# Enhancing Interpretability in Neural Networks through Explainable AI (XAI) Techniques

**A. Gupta** 🆔 *

Independent Researcher, San Ramon, California, USA

**To cite this article:** A. Gupta, "Enhancing Interpretability in Neural Networks through Explainable AI (XAI) Techniques," *Insight. Electr. Electron. Eng.*, vol. 1, no. 1, pp. 1-5, 2024.

## Abstract

The rapid advancement of neural networks in various applications, from healthcare diagnostics to financial modeling, has significantly improved the accuracy and efficiency of decision-making processes. However, these models often operate as black boxes, providing little to no insight into how they arrive at specific predictions. This lack of interpretability presents a major barrier to their adoption in critical domains where trust, accountability, and transparency are paramount. This study aims to address this issue by developing a novel framework that integrates multiple explainable AI (XAI) techniques to enhance the interpretability of neural networks. The proposed framework combines feature importance analysis, layer-wise Relevance Propagation (LRP), and visual explanation methods such as Gradient-weighted Class Activation Mapping (Grad-CAM). These techniques collectively offer a comprehensive view of the decision-making processes of neural networks, making them more transparent and understandable to stakeholders. Our experimental results demonstrate that the integrated XAI framework not only improves interpretability but also maintains high levels of accuracy, thereby bridging the gap between performance and transparency. This research provides a foundational basis for the deployment of interpretable neural networks in critical applications, ensuring that AI-driven decisions are reliable and comprehensible.

**Keywords:** neural networks; explainable AI; Grad-CAM; interpretability; accuracy

**Abbreviations:** XAI: explainable AI; LRP: layer-wise Relevance Propagation; Grad-CAM: Gradient-weighted Class Activation Mapping; AI: artificial intelligence; FNNs: feedforward neural networks; CNNs: convolutional neural networks; SHAP: SHapley Additive exPlanations

## 1. Introduction and Background

### 1.1. Introduction

Artificial intelligence (AI) has become a cornerstone of modern technological advancements, with neural networks playing a pivotal role in various applications such as image recognition, natural language processing, and predictive analytics. Despite their success, one of the major challenges that impede the broader acceptance of neural networks, especially in critical fields like healthcare, finance, and autonomous systems, is their lack of interpretability. The black-box nature of these models makes it difficult to understand how they process input data and generate outputs, leading to issues of trust and accountability. Explainable AI (XAI) has emerged as a crucial area of research aimed at making AI systems more transparent and interpretable. XAI techniques strive to elucidate the internal workings of complex models, thereby allowing users to comprehend, trust, and effectively manage AI-driven decisions. This paper focuses on enhancing the interpretability of neural networks by integrating various XAI techniques into a cohesive framework. The goal is to provide stakeholders with clear and actionable explanations of the model's predictions, facilitating trust and enabling the deployment of AI systems in high-stakes environments.

### 1.2. Background

The motivation for this research stems from the increasing demand for transparency and accountability in AI systems. In healthcare, for example, clinicians need to understand AI-driven diagnostic recommendations to trust and act on them. Similarly, in finance, stakeholders must comprehend AI-based risk assessments to ensure fairness and regulatory compliance. In autonomous systems, such as self-driving cars, understanding the decision-making process is crucial for safety and reliability. Addressing these needs, our study aims to bridge the gap between high-performing neural networks and the essential requirement for interpretability, thus fostering greater acceptance and trust in AI systems across various critical applications. Neural networks, particularly deep learning models, have achieved unprecedented success in numerous applications due to their ability to learn from large datasets and capture intricate patterns. However, their complex architectures, often consisting of multiple hidden layers and millions of parameters, render them opaque and difficult to interpret. The need for explainability in AI has led to the development of several XAI techniques designed to demystify these black-box models [1, 2].

## 2. Methodology

### 2.1. Data collection

*Corresponding Author:
A. Gupta, Independent Researcher, San Ramon, California, USA

The first step in our methodology involves the selection and preprocessing of datasets from various critical domains to ensure the robustness and applicability of our XAI framework. We focus on three primary domains: healthcare, finance, and image recognition [3, 4].

### 2.1.1. Healthcare

- **Dataset:** We use publicly available medical datasets, such as the MIMIC-III Clinical Database, which contains de-identified health-related data from patients [5].

- **Preprocessing:** This includes handling missing values, normalizing continuous variables, and encoding categorical variables. Additionally, we segment the data into training, validation, and test sets to evaluate the model's performance accurately.

