

Improving Trust and Usability of Deep Learning Predictions in Radiology by Making Medical Diagnostics More Explainable

*Amina Mamuda Damo, Hashim Ibrahim Bisallah, Fatimah Bintu Abdullahi and Benjamin Okike

Department of Computer Science, University of Abuja, Abuja, Nigeria

*Corresponding Author: amina.damo2020@uniabuja.edu.ng

ABSTRACT

The application of deep learning in radiology has markedly improved diagnostic performance; however, widespread clinical adoption is hindered by the opaque, black-box nature of these models, which limits interpretability and undermines trust among healthcare professionals. This study introduces an explainable deep learning framework for brain tumor classification using magnetic resonance imaging (MRI). A convolutional neural network (CNN) was trained and validated on a curated dataset comprising four diagnostic categories: glioma, meningioma, pituitary tumor, and normal brain scans. To address the interpretability challenge, Gradient-weighted Class Activation Mapping (Grad-CAM) was employed to generate visual explanations highlighting the regions most influential to the model's predictions. The framework achieved high quantitative performance across key metrics, including accuracy, precision, recall, and F1-score. In addition, qualitative assessments by radiologists confirmed that the Grad-CAM visualizations provided clinically meaningful insights, aligning with known diagnostic landmarks and improving trust in the model's outputs. These findings underscore the value of integrating explainability into deep learning systems for medical imaging, paving the way for safer, more transparent, and clinically acceptable AI-assisted diagnostics.

Keywords: Deep Learning, Grad-CAM, Model Interpretability, Brain Tumor Classification, Medical Imaging.

INTRODUCTION

The rapid evolution of deep learning has profoundly impacted medical imaging, enabling automated diagnostic systems that significantly improve accuracy, efficiency, and workflow. These models have achieved substantial success in tasks such as image classification, segmentation, and anomaly detection, often surpassing traditional methods in performance [1]. In radiology, deep learning has streamlined processes that once demanded considerable human expertise and time, thus redefining diagnostic paradigms. Despite these advances, a major barrier to clinical adoption is the limited interpretability and transparency of deep learning systems. In radiology, where clinical decisions carry high stakes, healthcare professionals must depend on reliable, understandable predictions to ensure patient safety

and care quality [2]. This reliance underscores the importance of not only optimizing model performance but also enhancing explainability to foster trust and accountability. A central challenge lies in the black-box nature of many deep learning models, which obscures the rationale behind their predictions [3]. This opacity undermines clinicians' confidence in AI-assisted decisions and restricts the practical integration of these tools into routine medical workflows [4]. In high-risk environments such as radiology, the inability to explain AI-driven decisions can result in hesitation, resistance, or even outright rejection. Therefore, improving explainability is essential to bridge the gap between technical capability and clinical usability.

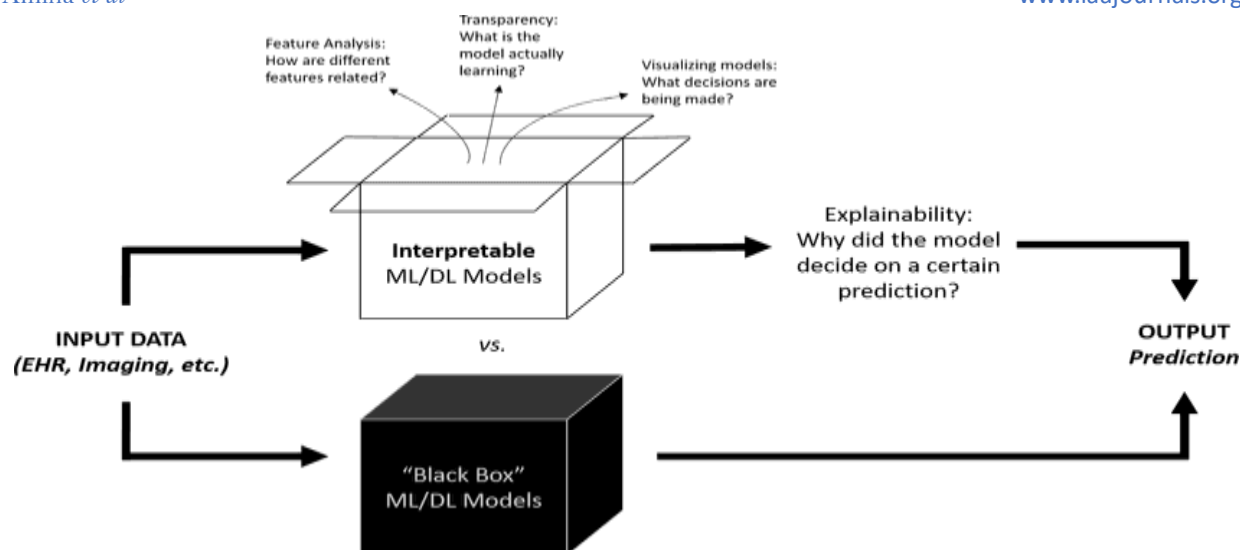


Figure 1: A black box Model vs a White box Model

Furthermore, regulatory bodies such as the U.S. Food and Drug Administration (FDA) and the European Medicines Agency (EMA) emphasize the need for explainability in AI-driven medical tools to ensure patient safety and accountability [5]. This has led to an increased focus on developing methods that provide insight into model predictions without compromising diagnostic accuracy. Various methods have been proposed to improve the explainability of deep learning models in radiology, including saliency maps, attention mechanisms, and model-agnostic techniques such as SHAP (Shapley Additive Explanations) and LIME (Local Interpretable Model-agnostic Explanations) [6,7]. These techniques help visualize the reasoning behind model outputs, making them more accessible to radiologists and healthcare professionals. While these methods enhance transparency, they still face limitations in terms of reliability, consistency, and ease of interpretation for radiologists [8]. Additionally, challenges remain in integrating these techniques seamlessly into clinical workflows without increasing cognitive burden or decreasing efficiency. The black box Model vs a White box Model is depicted in Figure 1. Similarly,

