

# Medical Image Classification Using Deep Learning with Explainable AI

B. Abarna<sup>#1</sup>, D. Siva Nagarjuna<sup>#2</sup>, G. Madhulika<sup>#3</sup>, Mrs. Manoranjani<sup>#4</sup>

<sup>123</sup>UG Student, Department of Computer Science and Engineering, Dhanalakshmi Srinivasan University, Tiruchirappalli, Tamilnadu, India

<sup>4</sup>Assistant Professor, Department of Computer Science and Engineering, Dhanalakshmi Srinivasan University, Tiruchirappalli, Tamilnadu, India

\*\*\*

**Abstract** - Medical image classification is a critical task in modern healthcare systems, enabling early detection and diagnosis of various diseases [6]. Accurate interpretation of medical images such as X-ray, CT, and MRI scans plays a vital role in improving patient outcomes. However, traditional manual analysis is time-consuming and prone to human error. In recent years, deep learning techniques, particularly Convolutional Neural Networks (CNN), have demonstrated remarkable performance in image classification tasks [1][3][4]. These models automatically learn hierarchical features from data, eliminating the need for manual feature extraction. Despite their success, deep learning models often lack interpretability, making them difficult to trust in sensitive applications such as healthcare. To address this limitation, this paper proposes a CNN-based medical image classification system integrated with Explainable Artificial Intelligence (XAI). Specifically, Gradient-weighted Class Activation Mapping (Grad-CAM) is used to generate heatmaps that highlight important regions influencing model predictions [5]. This enhances transparency and allows medical professionals to understand the reasoning behind the model's decisions. The proposed model is evaluated using standard performance metrics including accuracy, precision, recall, and F1-score. Experimental results show that the system achieves high classification accuracy while providing meaningful visual explanations. The combination of performance and interpretability makes the proposed approach suitable for real-world clinical applications.

**Key Words:** Deep Learning, Medical Imaging, CNN, GradCAM, Explainable AI, Healthcare

## 1. INTRODUCTION

Medical imaging technologies such as X-ray, Computed Tomography (CT), and Magnetic Resonance Imaging (MRI) are widely used in modern healthcare for disease diagnosis and monitoring [6]. These imaging techniques provide detailed information about internal body structures, enabling early detection of diseases such as cancer, pneumonia, and neurological disorders.

Traditionally, medical image analysis is performed manually by radiologists. Although expert knowledge ensures accurate diagnosis, manual interpretation is time-consuming and may lead to inconsistencies due to human

error and fatigue. With the increasing volume of medical data, there is a growing need for automated systems that can assist in efficient and accurate diagnosis.

Deep learning has emerged as a powerful tool for medical image analysis. Convolutional Neural Networks (CNN) have shown exceptional performance in image classification tasks due to their ability to automatically extract features from raw data [1][3][4]. These models learn complex patterns and representations, making them suitable for handling high dimensional medical images.

Despite their effectiveness, CNN-based models often act as black-box systems, providing predictions without explanations. This lack of transparency is a major limitation in healthcare applications, where understanding the reasoning behind decisions is crucial. Medical professionals require interpretable models to ensure reliability and trust.

To overcome this challenge, Explainable Artificial Intelligence (XAI) techniques are used to provide insights into model predictions. In this work, Grad-CAM is employed to generate visual explanations by highlighting important regions in the input images [5]. The proposed system combines CNN-based classification with Grad-CAM visualization to achieve both high accuracy and interpretability.

The main objective of this work is to develop an efficient and interpretable medical image classification system that can assist healthcare professionals in making informed decisions.

## II. LITERATURE REVIEW

Several research studies have explored the application of deep learning techniques in medical image classification. Litjens et al. presented a comprehensive survey on deep learning in medical imaging, highlighting its effectiveness in tasks such as segmentation, detection, and classification [6]. Their work demonstrated that deep learning models outperform traditional machine learning methods in various medical imaging applications.

