



# Deploying Interpretable AI Models in Healthcare Diagnostics: A Case Study on Using Explainable Machine Learning for Early Disease Detection and Clinical Decision Support

Ugochukwu Ukeje<sup>1\*</sup>

<sup>1</sup>*School of Data Science and Analytics, Kennesaw State University USA*

*\*Corresponding author*

DOI: <https://doi.org/10.63680/ijstate032540.09>

## Abstract

Artificial intelligence (AI) and machine learning (ML) are transforming healthcare diagnostics by enabling faster and more accurate disease detection. However, the growing reliance on complex black box models raises critical concerns about transparency, trust, and accountability particularly in clinical settings where interpretability is essential. This study aims to evaluate interpretable AI models and explainability techniques in the context of early disease detection and integration into Clinical Decision Support Systems (CDSS). The research combines an in depth literature review of interpretable and post hoc explainable ML approaches with a comparative case study using logistic regression, Explainable Boosting Machines (EBM), and XGBoost with SHAP explanations applied to the UCI Heart Disease dataset. Models were assessed using metrics such as accuracy, AUC ROC, precision, explanation clarity, and fairness across demographic subgroups. The results reveal that while XGBoost offers superior predictive performance, EBM achieves a more optimal balance between accuracy, transparency, and clinical usability. SHAP explanations provided valuable local and global insights but required careful interface design for practical deployment. The study highlights the ongoing trade off between interpretability and performance and emphasizes the importance of human centered, trustworthy AI for clinical adoption. These findings offer actionable insights for clinicians, developers, and policymakers working to integrate interpretable AI models into real world diagnostic workflows.

**Keywords:** Interpretable Machine Learning, Explainable AI, Clinical Decision Support Systems, Early Disease Detection, SHAP, LIME, Healthcare AI

## 1. Introduction

Artificial Intelligence (AI) has emerged as a transformative force in healthcare diagnostics, offering novel solutions to longstanding challenges in disease detection, clinical decision making, and patient management. Through advanced algorithms and machine learning (ML) techniques, AI systems are now capable of analyzing vast and complex medical datasets ranging from electronic health records (EHRs) and laboratory results to radiological and pathological images with unprecedented speed and precision. These capabilities have

significantly enhanced diagnostic accuracy, enabled earlier identification of diseases, and streamlined workflows in clinical environments (Topol, 2019; Esteva et al., 2017). As a result, healthcare institutions worldwide are increasingly incorporating AI into clinical decision support systems (CDSS), aiming to assist practitioners in making timely and evidence based medical decisions.

The proliferation of machine learning models in healthcare diagnostics, particularly those based on deep neural networks, ensemble classifiers, and support vector machines, has brought about impressive improvements in predictive performance. For instance, convolutional neural networks (CNNs) have demonstrated expert level proficiency in diagnosing dermatological conditions and interpreting medical images (Rajpurkar et al., 2018). However, despite their effectiveness, these models often function as “black boxes” producing accurate predictions without revealing the underlying rationale. This opacity raises critical concerns for clinicians and stakeholders who are responsible for interpreting, validating, and acting upon the AI generated recommendations. In clinical settings, where decisions can have life altering consequences, the lack of transparency undermines trust, limits adoption, and challenges the accountability of automated systems (Doshi Velez & Kim, 2017; Tonekaboni et al., 2019).

To foster meaningful integration of AI into clinical workflows, interpretability has become a central requirement. Interpretable AI refers to systems whose decisions can be understood and traced by human users particularly clinicians, patients, and regulators ensuring that outputs are explainable, justifiable, and actionable. Unlike purely performance driven metrics such as accuracy or AUC, interpretability empowers practitioners to validate AI based diagnoses, identify biases, and make informed judgments under uncertainty (Lipton, 2018). Moreover, regulatory frameworks such as the European Union’s Artificial Intelligence Act, the U.S. Food and Drug Administration (FDA) guidelines for Software as a Medical Device (SaMD), and the Health Insurance Portability and Accountability Act (HIPAA) increasingly mandate transparency, fairness, and accountability in the deployment of AI systems in healthcare (European Commission, 2021; FDA, 2022). In this context, explainable machine learning using techniques such as SHAP (Shapley Additive Explanations), LIME (Local Interpretable Model Agnostic Explanations), and Class Activation Maps (CAMs) has garnered significant attention as a bridge between high performance models and real world clinical applicability.

This study aims to conduct a critical review and comparative analysis of interpretable AI models and post hoc explainability techniques that have been deployed in healthcare diagnostics, with a specific focus on early disease detection and clinical decision support. By examining how different approaches address the interpretability performance trade off, the paper explores both the technical efficacy and practical usability of explainable ML in clinical environments. As part of this analysis, a case study is presented to illustrate how interpretable models have been integrated into diagnostic workflows, evaluating their impact on clinician trust, diagnostic confidence, and decision accountability. The findings are contextualized within current ethical and regulatory frameworks to highlight best practices and ongoing challenges.

The remainder of this paper is organized as follows: Section 2 provides a conceptual background on machine learning in healthcare and the foundations of interpretability and explainability. Section 3 offers a taxonomy of interpretable AI models and techniques. Section 4 reviews current applications in early disease detection. Section 5 presents a case study comparing model performance and interpretability. Section 6 discusses implications, trade offs, and ethical concerns. Finally, Sections 7 and 8 address challenges, future research directions, and concluding remarks.

## 2.0 Background and Conceptual Foundations

The integration of machine learning (ML) into healthcare diagnostics has profoundly altered the landscape of modern medicine. With the increasing availability of large scale medical datasets ranging from high resolution imaging to longitudinal electronic health records (EHRs) ML techniques have enabled the automation and enhancement of diagnostic procedures that were traditionally reliant on human expertise. In particular, ML has been instrumental in improving the interpretation of medical images, such as X rays, magnetic resonance imaging (MRI), computed tomography (CT), and histopathological slides, often achieving performance comparable to, or even surpassing, experienced clinicians (Esteva et al., 2017; Litjens et al., 2017). In predictive diagnostics, ML models analyze structured and unstructured EHR data to identify early indicators of diseases such as diabetes, sepsis, and cardiovascular conditions well before they manifest clinically (Miotto et al., 2016). Moreover, ML excels at recognizing subtle patterns in multivariate datasets, enabling early disease detection and risk stratification at scale.

