

# MRI-Based Early Detection and Classification of Alzheimer's Disease with Grad-CAM Explainability

Pratham Agarwal\*, Pranshu Saini\*, Debjani Ghosh\*, Vimal Kumar\*, Punit Gupta†

\*School of Computer Science Engineering and Technology,

Bennett University, Greater Noida, India

Emails: e22cseu1308@bennett.edu.in, e22cseu1300@bennett.edu.in, debjani.ghosh@bennett.edu.in, vimal.kumar@bennett.edu.in

†National College of Ireland, Dublin, Ireland

Email: pgupta@staff.ncirl.ie

**Abstract**—Alzheimer's disease is a progressive neurodegenerative disorder that profoundly affects memory and cognitive abilities. For prompt intervention and care planning, an accurate and early-stage classification is essential. From conventional machine learning models like K-Nearest Neighbors (KNN) and Ensemble methods, to sophisticated deep learning techniques like a custom Dual-Branch Convolutional Neural Network (CNN) and a fine-tuned CNN model, this study offers a thorough assessment of several classification approaches. Among the conventional models, the CNN with Fine-Tuning performed the best in classification with F1-score, recall and accuracy of 96%, 94% and 95% respectively, followed by the KNN. To evaluate alignment with medically relevant brain regions and interpret model predictions, Grad-CAM heatmap visualizations were utilized. The results highlight how crucial it is for clinical decision support systems to diagnose Alzheimer's disease to be both accurate and interpretable.

**Keywords**- Convolutional Neural Network (CNN), Alzheimer's Disease, MRI, Deep Learning, Explainability, Grad-CAM, KNN, Random Forest.

## 1. INTRODUCTION

**A**LZHEIMER'S disease is a progressive neurodegenerative condition that predominantly impairs memory, cognitive functions, and behavior, severely impacting the quality of life of those affected. Improving patient outcomes and limiting the disease's course depend on an early and precise diagnosis. Magnetic resonance imaging (MRI) is a crucial non-invasive technique for identifying structural brain changes linked to AD. In recent years, there has been a lot of interest in using deep learning and machine learning approaches to categorize various phases of Alzheimer's disease from MRI data. In order to better distinguish between AD stages and enhance diagnostic performance, this study examines and contrasts a number of models, including custom convolutional neural networks (CNNs) and conventional classifiers.

Despite advances in deep learning, building models that are both highly accurate and interpretable remains challenging. Many existing methods lack either precision or transparency, limiting their usefulness for clinicians, especially early-career professionals. Moreover, identifying specific brain regions influencing predictions is often neglected. The classification of Alzheimer's disease into four stages, like non-demented,

very mild, mild and moderate demented, adds further complexity due to subtle anatomical differences. This research is motivated by the need for accurate, interpretable models to support early diagnosis using brain MRI.

This study presents a comparative analysis of several models—including CNN with Fine-Tuning, K-Nearest Neighbors (KNN), Ensemble Learning and a Dual-Branch CNN approach—for classifying Alzheimer's stages using MRI. A custom Dual-Branch CNN was proposed to enhance spatial feature extraction, though its attention maps remained inconsistent. Grad-CAM visualizations were used to interpret model focus, revealing that the Fine-Tuned CNN not only achieved the highest accuracy (95%) but also correctly attended to clinically important brain regions. KNN was found to be the most effective among traditional models, ranking second overall in accuracy.

The subsequent sections of the paper are arranged as follows: Section 2 reviews related work on Alzheimer's classification. Section 3 details the methodology, including the dataset and preprocessing. Section 4 presents results, evaluation metrics, comparative analysis, and Grad-CAM visualizations. Section 5 discusses key findings, highlights limitations, and concludes the study.

## 2. RELATED WORK

Several studies have explored traditional machine learning approaches for Alzheimer's classification using MRI. In [1], an automated machine learning technique was applied to 1,167 MRI scans across four categories, achieving 75% accuracy via 10-fold cross-validation. The work in [2] employed kernel PCA for gray matter MRI feature extraction and classified subjects using AdaBoost and SVM, achieving 84% accuracy—outperforming standard PCA-based methods. Similarly, [3] integrated rs-fMRI and sMRI features with Random Subspace Feature Selection and machine learning classifiers to differentiate MCI converters from non-converters, reaching up to 89.80% accuracy.

