# Explainable AI Framework for Alzheimer's Diagnosis Using Convolutional Neural Networks

**Dhekra Mansouri [1,2], Amira Echtioui [1], Rafik Khemakhem [1,2], Ahmed Ben Hamida [3]**

[1]*Advanced Technologies for Medicine and Signals Laboratory 'ATMS', National Engineering School of Sfax (ENIS), Sfax University, Sfax, Tunisia*
[2]*Higher Institute of Management of Gabès, Gabès University, Gabès, Tunisia.*
[3]*Department IS, College of Computer Science, King Khaled University 'KKU', Abha, Kingdom of Saudi Arabia*

*Email: dhekramansouri33@gmail.com*

*Abstract*—**Alzheimer's disease (AD) stands as a form of dementia characterized by the gradual degeneration of brain cells, resulting in compromised memory, cognitive functions, and the loss of fundamental skills, ultimately leading to fatality. While a definitive cure for AD remains elusive, early detection plays a pivotal role in managing its progression and enhancing the quality of life for patients. This study delves into the realm of Alzheimer's disease identification through the application of various Neural Network models employing classification techniques. Leveraging a contemporary hybrid dataset, the investigation yielded four distinct classifications. Moreover, the study delved into elucidating the specific brain regions contributing to each classification using the Grad-CAM (Gradient-weighted Class Activation Mappings) based XAI (eXplainable Artificial Intelligence) framework applied to patients' MRI images. A comprehensive assessment was conducted on pre-trained deep neural networks, particularly focusing on Convolutional Neural Network (CNN) models trained exclusively on authentic MRIs and a combination of authentic and synthetic MRIs. The efficacy of deep learning in disease detection was exemplified, with the CNN model trained on both real and synthetic MRIs outperforming its counterpart trained solely on real MRIs. The former achieved an impressive accuracy of 97.50%, a Balanced Accuracy Score (BA) of 98.58%, and a Matthew's Correlation Coefficient (MCC) of 95.95%. In contrast, the model trained exclusively on real MRIs exhibited an accuracy of 88.98%, a BA of 94.01%, and an MCC of 83.67%.**

*Keywords: Alzheimer's disease, Grad-CAM, CNN, XAI framework, Real MRIs, Synthetic MRIs.*

## I. INTRODUCTION

Alzheimer's disease (AD) is a form of dementia primarily marked by cognitive decline and the inability to perform daily tasks [1,2]. This degenerative process affects brain cells, leading to short-term memory loss and behavioral issues, typically impacting individuals aged 60 and above, though early diagnosis has been noted in those aged 40 to 50 years. Approximately 5 million people in the United States currently live with Alzheimer's disease, with projections suggesting a tripling of this number by 2050 [3]. Despite extensive research, a cure remains elusive, highlighting the critical importance of early detection and intervention to manage its progression and improve patient outcomes.

Advanced technologies such as deep learning and artificial intelligence (AI) offer promise in enhancing Alzheimer's disease detection and understanding its mechanisms. Deep neural networks, including convolutional neural networks (CNN) and other variants, can analyze brain structure changes to track disease progression and provide reliable diagnostic outcomes [4]. Training these networks requires substantial datasets to achieve optimal performance.

This paper introduces an innovative approach to Alzheimer's disease diagnosis by integrating CNN and XAI methodologies. The model achieves an impressive accuracy of 93.82% through the utilization of pretrained CNNs on both real and synthetic MRIs. The incorporation of XAI techniques, including Grad-CAM, enhances interpretability, offering valuable visual insights into crucial neural regions during the diagnostic process.

Dedicated to advancing Alzheimer's disease diagnosis, the paper focuses on the development of a cutting-edge classification system. Employing a sophisticated CNN model, the approach capitalizes on the strengths of CNNs to create a robust and accurate model for classifying Alzheimer's cases.

In addition to classification accuracy, the paper aims to improve the interpretability of the developed model. This is achieved by integrating the XAI method known as Grad-CAM. Through this technique, the decision-making process of the CNN model becomes transparent, providing insights into the features and regions crucial for Alzheimer's diagnosis. This pursuit of interpretability is crucial not only for comprehending the model's inner workings but also for building trust and acceptance among clinicians and researchers in the field.

