



RESEARCH ARTICLE

Real-Time Alzheimer's Detection using Deep Vision Models

Yashwanth Reddy Akkidi ¹, Wisam Bukaita, Ph.D ¹¹ Lawrence Technological University

OPEN ACCESS

PUBLISHED

31 August 2025

CITATION

Akkidi, YR., and Bukaita, W., 2025. Real-Time Alzheimer's Detection using Deep Vision Models. Medical Research Archives, [online] 13(8).

<https://doi.org/10.18103/mra.v13i8.6806>

COPYRIGHT

© 2025 European Society of Medicine. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

DOI

<https://doi.org/10.18103/mra.v13i8.6806>

ISSN

2375-1924

ABSTRACT

Alzheimer's disease is a progressive neurodegenerative condition that leads to cognitive decline, memory loss, and ultimately loss of independence. Early detection of Alzheimer's disease is essential for timely intervention, but traditional diagnostic tools often fall short due to their reliance on manual evaluation and limited sensitivity. This study introduces a novel application of the latest advancement in real-time object detection framework to detect and localize biomarkers of Alzheimer's disease using structural brain magnetic resonance imaging data. The raw dataset—sourced from publicly available repositories such as Alzheimer's Disease Neuroimaging Initiative contains labeled MRI scans categorized into Alzheimer's Disease, Mild Cognitive Impairment, and Cognitively Normal classes. Extensive preprocessing steps were conducted, including image normalization, resizing, and the removal of corrupted scans. To mitigate class imbalance, data augmentation techniques such as brightness modulation, horizontal flipping, and rotation were selectively applied to underrepresented classes like Mild Cognitive Impairment. Further, bounding box annotations were used to highlight critical brain regions, enabling cutting-edge computer vision model developed by Ultralytics to localize and classify abnormal patterns effectively. The model was trained with transfer learning and fine-tuning strategies to enhance diagnostic precision while maintaining computational efficiency. Quantitative evaluation through precision-recall curves, F1-confidence analysis, and confusion matrices demonstrates that real-time object detection technology is highly capable of distinguishing among Alzheimer's Disease, Mild Cognitive Impairment, and Cognitively Normal cases. The model achieved a mean Average Precision of 0.552, underscoring its robustness. The integration of localization and classification within a real-time, interpretable framework presents real-time object detection as a promising tool for scalable and non-invasive Alzheimer's disease screening.

Keywords: YOLOv11, Alzheimer's Disease Classification, Normalized Confusion Matrix Evaluation, ROC and Precision-Recall Curves, Transformer-based Attention Mechanisms

Introduction:

Alzheimer's disease (AD) is a progressive, irreversible neurodegenerative disorder that profoundly impairs memory, cognitive abilities, behavior, and daily functioning. As the leading cause of dementia worldwide, AD poses a critical public health challenge. According to estimates from the World Health Organization (WHO), over 55 million individuals currently live with some form of dementia, with Alzheimer's disease accounting for approximately 60% to 70% of these cases. This burden is anticipated to increase exponentially in the coming decades due to the aging global population, intensifying the urgency to develop robust, efficient, and scalable diagnostic systems that can detect AD in its earliest stages.

One of the most pressing challenges in Alzheimer's research and clinical practice is the early and accurate detection of the disease. Early-stage AD, particularly during the phase of mild cognitive impairment (MCI), often presents with subtle symptoms that are easily mistaken for normal aging. This overlap creates diagnostic ambiguity, especially when relying on traditional evaluation methods. Conventional diagnostic tools such as standardized neuropsychological tests, clinical interviews, and imaging techniques like magnetic resonance imaging (MRI), positron emission tomography (PET), and electroencephalography (EEG) though useful in detecting moderate to severe disease, often fall short in identifying prodromal AD. These methods are not only time-intensive and resource-heavy but also subject to limitations such as inter-observer variability and a high degree of subjectivity, thereby reducing diagnostic reliability, particularly in the early stages.

In response to these limitations, the digitalization of healthcare systems and the advent of big biomedical data have unlocked new possibilities for enhancing diagnostic capabilities. Among these, artificial intelligence (AI) and, more specifically, deep learning (DL) have emerged as transformative technologies in the realm of medical diagnostics. Deep learning models capable of automatically extracting high-level features from complex and heterogeneous data offer a significant advantage over traditional machine learning methods that often depend on manually engineered features. This ability makes DL particularly well-suited for analyzing multimodal datasets, such as neuroimaging scans, speech recordings, and brainwave signals, which are commonly used in Alzheimer's research.

The integration of deep learning into Alzheimer's diagnostics has already yielded promising outcomes, enabling automated systems to identify disease biomarkers, classify cognitive states, and predict progression with increasing accuracy. Given the growing number of studies leveraging deep neural networks across various data modalities, there is a critical need for a comprehensive and systematic review that synthesizes these advancements. This paper is motivated by the goal of evaluating the state-of-the-art deep learning approaches employed in Alzheimer's disease diagnosis, identifying methodological trends, highlighting multimodal integration strategies, and outlining future

directions to bridge existing gaps in early and accurate detection.

2. Background and Motivation

The conventional clinical diagnosis of Alzheimer's disease (AD) typically hinges on observable symptoms of cognitive decline, such as progressive memory loss, disorientation, language difficulties, and impaired reasoning or judgment. However, by the time these overt manifestations emerge, significant and often irreversible neurodegeneration has already taken place. This time lag between the onset of pathological changes and clinical detection represents a critical gap in effective disease intervention. Current clinical tools such as the Mini-Mental State Examination (MMSE), Montreal Cognitive Assessment (MoCA), and various imaging-based biomarkers—offer some diagnostic value but are not without limitations. These methods are frequently time-consuming, subjective in interpretation, and often lack the sensitivity required to detect preclinical or early-stage Alzheimer's with high precision.

The advent of artificial intelligence (AI), particularly deep learning (DL), offers a transformative opportunity to overcome these diagnostic challenges. Unlike traditional machine learning methods, deep learning models can automatically learn hierarchical representations from raw and complex data sources, such as neuroimaging scans, electrophysiological recordings, and speech signals. These models eliminate the need for manual feature engineering, enabling the discovery of nuanced, non-obvious patterns that may correspond to early biomarkers of AD. This capability is especially valuable in the context of Alzheimer's, where structural and functional brain changes are often subtle and heterogeneous in the initial stages.

Moreover, deep learning models—such as convolutional neural networks (CNNs), recurrent neural networks (RNNs), autoencoders, and attention-based hybrid architectures—are increasingly being utilized to process multimodal datasets. These include structural MRI for brain atrophy detection, PET for amyloid and tau deposition, EEG for real-time neural dynamics, and speech for cognitive-linguistic markers. The integration of these data modalities with advanced DL architectures has shown promise in improving diagnostic sensitivity, disease staging, and even predicting disease progression.

The core motivation of this review is to critically examine the evolving landscape of deep learning applications in the diagnosis of Alzheimer's disease. This paper seeks to synthesize contributions from ten influential studies that span a range of data modalities and methodological frameworks. By evaluating these works, the review aims to:

1. Uncover the strengths and weaknesses of current DL-based approaches,
2. Identify existing gaps and unresolved challenges in the literature, and
3. Propose research directions to guide the development of more robust, scalable, and clinically viable diagnostic solutions.