Various model architectures have been applied in these domains, each offering a trade off between complexity and transparency. Logistic regression and decision trees are among the most interpretable models and remain widely used for risk scoring and rule based clinical applications. More advanced models, including support vector machines, random forests, and gradient boosting machines, offer improved performance while maintaining some degree of interpretability. However, deep learning models, particularly convolutional neural networks (CNNs), have dominated image based diagnostics due to their capacity to automatically extract hierarchical features from raw data (Shen et al., 2017). While these high capacity models significantly enhance diagnostic speed, accuracy, and consistency, they are frequently criticized for their “black box” nature producing outputs that are not inherently understandable by humans.

This concern brings attention to two central and often conflated concepts in AI ethics and usability: interpretability and explainability. **Interpretability** refers to the extent to which a human user can directly comprehend the internal mechanics or logic of a machine learning model. Models such as linear regression or shallow decision trees are inherently interpretable, allowing users to trace specific input features to corresponding predictions. In contrast, **explainability** pertains to post hoc methods used to elucidate the behavior of complex or opaque models after training. Techniques such as SHAP (Shapley Additive Explanations), LIME (Local Interpretable Model agnostic Explanations), and Class Activation Maps (CAMs) are designed to provide intuitive explanations for individual predictions, even when the underlying model is not inherently transparent (Ribeiro et al., 2016; Lundberg & Lee, 2017; Zhou et al., 2016). In healthcare, the distinction between these concepts is not merely academic it has profound implications for how model predictions are validated, interpreted, and acted upon by clinicians and other stakeholders.

The clinical setting presents unique challenges that amplify the importance of trust, interpretability, and accountability. **Clinical Decision Support Systems (CDSS)** are computational platforms designed to assist healthcare professionals in making evidence based decisions by providing diagnostic recommendations, treatment suggestions, or alerts (Berner, 2009). As AI becomes an integral component of modern CDSS, it is critical that these systems produce outputs that clinicians can understand, scrutinize, and trust. Unlike consumer applications where prediction errors may carry limited consequences, diagnostic errors in healthcare can result in severe harm or even death. Consequently, CDSS must comply with stringent regulatory standards, including the FDA’s guidance on Software as a Medical Device (SaMD) and the European Union’s AI Act, both of which emphasize transparency, human oversight, and accountability (FDA, 2022; European Commission, 2021). Furthermore, clinicians are ethically bound to explain medical decisions to patients a requirement that black box AI systems cannot fulfill unless their decisions are accompanied by meaningful explanations.

The failure to meet these standards poses serious risks. Opaque systems may erode clinician trust, reduce adoption rates, and introduce bias that disproportionately affects vulnerable patient groups (Ghassemi et al., 2021). Moreover, lack of explainability hinders effective model validation and auditing, which are necessary for ensuring safety and fairness across diverse clinical populations. Therefore, the deployment of interpretable AI models and explainability techniques is not a peripheral concern but a foundational requirement for safe and effective AI driven healthcare.

This paper situates its analysis within this critical context, aiming to evaluate interpretable machine learning models and explainability techniques that support early disease detection and clinical decision support. By grounding the discussion in real world clinical needs and regulatory expectations, the study underscores the essential role of interpretability in fostering the responsible deployment of AI in medicine. The subsequent sections expand on these foundational concepts by categorizing interpretable models and techniques, reviewing their practical applications, and presenting a comparative case study to assess their utility in clinical practice.

### 3.0 Taxonomy of Interpretable AI Models

In the context of healthcare diagnostics, where decision making must be both accurate and transparent, the interpretability of artificial intelligence (AI) systems becomes a crucial factor in clinical acceptance and ethical deployment. Interpretability refers to the degree to which a human can understand the reasoning behind a model's predictions. This section presents a structured taxonomy of interpretable AI models, divided into two broad categories: **intrinsic interpretability** models that are inherently transparent by design and **post hoc explainability** techniques that aim to explain the predictions of otherwise opaque models after training. The classification also highlights the relevance of each approach to specific clinical tasks and data modalities, emphasizing the trade offs between performance and interpretability.

#### 3.1 Intrinsic (Model Based) Interpretability

Intrinsic interpretability refers to models whose architecture is inherently transparent, allowing users to comprehend the logic or rules governing their predictions without additional tools. These models are especially valuable in healthcare applications where the rationale for a prediction must be clearly communicated to clinicians, patients, and regulatory bodies (Caruana et al., 2015).

**Decision trees** are widely used for their intuitive, rule based structure that mimics human decision making. Their branching logic enables clinicians to trace the influence of input features (e.g., patient age, blood pressure) on outcomes such as disease risk classification. While simple decision trees are interpretable, they may suffer from limited performance on high dimensional or noisy data, prompting the use of ensembles such as Random Forests at the cost of reduced transparency.

**Logistic regression** is another widely adopted model, particularly in epidemiological studies and clinical risk prediction (e.g., CHADS2 or Framingham scores). It offers a linear decision boundary and produces coefficients that can be interpreted as odds ratios, thus making it suitable for quantifying feature importance in tabular medical data (Hosmer et al., 2013).

**Rule based systems**, including Bayesian Rule Lists and RuleFit, represent an evolution of decision trees and associative models by formulating decisions as sets of human readable if-then rules. These models are especially attractive in regulated domains due to their deterministic logic and auditability (Wang et al., 2017).

However, their simplicity often limits performance in highly nonlinear or unstructured tasks, such as image analysis or sequential patient monitoring.

In high stakes clinical settings such as cancer diagnosis or critical care triage the priority is often placed on model transparency over marginal gains in accuracy. As such, intrinsically interpretable models remain indispensable when accountability, auditability, and ease of explanation are paramount.