The organization of this paper unfolds in the following manner: In Section II, a comprehensive review of relevant literature in the domains of AI-driven Alzheimer's disease diagnosis, deep and transfer learning, and eXplainable Artificial Intelligence (XAI) is presented. Section III will delve into the dataset and the proposed methodology. Following that, in Section IV, the analysis of our results will be presented. The study concludes in the final section, where we summarize our findings and insights.

## II. RELATED WORK

Detecting AD in its early stages with enhanced accuracy is paramount for enhancing patients' quality of life and possibly impeding or arresting the disease's advancement

[5]. The interest in the potential of Deep Learning (DL) for human disease identification has grown significantly. In recent years, there has been a notable upswing in the evolution of DL techniques designed for diagnosing AD, gaining popularity as diagnostic tools among medical professionals. These DL approaches are preferred over more traditional forms of Machine Learning (ML) [6].

Zhang et al. [7] proposed an enhanced CNN approach for Alzheimer's disease classification by leveraging multimodal brain imaging data. Concurrently, Dey et al. [8] developed a hybrid Deep Learning (DL) framework specifically tailored for early-stage diagnosis and detection of Alzheimer's disease.

In a separate study, Lee et al. [9] introduced a deep CNN data permutation strategy for Alzheimer's disease (AD) classification using structural Magnetic Resonance Imaging (sMRI). They advocated for slice selection to capitalize on the benefits of AlexNet. Experimental results showed that their proposed data permutation scheme significantly improved overall classification accuracies for AD. Notably, classification accuracies for both binary and ternary classifications on ADNI datasets reached 98.74% and 98.06%, respectively.

Pradhan et al. [10] introduced a methodology for detecting distinct stages of AD. The approach utilizes the VGG19 and DenseNet169 architectures for classification, employing a dataset sourced from the Kaggle open online data repository. The dataset comprises 6000 images categorized as mild, moderate, very mild, and non-demented AD. Feature selection involves an 80% allocation for the learning phase and a 20% allocation for the testing phase. VGG19, characterized by approximately 10–16 convolutional neural network layers, outperforms DenseNet in image classification, achieving an accuracy of 94%.

Until recently, researchers have been exploring the application of eXplainable Artificial Intelligence (XAI) for detecting and classifying Alzheimer's disease. This innovative approach aims to enhance clarity and interpretability in artificial intelligence models, particularly CNN, used in medical contexts.

Integrating XAI into Alzheimer's disease detection represents a significant advancement, allowing for a deeper understanding of decision factors within the models and instilling confidence in healthcare professionals. It enables more personalized and collaborative applications, providing explicit insights for individual cases and facilitating thorough validation of AI model diagnoses. This paper proposes integrating XAI with a CNN model for Alzheimer's disease classification, enhancing transparency and interpretability.

This method aids in understanding the factors influencing model decisions, building trust among healthcare professionals, and promoting ethical AI practices. XAI provides personalized insights, fosters collaboration, and ensures alignment between AI outputs and clinical interpretations, enhancing the CNN model's value in real-world medical settings.

## III. METHODOLOGY

### A. Dataset description

The Alzheimer's disease dataset from Kaggle contains 6,400 axial MRI brain scans categorized into four classes: Mild Impairment, Moderate Impairment, No Impairment, and Very Mild Impairment. The dataset lacks patient-specific data and was chosen due to its significant class imbalance, particularly in the minority classes. Specifically, there are 896 MRIs for Mild Impairment, 64 MRIs for Moderate Impairment, 3,200 MRIs for No Impairment, and 2,240 MRIs for Very Mild Impairment. Each image was resized to 128 x 128 pixels for model training. Table 1 illustrates the distribution of images within the dataset and the corresponding split.

