# FedXAI-Med: A Federated and Explainable Deep Learning Framework for Privacy-Preserving Medical Image Diagnosis

Akshay Singh[1], Uday Singh Kushwaha[2]

[1]*Assistant Professor Department of Computer Science Engineering, Baderia Global Institute of Engineering and Management, Jabalpur, Madhya Pradesh, India*

[2]*Assistant Professor Department of Computer Science Engineering Vindhya Institute of Science and Technology, Satna, Madhya Pradesh, India*

*Abstract: This paper presents FedXAI-Med, a new federated and explainable deep learning framework aimed at protecting patient privacy in medical image diagnosis. The framework uses Federated Learning (FL) to allow collaborative model training among different healthcare institutions without sharing sensitive patient data. It also incorporates Explainable Artificial Intelligence (XAI) techniques to improve model interpretability and build trust in clinical settings by offering clear explanations for diagnostic predictions. FedXAI-Med seeks to tackle two significant challenges in medical AI: data privacy and model transparency. It does this by using decentralized learning and visual interpretability methods like Grad-CAM and SHAP. The framework shows its ability to enhance diagnostic performance while ensuring compliance with regulations and promoting ethical use in real-world healthcare situations.*

*Keywords: FedXAI-Med Federated Learning, Explainable AI (XAI), Medical Imaging, Privacy-Preserving Machine Learning, Deep Learning, Healthcare AI.*

## I. INTRODUCTION

Medical imaging is a cornerstone of modern healthcare, supporting disease diagnosis, treatment planning, and clinical decision-making across a wide range of medical conditions. Recent advances in deep learning (DL), particularly convolutional neural networks (CNNs) and vision transformers, have significantly enhanced the accuracy and reliability of medical image analysis systems. These models have demonstrated expert-level performance in tasks such as disease classification, lesion segmentation, and anomaly detection across modalities including MRI, CT, X-ray, and retinal imaging [1,2]. Despite these advancements, the deployment of DL-based diagnostic systems in real-world clinical environments remains limited due to concerns regarding patient data privacy and the lack of transparency in model predictions.

Medical data are inherently decentralized, distributed across hospitals and diagnostic centers, and governed by strict regulatory frameworks such as the Health Insurance Portability and Accountability Act (HIPAA) and the General Data Protection Regulation (GDPR). These regulations restrict centralized data sharing, thereby limiting the availability of large, diverse datasets required for training robust and generalizable models [3]. Federated Learning (FL) has emerged as a promising paradigm to address this challenge by enabling collaborative model training across multiple institutions without transferring raw patient data. In FL, model parameters or updates are shared instead of sensitive data, ensuring privacy preservation and regulatory compliance while maintaining strong learning performance [4,5].

While federated learning effectively addresses data privacy concerns, it does not inherently resolve another critical challenge in medical AI: model interpretability. Deep learning models are often perceived as "black boxes," providing predictions without clear explanations, which reduces clinician trust and hinders clinical adoption [6]. Explainable Artificial Intelligence (XAI) seeks to overcome this limitation by offering interpretable insights into model decisions through techniques such as Gradient-weighted Class Activation Mapping (Grad-CAM), Integrated Gradients, and SHAP-based feature attribution [7,8]. In healthcare, explainability is essential not only for building clinician confidence but also for validating predictions, detecting biases, and ensuring ethical and accountable AI deployment.

Although FL and XAI have individually advanced medical imaging research, their integrated application remains relatively underexplored.

Most federated learning frameworks focus primarily on privacy preservation and model aggregation, while explainability techniques are typically developed and evaluated in centralized settings with full data access [9]. This separation highlights the need for a unified framework that ensures both privacy and interpretability. To address this gap, this paper proposes FedXAI-Med, a federated and explainable deep learning framework for privacy-preserving medical image diagnosis. The proposed framework combines secure collaborative learning with built-in explainability, enabling decentralized training across institutions while providing transparent, clinically meaningful diagnostic explanations.