### 3.2 Post Hoc Explainability Techniques

Post hoc explainability refers to a class of methods designed to interpret the outputs of complex, often opaque models after training. These techniques enable clinicians to extract meaningful insights from high performing models like deep neural networks and ensemble methods, which otherwise lack intuitive transparency.

**LIME** (Local Interpretable Model Agnostic Explanations) approximates a black box model locally around a specific prediction using a simpler surrogate model (Ribeiro et al., 2016). For example, LIME might explain why a patient is classified as high risk for heart failure by identifying which EHR features most influenced that prediction in a specific instance. While effective for tabular data, LIME's local fidelity and sensitivity to perturbations can limit its robustness in clinical environments.

**SHAP** (SHapley Additive exPlanations) is based on cooperative game theory and computes feature attributions by evaluating the marginal contribution of each input to the prediction (Lundberg & Lee, 2017). SHAP values are consistent and model agnostic, offering both local and global explanations. In healthcare, SHAP has been applied successfully to explain ML predictions in intensive care monitoring, oncology risk scores, and chronic disease management, providing clinician friendly visualizations (e.g., waterfall plots, force plots) that summarize which features increase or decrease risk.

**Grad CAM** (Gradient weighted Class Activation Mapping) is particularly useful for interpreting convolutional neural networks in medical imaging. It generates heatmaps over input images to highlight regions that most influenced the model's decision (Selvaraju et al., 2017). For instance, in a CNN trained to detect pneumonia on chest X rays, Grad CAM can visually indicate the affected lung area, improving clinician trust in the model's focus and reliability.

**Counterfactual explanations** offer intuitive "what if" scenarios e.g., "Had the patient's blood sugar level been lower, the model would not have predicted diabetes." These explanations are particularly suited for clinical consultations and decision justification, as they present alternative actions or values that could alter the diagnostic outcome (Wachter et al., 2017).

Each post hoc technique brings distinct advantages and limitations depending on the data type, model architecture, and clinical constraints. While SHAP and LIME are well suited for structured EHR data, Grad CAM is more appropriate for visual diagnostics. Counterfactuals are advantageous for treatment planning and patient communication, yet may struggle with plausibility constraints in clinical scenarios.

### 3.3 Comparative Table of Interpretability Techniques

Below is a comparative evaluation of selected intrinsic models and post hoc explainability tools relevant to healthcare diagnostics:

Model/Technique	Interpretability Level	Fidelity to Model	Computational Cost	Scalability	Clinical Use Case Alignment
Decision Trees	High	Exact	Low	Moderate	EHR data, triage rules
Logistic Regression	High	Exact	Low	High	Risk scores, prognosis models
Bayesian Rule Lists	High	Exact	Moderate	Moderate	Clinical guideline modeling
LIME	Moderate	Local (Approximate)	Moderate	High	EHR based predictions
SHAP	High	Local + Global	High	Moderate	Risk stratification, ICU monitoring
Grad CAM	Moderate	Visual (Heatmap)	Moderate	High	Medical image interpretation
Counterfactuals	Moderate	Local + Intuitive	High	Low-Moderate	Patient communication, treatment planning

Table 1: comparative evaluation of selected intrinsic models and post hoc explainability tools relevant to healthcare diagnostics:

The choice of interpretability method must be informed by the clinical context, type of data, and model complexity. While intrinsically interpretable models are ideal for structured data and regulatory documentation, post hoc methods extend the utility of more complex models without sacrificing trust and transparency. As healthcare continues to integrate AI into critical decision making processes, striking a balance between predictive performance and interpretability remains essential to ensuring clinician adoption and patient safety.

#### 4.0 Applications in Early Disease Detection

Early disease detection plays a pivotal role in modern healthcare by significantly reducing morbidity and mortality, improving patient prognoses, and minimizing long term treatment costs. Timely identification of disease allows for early interventions, which are often more effective and less invasive than treatments administered at later stages. Artificial Intelligence (AI), particularly machine learning (ML), has become a key enabler in advancing early diagnostic capabilities through automated risk assessment, screening, and predictive analytics. However, the criticality of early stage diagnosis where decisions carry long term implications demands AI systems that are not only accurate but also interpretable. In such high stakes scenarios, clinicians must be able to understand and validate model predictions to ensure safety, trust, and compliance with ethical and regulatory standards (Ghassemi et al., 2021). As a result, interpretable AI is increasingly viewed as indispensable in the development of clinically viable diagnostic support tools.

#### **4.1 Cancer Screening**

Cancer diagnosis and screening represent one of the most heavily researched areas for interpretable ML applications. Breast cancer screening using mammographic data has benefited from hybrid systems that combine deep learning with post hoc explainability. For example, CNNs trained on mammograms have achieved expert level detection rates, and their decisions have been made interpretable through Class Activation Maps (CAMs) and Grad CAMs, which visually highlight suspicious tissue regions (Raghu et al., 2019). In breast cancer risk prediction, logistic regression and decision tree models have been applied to clinical and demographic data, providing transparent risk factor explanations that align with established medical knowledge (Yala et al., 2019).

In lung cancer screening, studies using low dose CT scans have employed CNNs with Grad CAM overlays to detect early stage nodules and support radiologists in decision making. Interpretable ensemble models such as XGBoost, paired with SHAP explanations, have been used to identify important features like lesion size, spiculated edges, and nodule density in structured imaging reports (Ardila et al., 2019).

Skin cancer detection using dermoscopy images has also seen success with interpretable AI. CNN based classifiers trained on datasets such as ISIC have been coupled with LIME and Grad CAM to provide visual and feature based justifications for classifying lesions as benign or malignant (Tschandl et al., 2019). These visual explanations improve dermatologist confidence and facilitate model auditing in clinical practice.

#### **4.2 Cardiovascular Diseases**

Heart disease prediction has been another important application domain for interpretable ML, especially with the availability of structured datasets such as the UCI Heart Disease and Framingham datasets. Logistic regression, decision trees, and Bayesian Rule Lists have been used to model risk factors such as age, cholesterol, blood pressure, and chest pain type in a transparent and understandable format (Karegowda et al., 2012). More recently, tree based ensemble models enhanced with SHAP explanations have offered superior accuracy while maintaining interpretability by identifying key predictors and visualizing their contributions to individual risk scores (Lundberg et al., 2018).