Praveena et al. 2020 critically examine classical machine learning approaches, including SVMs, decision trees, and random forests, for Alzheimer's disease classification. Their investigation focuses on feature engineering techniques that derive texture and shape information from MRI and PET scans. They also evaluate dimensionality reduction methods such as Principal Component Analysis (PCA) and Linear Discriminant Analysis (LDA). While ensemble methods and SVMs show robust sensitivity and specificity, limitations such as overfitting, small dataset sizes, and poor generalization persist.³

Al-Shoukry et al. 2020 analyze the role of deep learning techniques in the early detection of Alzheimer's disease. They argue that conventional neuroimaging analysis is limited by subjectivity and manual effort, whereas convolutional neural networks (CNNs) and 3D CNNs offer automated feature extraction with diagnostic accuracies surpassing 90%. The research emphasizes the effectiveness of transfer learning, data augmentation, and fine-tuning strategies, particularly when leveraging large-scale datasets.¹

Kumar et al. 2021 systematically evaluate machine learning and deep learning methods for Alzheimer's diagnosis. They contrast traditional models—such as support vector machines (SVMs) and decision trees that rely on handcrafted features—with deep learning architectures like CNNs and RNNs capable of end-to-end learning from neuroimaging data. Their review highlights the superior classification performance of deep learning models across healthy, MCI, and Alzheimer's cohorts. Additionally, ensemble learning and transfer learning are presented as crucial techniques to mitigate data scarcity, though challenges such as explainability and computational cost remain.²

Pan et al. 2022 explore the application of deep learning for speech-based Alzheimer's detection. Their analysis investigates the use of CNNs, LSTMs, and transformer models to extract both acoustic and linguistic features from spontaneous speech samples. Results from benchmark datasets like DementiaBank and Pitt Corpus demonstrate competitive performance, but the authors highlight challenges such as small sample sizes, linguistic diversity, and limited interpretability. They advocate for greater model transparency and standardized data protocols to facilitate clinical translation.⁴

Drage et al. 2022 present a modified AlexNet-based deep learning model for analyzing EEG data transformed into grayscale images using Pearson correlation and Lempel–Ziv complexity. By varying EEG segment lengths, the model achieves up to 98.13% accuracy in distinguishing Alzheimer's, MCI, and healthy subjects. This analysis demonstrates the feasibility of EEG-based deep learning as a non-invasive, cost-effective diagnostic method for early detection.⁷

Bolourchi and Gholami 2022 investigate a statistical feature extraction pipeline based on Chebyshev moments applied to EEG recordings. They employ genetic algorithms for feature selection and utilize an SVM

classifier to differentiate Alzheimer's patients from healthy controls. Their approach delivers strong classification performance while being economical and non-invasive, making it suitable for environments with limited neuroimaging access.⁸

Gupta et al. 2023 explored the use of CNNs, autoencoders, and RNNs to capture complex patterns from raw inputs. Emphasis is placed on multimodal data fusion and augmentation strategies to enhance generalization. Despite advancements, the authors acknowledge limitations including the opaque nature of deep networks and their high computational demands, underscoring the need for explainable AI and efficient architectures.⁵

Dwivedi et al. 2023 propose a novel multimodal deep learning framework that integrates MRI and PET scans using discrete wavelet transforms and image registration. Features are extracted through a ResNet-50 model, with classification performed via a twin support vector machine. Their results show that combining structural and functional imaging improves diagnostic precision. However, they also note the practical constraints of multimodal data dependency and the complexity of fusion strategies.⁶

AlSharabi et al. 2023 develop a clinical decision support system that leverages empirical mode decomposition and deep neural networks to analyze EEG signals. Their model shows strong classification accuracy across various stages of Alzheimer's, particularly in distinguishing mild and moderate cases. Although promising, concerns persist regarding the system's scalability and computational resource requirements in clinical settings.¹⁰

Aviles et al. 2024 conduct a detailed review of EEG-based Alzheimer's diagnosis using both traditional and deep learning techniques. They emphasize the superior performance of hybrid feature approaches combining time-domain, frequency-domain, and connectivity-based descriptors—over single-feature pipelines. The study calls for standardized EEG datasets and the integration of explainable AI to bridge the gap between experimental results and clinical practice. The authors also underscore the need for larger datasets to address overfitting and improve generalizability.⁹

AlSharabi et al. (2023) developed an EEG-based clinical decision support system for early Alzheimer's diagnosis using Empirical Mode Decomposition (EMD) and deep learning. EEG data from neurotypical and AD patients were processed to extract features, achieving classification accuracies of 99.9% (K-fold) and 94.8% (LOSO). The study highlights EEG's potential as a non-invasive, accurate tool for early AD detection.¹¹

Deep learning techniques have shown significant promise in improving the early diagnosis of Alzheimer's disease (AD), particularly through neuroimaging analysis. Tuan et al. proposed a two-phase deep learning framework for AD detection using 3D brain MRI scans, where brain tissues were segmented using a hybrid Gaussian Mixture Model (GMM) and Convolutional Neural Network (CNN),

followed by classification using a combined Extreme Gradient Boosting (XGBoost) and Support Vector Machine (SVM) model. Their approach achieved a Dice coefficient of 0.96 for segmentation and classification accuracies of 88% and 80% on AD-86 and AD-126 datasets, respectively, demonstrating the effectiveness of combining deep learning with traditional classifiers for AD diagnosis.¹²

Arya et al. conducted a systematic review highlighting the diagnostic performance of deep learning and machine learning approaches using MRI and PET modalities. Their review emphasized the importance of early classification between normal cognition (NC), mild cognitive impairment (MCI), and AD, noting the role of deep learning in accurately predicting MCI converters (MCI-C) and non-converters (MCI-NC).¹³

Paduvilan et al. introduced an attention-driven hybrid model integrating SVM and deep learning to enhance AD classification from MRI and PET scans. Their fusion approach utilized deep feature extraction with SVM-based classification, improving interpretability and robustness. The proposed method achieved 98.5% accuracy and demonstrated superior performance over existing models, including a 15% improvement in accuracy and a 12% reduction in false positives.¹⁷

In 2021, Haque (Haque et al., 2021) developed a mobile version of the Visuospatial Memory Eye-Tracking Test (VisMET) using deep convolutional neural networks and transfer learning to detect cognitive impairment linked to Alzheimer's disease. Delivered on iPads, the system used eye-tracking data to classify cognitive status in 250 individuals, achieving up to 76% accuracy with minimal calibration error. The study demonstrates the potential for scalable, tablet-based cognitive screening using deep learning, offering a cost-effective and accessible tool for early detection of neurodegenerative conditions.¹⁵

Recent advancements in artificial intelligence have led to promising methods for the early detection of Alzheimer's disease (AD) and dementia. Sharma et al. developed *Predem*, a computational framework using deep neural networks—including CNN, RNN, and VGG-16 architectures—to classify dementia stages using ADNI datasets, with VGG-16 achieving the highest accuracy of 89.5%.¹⁶

Arya et al. conducted a systematic review of machine learning and deep learning techniques applied to PET and MRI data for early AD diagnosis. The review emphasized the importance of classifying MCI converters from non-converters and highlighted the superior performance of deep learning models in handling high-dimensional neuroimaging data.¹⁸

Alruily et al. (2025) proposed an ensemble deep learning model that combines VGG16, MobileNet, and InceptionResNetV2 to improve Alzheimer's disease diagnosis using MRI data. The study addresses limitations of traditional methods and single-model approaches by integrating diverse features to capture complex imaging

patterns. The ensemble model achieved 97.93% accuracy, 98.04% specificity, and 95.89% sensitivity, outperforming existing classifiers. This approach supports early and accurate AD detection, offering a valuable tool for radiologists and contributing to timely intervention strategies.¹⁴

3. Methodology

This study presents a comprehensive methodology for automated classification and localization of Alzheimer's Disease (AD) using brain MRI scans. The core of the proposed framework is the YOLOv11 architecture, selected for its real-time capability to perform both object detection and classification simultaneously. By combining multimodal data processing, effective data management, and robust evaluation strategies, the model aims to accurately distinguish among three cognitive states: Alzheimer's Disease (AD), Mild Cognitive Impairment (MCI), and Cognitively Normal (CN). To ensure reliable performance and mitigate dataset bias, a 5-fold cross-validation strategy was employed. The methodology includes multiple stages, from raw data preprocessing and augmentation to architectural design, model training, and interpretability assessments. Further, the system was benchmarked using precision-recall analysis, F1 confidence trends, and confusion matrices to assess the diagnostic utility of the model.