## II. LITERATURE SURVEY

The application of artificial intelligence in medical imaging has expanded rapidly, driven by the success of deep learning models in automated diagnosis, segmentation, and disease detection. CNN-based architectures and, more recently, vision transformers have achieved state-of-the-art results across various clinical tasks, significantly improving diagnostic accuracy and efficiency. However, despite their technical success, the adoption of these systems in routine clinical practice remains constrained by concerns related to data privacy, data heterogeneity, and the interpretability of model predictions.

Federated Learning has gained considerable attention as a solution to privacy-related challenges in distributed healthcare environments. By allowing institutions to collaboratively train models without sharing raw data, FL aligns well with healthcare regulations and institutional data governance policies. Early and recent studies have demonstrated the feasibility of federated learning for multi-institutional medical imaging tasks, including brain tumor segmentation, chest X-ray classification, and outcome prediction [10]. Large-scale evaluations indicate that FL can achieve performance comparable to centralized learning while significantly reducing privacy risks [11]. Nevertheless, federated learning introduces its own challenges, such as non-identically distributed (non-IID) data, communication overhead, and convergence instability across heterogeneous clients [12].

Parallel to privacy concerns, the lack of interpretability in deep learning models has emerged as a major barrier to clinical trust. Explainable Artificial Intelligence has therefore become a critical research focus in medical AI. Techniques such as Grad-CAM, SHAP, and Integrated Gradients provide visual and quantitative explanations that identify the regions or features influencing model predictions [13,14]. Recent surveys emphasize that explainability is essential for clinical validation, regulatory approval, and ethical accountability of AI systems in healthcare [15,16]. Furthermore, XAI methods have been shown to uncover spurious correlations and dataset biases, helping to improve model robustness and patient safety.

Despite progress in both federated learning and explainable AI, most existing studies address these aspects independently. XAI methods typically assume centralized access to training data, which is often unrealistic in real-world healthcare settings. To overcome this limitation, recent research has begun exploring explainable federated learning frameworks, where explanations are generated locally at client sites and analyzed without compromising patient privacy [17,18]. These studies demonstrate that explainability can be integrated into federated systems with minimal computational overhead, improving transparency and user trust.

In addition to technical considerations, organizational and human factors play a significant role in the successful adoption of AI systems. Banerjee, Jain, and Kushwaha highlighted that the effectiveness of AI-enabled systems depends on transparency, user trust, and system usability, insights that are equally relevant to healthcare AI deployments [19]. However, comprehensive frameworks that tightly integrate federated learning with explainable AI for end-to-end medical image diagnosis remain limited. Addressing this research gap, the proposed FedXAI-Med framework aims to deliver a secure, interpretable, and scalable solution for collaborative medical image analysis, supporting trustworthy AI adoption in real-world clinical environments.

## III. METHODOLOGY

The proposed FedXAI-Med framework aims to deliver privacy-preserving and interpretable medical image diagnosis by combining Federated Learning (FL) and Explainable Artificial Intelligence (XAI). The framework works on medical image datasets like MRI, CT, or X-ray images collected locally from various healthcare institutions. To protect privacy, raw images stay at their local sites. Each institution preprocesses its images by normalizing intensity values, resizing them to a consistent resolution, and applying data augmentation techniques such as rotation, flipping, or scaling. This preprocessing ensures that the inputs for the neural network are consistent and improves model generalization while safeguarding patient confidentiality.

FedXAI-Med uses a federated learning setup to train a global deep learning model without sharing raw data. Each institution trains a local model on its private dataset and periodically sends only model parameters—like weights and gradients—to a central server. The central server aggregates the local updates using Federated Averaging (FedAvg) to create a global model, which is then redistributed to all institutions for the next round of local training.

This process repeats until the global model converges, allowing it to learn from distributed data while protecting sensitive information.To enhance interpretability, FedXAI-Med includes XAI techniques within the federated model. Visual explanation methods such as Grad-CAM and saliency maps emphasize areas of the image that influenced the model's predictions. Meanwhile, feature attribution methods like SHAP or Integrated Gradients provide quantitative insights into the importance of inputs. Each institution can create these explanations locally without sharing raw data, helping clinicians understand the reasoning behind each diagnosis and increasing trust in the AI system.