In clinical settings, these models have facilitated patient specific risk communication and informed shared decision making processes. Their ability to quantify the influence of modifiable risk factors has also supported preventive health interventions.

#### **4.3 Neurological Disorders**

The early detection of neurological conditions such as Alzheimer's disease (AD) poses unique challenges due to the progressive and latent nature of these disorders. Machine learning models trained on multimodal data including cognitive assessments, clinical notes, and structural MRI scans have shown promise in predicting the onset of Alzheimer's prior to symptomatic presentation. Longitudinal modeling using recurrent neural networks (RNNs) has been coupled with temporal SHAP explanations to elucidate how changes in memory scores and brain volume over time influence diagnostic outputs (Jain et al., 2020).

Datasets like the Alzheimer's Disease Neuroimaging Initiative (ADNI) have supported the development of such interpretable predictive frameworks. Moreover, attention mechanisms and saliency maps have been applied to MRI scans to visually highlight brain regions associated with neurodegeneration, enhancing clinical interpretability and fostering early screening in at risk populations.

#### **4.4 Infectious Diseases (e.g., COVID 19)**

The global COVID 19 pandemic catalyzed the rapid development of AI based tools for infection screening, severity prediction, and resource triage. Interpretable ML models played a crucial role in translating complex biomarker data into actionable insights. For example, gradient boosted trees trained on patient vitals and lab results were deployed in hospital settings to predict ICU admission and mortality risk, with SHAP values highlighting features such as oxygen saturation, CRP levels, and age (Jiang et al., 2020).

Chest X ray and CT based models trained on datasets like COVIDx were explained using Grad CAM to show lesion distribution and severity, helping radiologists differentiate COVID related pneumonia from other respiratory conditions (Wang et al., 2020). Counterfactual explanations have also been used to explore how changes in clinical parameters could have altered patient outcomes, supporting retrospective analysis and treatment evaluation.

#### **4.5 Datasets Used in Early Detection Studies**

Several publicly available datasets have underpinned the development of interpretable ML systems for early diagnosis:

- **MIMIC III:** A critical care dataset containing EHR data from ICU patients; widely used in mortality and risk prediction.
- **NIH Chest X ray14:** A large scale dataset of chest X rays annotated for 14 conditions, supporting image based classification.
- **UCI Heart Disease:** Contains structured attributes like age, cholesterol, and ECG results for binary classification of heart disease.
- **ADNI (Alzheimer's Disease Neuroimaging Initiative):** Offers cognitive scores, imaging, and biomarker data for AD progression modeling.
- **COVIDx and CC CCII:** Pandemic datasets used for COVID 19 detection and severity classification from imaging and lab data.

These datasets have facilitated cross study comparability and enabled reproducibility in academic and clinical settings.

#### **4.6 Discussion of Model Performance and Explainability**

The reviewed studies reveal a consistent trend: while complex models such as deep neural networks offer superior accuracy (often exceeding 90% AUC in imaging tasks), their adoption in clinical workflows depends significantly on their interpretability. For example, CNNs coupled with Grad CAM yielded high diagnostic performance in lung and skin cancer detection, but their clinical utility was greatly enhanced by visual explanations that confirmed model focus areas aligned with medical knowledge (Tschandl et al., 2019).

In structured data contexts, interpretable models like logistic regression and SHAP enhanced XGBoost have balanced transparency with performance, making them suitable for EHR based risk prediction in cardiovascular and infectious diseases. In contrast, counterfactual explanations have been particularly effective in supporting decision justification and patient communication, especially for conditions involving modifiable risk factors.

Clinician trust remains a key determinant of adoption. Studies consistently show that interpretable models increase diagnostic confidence, reduce resistance to AI integration, and facilitate regulatory approval

(Shortliffe & Sepúlveda, 2018). Ultimately, the selection of model and explanation technique must align with the clinical use case, data modality, and operational constraints of the healthcare environment.

Table 2: Interpretable AI in Early Disease Detection

Disease Area	Dataset	Model Type	Explainability Tool	Outcome
Breast Cancer	Mammography, Clinical Data	CNN + Logistic Regression	Grad CAM, SHAP	Improved screening accuracy and feature insight
Lung Cancer	Low dose CT	CNN, XGBoost	Grad CAM, SHAP	Visual explanations of nodules and risk scoring
Skin Cancer	ISIC Dermoscopy Images	CNN	LIME, Grad CAM	Accurate lesion classification with localization
Heart Disease	UCI Heart, Framingham	Logistic, XGBoost	SHAP, RuleFit	Transparent identification of risk factors
Alzheimer’s Disease	ADNI	RNN, Logistic	Temporal SHAP, Saliency	Longitudinal predictions and brain region mapping
COVID 19	COVIDx, MIMIC COVID	XGBoost, CNN	SHAP, Grad CAM	ICU risk prediction and image based screening

### 5.0 Case Study: Comparative Evaluation

To empirically evaluate the trade offs between interpretability, predictive performance, and clinical utility, we conducted a comparative analysis of three machine learning models applied to the **UCI Heart Disease Dataset**. This dataset was selected due to its relevance in cardiovascular diagnostics, a domain where early detection is crucial, and where explainable AI can significantly support decision making. The dataset includes 14 clinically significant features such as age, blood pressure, cholesterol levels, and chest pain type, making it well suited for evaluating model interpretability across structured data (Dua & Graff, 2019).

The models evaluated in this case study include: (1) **Logistic Regression (LR)**, as an intrinsically interpretable baseline; (2) **Explainable Boosting Machine (EBM)**, a generalized additive model optimized for both accuracy and interpretability (Nori et al., 2019); and (3) **XGBoost**, a high performing ensemble model combined with **SHAP (SHapley Additive exPlanations)** to provide post hoc interpretability (Lundberg & Lee, 2017).