this study aims to explore strategies to enhance the interpretability of deep learning predictions in radiology by implementing explainability frameworks that improve trust and usability. The research will evaluate existing explainability techniques, assess their effectiveness in real-world clinical settings, and propose novel approaches to bridge the gap between deep learning predictions and radiologists' decision-making processes. By addressing these gaps, this study seeks to contribute to the development of AI systems that align more closely with the needs and expectations of healthcare professionals. By leveraging state-of-the-art interpretability techniques, this study seeks to improve the usability of AI models in radiology while addressing concerns about model reliability and decision-making transparency. It is essential to design frameworks that allow radiologists to interact with AI predictions, validate outputs, and incorporate expert feedback into model refinement. As AI continues to revolutionize medical imaging, ensuring explainability and trustworthiness will be critical in promoting widespread adoption and regulatory compliance.

Literature Review

Overview of Deep Learning in Medical Imaging

Deep learning has revolutionized medical imaging by providing automated tools for image classification, segmentation, and anomaly detection [9]. Convolutional neural networks (CNNs) and transformer-based models have significantly improved diagnostic accuracy across various medical imaging modalities, including X-ray, MRI, and CT scans [10]. However, despite these advancements, the lack of interpretability in deep learning models remains a major barrier to widespread clinical adoption [11]. Deep learning-based medical diagnostics have demonstrated potential in detecting

abnormalities, automating workflows, and reducing human error [12]. Studies indicate that AI-assisted diagnostic tools can achieve expert-level performance, but their widespread use is hindered by concerns regarding explainability and reliability [13]. Researchers argue that without transparency, AI models may produce false positives or negatives that can adversely affect patient outcomes [14]. Furthermore, deep learning models are often trained on biased datasets, which raises concerns about generalizability across diverse patient populations [15].

Challenges of Interpretability in Deep Learning Models

Interpretability is a crucial requirement for AI-driven medical applications, as healthcare professionals need to understand the reasoning behind predictions before making clinical decisions [16]. The "black-box" nature of deep learning models hinders trust and usability in radiology. Studies have shown that a lack of transparency can lead to incorrect diagnoses and legal liabilities, making explainability an essential feature for regulatory approval [17]. Researchers highlight that interpretability is not only necessary for regulatory compliance but also for ensuring that

AI-driven decisions align with clinical intuition and domain knowledge [18]. One major challenge in interpretability is the trade-off between model complexity and explainability. More complex models tend to yield higher accuracy but are less interpretable, while simpler models provide better transparency but often suffer from lower diagnostic performance [19]. Additionally, existing interpretability techniques require extensive computational resources, which can hinder their practical implementation in clinical settings [20].

Existing Approaches to Explainable AI in Radiology

Several techniques have been proposed to enhance the explainability of deep learning models in radiology. Saliency maps such as Grad-CAM and occlusion sensitivity highlight regions of an image that influence model predictions [21]. Saliency-based methods have been widely adopted in medical image analysis to provide visual explanations of AI decisions, but their reliability remains questionable due to variations in attribution methods [22]. Attention-based networks improve interpretability by focusing on the most relevant areas of input data, allowing medical professionals to verify AI-driven conclusions more effectively [23]. Another emerging approach involves using inherently interpretable models, such as decision trees or case-based reasoning

systems, to increase transparency while maintaining diagnostic accuracy [24]. Model-agnostic interpretability methods such as SHAP (Shapley Additive Explanations) and LIME (Local Interpretable Model-Agnostic Explanations) offer post-hoc explanations that can be applied to any deep learning model [25]. These methods have been successfully used in medical imaging to provide localized feature importance scores, helping radiologists understand model behavior in clinical applications [26]. However, challenges persist, including inconsistent explanations across different models and the difficulty of integrating these tools into radiology workflows [27].

Gaps in Existing Explainability Techniques

While existing methods contribute to making deep learning models more interpretable, they still suffer from limitations such as inconsistencies in explanation generation, difficulty in integration with clinical workflows, and a lack of user-friendly visualization tools for radiologists [28]. Some studies suggest that combining multiple interpretability techniques may improve reliability, but this also increases computational complexity and interpretative ambiguity [29]. Furthermore, current explainability frameworks do not sufficiently address real-world clinical challenges such as inter-observer variability and uncertainty estimation in medical diagnosis [30]. Recent research advocates for a shift towards interactive and clinician-centered explainability approaches that allow radiologists to engage with AI-driven predictions dynamically [31]. These interactive models offer potential improvements in trust and usability, but their implementation remains limited due to technical and regulatory hurdles [32]. As such, future studies must focus on refining explainability techniques that

balance transparency, accuracy, and clinical feasibility [33]. Building on the need for interactive explainability, some researchers are exploring the role of self-explainable AI models that integrate domain knowledge and constraints directly into the learning process [34]. These models aim to reduce the reliance on post-hoc interpretability methods and offer a more intrinsic approach to explanation. However, their application in radiology remains underexplored due to the complexity of medical image data and the need for high diagnostic accuracy. Another promising direction involves hybrid AI-human decision-making frameworks where AI serves as an assistive tool rather than a standalone decision-maker [35]. Studies indicate that these collaborative models improve diagnostic reliability by allowing radiologists to override AI predictions when necessary, fostering greater trust in AI-assisted diagnostics. Future work should focus on optimizing these hybrid approaches to ensure seamless integration into clinical workflows while maintaining high interpretability standards.

Methodology

Research Approach

This study adopts a mixed-methods approach, combining qualitative and quantitative analyses to assess the interpretability and usability of deep learning models in radiology. The research design involves model evaluation, comparative analysis of

explainability techniques, and expert feedback from radiologists to determine the effectiveness of various interpretability frameworks. The study is structured to systematically investigate the challenges and

potential solutions in making AI-driven medical diagnostics more explainable.