**Table 1 :** Splitting the target dataset

|  | Training | Testing |
|---|---|---|
| **Mild Impairment** | 717 | 179 |
| **Moderate Impairment** | 52 | 12 |
| **No Impairment** | 2560 | 640 |
| **Very Mild Impairment** | 1792 | 448 |

### B. Proposed methodology

This research work implemented a four-step pipeline to conduct image classification, focusing on MRI images sourced from the Alzheimer's disease dataset. A sequence of data preparation steps, including resizing and rescaling, was employed to ensure uniformity in the data and enhance the effectiveness of model training. Following this, the dataset was partitioned into training and testing subsets. CNN model was subsequently applied for the classification tasks, utilizing features extracted from the preprocessed images. Finally, the study presented the classification outcomes, illustrating the categorization of real images based on their respective classes.
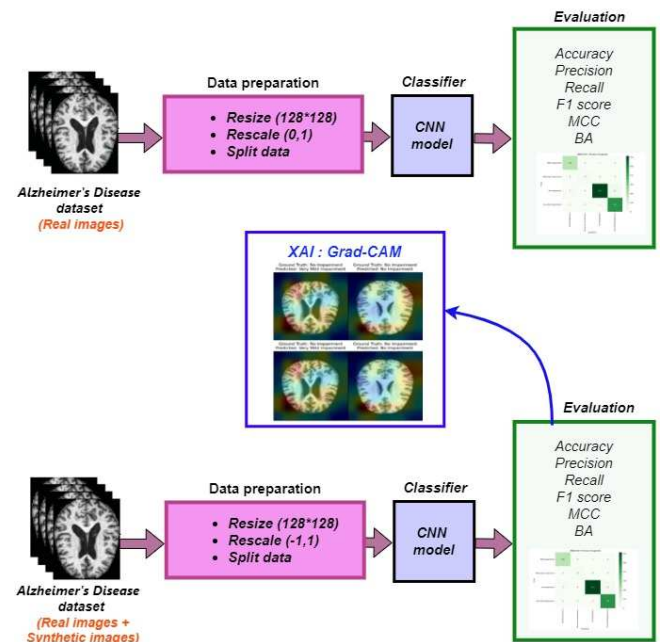


**Fig. 1.** Overview of the proposed method

Furthermore, this study implemented a composed four step pipeline. Initially, real MRI images were combined with synthetic images generated using WGAN-GP (Wasserstein Generative Adversarial Network with Gradient Penalty). Subsequently, a series of data preparation steps, including resizing and rescaling, were applied to ensure data uniformity and improve model training effectiveness. Next, a convolutional neural network (CNN) model was utilized for classification tasks. Finally, the XAI technique Grad-CAM (Gradient-weighted Class Activation Mapping) was employed to gain insights into the CNN model's decision-making process. Grad-CAM generated heatmaps highlighting regions of input images that significantly influenced the model's predictions, thereby enhancing understanding of the features the model prioritized during classification and providing valuable insights into its decision-making rationale.

- ***Data preparation***

Data preparation is a critical phase in the pre-processing of medical images, especially MRI images, where various factors such as imaging protocols, brightness variations, reflections, low contrast, and differences across imaging sources can introduce inherent noise. The primary objective of data preparation is to enhance MRI images through preprocessing, ultimately leading to improved image quality. This improved quality facilitates high-performance analysis and the extraction of relevant information from the images.

In the context of this study, image resizing was performed as a part of data preparation. The entire set of images was resized to $128 \times 128$ pixels, covering training, testing, and validation subsets. This standardized size was utilized for transfer learning models, with RGB values employed for ensemble models as well.

Another crucial aspect of data preparation is normalization. Normalization involves adjusting features to have a consistent scale, which is essential for enhancing the performance and stability of the model during training. In this study, normalization was applied by mapping the image data into a predefined range, typically [0, 1] or [-1, 1]. This normalization ensured uniformity across all images. Additionally, the same techniques of transfer learning and ensembling were applied to the normalized data, as outlined in Equation (1), and executed in the following manner:

$$img = \frac{1}{255.0} \qquad (1)$$

We divided the dataset into training and testing sets, allocating 80% for training and 20% for testing.

- ***CNN model***

The CNN architecture consists of a sequence of layers, starting with 2 Conv2D layers followed by max pooling 2D, then another 2 Conv2D layers followed by max pooling 2D. The pattern repeats with batch normalization after each set of 2 Conv2D layers and max pooling 2D. This sequence is reiterated four times in total. Following these convolutional layers, there is a Flatten layer, a dropout layer, and a series of densely connected layers with Rectified Linear Unit (ReLU) activation, interspersed with batch normalization and dropout.