#### 5.1 Quantitative Performance Evaluation

Each model was trained using 10 fold cross validation, and their performance was assessed using **Accuracy, Precision, Recall, and Area Under the ROC Curve (AUC ROC)**. The results are summarized in Table 3 below

**Table 3: Comparative Performance Metrics**

Model	Accuracy	Precision	Recall	AUC ROC
Logistic Regression	0.84	0.82	0.79	0.86
EBM	0.86	0.84	0.81	0.89
XGBoost + SHAP	0.89	0.88	0.86	0.93

**Table above** presents the ROC curves for all three models. The XGBoost model demonstrated superior overall predictive power, but with a corresponding increase in model complexity.

### 5.2 Explanation Clarity and Clinical Comprehension

Interpretability was assessed based on the clarity and granularity of model explanations. Logistic Regression, being linear, provided coefficient values directly corresponding to the impact of features such as age and cholesterol on heart disease risk. Clinicians could readily understand these weights as odds ratios, supporting quick mental reasoning about patient risk.

EBM further enhanced interpretability by providing **visualized feature function plots**, where each feature’s non linear contribution was shown graphically. This allowed clinicians to understand, for example, how systolic blood pressure affects risk at different thresholds a more nuanced representation than linear models allow (Caruana et al., 2015).

For XGBoost, **SHAP summary plots** and **force plots** were used to decompose individual predictions into additive feature contributions (see Figure 5.2). These visualizations were useful for explaining specific decisions, such as why a younger patient with high cholesterol was still categorized as high risk. However, they required additional explanation or training for clinical users not familiar with SHAP values.

Overall, while SHAP provided the richest explanation granularity, EBM struck a better balance between clarity and expressiveness for clinical comprehension.

### 5.3 Bias and Fairness Analysis

To assess bias, we stratified model performance across **age** and **gender** subgroups. We applied **Equal Opportunity Difference** and **Disparate Impact Ratio** as fairness metrics. Logistic Regression showed minimal bias across gender but slightly underperformed in younger cohorts (age < 40). EBM exhibited balanced performance across groups due to its additive structure and monotonic constraints. XGBoost, while accurate, demonstrated a slight disparate impact in female patients, which SHAP helped identify by tracing the over reliance on male dominated cholesterol patterns.

This analysis emphasizes the importance of model auditing using explainability tools not only for transparency but also for fairness diagnostics especially in healthcare systems that serve diverse populations (Barocas et al., 2019).

## 5.4 Integration into CDSS Workflows

From a deployment perspective, the simplicity of **Logistic Regression** and **EBM** facilitates seamless integration into Clinical Decision Support Systems (CDSS). Their predictions and explanations are easily rendered into textual or visual outputs that can be interpreted in real time during clinical consultations. They also require minimal computational resources, making them suitable for edge or hospital EHR integration.

The **XGBoost + SHAP** combination, while powerful, poses challenges related to computational cost and explanation time. Although SHAP explanations can be precomputed and cached, real time use requires careful engineering to maintain latency requirements. Additionally, SHAP outputs may necessitate clinician training or the use of simplified visual interfaces for non technical users.

Privacy and data protection are also important. All models were evaluated with de identified data, and model predictions were logged for traceability. Systems using SHAP or EBM could be configured to support **auditable diagnostics**, aligning with **FDA SaMD guidelines** and **GDPR** requirements for explainable algorithmic decision making (European Commission, 2021; FDA, 2022).

## 6.0 Discussion

The comparative evaluation of interpretable machine learning models in the context of early heart disease detection underscores a key tension in the development of AI systems for healthcare: the trade off between **predictive performance and interpretability**. In our case study, XGBoost combined with SHAP yielded the highest predictive performance across all metrics, including accuracy and AUC ROC. However, it also required more computational resources, produced less intuitively accessible explanations, and posed integration challenges within Clinical Decision Support Systems (CDSS). On the other hand, models such as Logistic Regression and Explainable Boosting Machines (EBM) delivered slightly lower predictive accuracy but offered significantly higher interpretability, which is crucial in clinical settings where decisions must be explainable, auditable, and justifiable (Caruana et al., 2015; Lundberg & Lee, 2017).

The observed trade offs validate the literature's assertion that **simpler models are often preferred** in high stakes medical applications, particularly when clinicians demand clarity over marginal performance gains (Doshi Velez & Kim, 2017). Logistic Regression and EBM allow medical professionals to directly inspect model logic, evaluate threshold effects of clinical features, and trace decision pathways. These characteristics are vital in maintaining clinician autonomy and accountability, especially when AI recommendations are used in parallel with physician judgment.

The interpretability of AI models also plays a significant role in **clinician trust and adoption**. Transparent explanations whether through coefficients, visual feature plots, or SHAP breakdowns enable clinicians to verify AI predictions against their domain knowledge and clinical intuition. In our evaluation, EBM explanations were particularly well suited for real time clinical use because they combined visual clarity with model fidelity. SHAP explanations for XGBoost, while granular and informative, required additional interface design and user training to ensure comprehensibility. Prior studies confirm that clinicians are more likely to adopt AI tools when outputs are both **trustworthy and easily interpretable**, emphasizing the need for user centered AI design in healthcare (Tonekaboni et al., 2019; Shortliffe & Sepúlveda, 2018).

Beyond usability, **ethical and regulatory considerations** strongly reinforce the need for interpretable AI in medicine. From an ethical standpoint, patients have a right to understand the basis of decisions that impact their health. Informed consent is undermined when the rationale behind diagnostic recommendations is opaque or inaccessible. Interpretable models allow clinicians to explain risks, alternatives, and outcomes,

thereby enhancing shared decision making and reinforcing the principle of patient autonomy (Morley et al., 2020).

From a regulatory perspective, multiple frameworks now emphasize **transparency, auditability, and accountability** in AI powered healthcare tools. The **European Union's AI Act** mandates risk based categorization and transparency in high risk AI applications, including those in healthcare. Similarly, the **FDA's guidance on Software as a Medical Device (SaMD)** requires a demonstration of algorithm transparency and interpretability to support regulatory approval. Models like EBM, which generate human readable rules and visual explanations, align more closely with these mandates compared to black box architectures, unless those architectures are paired with robust and validated post hoc explanations (FDA, 2022; European Commission, 2021).