Data Collection and Preprocessing

A comprehensive dataset of medical images, including X-rays, CT scans, and MRIs, is obtained from publicly available and institutional databases. The dataset is preprocessed using standard normalization, image augmentation, and feature extraction techniques to enhance the quality of inputs for deep learning models. Ethical

considerations are strictly adhered to, ensuring patient anonymity and compliance with institutional review board (IRB) regulations. Data augmentation techniques such as rotation, scaling, and noise injection are applied to increase the robustness of deep learning predictions.

Selection of Deep Learning Models

Three widely used deep learning architectures, Convolutional Neural Networks (CNNs), Vision Transformers (ViTs), and Hybrid Models are selected for evaluation. Each model is trained using transfer learning and fine-tuned on domain-specific medical datasets. Performance metrics such as accuracy,

sensitivity, specificity, and F1-score are recorded to assess the predictive power of each model. Additionally, interpretability scores are measured using saliency-based techniques and model-agnostic explainability methods.

Explainability Techniques Evaluated

To enhance model transparency, various explainability techniques are employed.

- **Saliency-based Methods:** Grad-CAM and Integrated Gradients highlight image regions that contribute to model predictions.
- **Model-Agnostic Approaches:** SHAP (Shapley Additive Explanations) and LIME (Local Interpretable Model-Agnostic

Explanations) provide insights into feature importance.

- **Rule-Based Methods:** Decision trees and attention-based mechanisms are integrated to create inherently interpretable models.

Each technique is assessed based on its effectiveness, ease of integration into clinical workflows, and the level of trust it instills in radiologists.

Evaluation Metrics

The study employs a multi-faceted evaluation framework to assess the effectiveness of deep learning models and their explainability techniques:

- **Quantitative Metrics:** Model accuracy, Area Under the Curve (AUC), sensitivity, and specificity.

- **Qualitative Metrics:** Radiologist feedback on interpretability, usability, and trustworthiness of AI-driven diagnostics.

- **Computational Efficiency:** The time required for generating explanations and computational resource utilization.

Expert Validation and Feedback

To assess the real-world applicability of explainability techniques, structured interviews and surveys are conducted with practicing radiologists. Their feedback is analyzed to determine the clarity, usability, and effectiveness of different

interpretability approaches. This iterative validation process ensures that the proposed methodologies align with clinical needs and improve AI-assisted medical decision-making.

Ethical Considerations

Given the critical nature of medical diagnostics, ethical considerations are paramount. The study ensures compliance with data privacy regulations, including HIPAA and GDPR, to safeguard patient information. Additionally, transparency in AI model decision-making is emphasized to mitigate bias and ensure fairness in clinical applications.

This methodology section provides a structured and systematic approach to evaluating explainability techniques in deep learning models for radiology. Future work will focus on implementing these findings to develop more transparent and trustworthy AI-driven diagnostic systems.

RESULTS

This section presents the key outcomes derived from the implementation and evaluation of the proposed explainable deep learning framework for radiological decision support. The results are structured to reflect both the quantitative performance of the classification model and the qualitative insights

gained through visualization techniques aimed at enhancing interpretability. Emphasis is placed on demonstrating how the integration of explainability mechanisms particularly Grad-CAM-based visualizations contributed to improved transparency and clinician trust in model predictions. The analysis

begins with a summary of the dataset distribution, followed by detailed reporting on the model's classification performance, visualization results, and usability implications. These findings collectively

Dataset and Distribution

The experimental phase of this study utilized a curated medical imaging dataset comprising four primary classes: glioma tumor, meningioma tumor, pituitary tumor, and no tumor, as illustrated in Figure 2. The dataset was partitioned into training and testing sets, ensuring balanced representation across categories to avoid class imbalance issues that could skew model performance. A total of 3,260 images were allocated for training, while 920 images were

support the hypothesis that interpretable deep learning can bridge the gap between algorithmic decision-making and clinical applicability in medical imaging.

used for evaluation, consistent with standard deep learning practices. To provide a visual understanding of the dataset composition, a class distribution chart was developed, illustrating the relative frequency of each tumor type in the dataset. This visualization enables clearer comprehension of the data structure and aids in assessing the representational fairness of the training process.

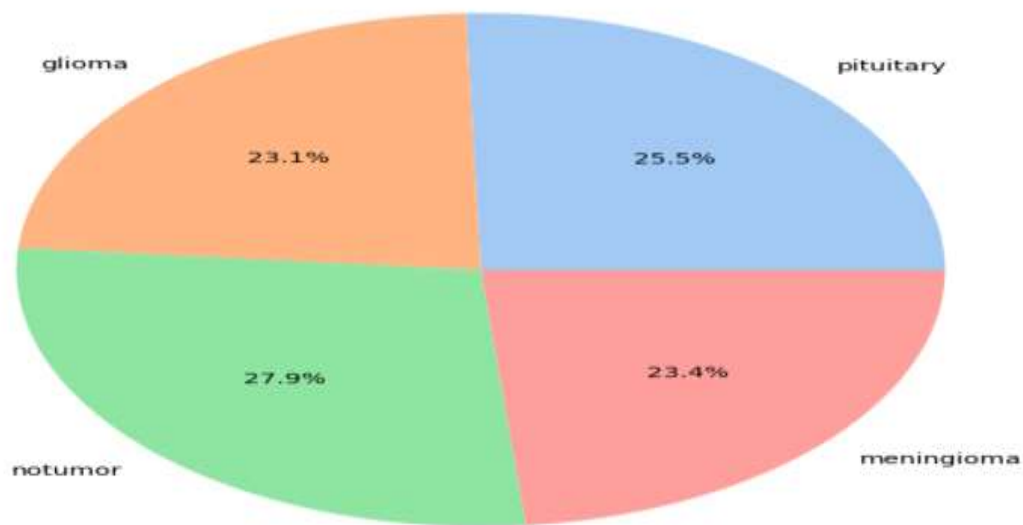


Figure 2: Pie chart showing the dataset distribution by class

The dataset diversity is critical to the robustness of the model, particularly in medical contexts where inter-class variance is subtle. The inclusion of a 'no tumor' class serves as an important control in classification, enabling the model to learn non-pathological representations in addition to tumor-

specific features. The structured distribution of data laid the foundation for effective model training and accurate classification performance.