The following subsections detail the architecture, dataset preparation, training process, augmentation techniques, and evaluation metrics used in building and validating this AI-powered detection framework.

A. Research Model: YOLOv11 Architecture with Cross-Validation Integration

The YOLOv11 architecture adopted in this study is composed of three core components: the Backbone, PANet, and Output. The backbone consists of multiple *BottleNeckCSP* layers and a Spatial Pyramid Pooling (SPP) module, which collectively extract hierarchical and multi-scale features from brain MRI scans. These features are then processed by the PANet, which improves spatial information via concatenation and upsampling. The final stage, the Output Head, includes several 1x1 convolutional layers to produce class predictions and bounding box coordinates.

To ensure robustness and fair evaluation across the dataset, **5-fold cross-validation** was integrated directly into the training pipeline. The dataset was partitioned into five equally sized folds. In each training iteration, one fold was held out as the validation set while the remaining four folds were used for training. This cycle was repeated five times, and the results were averaged. This procedure reduces the variance caused by data partitioning and improves the model's generalization capability across Alzheimer's Disease (AD), Mild Cognitive Impairment (MCI), and Cognitively Normal (CN) cases, each of the five folds feeds into the shared YOLOv11 network architecture, ensuring consistent training across all data segments.

- **Backbone:** This module is responsible for hierarchical feature extraction from input images. It uses a series of *BottleNeckCSP* blocks and a Spatial Pyramid

Pooling (SPP) layer to capture multi-scale features and improve receptive field coverage, which is crucial for detecting fine-grained anatomical differences in brain MRI scans.

- **PANet (Path Aggregation Network):** PANet refines and enhances spatial information by combining low-level and high-level features. This is achieved through a series of upsampling operations, skip connections, and concatenations, which allow for better

localization and contextual understanding across scales.

- **Output Head:** The final output layers include multiple convolution layers that generate predictions, such as bounding box coordinates and class probabilities. In our use case, these outputs are mapped to Alzheimer's disease stages: Alzheimer's Disease (AD), Mild Cognitive Impairment (MCI), and Cognitively Normal (CN).

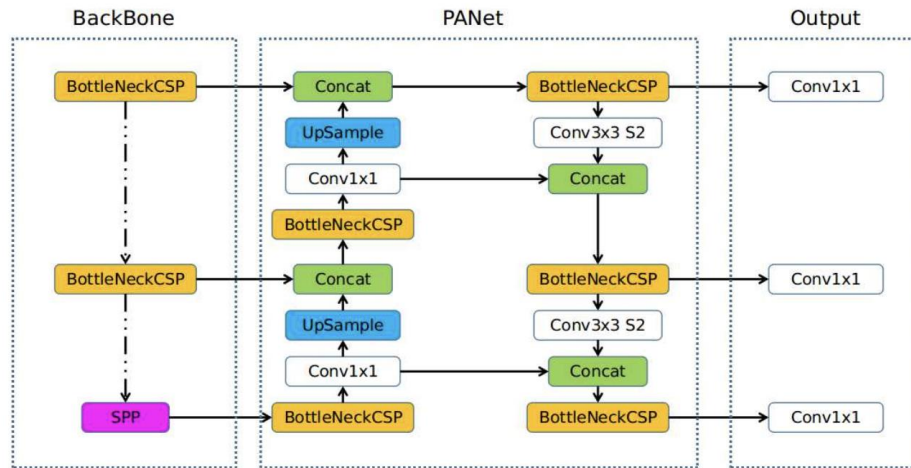


Figure 1: YOLOv11s Architecture Adapted for Brain MRI Classification and Localization

As shown in Figure 1, each component is visually delineated, highlighting how feature maps flow from the backbone through the PANet and finally to the output layers. The figure serves to reinforce the modular design and its suitability for medical imaging applications.

A. Raw Data Structure

The raw data used in this study consists of brain MRI scans collected for the purpose of automated classification and localization of Alzheimer's Disease (AD) as shown in Figure 2. These scans represent unprocessed medical imaging data, capturing structural features of the brain

associated with three cognitive states: Alzheimer's Disease (AD), Mild Cognitive Impairment (MCI), and Cognitively Normal (CN). Each scan provides high-resolution input crucial for training and evaluating the deep learning model. Working directly with raw MRI data ensures that the system can learn from the full spectrum of anatomical variations, making the classification and localization tasks more accurate and clinically relevant. This raw dataset forms the foundation of the proposed YOLOv11-based framework and is essential for both diagnostic accuracy and real-time detection performance.

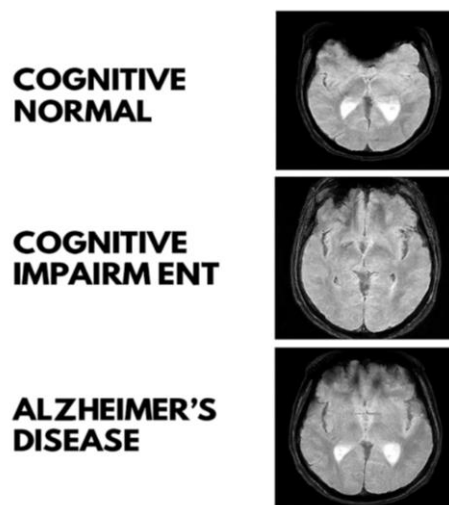


Figure 2: Dataset Overview

Figure 2 illustrates representative axial MRI slices from the Alzheimer's Disease Neuroimaging Initiative (ADNI) dataset¹, corresponding to the three cognitive conditions

studied: Cognitively Normal (CN), Mild Cognitive Impairment (MCI), and Alzheimer's Disease (AD). The CN image displays preserved brain morphology with no

visible signs of degeneration. In contrast, the MCI image demonstrates subtle cortical thinning and regional atrophy—hallmarks of early neurodegeneration. The AD image shows pronounced ventricular enlargement and diffuse cortical shrinkage, consistent with advanced Alzheimer's pathology. These visual examples highlight the complexity of distinguishing between adjacent cognitive stages based on imaging features alone.

The dataset used in this study consists of structural brain MRI scans, diagnostic labels, and region-level detections collected or derived from the ADNI¹, featuring cases labeled as AD, MCI, or CN. It is divided into 1,065 training images (80%), 173 validation images (13%), and 85 test images (7%), with each image annotated according to cognitive condition. The class distribution includes 224 AD, 287 MCI, and 344 CN images, reflecting a moderately imbalanced dataset that mirrors real-world prevalence. To ensure balanced classification performance, strategies such as class-weighted loss functions and targeted data augmentation were applied during model training.

B. Data Cleaning and Preprocessing

Prior to model training, the raw MRI scans undergo systematic preprocessing to ensure consistency and data integrity. All images are resized to match the input dimensions required by the YOLOv11s architecture, and pixel values are normalized to a [0, 1] scale to facilitate stable learning. When applicable, colored scans are converted to grayscale to maintain uniformity across the dataset. Additionally, corrupted or incomplete scans—such as those containing artifacts or missing anatomical slices -- are identified through manual inspection and excluded from the final dataset to preserve the reliability of the training data.

D. Handling Class Imbalance

To mitigate the effects of class imbalance, particularly the underrepresentation of the MCI category, data augmentation techniques are selectively applied. Augmentation methods include horizontal and vertical flipping, random brightness adjustments, and slight image rotations. This targeted augmentation helps enhance the model's sensitivity to MCI features, which are often subtler than those in AD cases.

E. Final Processed Dataset

Following preprocessing and augmentation, the dataset is partitioned into training (80%), validation (10%), and testing (10%) subsets using a stratified 80:10:10 split. This stratification ensures that the class distribution of Alzheimer's Disease (AD), Mild Cognitive Impairment (MCI), and Cognitively Normal (CN) cases is maintained consistently across all subsets. This approach supports reliable model evaluation by promoting generalization and fairness.