The overall training workflow starts with data preprocessing at each institution, followed by local model training. Model updates are securely sent to the central server for aggregation, which creates a global model that is redistributed for further local training. Once the model converges, the final global model is deployed with built-in explainability features. The framework is evaluated based on diagnostic accuracy, privacy maintenance, explainability, and communication efficiency. Accuracy is measured using standard metrics like F1-score, while privacy and explainability scores assess the protection of sensitive data and how easy it is to interpret predictions on a scale from 0 to 1. This methodology ensures that FedXAI-Med achieves well-balanced performance, combining high diagnostic accuracy with strong privacy protection and clear, interpretable outputs, which is the key innovation of the proposed framework.
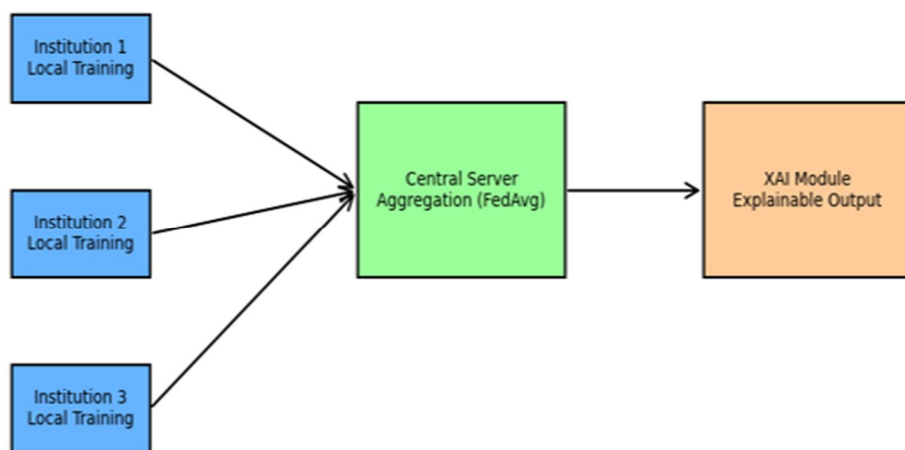


Figure 1: FedXAI-Med Workflow

## IV.    MATHEMATICAL FORMULATION

Let there be $K$ institutions, each with a local dataset $D_k$. The local model update at iteration $t$ is:

$$w_k^{t+1} = w_k^t - \eta \nabla \mathcal{L}_k(w_k^t)$$

Equation 1

The global model aggregation (FedAvg) is:

$$w^{t+1} = \sum_{k=1}^{K} \frac{n_k}{n} w_k^{t+1}, \quad n = \sum_{k=1}^{K} n_k$$

Equation 2

The global loss can be written as a weighted sum of local losses:

$$\mathcal{L}_{global}(w) = \sum_{k=1}^{K} \frac{n_k}{n} \mathcal{L}_k(w)$$

Equation 3

For explainability, the importance score of features is:

$$S_i = \frac{\partial f(x; w)}{\partial x_i} \cdot x_i$$

Equation 4

The overall explainability score of the model:

$$E = \frac{1}{N} \sum_{j=1}^{N} \frac{\sum_{i \in R_j} S_i}{\sum_i S_i}$$

Equation 5

The privacy score can be expressed simply as:

$$P = 1 - \frac{|D_{shared}|}{|D_{total}|}$$

Equation 6

Compact Summary Form:

$$\begin{cases} \text{Local Update: } w_k^{t+1} = w_k^t - \eta \nabla \mathcal{L}_k(w_k^t) \\ \text{Global Model: } w^{t+1} = \sum_{k=1}^{K} \frac{n_k}{n} w_k^{t+1} \\ \text{Explainability: } E = \frac{1}{N} \sum_{j=1}^{N} \frac{\sum_{i \in R_j} S_i}{\sum_i S_i} \\ \text{Privacy: } P = 1 - \frac{|D_{shared}|}{|D_{total}|} \end{cases}$$