On the deep learning front, various convolutional neural network (CNN) architectures have shown promising results. The authors of [4] proposed a Siamese CNN based on VGG-16

and achieved 99.05% accuracy on the OASIS dataset. In [5], a hybrid model combining 2D/3D CNNs and RNNs, enhanced via transfer learning, yielded 96.8% accuracy with a voxel-based 3D CNN. Transfer learning was further explored in [6], where MobileNet-based CNN classifiers achieved 96.6% accuracy on a five-class dataset. Ramzan et al. [7] used a fine-tuned ResNet-18 model on ADNI MRI images, achieving 97.88% accuracy. Parmar et al. [8] applied a 3D CNN on the same dataset, attaining 93.00% accuracy by leveraging volumetric features. Akter et al. [9] employed Inception V3 with fine-tuning and achieved 98.68% accuracy on the ADNI dataset. Chahd M. Chabib [10] introduced a Curvelet-based CNN model, DeepCurvMRI, reaching 98.62% accuracy using LOGO cross-validation. Öznur Özaltın [11] benchmarked five pre-trained models on a Kaggle dataset, where DenseNet-201 performed best at 82.11%, and a ResNet50-FTNCA-KNN pipeline achieved perfect accuracy after feature reduction. Lastly, in [12], a 3D CNN was trained on both synthetic and real MRI scans, achieving an AUC of 0.850 on real test data, demonstrating the potential of diffusion model-generated data for training.

Unlike many prior studies limited to binary classification or lacking interpretability, our work addresses multi-class Alzheimer's classification using a diverse set of models, including a novel Dual-Branch CNN. We further enhance clinical relevance by employing Grad-CAM to highlight key brain regions, offering both high accuracy (95%) and greater transparency.

### 3. METHODOLOGY

In this work, we present a comprehensive assessment framework for Alzheimer's disease detection using brain MRI images, employing four classification methods, including traditional machine learning models (KNN, Ensemble) and deep learning architectures (a combined CNN and a fine-tuned VGG16). To enhance model interpretability, we apply Grad-CAM visualizations on both the dual-branch CNN and the fine-tuned CNN, highlighting the key MRI regions that most influenced prediction outcomes.

#### 3.1. Dataset

TABLE I: Number of images in each training and validation dataset.

Class	Training Samples	Testing Samples	Total
Non-Demented	2,560	640	3,200
Very Mild Demented	1,792	448	2,240
Mild Demented	717	179	896
Moderate Demented	52	12	64
<b>Total</b>	<b>5,121</b>	<b>1,279</b>	<b>6,400</b>

There are 6,400 grayscale MRI images in the Alzheimer's MRI dataset used in this investigation [13]. The data set is categorized into four classes:

- **Non-Demented:** MRI scans, illustrated in Fig. 1, showing no apparent signs of Alzheimer's disease. This serves as the control group for model training.

- **Very Mild Demented:** Fig. 2 represents the earliest stage of Alzheimer's disease with subtle signs, often difficult to identify visually or computationally.
- **Mild Demented:** Fig. 3 comprises MRI scans displaying clear structural alterations typical of early-stage Alzheimer's disease, providing more distinguishable features that aid in accurate classification.
- **Moderate Demented:** Contains scans from individuals, shown in Fig. 4, in advanced Alzheimer's stages, contributing to class imbalance due to underrepresentation.

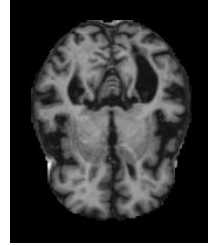


Fig. 1:  
MildDemented

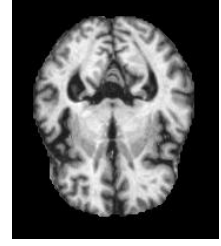


Fig. 2: Moderat-  
eDemented

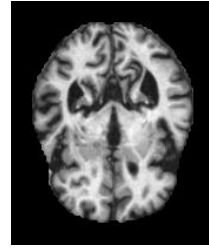


Fig. 3:  
NonDemented

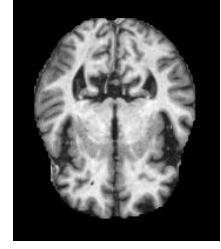


Fig. 4: VeryMild-  
Demented

*Data Splitting Strategy:* To promote fair evaluation and effective training, the dataset was divided into training and testing sets following an 80:20 stratified splitting approach presented in Table I. This stratification process ensured that the distribution of each class remained proportionally consistent across both subsets, thereby maintaining class balance and reducing potential sampling bias.