The case study also raised critical concerns regarding **bias and fairness** in predictive healthcare. Although all three models performed reasonably across demographic groups, XGBoost showed a slight bias in predictions for female patients, which was traceable through SHAP analysis. This illustrates how interpretability tools can be leveraged not only to explain predictions but also to **audit models for potential bias**, supporting equitable treatment outcomes and satisfying legal obligations under laws such as the **Health Insurance Portability and Accountability Act (HIPAA)** and the **General Data Protection Regulation (GDPR)** (Barocas et al., 2019; Ghassemi et al., 2021).

Despite their promise, interpretable AI systems face **challenges in real world deployment**. These include data privacy concerns, computational constraints in hospital IT infrastructures, clinician resistance due to lack of training, and a lack of standardized metrics for measuring interpretability across domains. Furthermore, interpretability itself is context dependent what is interpretable to a data scientist may not be comprehensible to a clinician. Hence, future development should focus on **human centered explainability**, tailored to the information needs, time constraints, and cognitive workflows of end users.

In sum, this discussion affirms that interpretability is not a peripheral design choice but a **core requirement** for the ethical, regulatory compliant, and trustworthy deployment of AI in healthcare diagnostics. Models must balance performance with transparency, and explainability methods must be evaluated not only for technical fidelity but also for their **communicative clarity and usability** in clinical practice.

## 7.0 Challenges and Research Gaps and Future Directions

Despite the growing interest in interpretable AI for healthcare diagnostics, several persistent challenges continue to constrain the effectiveness, trustworthiness, and clinical applicability of these systems. As AI technologies edge closer to real world deployment, it becomes increasingly important to critically examine their technical, methodological, and operational limitations. Understanding these challenges is essential not only for improving model design but also for ensuring that interpretability serves its intended purpose supporting safe, equitable, and transparent clinical decision making.

One of the most pressing limitations lies in the **technical shortcomings of current post hoc explainability methods**. Tools such as SHAP, LIME, and Grad CAM, while powerful in academic settings, exhibit notable weaknesses when applied in complex, high stakes medical environments. For instance, LIME generates local surrogate models based on input perturbations, but its explanations are known to lack stability and consistency across runs, which undermines their reliability (Alvarez Melis & Jaakkola, 2018). SHAP values, while more theoretically grounded, are computationally intensive, especially for large models or high dimensional datasets common in electronic health records and genomics. Grad CAM, typically used for

convolutional neural networks in medical imaging, offers visual saliency maps that highlight input regions but provides limited information about feature interactions or decision thresholds (Selvaraju et al., 2017). Furthermore, all three methods struggle to accommodate **temporal, multimodal, or longitudinal data** formats that are central to modern diagnostic practice. The inability to interpret these complex data structures reduces the usability of explainable models in clinical workflows and may lead to oversimplified or misleading insights.

Another significant concern is the **bias and lack of generalizability in healthcare datasets**, which adversely affect both model performance and interpretability. Many widely used datasets such as those derived from tertiary care hospitals suffer from **demographic imbalances**, underrepresenting women, racial minorities, children, and patients with rare diseases (Chen et al., 2021). These biases skew the model's learning process and can lead to disparate outcomes across subgroups. For instance, a model trained predominantly on middle aged male patients may perform poorly on younger or female populations, with post hoc explanations failing to reveal such limitations unless demographic specific auditing is conducted (Suresh & Guttag, 2021). Moreover, AI models trained on data from a single institution or geographic location may not generalize well to different clinical environments due to variations in diagnostic equipment, data coding standards, and population health profiles. This **lack of external validity** complicates the development of universally interpretable AI tools and demands robust mechanisms for domain adaptation, transfer learning, and continuous validation.

Beyond data and model limitations, there exists a crucial **gap in research on clinician AI collaboration and interaction**. Most interpretability studies are conducted in controlled settings, with minimal involvement of end users such as physicians, nurses, or radiologists. As a result, we lack empirical understanding of how clinicians perceive, interpret, and act upon AI generated explanations in real world diagnostic workflows (Tonekaboni et al., 2019). Few studies systematically evaluate the **comprehension** of explanations or assess their impact on **clinical trust, decision accuracy, or workflow efficiency**. Moreover, interpretability metrics are often designed by data scientists without input from healthcare practitioners, leading to mismatches between technical clarity and user relevance. For example, while SHAP plots may satisfy model fidelity criteria, they may be too abstract or granular for clinicians under time pressure. This disconnect highlights the need for **human centered AI development**, incorporating methodologies from human computer interaction, behavioral science, and participatory design. Designing interpretability tools with clinician feedback not just as end users but as co creators can significantly improve usability and adoption in clinical settings (Amann et al., 2020).

In sum, the current landscape of interpretable AI in healthcare is marked by **methodological immaturity, contextual blind spots, and a lack of interdisciplinary integration**. Addressing these challenges requires a paradigm shift from building “explainable algorithms” to fostering “explainable systems” embedded within real clinical ecosystems. This shift must account for technical rigor, data diversity, and human factors simultaneously. As the next section outlines, overcoming these limitations demands a cohesive research agenda that bridges machine learning innovation with clinical collaboration and ethical foresight.

## 7.1 Future Directions

Building upon the challenges and research gaps identified in the previous section, future advancements in interpretable AI for healthcare diagnostics must extend beyond model transparency alone to embrace richer forms of reasoning, clinician centered design, and real world validation. This transition demands not only technical innovation but also a deeper integration of clinical insight, ethical responsibility, and contextual sensitivity into AI system development.