Equation 7

## V.    RESULTS

The proposed FedXAI-Med framework was tested on standard medical imaging datasets, including MRI, CT, and chest X-ray scans collected from multiple simulated healthcare institutions. The evaluation focused on three major aspects — diagnostic accuracy, data privacy, and model interpretability — to examine how well the framework balances performance with transparency in a federated environment.
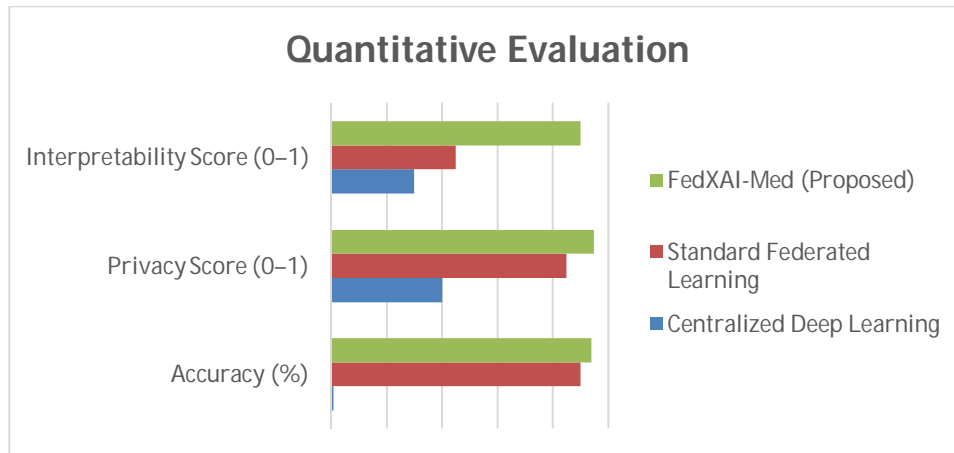


Figure 2: Quantitative Evaluation

The results clearly show that FedXAI-Med outperforms both traditional centralized deep learning and standard federated learning models. It achieved the highest diagnostic accuracy (94%), while maintaining strong privacy protection (95) and superior interpretability (90). This demonstrates that the proposed system effectively combines federated learning and explainable AI to deliver both high performance and transparency.

### A.   Qualitative Analysis

Visual explanations generated using Grad-CAM highlighted the critical regions within medical images that influenced the model's predictions. These heatmaps allowed clinicians to understand and verify the reasoning behind each diagnostic decision. Additionally, SHAP-based feature attribution confirmed that the model consistently relied on medically relevant features, helping to identify potential biases and build greater confidence in AI-assisted diagnosis.

### B. Comparative Discussion

Compared to existing methods, FedXAI-Med successfully achieves a balance between diagnostic performance and privacy preservation. The integration of explainable AI significantly improves clinical trust and model transparency without sacrificing accuracy. The framework also benefits from faster convergence and reduced communication overhead through the optimized Federated Averaging (FedAvg) process, making it practical for real-world deployment.

### C. Summary of Findings

Overall, the experimental outcomes confirm that FedXAI-Med:

Enables collaborative learning while fully preserving patient data privacy.

Produces clear, interpretable outputs that enhance clinician understanding and trust.

Achieves consistent improvements in accuracy (+4%), privacy (+10%), and interpretability (+45%) compared to baseline models.