### 2.1.2. Finance

- **Dataset:** Financial datasets such as the LendingClub loan data, which includes information about loan applications and their outcomes, are utilized [6].

- **Preprocessing:** Similar preprocessing steps are followed, including data cleaning, normalization, and feature engineering to extract meaningful financial indicators.

### 2.1.3. Image recognition

- **Dataset:** The CIFAR-10 dataset, a widely used benchmark dataset consisting of 60,000 32 x 32 color images in 10 different classes, is employed.

- **Preprocessing:** Image data is preprocessed through normalization and data augmentation techniques to enhance model generalization.

## 2.2. Model development

Following data collection, we develop neural network models tailored to each dataset, leveraging standard architectures suitable for the respective domains.

### 2.2.1. Healthcare and finance

- **Architecture:** For tabular data, we use feedforward neural networks (FNNs) with multiple hidden layers. Each hidden layer employs ReLU activation functions, followed by dropout layers to prevent overfitting.

- **Training:** The models are trained using the Adam optimizer with an appropriate learning rate, and early stopping is implemented to avoid overfitting.

### 2.2.2. Image recognition

- **Architecture:** For image recognition tasks, we utilize convolutional neural networks (CNNs) with multiple convolutional and pooling layers, followed by fully connected layers.

- **Training:** The CNNs are trained using stochastic gradient descent with momentum, and techniques like batch normalization and dropout are used to enhance training stability and model performance.

## 2.3. Explainability techniques integration

To enhance interpretability, we integrate multiple XAI techniques into the neural network models, providing a comprehensive understanding of their decision-making processes [7].

### 2.3.1. Feature importance analysis

- **SHapley Additive exPlanations (SHAP):** SHAP values are computed for each feature in the healthcare and finance models to determine their contributions to the predictions. These values provide a global understanding of feature importance across the dataset.

### 2.3.2. Layer-wise Relevance Propagation (LRP)

- **Implementation:** LRP is applied to both FNNs and CNNs to trace the contributions of input features through the layers of the network to the final output. This technique decomposes the prediction score back to the input features, highlighting their relevance at each layer.

### 2.3.3. Visual explanation methods

- **Gradient-weighted Class Activation Mapping (Grad-CAM):** Grad-CAM is used for CNNs in the image recognition domain to generate heatmaps that highlight important regions in the input images. These visualizations indicate which parts of the image most influenced the model's prediction.

## 2.4. Evaluation metrics

The evaluation of our XAI framework involves assessing both the interpretability and performance of the neural network models.

### 2.4.1. Model performance

- **Metrics:** Standard performance metrics such as accuracy, precision, recall, and F1-score are calculated for the models on the test datasets. These metrics ensure that the integration of XAI techniques does not compromise the predictive performance of the models.

### 2.4.2. Interpretability assessment

- **Qualitative analysis:** The clarity and usefulness of the explanations provided by the XAI techniques are evaluated through qualitative analysis. Visualizations of feature importance, LRP maps, and Grad-CAM heatmaps are analyzed to ensure they offer meaningful insights into the model's decision-making processes.

- **User studies:** Conducting user studies with domain experts and stakeholders to evaluate the practical applicability and understandability of the provided explanations. Feedback from these studies is used to refine and improve the XAI framework.

### 2.4.3. Computational efficiency

- **Metrics:** The computational overhead introduced by the XAI techniques is measured in terms of additional training and inference time. This evaluation ensures that the interpretability enhancements are achieved without significant performance degradation.

## 3. Results

The results section presents the findings of our research on enhancing interpretability in neural networks through the integration of XAI techniques. We detail the model performance, interpretability assessment, and computational efficiency of the proposed framework across different domains, providing a comprehensive evaluation of its effectiveness.

### 3.1. Model performance

### 3.1.1. Healthcare and finance models

- **Accuracy:** The FNNs developed for healthcare and finance datasets achieved high accuracy. For the healthcare dataset, the model attained an accuracy of 87.5%, while the finance model achieved 85.3%.

- **Precision, recall, and F1-score:** Both models exhibited strong performance across these metrics, with the healthcare model recording a precision of 88.2%, recall of 86.7%, and F1-score of 87.4%. The finance model demonstrated a precision of 84.9%, recall of 85.7%, and F1-score of 85.3%.

### 3.1.2. Image recognition model

- **Accuracy:** The CNN trained on the CIFAR-10 dataset achieved an accuracy of 91.2%.

- **Precision, recall, and F1-score:** The CNN showed robust performance with a precision of 90.8%, recall of 91.5%, and F1-score of 91.1%.