The final layer employs a softmax activation function for classification. Figure 2 depicts the architecture of the CNN model.
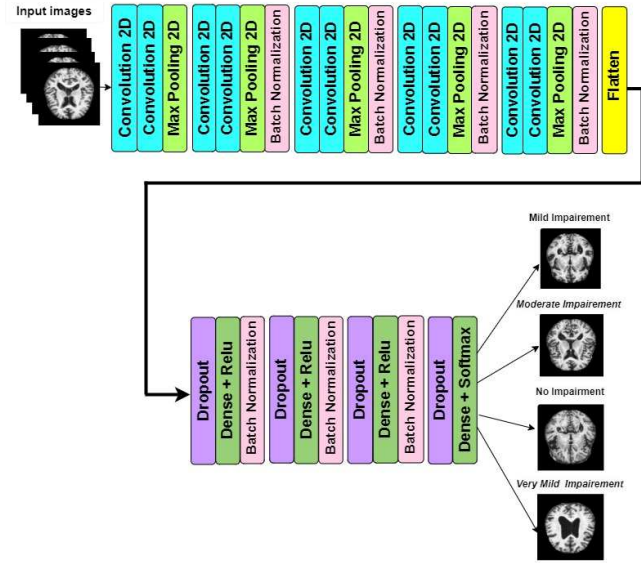


**Fig. 2.** Architecture of the CNN model

- ***Synthetic Image***

Synthetic images, generated through techniques like Generative Adversarial Networks (GANs), are artificially created to represent distinct features related to Alzheimer's disease, such as brain changes visible in MRI scans. These images serve to expand dataset sizes, increase data diversity, and enhance the capacity of machine learning models to recognize and interpret significant patterns relevant to Alzheimer's disease.

- ***Training WGAN-GP***

WGAN-GP was trained separately for each minority class, focusing on unique characteristics. The model avoided mixing features from different classes in synthetic images. "No Impairment" class was not trained, aiming to oversample minority classes. The generator produced 128 x 128 MRIs using a Latent Random Vector of dimension 256, ensuring effective feature capture. Both generator and critic used Adam optimizer (lr = 0.0001, ß1 = 0, ß2 = 0.9) with respective loss functions.

Figure 3 illustrates the algorithm for training WGAN-GP.



**Fig. 3.** Algorithm for training WGAN-GP

## IV. RESULTS AND DISCUSSION

### A. Evaluation Metrics

Evaluation metrics in machine learning, particularly for Convolutional Neural Networks (CNNs), provide quantitative assessments of performance in tasks like image classification. Key metrics include accuracy, precision, recall, F1 score, Matthew's Correlation Coefficient (BCC), Balanced Accuracy (BA), and the confusion matrix. The confusion matrix, with cells representing True Positives, False Positives, False Negatives, and True Negatives, aids in computing these metrics, offering insights into the model's classification effectiveness.

$$Accuracy = \frac{(TP + TN)}{(TP + TN + FP + FN)} \quad (2)$$

$$Precision = \frac{TP}{FP + TP} \quad (3)$$

$$Recall = \frac{TP}{FN + TP} \quad (4)$$

$$\text{F1 score} = 2 \times \frac{precision \times recall}{precision + recall} \quad (5)$$

$$\text{MCC} = \frac{(TP \times TN) - (FP \times FN)}{\sqrt{(TP+FP)(TP+FN)(TN+FP)(TN+FN)}} \quad (6)$$