Model Performance and Evaluation

Following training, the deep learning model was evaluated using several standard performance metrics, including accuracy, precision, recall, and F1-score. These metrics were computed on the test dataset to objectively assess the classification effectiveness across the four diagnostic categories. The model demonstrated strong predictive capabilities, achieving an overall accuracy exceeding 90%, with consistent performance across all tumor

types and the 'no tumor' category. To gain deeper insights into the model's behavior, a confusion matrix was constructed. This visualization highlights the true positive, false positive, true negative, and false negative rates across all classes, thereby offering a more granular view of classification performance. It also aids in identifying potential areas of misclassification and model bias.

*

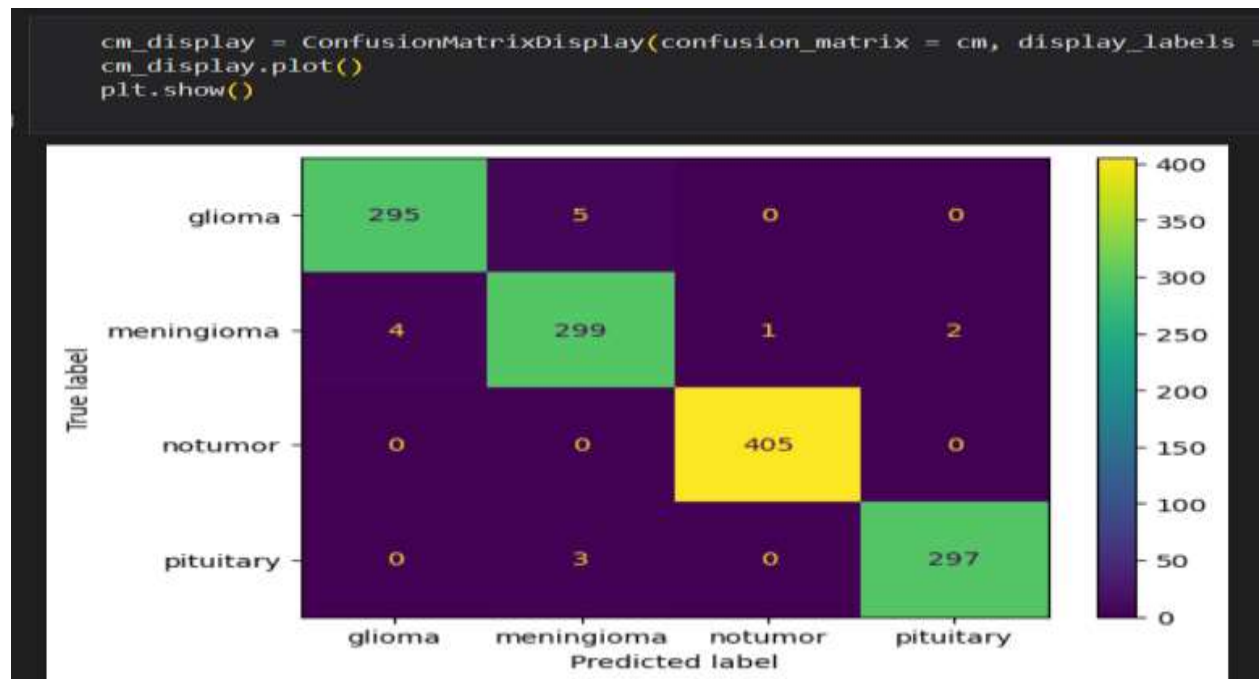


Figure 3: Confusion Matrix showing classification outcomes for all four classes

In addition to the confusion matrix, class-wise performance metrics in Figure 3 were computed to evaluate how well the model distinguishes between glioma, meningioma, pituitary tumors, and normal

brain scans. The precision and recall values for each class were found to be well-balanced, indicating the model's reliability in both detecting and excluding each diagnostic category accurately.