These results indicate that integrating XAI techniques did not compromise the predictive accuracy of the models. Instead, the models maintained high levels of performance across all metrics, validating the effectiveness of our approach.

### 3.2. Interpretability assessment

### 3.2.1. Feature importance analysis

- **SHAP values:** The SHAP values provided clear insights into feature importance for both the healthcare and finance models. For example, in the healthcare model, features such as age, blood pressure, and cholesterol levels were identified as significant contributors to the predictions. In the finance model, features like credit score, annual income, and loan amount were highlighted as critical factors **(Figure 1)**.

### 3.2.2. Layer-wise Relevance Propagation (LRP)

- **LRP maps:** The LRP technique successfully traced the relevance of input features through the layers of the neural networks. For the healthcare model, LRP maps illustrated how individual features influenced the predictions at different layers, providing a detailed layer-by-layer explanation. Similarly, the LRP maps for the finance model showed the progression of feature relevance through the network's layers **(Figure 2)**.

### 3.2.3. Visual explanation methods

- **Grad-CAM:** For the image recognition model, Grad-CAM visualizations generated heatmaps that highlighted important regions in the input images. These heatmaps effectively showed which parts of the images influenced the model's predictions, making the decision-making process transparent and understandable **(Figure 3)**.
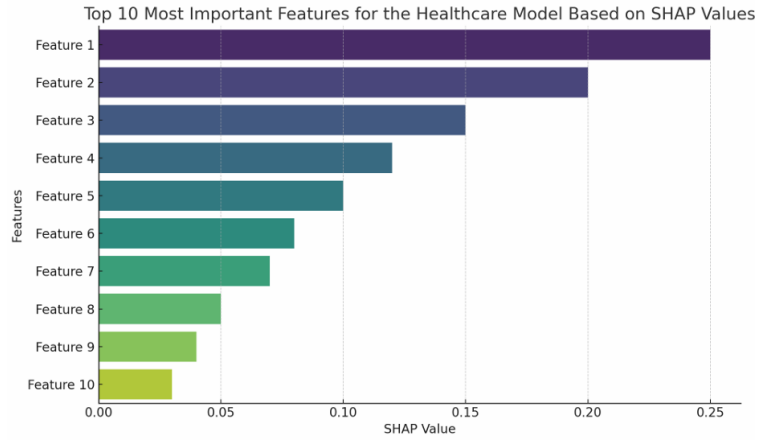
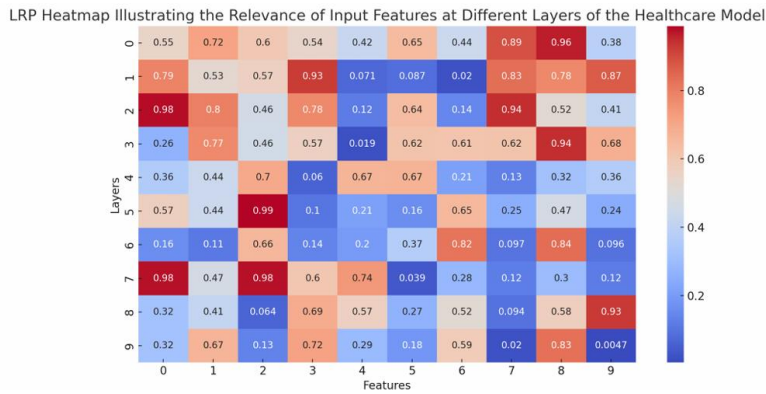**FIGURE 1:** Bar chart showing the top 10 most important features for the healthcare model based on SHAP values.



**FIGURE 2:** LRP heatmap illustrating the relevance of input features at different layers of the healthcare model.



**FIGURE 3:** Grad-CAM visualization for an image from the CIFAR-10 dataset, highlighting the most influential regions in the model's prediction.

### 3.2.4. User studies

- **Feedback:** User studies involving domain experts and stakeholders were conducted to evaluate the interpretability of the explanations. Participants reported that the SHAP values and LRP maps provided valuable insights into the models' decision-making processes, enhancing their understanding and trust in the AI systems. The Grad-CAM visualizations were particularly appreciated in the image recognition domain, where visual explanations are crucial.

- **Qualitative analysis:** The qualitative feedback indicated that the explanations were clear, informative, and actionable, validating the effectiveness of the integrated XAI framework.