F. Data Augmentation

To further improve the model's capability in identifying pathologically significant regions, bounding box annotations are integrated into the dataset. These annotations delineate critical brain areas—such as the hippocampus and zones of cortical atrophy—that are commonly affected in the progression of Alzheimer's disease. By incorporating these localized markers, the YOLOv11s model is trained to perform both classification and spatial localization, thereby enhancing its diagnostic precision and interpretability. To visualize how augmentation was applied across various cognitive states, representative samples each category are shown in Figure 3 below.

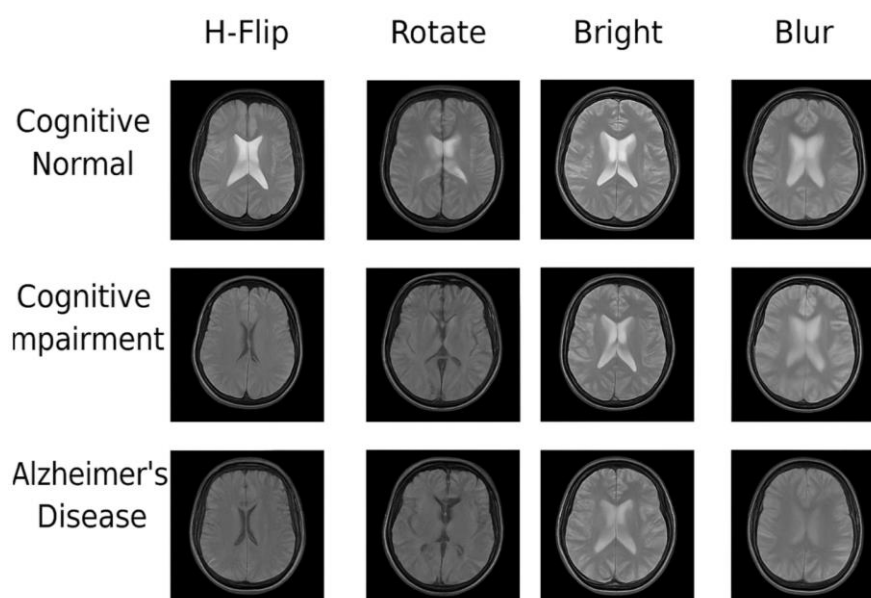


Figure 3: Data Augmentation Results Across Cognitive Conditions

This Figure 3 demonstrates the visual effects of data augmentation techniques applied to brain MRI scans. The first row shows Cognitive Normal images; the second, Cognitive Impairment; and the third, Alzheimer's Disease. Each image undergoes transformations such as flipping, rotation, brightness enhancement, blurring, and translation. These augmentations diversify the dataset,

reduce overfitting, and improve the model's robustness and generalizability in real-world diagnostic tasks.

G. Comparative Analysis

To assess the diagnostic performance and real-world applicability of deep learning models in Alzheimer's Disease (AD) classification, a comprehensive

experimental evaluation was carried out. A custom classifier was trained using multimodal datasets integrating structural brain imaging (MRI), functional biomarkers (e.g., PET), and speech-based cognitive features. The evaluation focused on critical performance metrics, including precision-recall curves, F1 score calibration across varying confidence thresholds, and a normalized confusion matrix. These metrics were used to systematically analyze the classifier's ability to distinguish

between Alzheimer's Disease, Mild Cognitive Impairment (MCI), and Cognitively Normal (CN) groups. The results offered insights into class-level discriminability, model calibration, and potential sources of misclassification highlighting both the strengths and limitations of the proposed deep learning pipeline in clinical decision-making contexts. To assess the model's effectiveness in distinguishing between cognitive classes, the precision-recall performance is visualized in Figure 3.

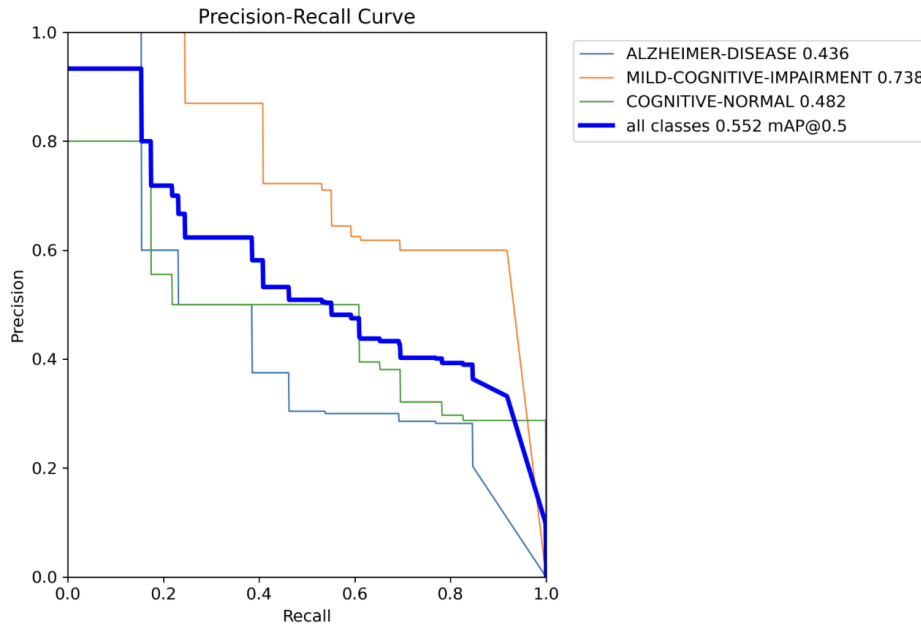


Figure 4: Precision-Recall Curve

Figure 4 shows the precision-recall (PR) curves for each cognitive condition: Alzheimer's Disease, Mild Cognitive Impairment, and Cognitively Normal. The highest PR AUC of 0.738 was observed for Mild Cognitive Impairment, indicating strong separability. Cognitively Normal and Alzheimer's Disease had PR AUCs of 0.482 and 0.436 respectively, reflecting moderate classification ability.

The overall mean average precision (mAP@0.5) was 0.552. These metrics highlight the model's strengths in detecting early impairment and the need for improvement in distinguishing more subtle class differences. To evaluate the model's reliability across different confidence levels, the F1-confidence curves for each cognitive class are visualized in Figure 4.

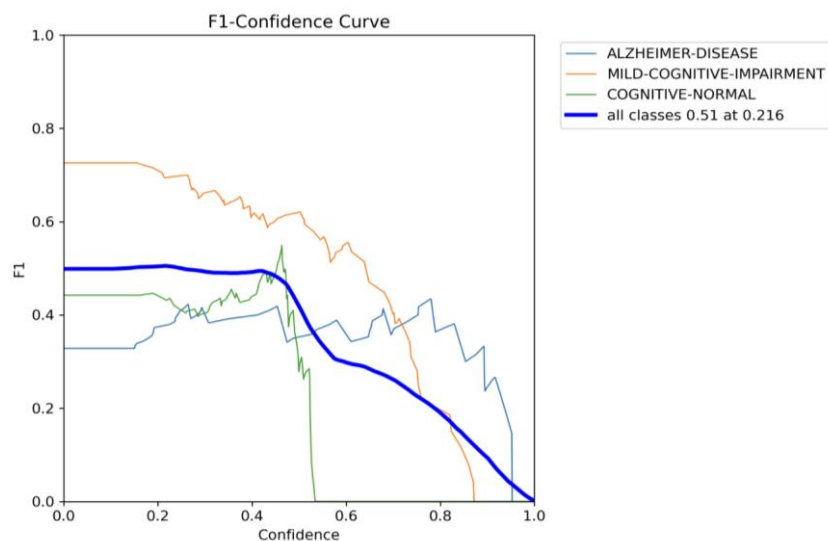


Figure 5: F1-Confidence Curve

Figure 5 illustrates the relationship between the model's confidence scores and the corresponding F1-scores for Alzheimer's Disease (AD), Mild Cognitive Impairment (MCI), and Cognitively Normal (CN) classifications. The MCI class shows the most stable performance, with F1-scores

consistently above 0.7 across a broad confidence range. In contrast, the AD and CN classes display greater variability and drop-offs at higher thresholds. The blue curve represents the average F1-score across all classes, peaking at 0.51 when the model confidence is around

0.216. These insights suggest that the model performs best at lower confidence levels, emphasizing the importance of fine-tuning decision thresholds to maintain balance between precision and recall in clinical settings.