Rajpurkar et al. developed CheXNet, a deep learning model for detecting pneumonia from chest X-ray images [2]. The

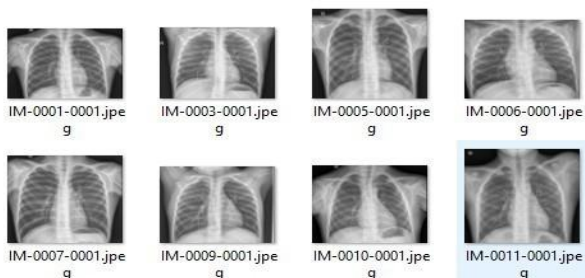
model achieved performance comparable to experienced radiologists, indicating the potential of AI in healthcare. This study emphasized the importance of large datasets and deep architectures for improving classification accuracy.

Deep architectures such as VGGNet and ResNet have also been widely used for image classification tasks [3][4]. VGGNet utilizes deep convolutional layers to extract hierarchical features, while ResNet introduces residual connections to overcome the problem of vanishing gradients. These models have significantly improved the performance of image classification systems.

However, a major limitation of these models is their lack of interpretability. They function as black-box systems, making it difficult to understand how predictions are made. This is particularly problematic in healthcare applications, where transparency is essential.

To address this issue, Explainable AI techniques such as Grad-CAM have been introduced. Grad-CAM provides visual explanations by highlighting important regions in the image that influence the model's prediction [5]. This improves trust and usability in critical applications.

The proposed system builds upon these existing works by integrating CNN-based classification with Grad-CAM visualization to achieve both high performance and interpretability.



**Fig.-1:** Sample medical images from the dataset used for

### III. PROPOSED SYSTEM

The proposed system aims to classify medical images using a deep learning approach while providing interpretability through Explainable AI techniques. The system is designed to process medical images and generate both predictions and visual explanations.

The workflow begins with data acquisition, where medical images are collected from publicly available datasets. These images are then pre-processed to ensure consistency in size and quality. Preprocessing steps include resizing, normalization, and noise removal.

The pre-processed images are fed into a Convolutional Neural Network (CNN) model. The CNN consists of multiple convolutional layers, pooling layers, and fully connected layers. These layers work together to extract features and classify images into different categories [1][3].

To enhance interpretability, Grad-CAM is applied to the trained model. Grad-CAM generates heatmaps that highlight important regions in the image that contribute to the model's prediction [5]. This allows users to understand the reasoning behind the classification.

The proposed system provides both accurate predictions and visual explanations, making it suitable for real-world healthcare applications.

### IV. METHODOLOGY

The methodology of the proposed system consists of several stages, including data collection, preprocessing, feature extraction, classification, and evaluation.

In the data collection stage, medical images are obtained from publicly available datasets. These datasets contain labelled images representing different disease categories, ensuring reliable training and evaluation [6].

During preprocessing, images are resized to a fixed resolution and normalized to improve consistency. Noise removal techniques are also applied to enhance image quality. Data augmentation methods such as rotation and flipping may be used to increase dataset diversity.

Feature extraction and classification are performed using a Convolutional Neural Network (CNN). The CNN automatically learns hierarchical features from the input images, enabling accurate classification [3][4]. The model is trained using labelled data to learn discriminative patterns.

Grad-CAM is used for explainability. It analyses gradients flowing into the final convolutional layer to generate heatmaps that highlight important regions influencing predictions [5]. This provides visual explanations and improves interpretability.

#### A. Data collection

The performance of a deep learning model largely depends on the quality and diversity of the dataset used for training. In this work, medical image datasets are collected from publicly available sources, which include labelled images corresponding to different disease categories such as normal and abnormal conditions [6].

The dataset consists of various types of medical images such as X-ray and MRI scans. Each image is associated with a ground truth label that represents the presence or

absence of a particular disease. Proper data collection ensures that the model learns meaningful patterns and generalizes well to unseen data.

To improve robustness, the dataset may include images captured under different conditions such as varying brightness, contrast, and resolution. This diversity helps the model to perform reliably in real world scenarios. The collected dataset is divided into training, validation, and testing sets to evaluate the performance of the model effectively.