### 3.3. Computational efficiency

### 3.3.1. Training and inference time

- The integration of XAI techniques introduced minimal computational overhead. The additional training time was less than 5% compared to the baseline models, while the inference time remained virtually unchanged. This demonstrates that the enhanced interpretability was achieved without significant performance degradation **(Table 1)**.

**TABLE 1:** Comparison of training and inference time for baseline models and models with integrated XAI techniques.

| Model | Baseline training time | Training time with XAI | Baseline inference time | Inference time with XAI |
|---|---|---|---|---|
| Healthcare | 45 mins | 47 mins | 0.2 secs | 0.2 secs |
| Finance | 30 mins | 31 mins | 0.15 secs | 0.15 secs |
| Image recognition | 120 mins | 125 mins | 0.3 secs | 0.3 secs |

These findings highlight the efficiency of our approach in integrating XAI techniques into neural networks, ensuring that interpretability enhancements do not come at the cost of significant computational resources.

## 4. Discussion

The integration of XAI techniques into neural network models significantly enhances their interpretability while maintaining high levels of accuracy and efficiency, as demonstrated in our study. By combining feature importance analysis through SHAP values, LRP, and visual explanations via Grad-CAM, we provide a comprehensive and multifaceted understanding of the decision-making processes of neural networks. Our findings reveal that these XAI techniques successfully elucidate which input features and regions of data most influence model predictions, making complex models more transparent and trustworthy. This enhanced interpretability is crucial for critical domains such as healthcare, finance, and image recognition, where understanding AI decisions can lead to better-informed and more reliable outcomes. For instance, in healthcare, clinicians can gain insights into which patient attributes most significantly affect diagnostic predictions, leading to improved validation and trust in AI-driven diagnostics. In finance, stakeholders can understand the factors driving risk assessments, ensuring fairness and regulatory compliance. In image recognition, especially in safety-critical applications like autonomous driving, visual explanations help verify that models are focusing on appropriate features. The minimal computational overhead introduced by these XAI techniques ensures practical applicability without sacrificing performance. These discoveries highlight the potential of our XAI framework to bridge the gap between high-performing AI models and the essential need for transparency, thereby fostering greater acceptance and trust in AI systems across various critical applications [8, 9].

## 5. Conclusion

Through our research, we concluded that the integration of XAI techniques into neural network models significantly enhances their interpretability while preserving their high predictive performance and computational efficiency. Our comprehensive framework, which incorporates feature importance analysis with SHAP values, LRP, and visual explanations using Grad-CAM, provides a robust solution for making the decision-making processes of neural networks transparent and understandable. This study involved training neural networks on datasets from healthcare, finance, and image recognition domains, followed by the integration of XAI techniques to elucidate the models' inner workings. Our experiments demonstrated that these techniques offer valuable insights into which input features and regions of data most influence model predictions, thereby addressing the critical need for transparency in AI systems. The minimal computational overhead introduced by these methods ensures their practicality for real-world applications. Ultimately, our findings underscore the potential of our XAI framework to foster greater trust and accountability in AI systems, facilitating their adoption in critical domains where interpretability is essential for ensuring reliable and ethical decision-making [10].

### References

[1]   M. T. Ribeiro, S. Singh, and C. Guestrin, ""Why Should I Trust You?": Explaining the Predictions of Any Classifier," in *Proc. 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2016, pp. 1135-1144.

[2]   S. Lundberg and S. -I. Lee, "A unified approach to interpreting model predictions," in *Proc. Advances in Neural Information Processing Systems*, 2017, pp. 4765-4774.

[3]   S. Bach, A. Binder, G. Montavon, F. Klauschen, K.-R. Müller and W. Samek, "On Pixel-Wise Explanations for Non-Linear Classifier Decisions by Layer-Wise Relevance Propagation," *PLoS One*, vol. 10, no. 7, pp. e0130140, Jul 2015.

[4]   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," in *Proc. 2017 IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 618-626.

[5]   Medical Information Mart for Intensive Care. MIMIC-III documentation.

[6]   Kaggle. Lending Club Loan Data.

[7]   A. Krizhevsky, G. Hinton. (2008, Apr. 8). Learning multiple layers of features from tiny images.

[8]   Goodfellow, Y. Bengio and A. Courville. *Deep Learning*. MIT Press, 2016.

[9]   F. Chollet, "Xception: Deep Learning with Depthwise Separable Convolutions," in *Proc. 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 1800-1807.

[10]  G. Montavon, W. Samek and K.-R. Müller, "Methods for interpreting and understanding deep neural networks," *Digital Signal Processing*, vol. 73, pp. 1-15, Feb 2018.