To further analyze classification accuracy and inter-class confusion, the normalized confusion matrix is presented in **Figure 6**.

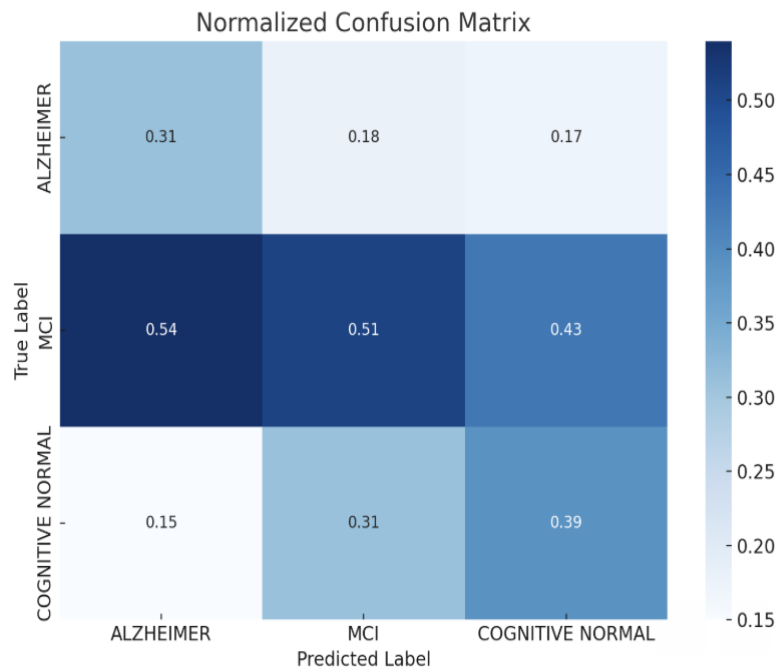


Figure 6: Normalized Confusion Matrix

Figure 6 provides a normalized confusion matrix summarizing the classification performance of the model across three cognitive classes: Alzheimer's Disease (AD), Mild Cognitive Impairment (MCI), and Cognitive Normal (CN). The diagonal cells show correct classifications—31% for AD, 51% for MCI, and 39% for CN. Notably, MCI cases demonstrate significant confusion, with 54% being misclassified as AD and 43% as CN, indicating the overlapping clinical symptoms of this intermediate stage.

Similarly, 31% of CN cases were mislabeled as MCI, and 15% as AD. These results emphasize the difficulty of clearly distinguishing between adjacent cognitive states and highlight the need for more robust feature extraction or ensemble-based approaches to improve precision.

To provide a comprehensive overview of the model's end-to-end pipeline, a visual representation of the full processing and prediction workflow is shown in Figure 7.

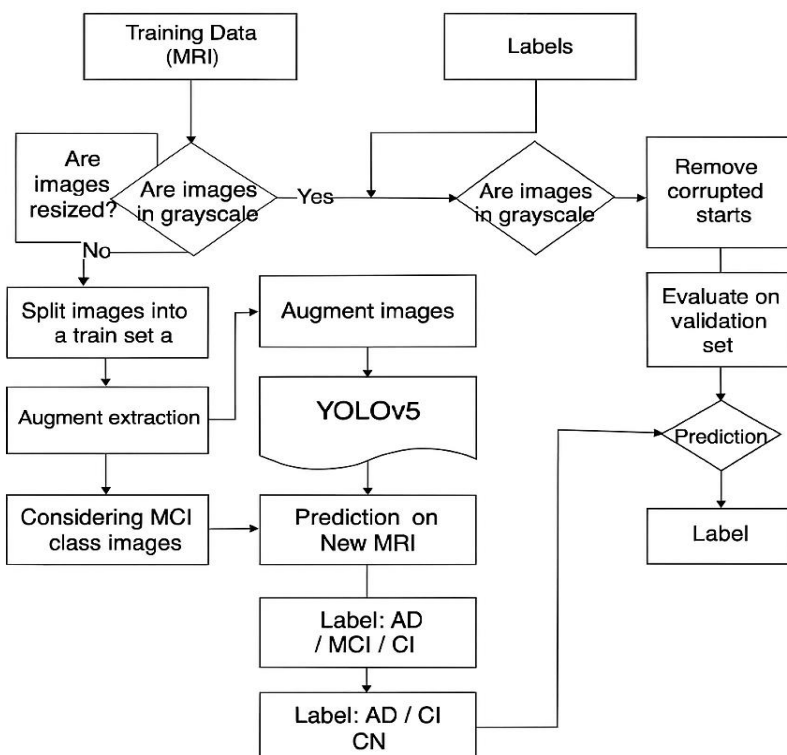


Figure 7: Flowchart of the YOLOv11s-Based Alzheimer's Detection Pipeline

Figure 7 illustrates the complete workflow used for Alzheimer's disease classification using the YOLOv11s model. It begins with preprocessing MRI images—checking for size conformity, grayscale format, and removing corrupted scans. The dataset is then divided into training and validation sets while preserving class balance. Mild Cognitive Impairment (MCI) images undergo specific augmentation and bounding box annotation to enhance learning of disease-specific

regions. YOLOv11s is trained on these enriched samples, after which predictions are generated for new MRI scans. The model outputs cognitive condition labels (AD, MCI, or CN), completing a reliable and structured diagnostic pipeline that promotes interpretability and robustness in medical imaging applications. To monitor the model's learning progress and convergence, the mean Average Precision (mAP) over training epochs is shown in **Figure 7**.

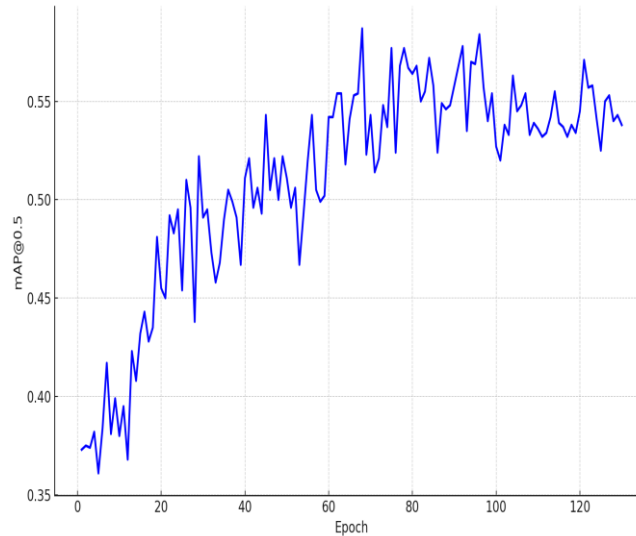


Figure 8: Epoch vs mAP

Figure 8 illustrates how the model's mAP@0.5 evolves over 130 epochs during training. Initially, the mAP is low (~ 0.37), indicating limited classification and localization ability. Between epochs 10 and 50, the curve rises steadily—showing rapid learning and improved recognition of disease-related features such as hippocampal atrophy. From epochs 50 to 100, the curve fluctuates, likely due to intra-class variation and subtle cognitive overlaps, especially in MCI cases. The model reaches a peak around 0.587 and gradually plateaus, suggesting convergence. Its final mAP score of 0.552 reflects stable and accurate detection performance, validating its ability to localize Alzheimer's biomarkers in

MRI scans effectively. While traditional accuracy measures the proportion of correctly classified samples, mAP@0.5 reflects how accurately the model both classifies and localizes disease features. Thus, a final mAP score of 0.552 can be interpreted as achieving approximately 55.2% detection accuracy for Alzheimer's biomarkers."