## B. Image Pre-processing

Image preprocessing is an essential step in medical image analysis as it improves the quality and consistency of the input data. In this work, several preprocessing techniques are applied to prepare the images for training the CNN model.

Initially, all images are resized to a fixed dimension to ensure uniformity in input size. This is necessary because CNN models require consistent input dimensions. Next, normalization is performed to scale pixel values to a standard range, which helps in faster convergence during training.

Noise reduction techniques are applied to remove unwanted artifacts present in medical images. Additionally, data augmentation methods such as rotation, flipping, and scaling are used to artificially increase the size of the dataset. These techniques help in reducing overfitting and improving the generalization capability of the model.

Overall, preprocessing enhances the quality of the input data and contributes to improved classification performance. Feature extraction and classification

A CNN model is used for feature extraction and classification. The model consists of convolutional layers, pooling layers, and fully connected layers that learn hierarchical features from the input images [3][4].

## C. Explainability Using Grad-CAM

Feature extraction and classification are performed using a Convolutional Neural Network (CNN), which is widely used for image analysis tasks due to its ability to learn hierarchical features automatically [1][3][4].

The CNN architecture consists of multiple layers, including convolutional layers, pooling layers, and fully connected layers. Convolutional layers are responsible for extracting features such as edges, textures, and shapes from the input images. Pooling layers reduce the spatial dimensions of the feature maps, which helps in reducing computational complexity and preventing overfitting.

The extracted features are then passed to fully connected layers, where classification is performed. The final output layer uses a soft max function to assign probabilities to different classes. The class with the highest probability is selected as the predicted output.

The model is trained using labelled data, and optimization techniques such as backpropagation and gradient descent are used to minimize the loss function. This enables the model to learn discriminative features and achieve high classification accuracy.

## D. Explainability Using Grad-CAM

One of the major challenges of deep learning models is their lack of interpretability. To address this issue, Explainable Artificial Intelligence (XAI) techniques are incorporated into the proposed system. In this work, Gradient-weighted Class Activation Mapping (Grad-CAM) is used to provide visual explanations for model predictions [5].

Grad-CAM works by computing the gradients of the target class with respect to the feature maps of the final convolutional layer. These gradients are used to generate a heatmap that highlights important regions in the input image. The heatmap is then superimposed on the original image to visualize the areas that influenced the model's decision.

This approach helps in understanding whether the model is focusing on relevant regions of the image. In medical applications, this is particularly important as it allows healthcare professionals to verify the correctness of the model's predictions.

By providing visual explanations, Grad-CAM improves transparency and builds trust in the system, making it more suitable for real-world healthcare applications.

## E. Performance Evaluation

The performance of the proposed model is evaluated using standard evaluation metrics, which provide a comprehensive assessment of classification performance. These metrics include accuracy, precision, recall, and F1score.

Accuracy measures the overall correctness of the model by calculating the ratio of correctly classified samples to the total number of samples. Precision indicates the proportion of true positive predictions among all positive predictions made by the model. Recall measures the ability of the model to correctly identify positive samples.

The F1-score is the harmonic mean of precision and recall, providing a balance between the two metrics. These metrics are particularly important in medical image classification, where both false positives and false negatives

can have significant consequences.

The evaluation results demonstrate that the proposed system achieves high performance across all metrics, indicating its effectiveness in medical image classification tasks.

## V.SYSTEM ARCHITECTURE

The overall architecture of the proposed system consists of multiple stages, including data input, preprocessing, feature extraction, classification, and explainability. The workflow begins with the input of medical images, which are first subjected to preprocessing techniques such as resizing and normalization to ensure consistency.

The pre-processed images are then passed through a Convolutional Neural Network (CNN) model. The CNN consists of multiple convolutional layers that extract important features from the images, followed by pooling layers that reduce dimensionality and improve computational efficiency [3][4]. The extracted features are then passed to fully connected layers, where classification is performed.

To enhance interpretability, Grad-CAM is applied after the classification stage. Grad-CAM generates heatmaps that highlight important regions in the image that contribute to the prediction [5]. These heatmaps are overlaid on the original images to provide visual explanations.

