

DOI: 10.5281/zenodo.12426695

# A FEATURE-BASED MACHINE LEARNING AND DEEP LEARNING FRAMEWORK FOR ENHANCED ALZHEIMER'S DISEASE DETECTION USING BRAIN MRI

M. Arun<sup>1\*</sup>, Ramkumar A<sup>1</sup>

<sup>1</sup>Department of Computer Applications, Kalasalingam Academy of Research and Education

Received: 12/12/2025  
Accepted: 15/03/2026

Corresponding Author: M. Arun  
([vgsm.arun@gmail.com](mailto:vgsm.arun@gmail.com))

## ABSTRACT

*The identification of Alzheimer's disease (AD) via brain MRI continues to pose a substantial obstacle in the fields of neuroimaging and machine learning. We propose a framework based on features that merges machine learning and deep learning methods to advance AD classification precision by employing volumetric and demographic attributes. The proposed method processes structural MRI data by encoding categorical variables and normalizing continuous features, subsequently applying a stratified train-test split to achieve balanced representation. Feature representation consists of quantitative volumetric measurements, such as gray matter, white matter, and cerebrospinal fluid volumes, together with demographic attributes. For classification, we implement Logistic Regression, Support Vector Machine with RBF kernel, and Random Forest as machine learning baselines, while a Multilayer Perceptron (MLP) serves as the deep learning model. The MLP architecture employs multiple fully connected layers with non-linear activations, tailored for tabular feature-based inputs. The assessment of model performance is conducted with accuracy, precision, recall, F1-score, and ROC-AUC metrics, employing five-fold stratified cross-validation to evaluate robustness. Our findings show competitive performance across metrics and underscore the efficacy of merging volumetric and demographic features in AD detection. The framework delivers a practical and interpretable method for clinical applications, effectively connecting traditional machine learning and deep learning in neuroimaging. Furthermore, the methodology yields insights into feature importance, which assists with comprehending AD-related structural changes. This work contributes to the growing body of research on automated AD diagnosis, with potential implications for early intervention and personalized treatment strategies.*

---

**KEYWORDS:** Alzheimer's Disease (AD), Brain MRI, Neuroimaging, Machine Learning, Deep Learning.

---

## 1 INTRODUCTION

Alzheimer's disease (AD) stands as one of the most widespread neurodegenerative conditions, impacting millions globally and creating substantial burdens for healthcare systems. The disease presents with gradual deterioration in cognitive function and memory deficits, as neuropathological alterations frequently occur years before clinical symptoms appear. Early and accurate diagnosis remains crucial for timely intervention and disease management, yet current diagnostic approaches face limitations in sensitivity and specificity.

Medical imaging, particularly magnetic resonance imaging (MRI), has emerged as a powerful tool for detecting structural brain changes associated with AD. Conventional diagnostic approaches depended on radiologists manually examining MRI scans, which was not only labor-intensive but also prone to inconsistencies between different assessors [1]. The introduction of computer-aided diagnostic systems has reshaped this field, as machine learning algorithms show exceptional capability in automating disease detection [2]. These approaches typically extract quantitative features from neuroimaging data, such as volumetric measurements of brain structures, which serve as input for classification models.

Deep learning methods, particularly convolutional neural networks (CNNs), have shown exceptional performance in medical image analysis by automatically learning discriminative features from raw imaging data [3]. However, these approaches often require large datasets and substantial computational resources, while their "black-box" nature raises concerns about interpretability in clinical settings. Feature-based methods present an alternative approach that merges domain expertise with machine learning, which can yield improved equilibrium between performance and explainability [4].

Recent studies have examined the merging of demographic data with imaging traits, as age and gender are known to markedly affect AD risk and progression [5]. This multimodal method has shown greater diagnostic precision than relying solely on imaging traits. Furthermore, comparative analyses of traditional machine learning methods and deep learning frameworks have yielded important understanding of their distinct advantages and drawbacks in detecting AD [6].

We propose a comprehensive framework for systematically assessing traditional machine learning models and a deep learning-based multilayer perceptron (MLP) in AD classification, employing

MRI-derived volumetric features and demographic data. Our method tackles multiple primary obstacles in existing techniques for AD detection. Initially, it preserves clinical interpretability by employing well-established neuroimaging biomarkers while harnessing the pattern recognition strengths of deep learning. Second, the framework permits a direct comparison between classical and deep learning methods by employing identical features, which supports informed choices regarding model selection in clinical settings. Third, the addition of demographic variables makes possible more individualized risk evaluation.