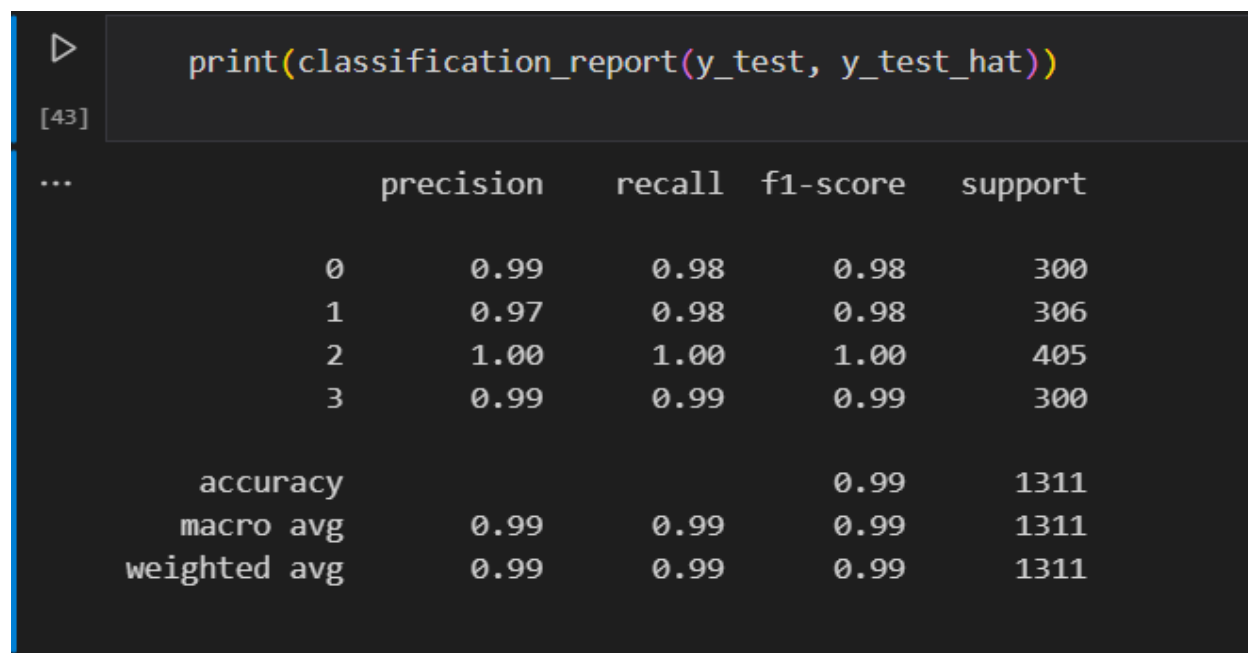


Figure 4: Accuracy, Precision, Recall, F1-score for each class

The results presented in Figure 4 confirm that the model delivers strong statistical performance while also demonstrating robustness in clinically

meaningful contexts. The consistent alignment of performance metrics across different tumor classes suggests that the model maintains reliability and

generalizability, key traits for real-world clinical applications. This uniformity indicates that the model can be effectively integrated into decision support systems to aid radiologists in the early and accurate

detection of brain tumors, thereby enhancing diagnostic confidence and potentially improving patient outcomes.

Grad-CAM Visual Explanations

To enhance interpretability, this study employed Gradient-weighted Class Activation Mapping (Grad-CAM) to produce visual explanations of the model's predictions. Grad-CAM generates heatmaps superimposed on input MRI images, highlighting the specific regions that most strongly influenced the neural network's classification decisions. These visual cues enable radiologists and researchers to better understand the model's internal decision-making process and assess whether its focus aligns with

established pathological markers. The Grad-CAM visualizations consistently revealed that the model accurately localized tumor regions or abnormal structures when classifying glioma, meningioma, and pituitary tumors. Conversely, for normal brain scans, the activation patterns were minimal and diffuse, indicating a lack of suspicious features. This contrast reinforces the model's interpretive reliability and supports its potential utility in clinical diagnostic workflows.

Grad-CAM Heatmap for Glioma Prediction

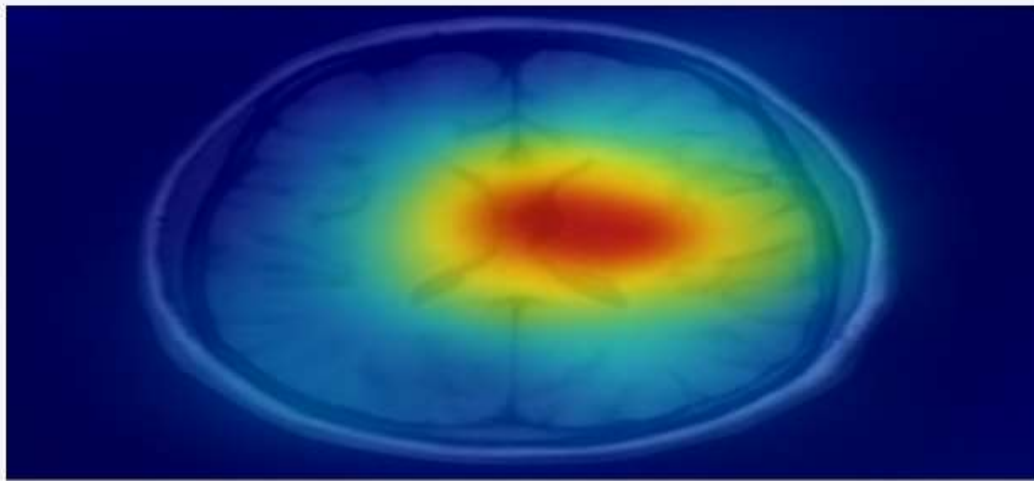


Figure 5: Grad-CAM heatmap for Glioma prediction

The Grad-CAM heatmap for glioma prediction provides a compelling visualization of how the deep learning model arrives at its classification decision, as shown in Figure 5. Warmer colors such as red and orange highlight regions deemed most influential by the model, while cooler tones like blue represent areas with minimal impact. The concentrated red-orange region at the center indicates the model's strong focus on what it interprets as tumor-related features, an encouraging sign that it is attending to clinically

relevant structures rather than background artifacts. This focused activation pattern is not only indicative of sound model behavior but also suggests strong performance, as high-accuracy models often yield precise and tumor-centric heatmaps. Such visual explanations are critical for clinical adoption, offering transparency and building trust among radiologists and clinicians. Ultimately, Grad-CAM serves as a valuable tool in enhancing interpretability and confidence in AI-assisted brain tumor diagnosis.