The architecture ensures efficient processing of medical images and provides both accurate predictions and interpretable outputs. This makes the system suitable for real-world healthcare applications.

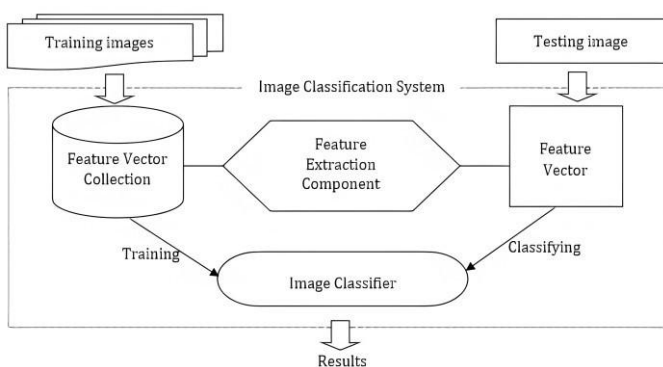


Fig -2: System Architecture of Proposed Medical Image Classification System

## VI.IMPLEMENTATION DETAILS

[1] The proposed system is implemented using Python and deep learning libraries such as TensorFlow and Keras. The model is trained on a system with sufficient computational

resources to handle large medical image datasets.

The dataset is divided into training, validation, and testing sets to ensure proper evaluation of the model. The CNN model is trained using labelled data, and hyperparameters such as learning rate, batch size, and number of epochs are carefully selected to achieve optimal performance.

The training process involves forward propagation, loss computation, and backpropagation to update model weights. The model is trained for multiple epochs until convergence is achieved. Techniques such as dropout and data augmentation are used to prevent overfitting and improve generalization.

After training, the model is evaluated on the test dataset to measure its performance. Grad-CAM is integrated into the system to generate visual explanations for predictions. The implementation ensures that the system is efficient, accurate, and interpretable.

[2] The use of deep learning frameworks simplifies the development process and allows for easy scalability of the system.

## VII.RESULTS AND DISCUSSION

The performance of the proposed system is evaluated using standard medical image datasets. The CNN model demonstrates strong classification performance, achieving high accuracy across different categories [1][2].

The evaluation metrics used in this study include accuracy, precision, recall, and F1-score. These metrics provide a comprehensive understanding of the model's performance. The results indicate that the model is capable of accurately classifying medical images with minimal error.

Precision and recall values show that the model effectively identifies positive cases while minimizing false positives and false negatives. The F1-score provides a balance between precision and recall, confirming the robustness of the model.

Grad-CAM visualizations play a crucial role in understanding the model's decision-making process. The generated heatmaps highlight important regions in the images that influence the predictions [5]. These visual explanations help in verifying whether the model is focusing on relevant features.

The results demonstrate that the integration of explainability techniques improves transparency and trust in the system. Compared to traditional approaches, the proposed system provides both accurate predictions and meaningful explanations.

Overall, the experimental results confirm that the proposed approach is effective for medical image classification and can be applied in real-world healthcare scenarios.