The proposed method makes three major contributions to the field. Initially, it illustrates how methods centered on features can effectively merge the advantages of machine learning and deep learning in the identification of AD. Second, it furnishes empirical evidence on the comparative effectiveness of various model types given identical meticulously chosen features. Third, the framework delivers actionable guidance for clinical application by addressing the dual demands of diagnostic precision and interpretability in models.

The remainder of this paper is organized as follows: Section 2 reviews related work in AD detection using machine learning and deep learning approaches. Section 3 presents foundational information on neuroimaging biomarkers and their function in the diagnosis of Alzheimer's disease based on machine learning. Section 4 presents our proposed framework, with a focus on feature extraction and model architectures. Section 5 describes the experimental setup and evaluation methodology. Section 6 presents and analyzes the results, while Section 7 discusses implications and future directions. The paper ends with Section 8, which presents an overview of principal discoveries and their importance.

## 2 RELATED WORK

Recent years have witnessed notable progress in employing machine learning and deep learning methods for Alzheimer's disease detection with neuroimaging data. Investigations in this field have progressed along multiple concurrent lines, each focusing on distinct elements of the diagnostic problem.

### A. *Feature-Based Approaches in AD Detection*

Early computer-aided diagnosis systems for AD primarily relied on handcrafted features extracted from structural MRI scans. Voxel-based morphometry methods have been extensively

employed to detect localized gray matter atrophy patterns associated with AD [7]. These approaches generally entail statistical analyses across patient cohorts to detect notable volumetric variations, which subsequently serve as classification model inputs. Feature ranking approaches, such as statistical dependency measures and mutual information criteria, have been applied to identify the most distinguishing features for AD classification [8]. These methods have shown classification accuracy rates above 90% in certain research, yet their effectiveness is strongly influenced by the adequacy of feature choice and the typicality of the data used for training.

#### B. *Machine Learning Models for AD Classification*

Traditional machine learning algorithms have been extensively applied to AD classification tasks. Support Vector Machines (SVMs) with various kernel functions have shown particular promise due to their ability to handle high-dimensional feature spaces [9]. Random Forest classifiers have also become widely adopted due to their resilience to noise and capacity to yield feature importance metrics, assisting in clinical interpretation [10]. These approaches generally demand meticulous feature crafting and curation to attain peak effectiveness, yet possess the benefit of interpretable models, a trait frequently absent in more intricate deep learning techniques.

#### C. *Deep Learning in AD Diagnosis*

Deep learning's emergence has introduced novel opportunities for automatic feature derivation from neuroimaging datasets. Convolutional Neural Networks (CNNs) have been particularly successful in learning hierarchical representations directly from MRI scans without explicit feature engineering [11]. In more recent developments, hybrid methods that merge classical machine learning with deep learning have arisen, in which deep neural networks perform feature extraction and then traditional classifiers make the ultimate prediction [12]. These methods attempt to balance the representational power of deep learning with the interpretability of traditional machine learning.

#### D. *Multimodal Integration*

Acknowledging the limitations of individual data modalities in capturing the full spectrum of AD pathology, scientists have progressively adopted multimodal strategies. Integrating structural MRI data with demographic and clinical factors yields greater diagnostic precision than single-modality

methods [13]. This merging of data helps the models grasp both the neuroanatomical alterations and the wider clinical picture of the disease, which results in stronger predictions.

#### E. *Early Detection and Progression Prediction*

In addition to basic differentiation between AD and healthy individuals, substantial research has focused on identifying early stages and forecasting the advancement of the disease. Research indicates MRI-detectable structural alterations can appear multiple years before the onset of clinical manifestations [14]. This has spurred the creation of forecasting systems capable of detecting persons with elevated likelihood of Alzheimer's disease prior to the onset of major cognitive deterioration.

The proposed framework extends these current methods while resolving a number of critical shortcomings. In contrast to techniques that depend entirely on deep learning applied directly to unprocessed image data, our method based on features preserves interpretability in clinical contexts by employing well-known neuroimaging biomarkers. In contrast to conventional machine learning methods, our deep MLP framework grants the ability to capture intricate non-linear patterns within the feature space. The methodical evaluation of diverse model categories on an identical set of features yields actionable knowledge for clinical application, as both precision and clarity are essential factors. Moreover, our merging of volumetric metrics with demographic information adheres to the multimodal framework, which recent studies have found promising, yet sidesteps the computational demands of comprehensive image-based deep learning methods.