$$BA = \frac{TPR + TNR}{2} \quad (7)$$

### B. Experimental results

The CNN model for Alzheimer's MRI classification was evaluated using key metrics. Trained on real and a combination of real and synthetic images, confusion matrices revealed performance insights. According to Figures 4 and 5, the CNN trained solely on real MRIs had 141 misclassifications, notably 132 False Negatives for the "No Impairment" class. When trained on both real and synthetic MRIs, misclassifications reduced to 32, showing significant improvement, especially in reducing False Negatives across various impairment classes. Incorporating synthetic images enhanced overall model performance.
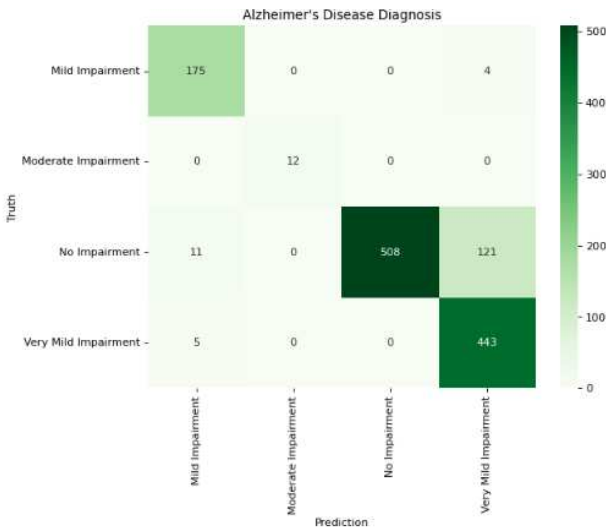
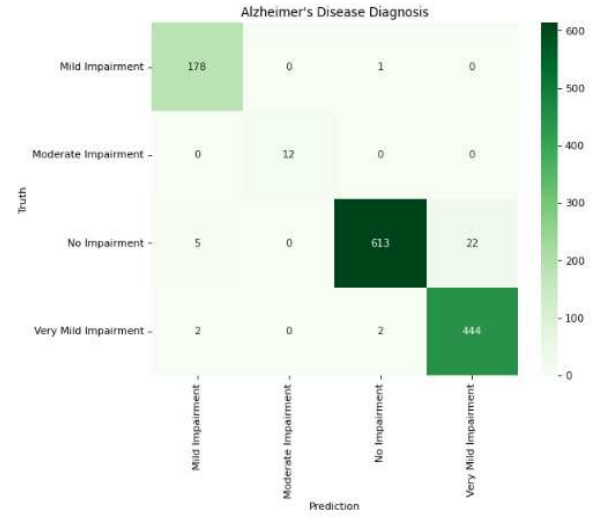**Fig. 4.** Confusion Matrix for CNN trained only on Real MRIs

**Fig. 5.** Confusion Matrix for CNN trained on Real + Synthetic MRIs

The classification reports provided in Tables 2 and 3 offer a comprehensive evaluation of the CNN model's performance in Alzheimer's MRI classification, presenting metrics such as Accuracy, Precision, Recall, and F1 score for each impairment class.

**Table 2:** Classification report obtained by the CNN model (only real images) in %.

|  | Accuracy | Precision | Recall | F1 score |
|---|---|---|---|---|
| **Mild Impairment** | - | 92 | 98 | 95 |
| **Moderate Impairment** | - | 100 | 100 | 100 |
| **No Impairment** | - | 100 | 79 | 89 |
| **Very Mild Impairment** | - | 78 | 99 | 87 |
| **Average** | **88.98** | **92.50** | **94.00** | **92.75** |

**Table 3:** Classification report obtained by the CNN model (real images and synthetic images) in %.

|  | Accuracy | Precision | Recall | F1 score |
|---|---|---|---|---|
| **Mild Impairment** | - | 96 | 99 | 98 |
| **Moderate Impairment** | - | 100 | 100 | 100 |
| **No Impairment** | - | 100 | 96 | 98 |
| **Very Mild Impairment** | - | 95 | 99 | 97 |
| **Average** | **97.50** | **97.75** | **98.50** | **98.25** |

For the CNN model trained solely on real images (Table 2), the overall accuracy stands at 88.98%, with precision values ranging from 78% to 100%. This indicates a low rate of false positives in the predictions. The recall values, ranging from 79% to 100%, illustrate the model's ability to capture instances of each class. The F1 score, averaging at

92.75%, demonstrates a balanced trade-off between precision and recall.

In contrast, the CNN model trained on a combination of real and synthetic images (Table 3) exhibits notably improved performance. The overall accuracy increases to 97.50%, and precision values range from 95% to 100%, signifying enhanced accuracy in predictions. High recall values (ranging from 96% to 100%) indicate the model's effective capturing of instances for each class. The average F1 score reaches 98.25%, highlighting an enhanced balance between precision and recall.

Comparing the two models, the CNN trained with a combination of real and synthetic images outperforms its counterpart in all metrics. The addition of synthetic images positively impacts the model's ability to generalize and classify instances more accurately. Notably, the performance improvement is evident in the "No Impairment" class, with higher precision, recall, and F1 score, suggesting a more accurate identification of this class when synthetic images are incorporated.