#### 3.2. Brain- MRI And Grad-CAM visualization

Understanding the decision-making mechanisms of deep learning models is crucial in medical imaging, especially for Alzheimer's classification. Gradient-weighted Class Activation Mapping (Grad-CAM) highlights the regions of the MRI image most influential to the model's prediction, offering visual explanations for CNN outputs.

Neurologists focus on regions like the entorhinal cortex, medial temporal lobe, and hippocampus, which show early signs of neurodegeneration [fig. 5]. These areas are key indicators of cognitive decline and memory loss.

#### 3.3. Preprocessing Techniques

To ensure effective training and evaluation of the Alzheimer's MRI dataset, a structured preprocessing pipeline

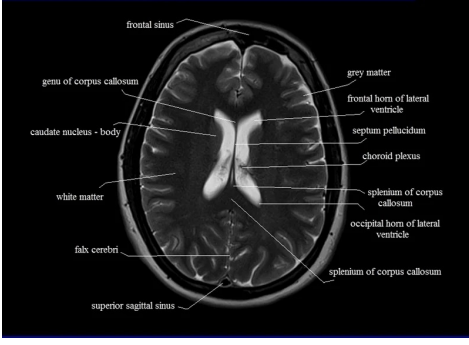


Fig. 5: Important parts in brain MRI

was implemented. Since several preprocessing steps were shared across different approaches, we first describe the common steps, followed by specific operations for each method.

Common Preprocessing Steps include image loading and resizing, normalization, label encoding and data augmentation. In image loading and resizing, all images were loaded and resized to a uniform shape of  $224 \times 224$  pixels. Next normalization will be performed where pixel values were scaled to the range  $[0, 1]$  by dividing by 255. Class labels were encoded into numeric form and, where required, converted to one-hot vectors in label encoding. Further, data augmentation (CNN-based approaches) has been applied using `ImageDataGenerator`, including random rotation, width/height shifts, zoom, horizontal flips, and rescaling.

Approaches 1–2 (KNN, and Ensemble) used grayscale images with handcrafted global features—one texture-based and one intensity-based—for classification. Approach 3 employed a fine-tuned VGG16 model, incorporating data shuffling to prevent order bias and training on augmented images with one-hot encoded labels. Approach 4 introduced a custom dual-branch CNN with spatial attention, utilizing a dual-input generator to enable parallel training of CNN1 and CNN2, allowing effective feature fusion.

### 3.4. Model Architectures and Approaches

This section presents the four methods investigated in this study. Each model was implemented using the respective preprocessing, feature extraction, and training techniques. A comprehensive evaluation was conducted to assess the performance of each model across classification metrics.

**3.4.1. K-Nearest Neighbors (KNN):** The KNN classifier with  $k = 3$  categorizes a test sample based on the majority label of its three nearest neighbors. Image data is transformed into 1D feature vectors and normalized to the range  $[0, 1]$  for use with distance metrics like Euclidean distance. Despite its simplicity, KNN serves as a baseline for evaluating raw image features.

**3.4.2. Ensemble Model:** Approach 2 integrates Random Forest and KNN using hard voting. Each model independently predicts the class label of an input image, and the final label is selected by majority voting. This fusion leverages both tree-based and distance-based perspectives, improving performance by reducing individual model weaknesses.

**3.4.3. CNN with Fine-Tuning:** This model uses **VGG16** (pre-trained on ImageNet) as a feature extractor, excluding its top layers. A series of layers are appended: **Global Average Pooling**, **Batch Normalization**, **Dropout (0.5)**, a **Dense layer (256 units, ReLU)** followed by another **Batch Normalization** and **Dropout (0.5)**. Finally, a **Dense(4, Softmax)** output layer performs classification. The model is optimized using Adam (initial learning rate 0.0001, fine-tuned to  $1e-5$ ) and Categorical Cross-Entropy loss, with callbacks like `EarlyStopping` and `ReduceLROnPlateau` for improved generalization.