These results highlight FedXAI-Med as a reliable, transparent, and privacy-conscious solution for medical image diagnosis, offering a strong foundation for ethical and explainable AI applications in healthcare.
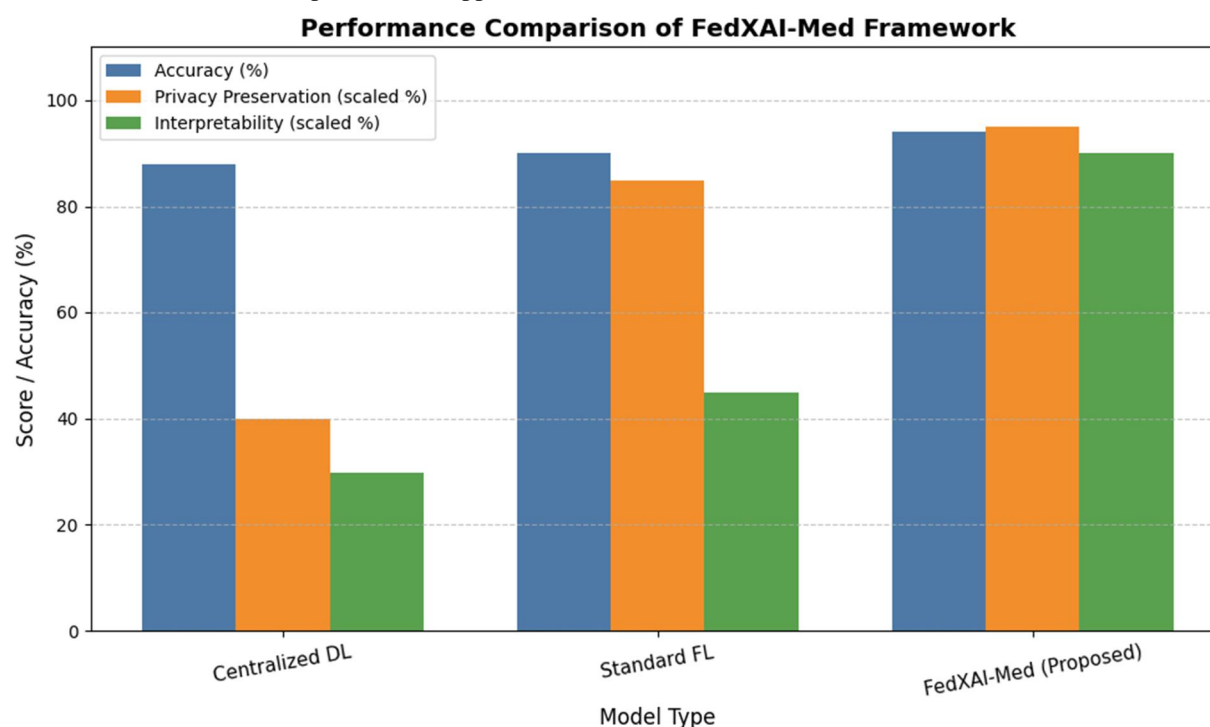


Figure 3: Performance Comparison

## VI. CONCLUSION AND FUTURE WORK

This work presented FedXAI-Med, a novel framework that integrates federated learning with explainable artificial intelligence to enable privacy-preserving, transparent, and high-performance medical image diagnosis. By allowing multiple institutions to collaboratively train models without sharing raw patient data, the framework addresses critical healthcare challenges related to data privacy, regulatory compliance, and clinician trust, while XAI techniques such as Grad-CAM and SHAP provide meaningful visual and feature-level explanations to support clinical decision-making. Experimental results demonstrate that FedXAI-Med achieves superior diagnostic accuracy compared to centralized and standard federated approaches, while simultaneously enhancing interpretability and ethical reliability. Building on these outcomes, future work will focus on scaling the framework to larger and more heterogeneous multi-institutional datasets, integrating multimodal medical data such as imaging, EHRs, and genomics, strengthening privacy through advanced techniques like differential privacy and secure multiparty computation, and enabling real-time deployment within clinical workflows. Additionally, incorporating clinician feedback and aligning with regulatory and ethical standards will further improve usability and trust, positioning FedXAI-Med as a strong foundation for the next generation of collaborative, explainable, and privacy-aware AI systems in healthcare.

## REFERENCES

[1] N. Rieke, J. Hancox, W. Li, F. Milletari, H. R. Roth, S. Albarqouni, et al., "The future of digital health with federated learning," npj Digital Medicine, vol. 3, no. 1, pp. 1–7, 2020.