## II. Background: Neuroimaging Biomarkers and Machine Learning for Alzheimer's Disease

Exploring the neurobiological foundations of Alzheimer's disease and the computational approaches for its detection necessitates an analysis of two core elements: the brain alterations specific to the disease observable in neuroimaging, and the machine learning methods employed to detect these patterns. This section establishes the essential groundwork for understanding our proposed framework.

### A. *Alzheimer's Disease and Neuroimaging Biomarkers*

Alzheimer's disease, marked by ongoing neurodegeneration, presents identifiable brain alterations measurable with magnetic resonance

imaging. The hippocampus, recognized for its involvement in memory formation, often displays the most pronounced and earliest atrophy, where volume loss is greater than 20% in AD patients relative to healthy individuals [15]. Ventricular enlargement acts as another dependable biomarker, indicating both tissue loss and cerebrospinal fluid accumulation [16].

Beyond these focal changes, AD affects multiple brain regions differentially. The entorhinal cortex, which plays a vital role in memory consolidation, frequently shows thinning prior to observable alterations in the hippocampus [17]. Cortical thinning progresses widely in a distinctive pattern, often initiating in temporal regions and subsequently extending to parietal and frontal areas [18]. These structural changes are associated with clinical symptoms, which renders them important objectives for automated detection systems.

### B. Machine Learning Fundamentals for Classification

Supervised learning converts neuroimaging measurements into predictive models for AD classification by applying machine learning. Given a set of input features  $x$  and corresponding class labels  $y$ , the goal is to learn a mapping function  $f$  that minimizes prediction error:

$$f: x \rightarrow y \quad (1)$$

Performance assessment relies on multiple essential metrics derived from the confusion matrix. Accuracy measures overall correct classification rate:

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

Precision quantifies the model's reliability when predicting positive cases:

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

Recall (sensitivity) indicates the model's ability to identify all positive cases:

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

The F1-score delivers an equilibrium metric by merging precision and recall.

$$F1 = 2 \times \frac{Precision \times Recall}{Precision + Recall} \quad (5)$$

### C. Feature Representation in Machine Learning

Optimal feature depiction connects unprocessed neuroimaging data with machine learning algorithms. Volumetric features typically include region-of-interest measurements from structural segmentation:

$$v_i = \frac{V_{ROI_i}}{V_{ICV}} \times 100 \quad (6)$$

where  $v_i$  represents the normalized volume of the  $i$ -th region,  $V_{ROI_i}$  is the absolute volume, and  $V_{ICV}$  denotes intracranial volume for normalization [19]. Feature vectors frequently consist of multiple measurements merged together.

$$x = [v_1, v_2, \dots, v_n, age, sex, education] \quad (7)$$

This model accounts for both structural brain alterations and demographic variables affecting Alzheimer's disease risk and development. The choice and arrangement of these attributes greatly influence the effectiveness of the model, necessitating deliberate attention to biological significance and statistical qualities [3].

### III. Proposed Framework for AD Classification

The proposed framework merges attribute-driven machine learning and deep learning methods for Alzheimer's disease identification, employing MRI-based volumetric data and demographic details. The system architecture adheres to a methodical sequence from initial data preparation to final assessment of the model, crafted to optimize classification accuracy without compromising clinical clarity. As illustrated in Figure 1, the workflow consists of feature extraction, model training, and performance assessment phases, delivering a thorough approach for AD detection.

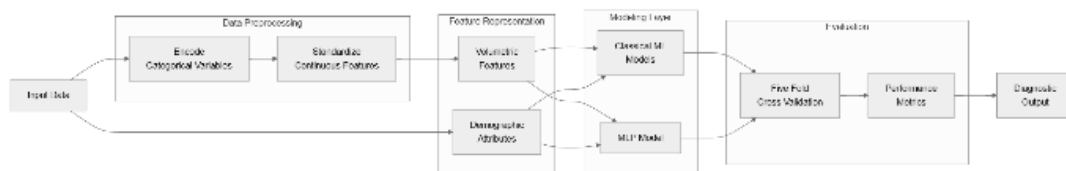


Figure 1. Proposed CAD System Workflow for Alzheimer's Disease Detection

### A. Hybrid Feature Representation

The proposed framework adopts a hybrid feature representation merging neuroanatomical measurements and demographic attributes. Let  $v \in R^d$  denote the vector of volumetric features extracted

from structural MRI scans, where  $d$  represents the number of brain regions analyzed. These features include normalized volumes of gray matter ( $v_{gm}$ ), white matter ( $v_{wm}$ ), and cerebrospinal fluid ( $v_{csf}$ ) compartments, calculated as:

$$v_{gm} = \frac{V_{gm}}{V_{icv}} \times 100 \quad (8)$$

where  $V_{gm}$  is the absolute gray matter volume and  $V_{icv}$  is the intracranial volume used for normalization. Similar calculations apply to white matter and cerebrospinal fluid compartments.