**3.4.4. Dual-Branch CNN (CNN1, CNN2, and Combined Model):** This approach merges two parallel networks—**CNN1** (standard Conv2D-based) and **CNN2** (SeparableConv2D-based with spatial attention)—into a unified classification model. Each branch starts with input of shape  $(224, 224, 3)$  and performs multiple convolutions and pooling operations to extract deep features. CNN1 uses regular convolutional layers, while CNN2 adds **Spatial Attention** to emphasize informative regions. Both branches end with **Global Average Pooling**, **Batch Normalization**, **Dropout**, and a **Dense(4, Softmax)** layer.

The **combined model** takes dual inputs, processes them through CNN1 and CNN2, then concatenates the resulting features. This merged feature vector passes through **Batch Normalization**, followed by **Dense layers (256 and 128 units)** with **Swish activation** and **Dropout (0.4 and 0.2)**. The final classification is done via a **Dense(4, Softmax)** output layer. This fusion captures diverse representations, enhancing robustness and accuracy.

## 4. RESULTS AND DISCUSSION

The performance of the model was assessed using recall, F1-score, and accuracy on the test dataset [Table I]. These metrics were chosen to comprehensively evaluate the balance between precision and recall, as well as the overall classification accuracy.

TABLE II: K-Nearest Neighbors (KNN) Classification Report

Class	Precision	Recall	F1-Score	Support
Mild Demented	0.57	0.61	0.59	179
Non Demented	0.79	0.82	0.80	640
Very Mild Demented	0.72	0.67	0.70	448
Moderate Demented	0.78	0.58	0.67	12
Overall Accuracy			0.73	1279
Weighted average	0.74	0.73	0.73	1279
Macro average	0.71	0.67	0.69	1279

TABLE III: Ensemble Model Classification Report

Class	Precision	Recall	F1-Score	Support
Mild Demented	0.58	0.64	0.61	179
Non Demented	0.68	0.90	0.77	640
Very Mild Demented	0.80	0.40	0.53	448
Moderate Demented	0.78	0.58	0.67	12
Overall Accuracy			0.69	1279
Weighted average	0.71	0.69	0.67	1279
Macro average	0.71	0.63	0.65	1279

TABLE IV: CNN with Fine-Tuning Model Classification Report

Class	Precision	Recall	F1-Score	Support
Mild Demented	1.00	0.85	0.92	179
Non Demented	0.93	1.00	0.97	640
Very Mild Demented	0.97	0.93	0.95	448
Moderate Demented	1.00	1.00	1.00	12
<b>Overall Accuracy</b>			<b>0.95</b>	1279
Weighted Average	0.95	0.95	0.95	1279
Macro Average	0.97	0.94	0.96	1279

TABLE V: Combined CNN Model Classification Report

Class	Precision	Recall	F1-Score	Support
Mild Demented	0.61	0.31	0.41	179
Non Demented	0.64	0.82	0.72	640
Very Mild Demented	0.61	0.50	0.55	448
Moderate Demented	0.00	0.00	0.00	12
<b>Overall Accuracy</b>			<b>0.63</b>	1279
Weighted average	0.62	0.63	0.61	1279
Macro average	0.47	0.41	0.42	1279

#### 4.1. Comparative Analysis

[Table VI] summarizes the performance of all models. The **CNN with Fine-Tuning** stands out as the top performer, achieving the highest **accuracy (0.95)**, **F1-score (0.96)**, and **recall (0.94)**, as detailed in [Table IV] and visualized in [Fig. 6, 7 and 9]. These results underscore the effectiveness of fine-tuning pre-trained networks for complex medical image classification tasks like Alzheimer's disease diagnosis [Fig. 8].

The Dual-Branch CNN (Combined Model) demonstrates incremental improvements over its individual components (CNN1 and CNN2), achieving an accuracy of 0.63 and F1-score of 0.42 [Table V]. While these results may appear modest, they reflect the specific characteristics and challenges of our dataset and approach. In contrast, El-Assy et al. [14] reported significantly higher accuracies—99.43%, 99.57%, and 99.13%—using a slightly different dual-branch CNN architecture on the ADNI dataset, supported by extensive class balancing and advanced preprocessing techniques.