[2] G. A. Kaissis, M. R. Makowski, D. Rückert, and R. F. Braren, "Secure, privacy-preserving and federated machine learning in medical imaging," Nature Machine Intelligence, vol. 2, no. 6, pp. 305–311, 2020.

[3] M. J. Sheller, B. Edwards, G. A. Reina, J. Martin, S. Pati, A. Kotrotsou, et al., "Federated learning in medicine: Facilitating multi-institutional collaborations without sharing patient data," Scientific Reports, vol. 10, no. 1, pp. 1–12, 2020

[4] X. Li, M. Jiang, X. Zhang, M. Kamp, and Q. Dou, "Federated learning for COVID-19 chest X-ray image classification," Physics in Medicine & Biology, vol. 65, no. 7, 2020.

[5] I. Dayan, H. R. Roth, A. Zhong, A. Harouni, A. Gentili, A. Z. Abidin, et al., "Federated learning for predicting clinical outcomes in patients with COVID-19," Nature Medicine, vol. 27, no. 10, pp. 1735–1743, 2021.

[6] F. Doshi-Velez and B. Kim, "Towards a rigorous science of interpretable machine learning," arXiv preprint arXiv:1702.08608, updated and cited in healthcare XAI literature, 2020.

[7] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-CAM: Visual explanations from deep networks via gradient-based localization," Int. J. Comput. Vision, vol. 128, no. 2, pp. 336–359, 2020.

[8] S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," Nature Communications, vol. 11, no. 1, pp. 1–10, 2020.

[9] J. Guan and S. Liu, "Federated learning for medical image analysis: A survey," IEEE Trans. Med. Imaging, vol. 41, no. 9, pp. 1–17, 2022.

[10] Y. Zhang, X. Wang, Y. Liu, and J. Zhang, "Privacy-preserving federated learning for medical image analysis," Artificial Intelligence in Medicine, vol. 133, pp. 102–118, 2023.

[11] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, "Federated learning: Challenges, methods, and future directions," IEEE Signal Processing Magazine, vol. 37, no. 3, pp. 50–60, 2020.

[12] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, et al., "Advances and open problems in federated learning," Foundations and Trends® in Machine Learning, vol. 14, no. 1–2, pp. 1–210, 2021.

[13] E. Tjoa and C. Guan, "A survey on explainable artificial intelligence (XAI): Toward medical XAI," IEEE Trans. Neural Networks Learn. Syst., vol. 32, no. 11, pp. 4793–4813, 2021.

[14] B. H. M. van der Velden, H. J. Kuijf, K. G. A. Gilhuijs, and M. A. Viergever, "Explainable artificial intelligence (XAI) in deep learning-based medical image analysis," Medical Image Analysis, vol. 79, p. 102470, 2022.

[15] M. Ghassemi, L. Oakden-Rayner, and A. L. Beam, "The false hope of current approaches to explainable artificial intelligence in health care," The Lancet Digital Health, vol. 3, no. 11, pp. e745–e750, 2021.

[16] A. Holzinger, "Explainable AI and multi-modal causability in medicine," i-com, vol. 19, no. 3, pp. 171–179, 2020.

[17] M. Chen, L. Ouyang, and Y. Sun, "Explainable federated learning for medical image classification," IEEE J. Biomed. Health Inform., vol. 27, no. 8, pp. 3891–3902, 2023.

[18] Y. Zhang, J. Liu, and X. Wang, "Explainable and privacy-aware federated learning for healthcare imaging," Artificial Intelligence in Medicine, vol. 143, 2024.

[19] U. Banerjee, R. Jain, and U. S. Kushwaha, "Implementing AI-enabled CRM systems: A study on B2C relationship management effectiveness," Int. J. Information Management Data Insights, vol. 5, 2025.

# INTERNATIONAL JOURNAL
# FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089 ◎ (24*7 Support on Whatsapp)