Speed SA, Mezzovico, Switzerland.
- [10] A. Karnik, “Performance of TCP congestion control with rate feedback: TCP/ABR and rate adaptive TCP/IP,” M. Eng. thesis, Indian Institute of Science, Bangalore, India, Jan. 1999.
- [11] J. Padhye, V. Firoiu, and D. Towsley, “A stochastic model of TCP Reno congestion avoidance and control,” Univ. of Massachusetts, Amherst, MA, CMPSCI Tech. Rep. 99-02, 1999.
- [12] Wireless LAN Medium Access Control (MAC) and Physical Layer (PHY) Specification, IEEE Std. 802.11, 1997.