In traditional machine learning models, KNN model achieved 2nd highest accuracy of 73% [Table II] and ensemble approach [Table III] yield moderate results, while in [15] they attained 93.18% using a modified KNN approach on the OASIS dataset. These comparisons underscore how dataset composition, class distribution, and preprocessing strategies influence model performance and highlight the exploratory nature of our study across different conditions.

TABLE VI: Evaluation Metrics for All Models

Model	Accuracy (Test)	Recall	F1-Score
K-Nearest Neighbors (KNN)	0.73	0.67	0.69
Ensemble Model	0.69	0.63	0.65
CNN with Fine-Tuning	<b>0.95</b>	0.94	0.96
CNN1	0.62	0.37	0.37
CNN2	0.56	0.37	0.38
Combined CNN Model	0.63	0.41	0.42

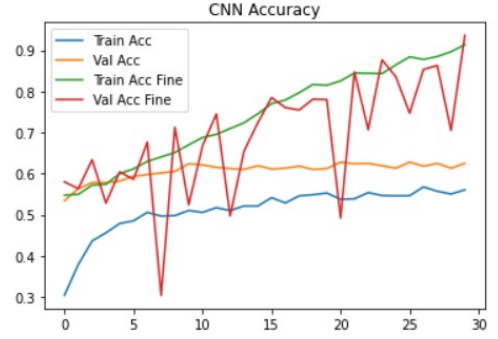


Fig. 6: Training and Validation Accuracy of CNN with Fine-Tuning

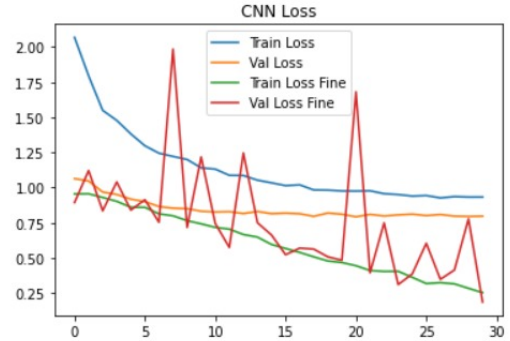


Fig. 7: Training and Validation Loss of CNN with Fine-Tuning

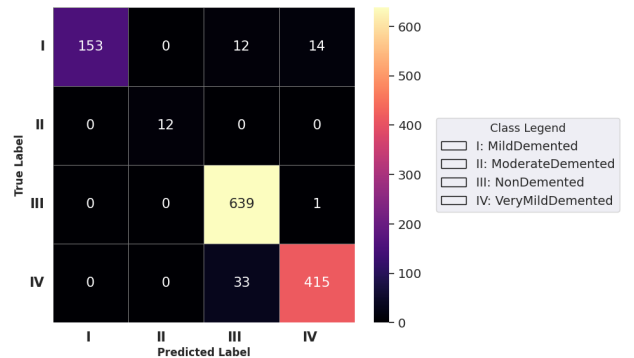


Fig. 8: Confusion matrix of the CNN with Fine-Tuning model



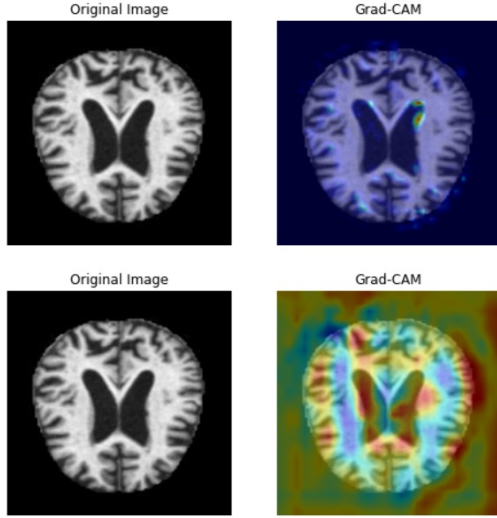


Fig. 9: Accuracy, recall and F1-score comparison across models

#### 4.2. Grad-CAM Visualization

In Grad-CAM heatmaps, red regions signify the most influential areas, while yellow, blue, and green indicate progressively lower relevance.

Grad-CAM visualization for the **Dual-Branch CNN** model (Fig. 10)—reveal several limitations in their attention mechanisms. **CNN1** focuses narrowly on a small portion of the frontal horn of the lateral ventricle, missing the broader neuroanatomical context required for accurate Alzheimer's classification. **CNN2** shows slightly better attention distribution, capturing parts of the frontal horn, grey matter, and splenium of the corpus callosum, but its focus remains scattered and misaligned with clinically relevant areas. The **Combined Model**, while aggregating features from both branches and generating a more structured heatmap, still fails to consistently emphasize clinically important brain regions.