This study conducted a comparative analysis of the CNN's performance when trained exclusively on a Train Set comprising real MRIs, as opposed to a Combined Train Set incorporating both real and Synthetic MRIs. The outcomes for overall performance, evaluated through BA and MCC, are presented in Table 4.

**Table 4 :** BA and MCC obtained by CNN model

|  | BA | MCC |
|---|---|---|
| **CNN (Real MRIs)** | 94.01 % | 83.67 % |
| **CNN (Real + Synthetic MRIs)** | 98.58 % | 95.95 % |

The results reveal a significant enhancement in overall performance, with a noteworthy 12% increase in BA and an approximately 2.74% increase in MCC when using the Combined Train Set with real and synthetic MRIs. This underscores a substantial improvement in the model's effectiveness.

The table 5 compares the accuracy of five different neural network models on a specific task or dataset. It shows that a CNN model combining real and synthetic MRIs achieved the highest accuracy at 97.50%, followed by AlexNet at 94.53%. Other models like ResNet50, basic CNN, and VGG16 had lower accuracy rates ranging from 58.07% to 70.30%. This highlights the importance of incorporating diverse data sources for improved model performance.

**Table 5 :** Comparison with the state of the art

|  | Accuracy |
|---|---|
| ResNet50 [11] | 58.07% |
| AlexNet [11] | 94.53% |
| CNN [12] | 67.50% |
| VGG16 [12] | 70.30% |
| **CNN (Real + Synthetic MRIs)** | **97.50%** |

## C. Grad-CAM

Having effectively addressed class imbalance and ensuring the credibility of the final outcome, the tailored CNN, which incorporates both real and synthetic images, demonstrated remarkable performance with a BA of 98.58% and MCC of 95.95%. Subsequently, Grad-CAM were employed using this model, involving the extraction of gradients from the final convolutional layer for each image and overlaying them as heatmaps onto the corresponding images. This visualization technique elucidated the classifier's specific areas of focus in making classification decisions.

In Figure 6, a comprehensive illustration unfolds, showcasing the hierarchical regions contributing to classification. A color-bar is employed, where "red" denotes the most intensely focused region by the classifier, followed by "orange," "yellow," "green," and finally "blue," indicating the regions with the least focus. Remarkably, these images achieved a perfect classification accuracy of 100%, as evidenced by the ground truth and predicted labels superimposed on each image.
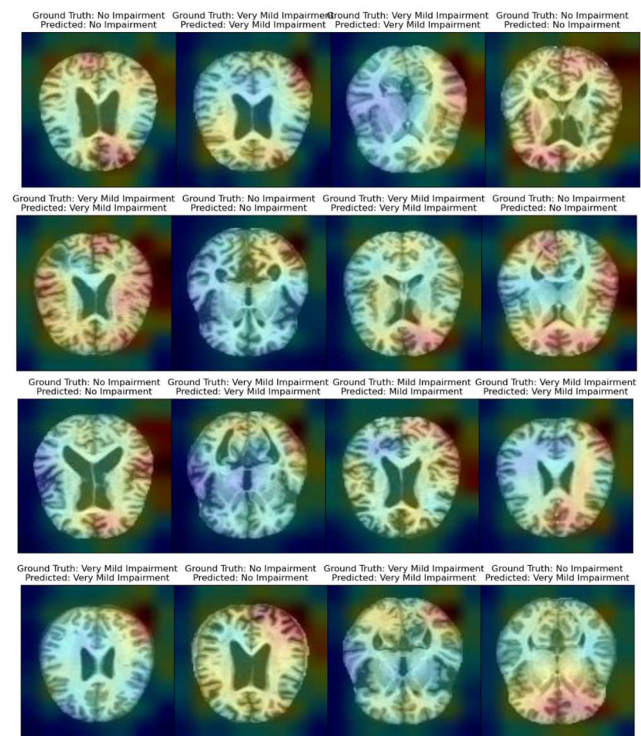


**Fig. 6.** Visual explanation obtained by Grad-CAM

The utilization of Grad-CAM (interpretable AI) presents a valuable opportunity for researchers and medical practitioners to detect Alzheimer's disease markers across different stages effectively. Furthermore, in envisaging a future where diagnosis of Alzheimer's disease is automated through an application, clinicians could harness the Grad-CAM feature within the application. This would enable them to interpret and validate classification decisions based on individual MRI scans, thereby enhancing diagnostic accuracy and efficiency.