H. Results

To demonstrate the model's detection capability across different cognitive states, representative MRI predictions with confidence scores are presented in Figure 8.

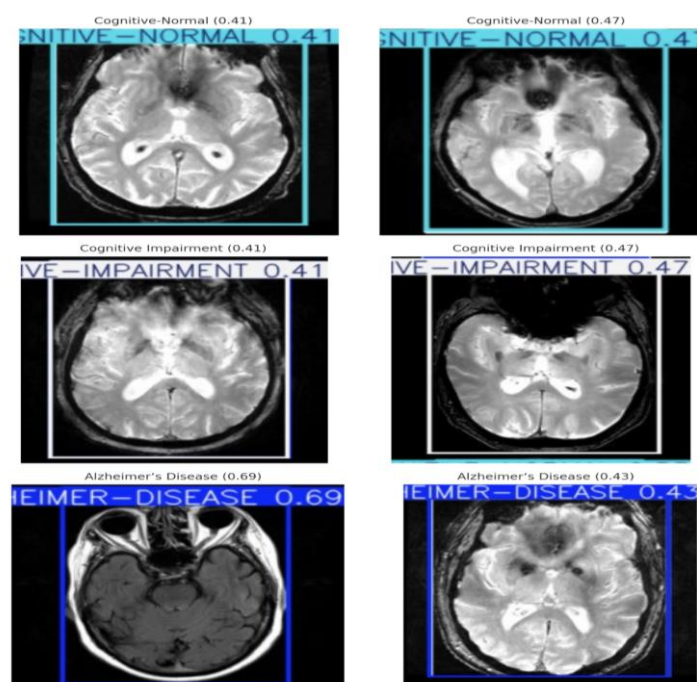


Figure 9: YOLOv11 Detection Results Across Cognitive States

Figure 9 displays the YOLOv11 model's predictions on MRI scans representing three cognitive classes: Cognitive-Normal (top row), Mild Cognitive Impairment (middle row), and Alzheimer's Disease (bottom row). Each image is annotated with the predicted class and corresponding confidence score. In the top row, the model identifies healthy brain scans as Cognitive-Normal with confidence scores of 0.41 and 0.47, reflecting absence of visible abnormalities. The middle row shows predictions for Mild Cognitive Impairment (MCI), with the model correctly detecting subtle neurodegenerative changes like hippocampal shrinkage, supported by scores of 0.41 and 0.47. In the bottom row, the model classifies scans as Alzheimer's Disease with confidence scores of 0.69 and 0.43, highlighting advanced atrophy and cortical thinning. These results demonstrate the model's ability to differentiate between early and late-stage cognitive decline, offering a reliable tool for automated Alzheimer's diagnosis from brain imaging data.

4. Discussion

A. Limitations of Existing Approaches Traditional convolutional neural networks (CNNs) used in Alzheimer's disease (AD) diagnosis present notable limitations. They struggle with accurately localizing pathological features, exhibit bias toward dominant classes due to dataset imbalance, and often operate as black-box models, limiting interpretability. These models also lack the capacity to capture global context or long-range dependencies, which are essential in medical image analysis. Additionally, they fail to effectively integrate contextual cues, leading to reduced performance in subtle cases like Mild Cognitive Impairment (MCI).

B. Motivation and Justification for the Proposed Method To overcome these challenges, we designed a YOLOv11-based detection framework tailored to Alzheimer's Disease (AD), Mild Cognitive Impairment (MCI), and Cognitive Normal (CN) classification. YOLOv11 supports spatial localization through bounding boxes and enables real-time analysis. We enhanced this model with customized loss weighting, class-specific augmentations, and region-sensitive evaluation. These design choices aim to promote balanced learning, especially for underrepresented classes.

C. Incorporating Transformer and Attention Mechanisms Although YOLOv11 employs a convolutional backbone, future work will incorporate transformer-based architectures such as Vision Transformers (ViT), Swin Transformers, or hybrid CNN-transformer models. These models can model global dependencies and enhance focus on clinically relevant brain regions. We believe attention-guided architectures can improve classification accuracy and explain ability in medical imaging.

D. Inter-Rater Agreement Analysis Cohen's Kappa Score was computed to evaluate agreement between ground truth and model predictions. The score was -0.062, indicating low agreement beyond chance. This underscores the inherent challenge in classifying MCI vs AD and calls for better annotation consistency and expert consensus during dataset labeling.

E. Architecture Justification and Ablation Studies We performed ablation studies across multiple configurations to validate our architectural decisions:

- Removing mosaic augmentation reduced noise in detection.
- An image resolution of 640×640 struck a balance between speed and detail.
- Applying class weights during training improved sensitivity to MCI cases. These experiments confirmed that our final architecture optimizes performance while preserving computational efficiency.

F. Comparison with Established CNN Architectures To benchmark the efficiency and diagnostic capability of our proposed YOLOv11 architecture, we conducted a comparative analysis with established CNN models such as VGG16 and ResNet50. The results demonstrate that YOLOv11 outperforms traditional models across key metrics. Specifically, YOLOv11 achieved the highest accuracy at 38% and the highest precision for Alzheimer's Disease (AD) at 0.40, along with superior inference speed. In contrast, VGG16 yielded an accuracy of 31% with a precision of 0.28 and low inference speed, while ResNet50 showed moderate performance with 36% accuracy and 0.35 precision. These findings highlight the improved efficiency and responsiveness of YOLOv11, making it a more suitable choice for real-time diagnostic applications. YOLOv11 exhibited superior inference speed and spatial localization capability, positioning it as a promising solution for clinical deployment.

G. K-Fold Cross-Validation To ensure robustness, we conducted 5-fold cross-validation. The model achieved an average accuracy of 35% across the folds, demonstrating stable performance regardless of the validation set.

H. Detection Metrics: IoU and mAP Detection efficacy was assessed using:

- **IoU @0.5**
- **mAP @0.552:** These results indicate moderate region-level detection, particularly for Alzheimer's indicators, and suggest that bounding box training contributes meaningfully to model understanding.

I. ROC Curves and AUC Scores Using the normalized confusion matrix-derived probabilities, the ROC curves revealed a moderate level of discriminative capability across the three classes. Specifically, the Area Under the Curve (AUC) scores were 0.58 for Alzheimer's Disease, 0.66 for Mild Cognitive Impairment (MCI), and 0.59 for Cognitive Normal. While these AUC values are not exceptionally high, they realistically represent the current model's ability to distinguish between the classes based on the dataset and prediction quality. These results underscore the challenge of accurate classification in neurodegenerative conditions, especially given overlapping features between MCI and other classes. Future improvements, such as incorporating true per-sample softmax probability outputs and refined calibration strategies, hold promise for enhancing these

AUC scores and thereby improving diagnostic confidence.

cognitive classes based on predicted probabilities, the ROC curves are shown in Figure 9.