These patterns suggest inadequate high-level feature extraction in the current Dual-Branch CNN architecture, which corresponds with its lower classification performance.

In the Grad-CAM visualization of the **fine-tuned CNN** model (Fig. 11) the most crucial areas of the MRI scans that are necessary for classifying Alzheimer's disease, like

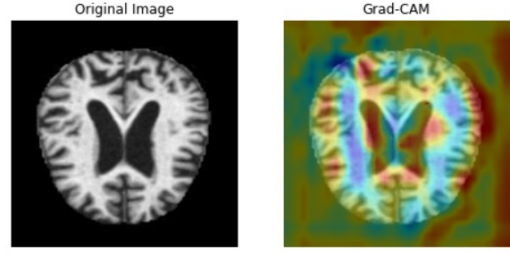


Fig. 10: Grad-CAM Analysis Results: Dual-Branch CNN (CNN1, CNN2, Combined Model)

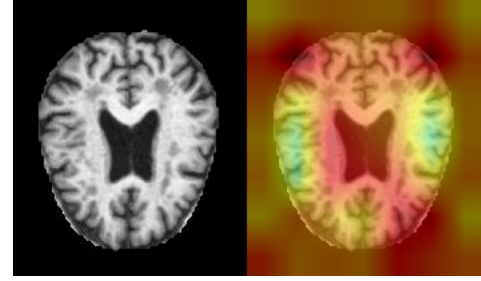


Fig. 11: Grad-CAM Analysis: CNN with Fine Tuning

the frontal horn and the splenium of the corpus callosum, are highlighted in red and yellow. White matter, which is represented by a very tiny blue-marked area, is thought to be less significant for the classification task .

#### 4.3. Discussion

Comparative analysis revealed notable differences in model performance for Alzheimer's MRI classification.

- **Fine-Tuned CNN:** Achieved the highest classification accuracy (95.0%) and strongest performance overall [Fig. 9]. Grad-CAM heatmaps (Fig. 11) showed focused attention on clinically relevant regions such as the entorhinal cortex, ventricles, and hippocampus.
- **K-Nearest Neighbors (KNN):** Secured the second-highest testing accuracy (73.0%) [Table VI], outperforming the ensemble model. Demonstrated the effectiveness of simple, distance-based classifiers on smaller datasets with informative features.
- **Ensemble Model (KNN + Random Forest):** Showed lower performance compared to KNN alone [Table III], indicating that complex feature pipelines may not always improve results on limited data.
- **Dual-Branch CNN (CNN1, CNN2, Combined):** Delivered suboptimal performance relative to the fine-tuned CNN. Among the three, the Combined model performed best but still lagged behind. Heatmaps (Fig. 10) revealed weak or scattered attention focus. CNN1 focused too narrowly on the frontal horn, CNN2 showed dispersed attention across brain areas, and the Combined model offered slight improvements but lacked consistent alignment with clinically significant regions.

## 5. CONCLUSION

This study conducted a comprehensive evaluation of four distinct approaches for classifying Alzheimer's disease using brain MRI scans, spanning traditional machine learning and modern deep learning architectures. Among the tested models, the CNN with fine-tuning demonstrated superior performance (0.95), benefiting from robust training strategies such as data augmentation and regularization. The K-Nearest Neighbors (KNN) model, despite its simplicity, emerged as a strong baseline, achieving the second-highest testing accuracy of 0.73, which reaffirms the value of classical models in small- to medium-scale medical imaging tasks. Grad-CAM visualizations further validated that high-performing models accurately focused on clinically relevant brain regions, enhancing interpretability. Conversely, while the Dual-Branch CNN showed potential, its attention mechanisms and classification performance indicated the need for architectural refinement.

Overall, our findings highlight the strength of deep learning—particularly fine-tuned CNNs—for reliable and interpretable Alzheimer's diagnosis, while emphasizing the continued relevance of traditional methods for benchmarking and complementary evaluation.