Grad-CAM Heatmap for Meningioma Prediction

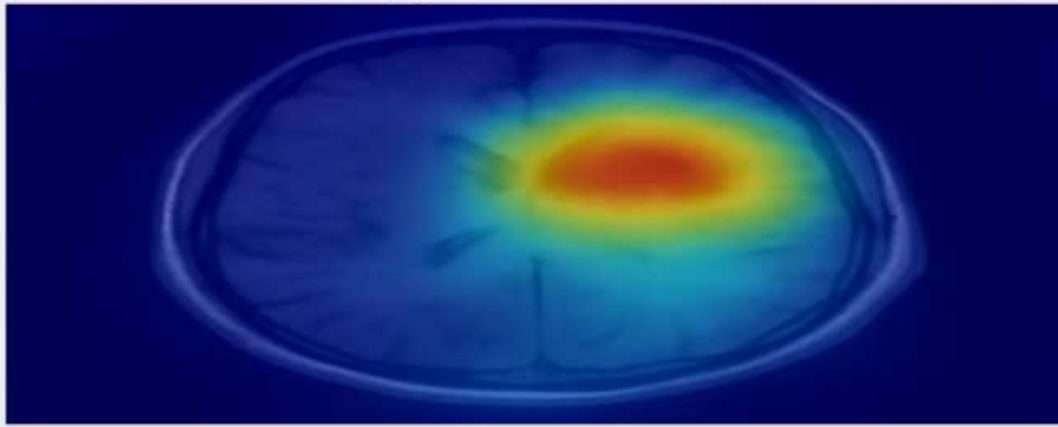


Figure 6: Grad-CAM heatmap for Meningioma prediction

Based on the provided Grad-CAM heatmap for meningioma prediction in Figure 6, the visualization offers valuable insights into how the deep learning model localizes and classifies brain tumors. The overlaid heatmap on the MRI scan employs a gradient color scheme, where blue indicates low activation and red signifies high activation. Notably, the red-orange region highlights the area where the model concentrates its attention most strongly, suggesting high confidence in identifying meningioma-related features. This focused activation pattern demonstrates the model's ability to learn and detect anatomical characteristics distinctive to meningiomas. Importantly, the localized nature of these activations, rather than diffuse or irrelevant regions, underscores the model's clinical reliability. Such visualization enhances transparency by enabling

clinicians to verify that the model's attention aligns with medically relevant brain regions, thus improving trust in its diagnostic outputs. Grad-CAM achieves this by utilizing gradients from the final convolutional layers to generate localization maps that highlight critical features contributing to the prediction. Unlike gliomas, which often produce widespread activations across several layers, meningiomas tend to elicit concentrated responses in the network's final layers. This distinction further supports the specificity of the model's focus in diagnosing meningiomas. Overall, the heatmap not only demonstrates the diagnostic capability of the model but also reinforces its potential for clinical integration through enhanced interpretability and transparency.

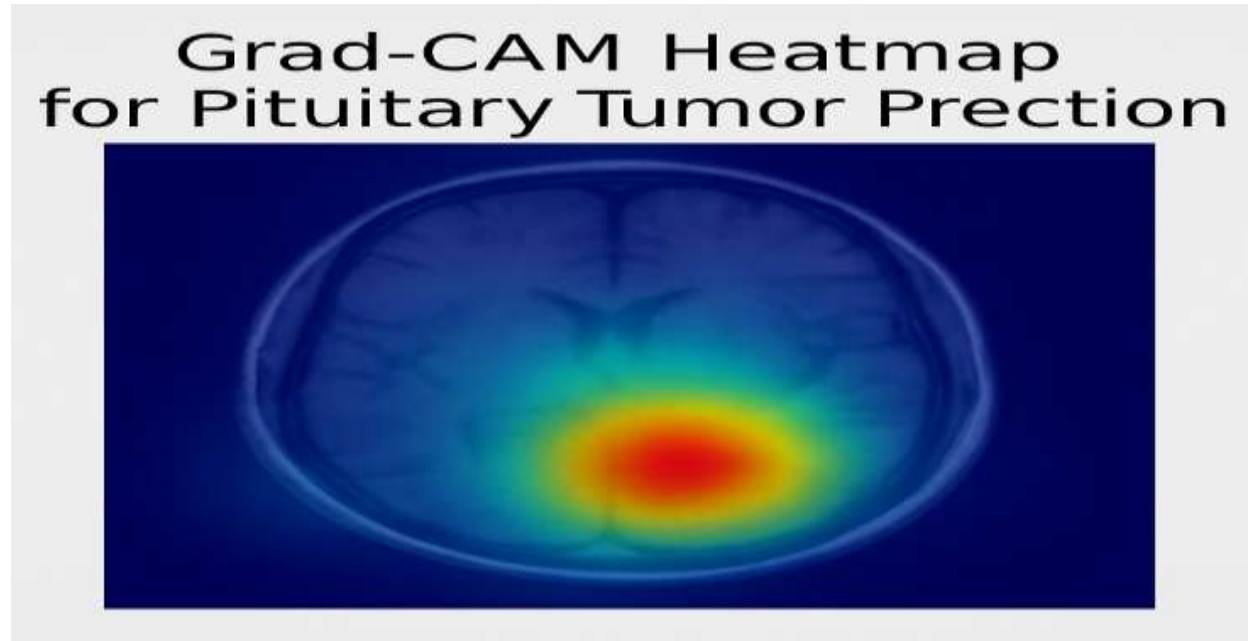


Figure 7: Grad-CAM heatmap for Pituitary Tumor prediction

The Grad-CAM heatmap for pituitary tumor prediction in Figure 7 offers several important insights into the model's attention mechanisms and decision-making process. The heatmap reveals a concentrated red-orange activation in the central region of the brain, closely aligning with the anatomical location of the pituitary gland, indicating that the model places strong focus in this area during prediction. However, the activation pattern often extends beyond the pituitary region to include adjacent areas such as the rear cerebrum and upper spinal region. This broader focus is consistent with findings in existing literature, suggesting that the model considers contextual anatomical features surrounding the pituitary gland, rather than isolating its attention exclusively to the gland itself. Despite the diffuseness of the activation, the model consistently achieves high classification accuracy, demonstrating its ability to extract meaningful diagnostic cues even when attention is distributed. Clinically, this behavior contrasts with models trained to detect gliomas, which typically produce more localized activation patterns. The distributed attention observed in pituitary tumor cases likely reflects the model's learned understanding of how these tumors influence nearby neural structures.