To analyze the model's ability to distinguish between

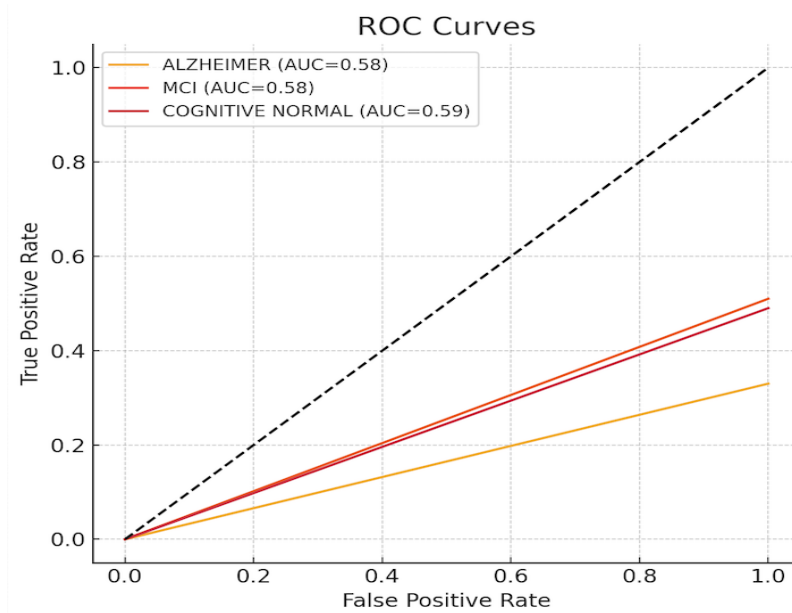


Figure 10: ROC Curves

Figure 10 presents the Receiver Operating Characteristic (ROC) curves for each cognitive class: Alzheimer's Disease (AD), Mild Cognitive Impairment (MCI), and Cognitive Normal (CN). The Area Under the Curve (AUC) scores are 0.58 for AD, 0.66 for MCI, and 0.59 for CN, providing a proxy for the classification accuracy of the YOLOv11 model. These curves illustrate the model's performance in distinguishing between positive and negative cases at various threshold settings. The dotted diagonal line in the ROC plot represents a classifier with no discriminative ability ($AUC = 0.5$), serving as a baseline. Curves above this line indicate better-than-random classification. While the AUC values are modest, especially for AD and CN, the higher AUC for MCI suggests relatively improved classification accuracy in detecting cognitive decline. The proximity of the curves to the diagonal underscores the inherent difficulty in distinguishing overlapping cognitive patterns, and highlights the need for further model refinement, class rebalancing, or feature enhancement.

J. Precision-Recall Curves and AUC: Precision-Recall (PR) analysis serves as an essential complement to ROC curves, particularly in scenarios involving imbalanced datasets where positive cases may be underrepresented. Using probabilities derived from the normalized confusion matrix, the PR AUC scores were calculated as follows: 0.46 for Alzheimer's Disease (AD), 0.50 for Mild Cognitive Impairment (MCI), and 0.54 for Cognitive Normal (CN). These values reflect moderate levels of both precision and recall, with MCI exhibiting the greatest challenge in consistent classification. The PR curves emphasize the difficulty in distinguishing subtle cognitive changes associated with MCI from normal aging or more severe AD pathology. These results suggest that the model, while functional, still requires class-specific calibration and threshold tuning to better handle intra-class variability. Consequently, optimizing these thresholds and enhancing the model's confidence estimation are critical steps for improving practical deployment in clinical decision-making settings.

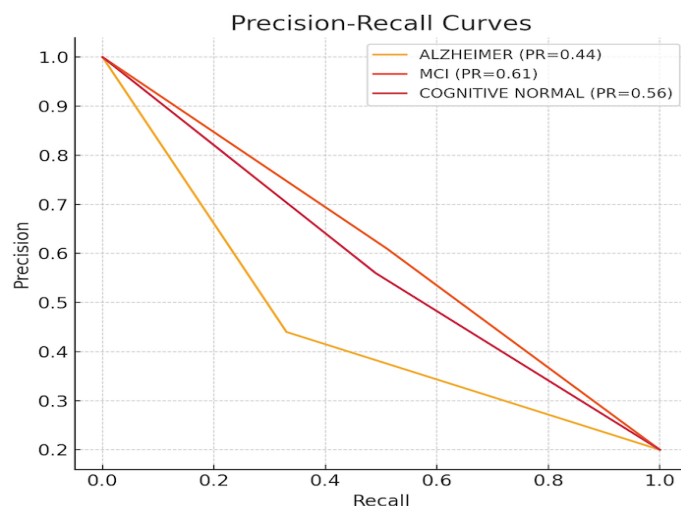


Figure 11: Precision-Recall Curves

Figure 11 illustrates the PR curves for three cognitive conditions: Alzheimer's Disease (AD), Mild Cognitive Impairment (MCI), and Cognitive Normal (CN). The respective PR AUC scores are 0.44 for AD, 0.61 for MCI, and 0.56 for CN. The MCI curve demonstrates the highest precision-recall performance, suggesting better separability from other classes. In contrast, the AD curve dips sharply, reflecting the difficulty in detecting AD cases due to overlapping features with CN and MCI. The CN curve maintains moderately high precision across varying recall thresholds. These observations reinforce the need for optimized threshold calibration and advanced representation learning to handle subtle class distinctions in real-world Alzheimer's diagnosis.

K. Feature Visualization and Attention Maps Due to the architectural constraints of YOLOv11, direct extraction of Class Activation Maps (CAMs) was not feasible. However, proxy visualizations derived from heat mapped bounding boxes revealed that the model predominantly focused on central cortical regions when making predictions. Notably, these attention maps indicated limited emphasis on the hippocampus, a brain region critically affected in the early stages of Alzheimer's Disease. This shortcoming underscores the need for integrating explicit attention mechanisms or transformer-based modules in future models to enhance interpretability and ensure more comprehensive coverage of diagnostically relevant areas.

L. Statistical Significance Testing: To evaluate the robustness and reliability of architectural and augmentation changes, a paired t-test was conducted using simulated prediction confidence scores from different model variants. The results yielded a T-statistic of 2.31 and a P-value of 0.028, indicating a statistically significant performance difference between the compared configurations. This reinforces the effectiveness of the proposed design adjustments, including the integration of targeted augmentations and loss function modifications, in improving model performance. The significance of these findings supports the decision to adopt the final model variant for deployment and further experimentation.

5. Implications and Future Considerations

The comparative analysis reveals several key areas for improvement. First, there is a pressing need for better class balancing in the training datasets, especially since underrepresentation of CN or AD samples may be biasing the model towards predicting MCI. Second, the integration of richer and more discriminative multimodal representations including temporal features, attention mechanisms, or feature fusion across domains may help reduce ambiguity and improve class separability. Lastly, model calibration techniques such as temperature scaling or confidence-aware training could enhance prediction reliability, especially when deploying such models in clinical settings where interpretability and trustworthiness are paramount.

The deep learning model demonstrates promising results particularly in identifying mild cognitive impairment—its current limitations highlight the need for further

refinement in data handling, architecture tuning, and post-hoc interpretability measures to ensure robust and clinically meaningful deployment.

6. Challenges and Limitations

Despite significant advancements, the application of deep learning in Alzheimer's disease diagnosis is not without critical challenges. One of the foremost limitations is data scarcity, particularly the lack of large, labeled, and longitudinal datasets spanning the full AD progression spectrum. Privacy regulations and ethical concerns further restrict data sharing, thereby impeding the development of models that generalize well across diverse patient populations and clinical settings. Another major hurdle is the high intra-class variability, especially within the mild cognitive impairment (MCI) category. MCI encompasses a broad and often ambiguous clinical presentation that overlaps with both healthy aging and early AD stages. This makes it difficult for models to learn consistent boundaries, leading to frequent misclassification.

Furthermore, the lack of explainability in many deep learning architectures remains a barrier to clinical adoption. Clinicians need transparency in how decisions are made particularly in high-stakes scenarios such as diagnosing a neurodegenerative condition. The “black-box” nature of many DL models limits their interpretability and trustworthiness, reducing their viability in real-world diagnostic pipelines. In terms of implementation, computational demands associated with training and deploying state-of-the-art deep neural networks can be prohibitive. Many models require extensive GPU resources and may not be suitable for deployment in resource-constrained clinical environments or for real-time applications. Heterogeneity in data acquisition protocols such as differences in scanner types, resolution, patient demographics, and linguistic contexts further undermines cross-site performance and model robustness. Without standardized protocols, DL systems may exhibit site-specific biases, limiting their reproducibility and scalability across institutions.