## REFERENCES

- [1] VP Subramanyam Rallabandi, Ketki Tulpule, Mahanandeeshwar Gattu, Alzheimer's Disease Neuroimaging Initiative, et al. Automatic classification of cognitively normal, mild cognitive impairment and alzheimer's disease using structural mri analysis. *Informatics in Medicine Unlocked*, 18:100305, 2020.
- [2] Yu Wang, Wen Zhou, Chongchong Yu, and Weijun Su. Assisted magnetic resonance imaging diagnosis for alzheimer's disease based on kernel principal component analysis and supervised classification schemes. *Journal of Information Processing Systems*, 17(1):178–190, 2021.
- [3] Tingting Zhang, Qian Liao, Danmei Zhang, Chao Zhang, Jing Yan, Ronald Ngetich, Junjun Zhang, Zhenlan Jin, and Ling Li. Predicting mci to ad conversation using integrated smri and rs-fmri: machine learning and graph theory approach. *Frontiers in Aging Neuroscience*, 13:688926, 2021.
- [4] Atif Mehmood, Muazzam Maqsood, Muzaffar Bashir, and Yang Shuyuan. A deep siamese convolution neural network for multi-class classification of alzheimer disease. *Brain sciences*, 10(2):84, 2020.
- [5] Amir Ebrahimi, Suhui Luo, and for the Alzheimer's Disease Neuroimaging Initiative. Convolutional neural networks for alzheimer's disease detection on mri images. *Journal of Medical Imaging*, 8(2):024503–024503, 2021.
- [6] Gowhar Mohi ud din dar, Avinash Bhagat, Syed Immamul Ansarullah, Mohamed Tahar Ben Othman, Yasir Hamid, Hend Khalid Alkahtani, Inam Ullah, and Habib Hamam. A novel framework for classification of different alzheimer's disease stages using cnn model. *Electronics*, 12(2):469, 2023.
- [7] Farheen Ramzan, Muhammad Usman Ghani Khan, Asim Rehmat, Sajid Iqbal, Tanzila Saba, Amjad Rehman, and Zahid Mehmood. A deep learning approach for automated diagnosis and multi-class classification of alzheimer's disease stages using resting-state fmri and residual neural networks. *Journal of medical systems*, 44:1–16, 2020.
- [8] Harshit Parmar, Brian Nutter, Rodney Long, Sameer Antani, and Sunanda Mitra. Spatiotemporal feature extraction and classification of alzheimer's disease using deep learning 3d-cnn for fmri data. *Journal of Medical Imaging*, 7(5):056001–056001, 2020.
- [9] FM Javed Mehedi Shamrat, Shamima Akter, Sami Azam, Asif Karim, Pronab Ghosh, Zarrin Tasnim, Khan Md Hasib, Friso De Boer, and Kawsar Ahmed. Alzheimernet: An effective deep learning based proposition for alzheimer's disease stages classification from functional brain changes in magnetic resonance images. *IEEE Access*, 11:16376–16395, 2023.
- [10] Chahd M Chabib, Leontios J Hadjileontiadis, and Aamna Al Shehhi. Deepcurvmri: Deep convolutional curvelet transform-based mri approach for early detection of alzheimer's disease. *IEEE Access*, 11:44650–44659, 2023.
- [11] Öznur Özaltn. Early detection of alzheimer's disease from mr images using fine-tuning neighborhood component analysis and convolutional neural networks. *Arabian Journal for Science and Engineering*, pages 1–20, 2025.
- [12] Nikhil J Dhinagar, Sophia I Thomopoulos, Emily Laltoo, and Paul M Thompson. Counterfactual mri generation with denoising diffusion models for interpretable alzheimer's disease effect detection. *2024 46th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, pages 1–6, 2024.
- [13] Marco Pinamonti. Alzheimer mri 4 classes dataset. Online dataset, 2022. Kaggle Dataset.
- [14] AM El-Assy, Hanan M Amer, HM Ibrahim, and MA Mohamed. A novel cnn architecture for accurate early detection and classification of alzheimer's disease using mri data. *Scientific Reports*, 14(1):3463, 2024.
- [15] Srinivasan Aruchamy, Veeramachaneni Mounya, and Ankit Verma. Alzheimer's disease classification in brain mri using modified knn algorithm. *IEEE International Symposium on Sustainable Energy, Signal Processing and Cyber Security (iSSSC)*, 2020.