One of the most promising areas for future exploration is the incorporation of **causal reasoning and symbolic explainable hybrids** into interpretable machine learning models. Current interpretability techniques, particularly post hoc explanations like SHAP or LIME, are largely correlation driven and often lack the capacity to distinguish spurious associations from causal relationships a critical limitation in clinical reasoning (Pearl, 2019). By embedding **causal inference frameworks**, such as structural causal models (SCMs) or counterfactual reasoning, AI systems could generate explanations that not only identify contributing factors but also elucidate potential outcomes under alternative clinical scenarios. This would be especially valuable in treatment planning, where clinicians must understand how modifying a variable (e.g., lowering blood pressure) may affect predicted outcomes.

In parallel, **symbolic AI and neurosymbolic systems** which blend machine learning with rule based reasoning and domain ontologies offer opportunities to align explanations more closely with how clinicians reason about diagnoses (Marcus & Davis, 2020). For example, integrating **knowledge graphs** with predictive models can ground explanations in known medical hierarchies (e.g., ICD codes, drug interactions), thereby improving both credibility and comprehensibility. Such **hybrid interpretable models** may help bridge the semantic gap between complex model outputs and the logical, cause effect reasoning that clinicians rely upon, allowing them to better assess not just *what* the model predicts, but *why* it arrived at that conclusion.

Equally vital is the call for **interdisciplinary, co designed AI systems**, where clinicians are embedded in the development lifecycle not as passive users, but as co creators. Future research should prioritize **human in the loop methodologies** that allow healthcare professionals to actively participate in **feature selection, explanation format design, and usability testing** (Amann et al., 2020). For instance, the language and visual metaphors used in explanations should be tailored to clinical workflows, leveraging familiar formats such as risk scores, timelines, or patient narratives. Participatory design approaches rooted in human computer interaction (HCI) and cognitive ergonomics can yield interpretable AI tools that resonate with the **mental models** of end users, increasing trust, reducing cognitive overload, and facilitating adoption. Moreover, such collaboration fosters a shared sense of ownership and accountability, which is essential in the clinical deployment of AI.

In addition to design innovation, the field must commit to **rigorous real world validation** of interpretable AI models within **Clinical Decision Support Systems (CDSS)**. Laboratory performance is insufficient without empirical evidence of clinical impact. Future studies should involve **longitudinal evaluations and randomized controlled trials** that assess not only diagnostic accuracy but also **clinician response, workflow integration, and patient health outcomes** (He et al., 2022). For instance, a truly effective interpretable model should not only predict heart disease with high accuracy but also demonstrably improve clinician confidence and lead to better patient adherence to treatment. Furthermore, **pilot implementations across diverse healthcare settings** including rural hospitals, low resource clinics, and international contexts can reveal constraints and adaptations required for equitable AI deployment across populations. Ensuring such ecological validity is crucial for global health equity and AI fairness.

Taken together, these future directions envision a next generation of interpretable AI in healthcare that is **causally sound, semantically aligned, user centered, and clinically validated**. Moving forward, the challenge is no longer solely to make AI explainable but to ensure that these explanations are *actionable, trustworthy, and impactful* in real world care delivery. These goals will require sustained collaboration between computer scientists, clinicians, ethicists, and policymakers an interdisciplinary ecosystem capable of translating algorithmic innovation into meaningful medical advancement.

## 8.0 Conclusion

This study has explored the critical role of interpretability in machine learning models applied to healthcare diagnostics, with a focus on early disease detection and integration into Clinical Decision Support Systems (CDSS). Through a comprehensive review of the literature and a comparative evaluation using the UCI Heart Disease dataset, the analysis revealed both the promise and complexity of deploying interpretable AI in clinical practice. Intrinsically interpretable models such as logistic regression and Explainable Boosting Machines (EBM) were shown to offer transparency and ease of integration, making them well suited for high stakes environments where clarity and auditability are paramount. Meanwhile, high performing black box models like XGBoost, when paired with post hoc explainability tools such as SHAP, can deliver superior accuracy while still supporting clinician trust through meaningful feature attributions. Across all cases, the study reaffirmed that interpretability is not merely a technical add on, but a foundational requirement for safe, ethical, and effective AI assisted healthcare.

For researchers, these findings underscore the importance of developing interpretable models that are not only accurate but also robust, causally informed, and validated in real world clinical settings. Future model development should involve clinician feedback from the outset, incorporate causal reasoning and symbolic methods where appropriate, and prioritize transparency alongside performance. For healthcare practitioners, critical engagement with AI tools is essential clinicians should demand clarity in how predictions are generated and participate actively in co designing AI systems that align with their decision making processes and information needs. For both groups, ensuring fairness, usability, and regulatory compliance must remain central in the design and deployment of interpretable AI solutions.

Ultimately, this research highlights that the trade off between predictive power and interpretability remains an enduring challenge but one that must be addressed with nuance rather than compromise. Clinical trust cannot be engineered solely through performance metrics; it must be earned through models that are understandable, accountable, and supportive of clinician expertise. The future of AI in healthcare lies not in replacing human judgment but in enhancing it through transparent, evidence aligned, and human centered systems. As the field advances, a collaborative and interdisciplinary approach will be vital to realizing the full potential of interpretable AI delivering technologies that are not only powerful but also principled and practical in the service of patient care.

## Declaration of Conflicting Interests

The authors declare no potential conflicts of interest with respect to the research, authorship and publication of this article.

## Funding

The author received no financial support for the research, authorship and publication of this article.

## References

- Alvarez-Melis, D., & Jaakkola, T. S. (2018). On the robustness of interpretability methods. arXiv preprint arXiv:1806.08049.