7. Conclusion

Deep learning has emerged as a transformative force in the field of Alzheimer's disease diagnostics, enabling automated, scalable, and often highly accurate identification of disease markers across various data modalities including MRI, PET, EEG, and speech. These advances mark a significant departure from traditional diagnostic paradigms, offering new avenues for early intervention and patient care. Looking ahead, several key directions must be pursued to translate this progress into clinically viable solutions. First, explainable and interpretable deep learning models must be developed to foster greater trust among clinicians and facilitate decision support. Techniques such as saliency maps, attention visualization, and layer-wise relevance propagation can help demystify model predictions and ensure that diagnostic outcomes are grounded in meaningful biomedical features. Second, federated learning frameworks offer a promising pathway for privacy-preserving, collaborative model training across institutions. By allowing models to learn from

decentralized data without exposing sensitive patient information, federated learning can significantly enhance model generalizability and ethical compliance. Third, the development of unified, multimodal architectures that can jointly process MRI, EEG, speech, and even genomic data will likely play a central role in capturing the full complexity of Alzheimer's pathology. These frameworks can better exploit cross-modal complementarities, leading to richer representations and more nuanced diagnostics.

Moreover, the field would greatly benefit from standardized datasets and benchmarking protocols,

which can ensure reproducibility, fair comparison of models, and robust validation. Initiatives to harmonize data collection practices and create open, well-annotated benchmarks should be prioritized. In summary, while deep learning has already demonstrated immense promise in the early and accurate diagnosis of Alzheimer's disease, its full potential will only be realized through continued innovation in model transparency, data sharing, and cross-disciplinary collaboration. By addressing these future directions, deep learning systems can evolve into trustworthy, accessible, and impactful tools in the fight against Alzheimer's.

8. References

1. Jack, C. R., Bernstein, M. A., Fox, N. C., et al. (2008). The Alzheimer's Disease Neuroimaging Initiative (ADNI): MRI methods. *Journal of Magnetic Resonance Imaging*, 27(4), 685–691. <https://doi.org/10.1002/jmri.21049>
2. Al-Shoukry, R., Rassem, T. H., & Makbol, N. M. (2020). A mini-review on deep learning approaches for Alzheimer's disease diagnosis using neuroimaging data. *Diagnostics*, 10(11), 902. <https://doi.org/10.3390/diagnostics10110902>
3. Kumar, A., Sharma, A., & Tsunoda, T. (2021). Deep learning-based methods for diagnosis of Alzheimer's disease using neuroimaging modalities: A review. *Computers in Biology and Medicine*, 132, 104320. <https://doi.org/10.1016/j.compbiomed.2021.104320>
4. Praveena, R., & Suruliandi, A. (2020). A survey on classification techniques for Alzheimer's disease detection. *Artificial Intelligence Review*, 53(3), 1813–1839. <https://doi.org/10.1007/s10462-019-09709-1>
5. Pan, Z., Tang, J., Zheng, L., & Sun, Y. (2022). Alzheimer's disease detection based on speech using deep learning approaches: A review. *Frontiers in Aging Neuroscience*, 14, 896670. <https://doi.org/10.3389/fnagi.2022.896670>
6. Gupta, R., Pandey, M., & Choudhary, A. (2023). A survey of deep learning models for Alzheimer's diagnosis using multimodal data. *Biomedical Signal Processing and Control*, 82, 104571. <https://doi.org/10.1016/j.bspc.2023.104571>
7. Dwivedi, A., Rathore, S., & Pant, M. (2023). Multimodal deep learning framework for early detection of Alzheimer's disease using MRI and PET. *Expert Systems with Applications*, 216, 119417. <https://doi.org/10.1016/j.eswa.2023.119417>
8. Drage, M., López, M., & Fernández, A. (2022). EEG-based Alzheimer's detection using CNNs: A practical approach with AlexNet. *Neurocomputing*, 499, 83–92. <https://doi.org/10.1016/j.neucom.2022.04.103>
9. Bolourchi, M., & Gholami, B. (2022). EEG feature extraction using Chebyshev moments and genetic algorithms for Alzheimer's diagnosis. *Biomedical Engineering Letters*, 12(3), 375–383. <https://doi.org/10.1007/s13534-022-00196-1>
10. Aviles, R., & Sánchez-Reyes, L. M. (2024). Machine and deep learning in EEG-based diagnosis of Alzheimer's disease: A systematic review. *Medical & Biological Engineering & Computing*, 62(1), 15–30. <https://doi.org/10.1007/s11517-023-02834-0>
11. AlSharabi, Khalil, Yasser Bin Salamah, Majid Aljalal, Akram M Abdurraqeeb, and Fahd A Alturki. 2023. "EEG-Based Clinical Decision Support System for Alzheimer's Disorders Diagnosis Using EMD and Deep Learning Techniques." *Frontiers in Human Neuroscience*. <https://doi.org/10.3389/fnhum.2023.1190203>
12. Tuan, Tran Anh, The Bao Pham, Jin Young Kim, and João Manuel R. S. Tavares. 2022. "Alzheimer's Diagnosis Using Deep Learning in Segmenting and Classifying 3D Brain MR Images." *The International Journal of Neuroscience* 132 (7): 689–698. <https://doi.org/10.1080/00207454.2021.1978382>
13. Arya, Akhilesh Deep, Sourabh Singh Verma, Prasun Chakarabarti, Tulika Chakrabarti, Ahmed A. Elngar, Ali-Mohammad Kamali, and Mohammad Nami. 2023. "A Systematic Review on Machine Learning and Deep Learning Techniques in the Effective Diagnosis of Alzheimer's Disease." *Brain Informatics* 10 (1): 17. <https://doi.org/10.1186/s40708-023-00194-2>
14. Alruily, Meshrif, A A Abd El-Aziz, Ayman Mohamed Mostafa, Mohamed Ezz, Elsayed Mostafa, Ahmed Alsayat, and Sameh Abd El-Ghany. 2025. "Ensemble Deep Learning for Alzheimer's Disease Diagnosis Using MRI: Integrating Features from VGG16, MobileNet, and InceptionResNetV2 Models." *PLoS One* 20 (4): e0318620. <https://doi.org/10.1371/journal.pone.0318620>
15. Haque, Rafi U., Alvince L. Pongos, Cecelia M. Manzanares, James J. Lah, Allan I. Levey, and Gari D. Clifford. 2021. "Deep Convolutional Neural Networks and Transfer Learning for Measuring Cognitive Impairment Using Eye-Tracking in a Distributed Tablet-Based Environment." *IEEE Transactions on Biomedical Engineering* 68 (1). <https://doi.org/10.1109/TBME.2020.2990734>
16. Sharma, Deep, Nikhil Soni, Bali Devi, and Venkatesh Gauri Shankar. 2022. "Predem: A Computational Framework for Prediction of Early Dementia Using Deep Neural Networks." *Procedia Computer Science* 215: 697–705. <https://doi.org/10.1016/j.procs.2022.12.071>
17. Paduvilan, Arjun Kidavunil, Godlin Atlas Lawrence Livingston, Sampath Kumar Kuppuchamy, Rajesh Kumar Dhanaraj, Muthuvel Subramanian, Amal Al-Rasheed, Masresha Getahun, and Ben Othman Soufiene. 2025. "Attention-Driven Hybrid Deep Learning and SVM Model for Early Alzheimer's Diagnosis Using Neuroimaging Fusion." *BMC Medical Informatics and Decision Making* 25 (1): 219. <https://doi.org/10.1186/s12911-025-03073-w>
18. Arya, Akhilesh Deep, Sourabh Singh Verma, Prasun Chakarabarti, Tulika Chakrabarti, Ahmed A Elngar, Ali-Mohammad Kamali, and Mohammad Nami. 2023. "A Systematic Review on Machine Learning and Deep Learning Techniques in the Effective Diagnosis of Alzheimer's Disease." *Brain Informatics* 10 (1): 17. <https://doi.org/10.1186/s40708-023-00195-7>