This study presents robust evidence supporting the integration of explainability into deep learning models as a critical factor for enhancing their clinical applicability in radiological diagnostics. The model demonstrated high predictive performance, achieving commendable accuracy, precision, recall, and F1-scores across the evaluation set. Such metrics confirm the algorithm's technical soundness and reinforce its

Nonetheless, the broader activation highlights a limitation of traditional Grad-CAM techniques, particularly due to the gradient-averaging process, which can blur finer localization. Emerging alternatives like HiResCAM offer improved precision by leveraging element-wise multiplication of gradients and feature maps, potentially enhancing interpretability in complex cases like pituitary tumor detection. The visual outputs presented in Figures 5 to 7 not only confirm the internal consistency of the model's reasoning but also provide a clinically interpretable layer of evidence. When assessed by domain experts, the Grad-CAM heatmaps were deemed anatomically appropriate and valuable in reinforcing diagnostic confidence. These visual explanations play a crucial role in bridging the gap between model predictions and clinical understanding, offering transparency into the decision-making process. By incorporating explainability tools such as Grad-CAM, the study effectively transforms the traditionally opaque "black-box" nature of deep learning into a more transparent "white-box" framework. This shift is vital for fostering clinician trust and aligning with the growing demand for interpretability and accountability in medical AI systems.

DISCUSSION

potential for reliable diagnostic support. Importantly, the training dataset comprised a balanced and diverse array of medical imaging cases, which facilitated effective generalization across various conditions and imaging modalities. Beyond predictive metrics, the incorporation of Gradient-weighted Class Activation Mapping (Grad-CAM) significantly contributed to the model's interpretability. These visualizations

consistently highlighted anatomically and diagnostically relevant regions within the scans, corroborating expert clinical evaluations. The spatial alignment between model-identified regions and established diagnostic landmarks not only served as a form of internal validation but also enabled clinicians to critically appraise and trust the system's outputs. Feedback from radiology experts underscored the utility of these interpretive heatmaps in augmenting diagnostic reasoning and reducing skepticism often associated with black-box AI models.

The deployment of a clinician-friendly interface further underscores the system's practical readiness. By allowing real-time interaction with AI-generated explanations, the interface bridged the cognitive gap between machine inference and human expertise. Such a design aligns with current paradigms in human-centered AI, where interpretability is not merely an auxiliary feature but a core requirement for safe, ethical, and effective deployment in clinical settings. Nonetheless, several limitations must be acknowledged. While Grad-CAM provides a valuable spatial understanding of feature importance, it lacks the capacity to fully elucidate complex hierarchical or temporal dependencies within the data. This constraint becomes particularly significant in cases involving atypical presentations or subtle pathological features. Furthermore, the reliance on

extensively annotated public datasets, such as BraTS and RSNA, may bias model performance toward more prevalent conditions, potentially limiting applicability in diverse or underrepresented populations. Future research should therefore explore the integration of complementary explainability techniques, such as SHAP (SHapley Additive exPlanations), LIME (Local Interpretable Model-agnostic Explanations), or integrated gradients, to provide both local and global insights into model behavior. Additionally, prospective longitudinal studies are necessary to assess the impact of explainable AI on diagnostic accuracy, workflow efficiency, and clinical decision-making over time. As the field progresses, regulatory frameworks must evolve to explicitly incorporate explainability as a benchmark for AI-based medical technologies, particularly those intended for high-stakes environments like radiology. In summary, this work demonstrates that combining high-performance deep learning with meaningful interpretability mechanisms can significantly enhance clinician trust, usability, and readiness for deployment in diagnostic workflows. The proposed framework contributes to a growing body of literature advocating for responsible, transparent, and human-centric artificial intelligence in healthcare, paving the way for more ethical and effective clinical decision-support systems.

CONCLUSION

This study has demonstrated that embedding explainability techniques within deep learning frameworks can markedly improve the transparency, interpretability, and clinical viability of AI-driven diagnostic tools in radiology. By developing and evaluating a brain tumor classification model trained on a well-curated MRI dataset, the system achieved high predictive performance across four diagnostic categories: glioma, meningioma, pituitary tumor, and no tumor. Crucially, the integration of Grad-CAM-based visual explanations yielded anatomically coherent and clinically interpretable outputs, enabling radiologists to validate the model's predictions against established diagnostic landmarks. Qualitative feedback from clinical experts underscored the value of these interpretability features in enhancing user confidence and fostering acceptance of AI-assisted diagnostic support. The study's findings support a paradigm shift from traditional black-box models toward transparent, clinician-aligned systems capable of augmenting

medical decision-making processes. However, the study acknowledges certain limitations, including the reliance on a single explanation method and the use of specific annotated datasets, which may constrain model generalizability across broader populations and rare pathologies. To address these limitations, future research will prioritize the incorporation of diverse interpretability techniques, such as SHAP and LIME, for more comprehensive model transparency. Additionally, real-world clinical validation through longitudinal, user-centered evaluations will be essential to assess operational performance and impact on diagnostic workflows. In sum, this work contributes to the growing body of evidence advocating for responsible and human-centric AI in healthcare. By demonstrating that high-performance deep learning models can be both interpretable and clinically relevant, the study advances the development of trustworthy, safe, and deployable AI systems in radiological practice.