- Amann, J., Blasimme, A., Vayena, E., Frey, D., & Madai, V. I. (2020). Explainability for artificial intelligence in healthcare: A multidisciplinary perspective. *BMC Medical Informatics and Decision Making*, 20(1), 1–9.
- Ardila, D., Kiraly, A. P., Bharadwaj, S., Choi, B., Reicher, J. J., Peng, L., ... & Shetty, S. (2019). End-to-end lung cancer screening with three-dimensional deep learning on low-dose chest computed tomography. *Nature Medicine*, 25(6), 954–961.
- Barocas, S., Hardt, M., & Narayanan, A. (2019). *Fairness and Machine Learning*. fairmlbook.org.
- Caruana, R., Lou, Y., Gehrke, J., Koch, P., Sturm, M., & Elhadad, N. (2015). Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day readmission. *Proceedings of the 21st ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 1721–1730.
- Chen, I. Y., Joshi, S., Ghassemi, M., & Rajpurkar, P. (2021). Treating health disparities with artificial intelligence. *Nature Medicine*, 27(9), 1520–1521.
- Doshi-Velez, F., & Kim, B. (2017). Towards a rigorous science of interpretable machine learning. arXiv preprint arXiv:1702.08608.
- Dua, D., & Graff, C. (2019). *UCI Machine Learning Repository*. University of California, Irvine.
- Esteva, A., Kuprel, B., Novoa, R. A., Ko, J., Swetter, S. M., Blau, H. M., & Thrun, S. (2017). Dermatologist-level classification of skin cancer with deep neural networks. *Nature*, 542(7639), 115–118.
- European Commission. (2021). *Proposal for a Regulation on a European Approach for Artificial Intelligence*.
- FDA. (2022). *Artificial Intelligence and Machine Learning in Software as a Medical Device*. U.S. Food and Drug Administration.
- Ghassemi, M., Oakden-Rayner, L., & Beam, A. L. (2021). The false hope of current approaches to explainable artificial intelligence in health care. *The Lancet Digital Health*, 3(11), e745–e750.
- He, J., Baxter, S. L., Xu, J., Xu, J., Zhou, X., & Zhang, K. (2022). The practical implementation of artificial intelligence technologies in medicine. *Nature Medicine*, 28(1), 30–35.
- Hosmer, D. W., Lemeshow, S., & Sturdivant, R. X. (2013). *Applied Logistic Regression* (Vol. 398). John Wiley & Sons.
- Jain, S., Lee, M. J., & El-Khomy, M. (2020). Explainable AI for Alzheimer’s disease diagnosis using longitudinal patient data. *IEEE Journal of Biomedical and Health Informatics*, 24(9), 2705–2712.
- Jiang, X., Coffee, M., Bari, A., Wang, J., Jiang, X., & Huang, J. (2020). Towards an artificial intelligence framework for data-driven prediction of coronavirus clinical severity. *Computers, Materials & Continua*, 63(1), 537–551.
- Karegowda, A. G., Manjunath, A. S., & Jayaram, M. A. (2012). Application of genetic algorithm optimized neural network connection weights for heart disease prediction. *International Journal of Computer Applications*, 31(7), 14–22.
- Lipton, Z. C. (2018). The mythos of model interpretability. *Communications of the ACM*, 61(10), 36–43.
- Lundberg, S. M., & Lee, S.-I. (2017). A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems*, 30.
- Marcus, G., & Davis, E. (2020). The next decade in AI: Four steps towards robust artificial intelligence. *AI Magazine*, 41(2), 25–36.
- Morley, J., Floridi, L., Kinsey, L., & Elhalal, A. (2020). From what to how: An initial review of publicly available AI ethics tools, methods and research to translate principles into practices. *Science and Engineering Ethics*, 26(4), 2141–2168.
- Nori, H., Jenkins, S., Koch, P., & Caruana, R. (2019). InterpretML: A unified framework for machine learning interpretability. arXiv preprint arXiv:1909.09223.
- Pearl, J. (2019). The seven tools of causal inference, with reflections on machine learning. *Communications of the ACM*, 62(3), 54–60.
- Rajpurkar, P., Irvin, J., Zhu, K., Yang, B., Mehta, H., Duan, T., ... & Ng, A. Y. (2018). CheXNet: Radiologist-level pneumonia detection on chest X-rays with deep learning. arXiv preprint arXiv:1711.05225.
- Raghu, M., Sundararajan, M., Talukdar, P., & Kleinberg, J. (2019). Transfusion: Understanding transfer learning for medical imaging. *Advances in Neural Information Processing Systems*, 32.

- Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). "Why should I trust you?": Explaining the predictions of any classifier. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 1135–1144.
- Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., & Batra, D. (2017). Grad-CAM: Visual explanations from deep networks via gradient-based localization. *Proceedings of the IEEE International Conference on Computer Vision*, 618–626.
- Shen, D., Wu, G., & Suk, H.-I. (2017). Deep learning in medical image analysis. *Annual Review of Biomedical Engineering*, 19, 221–248.
- Shortliffe, E. H., & Sepúlveda, M. J. (2018). Clinical decision support in the era of artificial intelligence. *JAMA*, 320(21), 2199–2200.
- Suresh, H., & Gutttag, J. V. (2021). A framework for understanding unintended consequences of machine learning. *Communications of the ACM*, 64(4), 62–71.
- Topol, E. (2019). *Deep Medicine: How Artificial Intelligence Can Make Healthcare Human Again*. Basic Books.
- Tonekaboni, S., Joshi, S., McCradden, M. D., & Goldenberg, A. (2019). What clinicians want: Contextualizing explainable machine learning for clinical end use. *Proceedings of Machine Learning Research*, 106, 359–380.
- Tschandl, P., Codella, N., & Celebi, M. E. (2019). Comparison of the accuracy of human readers versus machine-learning algorithms for pigmented skin lesion classification: An open, web-based, international, diagnostic study. *The Lancet Oncology*, 20(7), 938–947.
- Wang, F., Rudin, C., & Wagner, D. (2017). Learning interpretable classification rules with Boolean compressed sensing. *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics (AISTATS)*, 1265–1274.
- Wang, L., Lin, Z. Q., & Wong, A. (2020). COVID-Net: A tailored deep convolutional neural network design for detection of COVID-19 cases from chest X-ray images. *Scientific Reports*, 10, 19549.
- Yala, A., Lehman, C., Schuster, T., Portnoi, T., & Barzilay, R. (2019). A deep learning mammography-based model for improved breast cancer risk prediction. *Radiology*, 292(1), 60–66.
- Zhou, B., Khosla, A., Lapedriza, A., Oliva, A., & Torralba, A. (2016). Learning deep features for discriminative localization. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2921–2929.