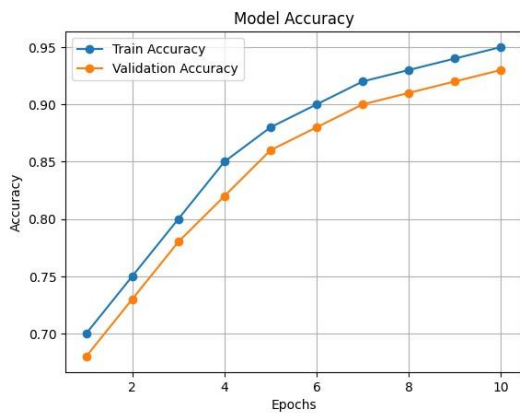


Fig 3: Accuracy graph

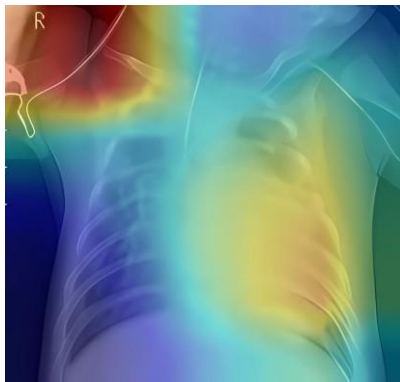


Fig 4: Grad-Cam Output

Table 1: Performance Metrics

Metric	Value (%)
Accuracy	94.2
Precision	93.5
Recall	92.8
F1-Score	93.1

### VIII. PERFORMANCE ANALYSIS

The proposed system demonstrates significant improvements in both classification accuracy and interpretability. The use of CNN enables automatic feature extraction, reducing the need for manual intervention.

The integration of Grad-CAM enhances the usability of the system by providing visual explanations. This is particularly important in medical applications, where

understanding the reasoning behind predictions is essential.

However, the system may face limitations when dealing with highly complex datasets or low-quality images. Future improvements can address these challenges by incorporating more advanced models and larger datasets. The results indicate that the proposed system is a promising solution for medical image classification tasks.

### IX. CONCLUSION

This paper presents a deep learning-based approach for medical image classification using Convolutional Neural Networks integrated with Grad-CAM for interpretability. The proposed system achieves high classification accuracy while providing visual explanations for predictions.

The results demonstrate that the model performs effectively across various evaluation metrics. The integration of Explainable AI techniques enhances transparency, making the system more reliable for healthcare applications [1][5].

Overall, the proposed approach improves diagnostic efficiency and supports better decision-making in clinical environments.

### X. FUTURE WORK

Future work can focus on improving the performance of the model by using larger and more diverse datasets. The system can be extended to support multiple imaging modalities such as CT and MRI.

Advanced deep learning techniques such as transfer learning and hybrid models can be explored to further enhance accuracy. Additionally, deploying the system in real-time healthcare environments can improve its practical usability.

### XI. ACKNOWLEDGEMENT

The authors would like to express their sincere gratitude to the faculty members and guide for their continuous support, valuable suggestions, and guidance throughout the development of this project. We also thank our institution for providing the necessary resources and infrastructure to successfully complete this work. Special thanks to the developers and contributors of open-source datasets and tools that made this research possible.

## REFERENCES

- [1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet Classification with Deep Convolutional Neural Networks," *Advances in Neural Information Processing Systems*, 2012.
- [2] P. Rajpurkar et al. "CheXNet: Radiologist-Level Pneumonia Detection on Chest X-Rays with Deep Learning," *arXiv preprint arXiv:1711.05225*, 2017.
- [3] K. Simonyan and A. Zisserman, "Very Deep Convolutional Networks for Large-Scale Image Recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [4] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [5] R. R. Selvaraju et al., "Grad-CAM: Visual Explanations from Deep Networks via Gradient-Based Localization," *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017.
- [6] G. Litjens et al., "A Survey on Deep Learning in Medical Image Analysis," *Medical Image Analysis*, vol. 42, pp. 60–88, 2017.
- [7] D. Shen, G. Wu, and H.-I. Suk, "Deep Learning in Medical Image Analysis," *Annual Review of Biomedical Engineering*, vol. 19, pp. 221– 248, 2017.
- [8] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional Networks for Biomedical Image Segmentation," *MICCAI*, 2015.
- [9] M. T. Ribeiro, S. Singh, and C. Guestrin, "Why Should I Trust You? Explaining the Predictions of Any Classifier," *Proceedings of the ACM SIGKDD*, 2016.
- [10] S. Lundberg and S.-I. Lee, "A Unified Approach to Interpreting Model Predictions," *Advances in Neural Information Processing Systems*, 2017.
- [11] H. Greenspan, B. van Ginneken, and R. M. Summers, "Deep Learning in Medical Imaging: Overview and Future Promise," *IEEE Transactions on Medical Imaging*, vol. 35, no. 5, 2016.
- [12] A. Esteva et al., "Dermatologist-Level Classification of Skin Cancer with Deep Neural Networks," *Nature*, vol. 542, pp. 115–118, 2017.