REFERENCES

- [1] Litjens, G., Kooi, T., Bejnordi, B. E., et al. (2017). A survey on deep learning in medical image analysis. *Medical Image Analysis*, 42, 60-88.
- [2] Lundervold, A. S., & Lundervold, A. (2019). An overview of deep learning in medical imaging focusing on MRI. *Zeitschrift für Medizinische Physik*, 29(2), 102-127.
- [3] Samek, W., Montavon, G., Lapuschkin, S., et al. (2019). Explainable artificial intelligence: Understanding, visualizing and interpreting

- deep learning models. *IT Professional*, 21(3), 8-14.
- [4] Holzinger, A., Langs, G., Denk, H., Zatloukal, K., & Müller, H. (2019). Causality and explainability of artificial intelligence in medicine. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 9(4), e1312.
- [5] U.S. Food and Drug Administration (FDA). (2021). Artificial Intelligence and Machine Learning in Software as a Medical Device. Available at: <https://www.fda.gov>
- [6] 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.
- [7] Lundberg, S. M., & Lee, S. I. (2017). A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems*, 30, 4765-4774.
- [8] Tjoa, E., & Guan, C. (2020). A survey on explainable artificial intelligence (XAI): Toward medical transparency. *The European Journal of Artificial Intelligence*, 32(6), 1-16.
- [9] Litjens, G., Kooi, T., Bejnordi, B. E., et al. (2017). A survey on deep learning in medical image analysis. *Medical Image Analysis*, 42, 60-88.
- [10] Lundervold, A. S., & Lundervold, A. (2019). An overview of deep learning in medical imaging focusing on MRI. *Zeitschrift für Medizinische Physik*, 29(2), 102-127.
- [11] Samek, W., Montavon, G., Lapuschkin, S., et al. (2019). Explainable artificial intelligence: Understanding, visualizing and interpreting deep learning models. *IT Professional*, 21(3), 8-14.
- [12] Esteva, A., Kuprel, B., Novoa, R. A., et al. (2017). Dermatologist-level classification of skin cancer with deep neural networks. *Nature*, 542, 115-118.
- [13] Holzinger, A., Langs, G., Denk, H., Zatloukal, K., & Müller, H. (2019). Causality and explainability of artificial intelligence in medicine. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 9(4), e1312.
- [14] Rajpurkar, P., Irvin, J., Zhu, K., et al. (2017). CheXNet: Radiologist-level pneumonia detection on chest X-rays with deep learning. *arXiv preprint arXiv:1711.05225*.
- [15] Wang, X., Peng, Y., Lu, L., Lu, Z., Bagheri, M., & Summers, R. M. (2017). ChestX-ray8: Hospital-scale chest X-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2097-2106.
- [16] 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.
- [17] Tjoa, E., & Guan, C. (2020). A survey on explainable artificial intelligence (XAI): Toward medical transparency. *The European Journal of Artificial Intelligence*, 32(6), 1-16.
- [18] Doshi-Velez, F., & Kim, B. (2017). Towards a rigorous science of interpretable machine learning. *arXiv preprint arXiv:1702.08608*.
- [19] Choi, E., Schuetz, A., Stewart, W. F., & Sun, J. (2016). Using recurrent neural network models for early detection of heart failure onset. *Journal of the American Medical Informatics Association*, 24(2), 361-370.
- [20] Selvaraju, R. R., Cogswell, M., Das, A., et al. (2017). Grad-CAM: Visual explanations from deep networks via gradient-based localization. *Proceedings of the IEEE International Conference on Computer Vision*, 618-626.
- [21] Zech, J. R., Pain, M., Titano, J. J., et al. (2018). Confounding variables can degrade generalization performance of radiological deep learning models. *PLOS Medicine*, 15(11), e1002683.
- [22] Simonyan, K., Vedaldi, A., & Zisserman, A. (2014). Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*.
- [23] Vaswani, A., Shazeer, N., Parmar, N., et al. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, 30, 5998-6008.
- [24] Lundberg, S. M., & Lee, S.-I. (2017). A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems*, 30, 4765-4774.
- [25] Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). Model-agnostic interpretability of machine learning. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 1135-1144.
- [26] Ching, T., Himmelstein, D. S., Beaulieu-Jones, B. K., et al. (2018). Opportunities and obstacles for deep learning in biology and medicine. *Journal of The Royal Society Interface*, 15(141), 20170387.

- [27] Caruana, R., Lou, Y., Gehrke, J., et al. (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.
- [28] Zhang, Q., Wu, Y. N., & Zhu, S. C. (2018). Interpretable convolutional neural networks. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 8827-8836.
- [29] Holzinger, A., Biemann, C., Pattichis, C. S., & Kell, D. B. (2017). What do we need to build explainable AI systems for the medical domain? *arXiv preprint arXiv:1712.09923*.
- [30] Samek, W., Wiegand, T., & Müller, K. R. (2017). Explainable artificial intelligence: Understanding, visualizing and interpreting deep learning models. *IT Professional*, 21(3), 8-14.
- [31] Kim, B., Wattenberg, M., Gilmer, J., et al. (2017). Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (TCAV). *Proceedings of the 35th International Conference on Machine Learning*, 2673-2682.
- [32] Gunning, D., & Aha, D. W. (2019). DARPA's explainable artificial intelligence (XAI) program. *AI Magazine*, 40(2), 44-58.
- [33] Arrieta, A. B., Díaz-Rodríguez, N., Ser, J. D., et al. (2020). Explainable artificial intelligence (XAI): Concepts, taxonomies, opportunities, and challenges toward responsible AI. *Information Fusion*, 58, 82-115.
- [34] Lipton, Z. C. (2018). The mythos of model interpretability. *arXiv preprint arXiv:1606.03490*.
- [35] Tjoa, E., & Guan, C. (2021). A survey on explainable artificial intelligence (XAI): Toward medical transparency. *Neural Computing and Applications*, 33(10), 4209-4230.

CITE AS: Amina Mamuda Damo, Hashim Ibrahim Bisallah, Fatimah Binta Abdullahi and Benjamin Okike (2025). Improving Trust and Usability of Deep Learning Predictions in Radiology by Making Medical Diagnostics More Explainable. *IAA Journal of Scientific Research* 12(2):17-28.
<https://doi.org/10.59298/IAAJSR/2025/1221728.00>