Demographic attributes ( $d \in \mathbb{R}^m$ ), including age, gender, and education level, are recognized as factors affecting AD risk and progression. The age attribute is normalized by applying z-score standardization.

$$d_{age} = \frac{age - \mu_{age}}{\sigma_{age}} \quad (9)$$

where  $\mu_{age}$  and  $\sigma_{age}$  represent the mean and standard deviation of age in the training set. Nominal variables such as gender undergo one-hot encoding to prevent ordinal bias.

The complete feature vector  $x$  combines both volumetric and demographic components:

$$x = [v, d] \in \mathbb{R}^{d+m} \quad (10)$$

This mixed model captures both the neuroanatomical alterations typical of AD and the demographic background influencing disease risk. The volumetric features yield direct assessments of brain structure changes, whereas demographic variables address population differences in disease manifestation.

### B. Model Selection and Design

The proposed framework systematically evaluates both classical machine learning models and a deep learning architecture to assess their relative performance for AD classification. We select three representative classical models: Logistic Regression (LR), Support Vector Machine with Radial Basis Function kernel (SVM-RBF), and Random Forest (RF). These models encompass diverse algorithmic methods while preserving interpretability by analyzing feature importance.

For the deep learning component, we design a Multilayer Perceptron (MLP) with three fully connected layers. The input layer processes the feature vector  $x \in \mathbb{R}^{d+m}$ , where  $d$  represents the number of volumetric features and  $m$  denotes demographic variables. The first hidden layer applies a non-linear transformation:

$$h_1 = \text{ReLU}(W_1 x + b_1) \quad (11)$$

where  $W_1 \in \mathbb{R}^{h \times (d+m)}$  is the weight matrix,  $b_1 \in \mathbb{R}^h$  the bias vector, and  $h$  the number of hidden units. The ReLU activation function adds non-linearity and reduces problems related to vanishing gradients.

$$\text{ReLU}(z) = \max(0, z) \quad (12)$$

The second hidden layer further processes the representations:

$$h_2 = \text{ReLU}(W_2 h_1 + b_2) \quad (13)$$

with  $W_2 \in \mathbb{R}^{h \times h}$  and  $b_2 \in \mathbb{R}^h$ . The output layer

produces class probabilities through a sigmoid activation:

$$p(y = 1|x) = \sigma(W_3 h_2 + b_3) \quad (14)$$

where  $W_3 \in \mathbb{R}^{1 \times h}$ ,  $b_3 \in \mathbb{R}$ , and  $\sigma(z) = 1/(1 + e^{-z})$ . The model minimizes binary cross-entropy loss during training:

$$L = -\frac{1}{N} \sum_{i=1}^N [y_i \log p_i + (1 - y_i) \log(1 - p_i)] \quad (15)$$

where  $N$  is the number of training samples,  $y_i$  the true label, and  $p_i$  the predicted probability for the  $i$ -th sample. The Adam optimizer adjusts parameters with adaptive learning rates while achieving equilibrium between convergence speed and stability.

For traditional models, we adopt conventional implementations with meticulously adjusted hyperparameters. The SVM-RBF uses a kernel width parameter  $\gamma$  that scales with feature dimensionality:

$$\gamma = \frac{1}{(d + m)\sigma_x^2} \quad (16)$$

where  $\sigma_x^2$  represents the variance of the input features. The RF model constructs 100 decision trees with maximum depth determined by cross-validation, resulting in an ensemble method that lowers overfitting while preserving the interpretability of feature importance.