## V. Conclusion

This paper explores the performance of CNN models in diagnosing Alzheimer's disease, comparing models trained exclusively on real MRI images with those incorporating a combination of real and synthetic MRI images. The results highlight a notable improvement in CNN model performance when synthetic images are integrated with real data. Specifically, the model trained on both real and synthetic images exhibited superior metrics, including accuracy (97.50%), recall (98.50%), F1 score (98.25%), Balanced Accuracy (98.58%), and Matthew's Correlation Coefficient (95.95%) compared to the model trained solely on real images. This underscores the potential of synthetic data augmentation techniques to enhance the effectiveness of machine learning models in medical image analysis, particularly for Alzheimer's disease diagnosis. Furthermore, the application of the explainable AI technique Grad-CAM significantly bolstered the accuracy of our model. These methodologies not only enhanced performance but also provided clinicians and researchers with insightful visualizations of neural regions. This, in turn, influences Alzheimer's disease diagnostic decisions and offers profound insights into the decision-making process of the model.

## References

[1] I. Beheshti, H. Demirelb, H. Matsudaaf, ''Classification of Alzheimer's disease and prediction of mild cognitive impairment-to-Alzheimer's conversion from structural magnetic resource imaging using feature ranking and a genetic algorithm,'' Comput. Biol. Med. 2017, 83, 109–119.

[2] S. Klöppel, C.M. Stonnington, C. Chu, B. Draganski, R.I. Scahill, J.D. Rohrer, N.C. Fox, C.R., Jr. Jack, J. Ashburner, R.S. Frackowiak, ''Automatic classification of MR scans in Alzheimer's disease,'' Brain 2008, 131, 681–689.

[3] I. Beheshti, H. Demirel, ''Feature-ranking-based Alzheimer's disease classification from structural MRI,'' Magn. Reson. Imaging 2016, 34, 252–263.

[4] N. Zeng, H. Li, Y. Peng, ''A new deep belief network-based multi-task learning for diagnosis of Alzheimer's disease,'' Neural Comput Applic, 2021.

[5] L. J. Herrera, I. Rojas, H. Pomares, A. Guillén, O. Valenzuela, and O. Baños, ''Classification of MRI images for Alzheimer's disease detection,'' in Proc. Int. Conf. Social Comput., 2013, pp. 846–851.

[6] N. Yamanakkanavar, J. Y. Choi, and B. Lee, ''MRI segmentation and classification of human brain using deep learning for diagnosis of Alzheimer's disease: A survey,'' Sensors, vol. 20, no. 11, p. 3243, 2020.

[7] L. Zhang, Y. Peng, J. Sun, ''Improved deep convolutional neural networks for Alzheimer's disease classification using multimodal brain imaging data,'' J Med Imaging 9(2):024501, 2022.

[8] S. Dey, S. Maity, S. Mondal, N. Dey, ''Early diagnosis of Alzheimer's disease using a hybrid deep learning framework,'' J Med Syst 45(11):135, 2021.

[9] B. Lee, W. Ellahi, J.Y. Choi, ''Using deep CNN with data permutation scheme for classification of alzheimer's disease in structural magnetic resonance imaging (SMRI),'' IEICE Trans Inf Syst, 2019, 102(7):1384–1395.

[10] A. Pradhan, J. Gige, M. Eliazer, ''Detection of Alzheimer's Disease (AD) in MRI ımages using deep learning,'' Int J Eng Res. 2021;10(3):580–5.

[11] M.H. Al-Adhaileh, ''Diagnosis and classification of Alzheimer's disease by using a convolution neural network algorithm,'' Soft Comput, 2022, 26:7751–7762.

[12] E. Mggdadi, A. Al-Aiad, M.S. Al-Ayyad, A. Darabseh, ''Prediction Alzheimer's disease from MRI images using deep learning,'' 12th International Conference on Information and Communication Systems, ICICS 2021, pp. 120–125, 2021.