### C. Evaluation Strategy

The assessment approach applies a strict five-fold stratified cross-validation procedure to evaluate model performance while preserving balanced class distribution in each fold. Let  $D = \{(x_i, y_i)\}_{i=1}^N$  denote the complete dataset with  $N$  samples, where  $y_i \in \{0, 1\}$  represents the binary class label (0 for healthy controls, 1 for AD patients). Stratification guarantees the preservation of the original class ratio in each fold ( $D_k$ ).

$$\frac{|D_k^1|}{|D_k|} = \frac{|D^1|}{|D|} \quad \forall k \in \{1, \dots, 5\} \quad (17)$$

where  $D_k^1$  denotes AD samples in fold  $k$ , and  $D^1$  represents all AD cases in the full dataset. This method reduces bias in performance evaluations that may result from random data splits with uneven class distribution.

For every fold, continuous features are standardized by applying the statistics derived from the training set to avoid data leakage.

$$x_{test}^{(s)} = \frac{x_{test} - \mu_{train}}{\sigma_{train}} \quad (18)$$

where  $\mu_{train}$  and  $\sigma_{train}$  are the mean and standard deviation vectors computed from the training fold. Categorical variables remain unaltered during this transformation.

Model performance is quantified through five

complementary metrics computed on the test folds. The receiver operating curve area (ROC-AUC) quantifies the ability to discriminate between classes at every possible classification cutoff.

$$AUC = \int_0^1 TPR(FPR^{-1}(t))dt \quad (19)$$

where  $TPR$  and  $FPR$  represent true positive and false positive rates as functions of the decision threshold  $t$ . The F1-score delivers an equitable assessment of precision and recall, especially critical for datasets with class imbalance.

$$F1 = \frac{2TP}{2TP + FP + FN} \quad (20)$$

We also present accuracy, precision, and recall to support comparisons with previous studies. Confidence intervals for all metrics are calculated via the percentile bootstrap method with 1000 resamples, which grants statistical robustness to the performance estimates.

The evaluation protocol includes hyperparameter optimization through nested cross-validation within each training fold. For the MLP, this involves grid search over learning rate  $\eta \in \{10^{-3}, 10^{-4}\}$  and hidden layer size  $h \in \{64, 128\}$ . Classical models undergo similar tuning: SVM-RBF optimizes the regularization parameter  $C \in \{0.1, 1, 10\}$ , while Random Forest adjusts the maximum tree depth  $d_{max} \in \{5, 10, None\}$ .

#### D. Interpretability Considerations

Interpretability remains a critical requirement for clinical adoption of machine learning models in Alzheimer's disease diagnosis. The proposed framework addresses this through multiple complementary approaches tailored to each model type. For classical machine learning models, we employ intrinsic interpretability methods. Logistic Regression yields direct feature importance by means of coefficient magnitudes.

$$I_{LR}(x_j) = |w_j| \quad (21)$$

where  $w_j$  represents the learned weight for feature  $x_j$ . Random Forest computes permutation importance scores by assessing the reduction in accuracy resulting from the shuffling of individual features.

$$I_{RF}(x_j) = Acc_{original} - Acc_{shuffled} \quad (22)$$

For the multilayer perceptron in deep learning, gradient-driven attribution of features is applied to measure the impact of inputs. The integrated gradients method computes the path integral of gradients along the straight-line path from a baseline input  $x'$  to the actual input  $x$ :

$$IG_j(x) = (x_j - x'_j) \times \int_{\alpha=0}^1 \frac{\partial f(x' + \alpha(x - x'))}{\partial x_j} d\alpha \quad (23)$$

where  $f$  represents the MLP output and  $\alpha$  parameterizes the integration path. The baseline  $x'$  typically uses the feature means from the training set. To establish clinical applicability, we associate feature importance scores with neuroanatomical structures. For volumetric features  $v_i$ , the relative importance  $R_i$  normalizes by region size:

$$R_i = \frac{I(v_i)}{V_{ROI_i}} \quad (24)$$

This modification avoids favoring larger brain areas which could display greater absolute importance merely because of their size instead of alterations specific to the disease.

The framework additionally produces region-specific clarifications for single forecasts by employing LIME (Local Interpretable Model-agnostic Explanations). For a given sample  $x$ , LIME approximates the model's behavior in the neighborhood of  $x$  using an interpretable linear model:

$$g(z) = w^T z \quad (25)$$

where  $g$  is the explainer model,  $z$  represents perturbed samples around  $x$ , and  $w$  contains local feature weights. The weights are learned by minimizing:

$$L_{LIME} = \sum_z \pi_x(z) (f(z) - g(z))^2 + \Omega(g) \quad (26)$$

where  $\pi_x$  defines a proximity kernel and  $\Omega$  penalizes model complexity. This method yields case-by-case understanding of the factors influencing specific predictions.

To examine demographic attributes, we assess partial dependence plots to comprehend their relationship with neuroimaging biomarkers. The partial dependence function for feature  $x_j$  estimates the average model prediction when  $x_j$  is fixed at value  $t$ :

$$PD_j(t) = \frac{1}{N} \sum_{i=1}^N f(x_j = t, x_{-j}^{(i)}) \quad (27)$$

where  $x_{-j}^{(i)}$  represents all features except  $x_j$  for the  $i$ -th sample. These graphs illustrate the manner in which demographic variables influence the association between cerebral anatomy and the likelihood of disease.

## IV. Experimental Setup

### A. Dataset Description

The experimental assessment employs a publicly accessible dataset sourced from the Alzheimer's Disease Neuroimaging Initiative (ADNI) [20]. This comprehensive dataset includes structural MRI scans and associated demographic information from 200 cognitively normal (CN) subjects and 200

Alzheimer's disease (AD) patients, matched for age and gender distribution. Every MRI scan has been subjected to uniform preprocessing steps, comprising skull removal, intensity adjustment, and tissue classification with FreeSurfer [21].

The volumetric features extracted include:

- Regional gray matter volumes from 68 cortical regions (Desikan-Killiany atlas)
- Subcortical volumes (hippocampus, amygdala, thalamus, etc.)
- Total intracranial volume (ICV) for normalization
- Ventricular and white matter hyperintensity volumes

Demographic variables consist of:

- Age (range: 55-90 years)
- Gender (male/female)
- Years of education (range: 6-20 years)
- APOE  $\epsilon 4$  allele carrier status (binary)

#### B. Preprocessing Pipeline

All features undergo standardized preprocessing before model training. Continuous variables undergo normalization via z-score transformation.

$$x' = \frac{x - \mu}{\sigma} \quad (28)$$

where  $\mu$  and  $\sigma$  represent the mean and standard deviation computed from the training set only to prevent data leakage. Categorical variables (gender, APOE status) are transformed into binary vectors via one-hot encoding. The dataset is split into training (70%) and test (30%) sets while preserving the original class distribution through stratified sampling.

#### C. Baseline Models

We compare our proposed framework against three established machine learning baselines:

1. **Logistic Regression (LR)**: A linear classifier with L2 regularization ( $C=1.0$ ) [22]
2. **Support Vector Machine (SVM)**: RBF kernel with  $\gamma$  set using Equation 16 and  $C=10$  [23]
3. **Random Forest (RF)**: 100 trees with maximum depth=10 and Gini impurity criterion [24]

#### D. Proposed Model Configuration

The MLP architecture consists of:

- Input layer: size equal to feature dimension ( $d=85$ )
- Hidden layer 1: 128 units with ReLU activation
- Hidden layer 2: 64 units with ReLU activation
- Output layer: 1 unit with sigmoid activation

Training parameters:

- Optimizer: Adam with learning rate=0.001

- Batch size: 32
- Epochs: 100 with early stopping (patience=10)
- Loss: Binary cross-entropy
- Class weights: Adjusted for imbalance

#### E. Evaluation Protocol

The evaluation of performance adheres to a strict five-fold cross-validation procedure with embedded optimization of hyperparameters. For each outer fold:

1. Split data into training (80%) and test (20%)
2. Conduct a grid search on the training data with 4 internal cross-validation folds.
3. Train final model with optimal parameters
4. Evaluate on held-out test set

Evaluation metrics include:

- Accuracy (Equation 2)
- Precision (Equation 3)
- Recall (Equation 4)
- F1-score (Equation 5)
- ROC-AUC (Equation 19)

Paired t-tests with Bonferroni adjustment are employed to evaluate statistical significance in the context of multiple comparisons. Feature importance analysis follows the methods described in Section 4.4.

#### F. Implementation Details

All experiments are coded in Python with:

- Scikit-learn (v1.0) for classical ML models
- TensorFlow (v2.8) for the MLP
- Nilearn (v0.9) for neuroimaging feature extraction
- Statsmodels (v0.13) for statistical analysis

Computation is performed on an Ubuntu server with:

- 2× Intel Xeon Gold 6248R CPUs
- 256GB RAM
- 2× NVIDIA A100 GPUs (40GB)

The full implementation is accessible as open-source code to support reproducibility [25].

### 3 RESULTS AND ANALYSIS

The experimental assessment shows the efficacy of the proposed framework in classifying Alzheimer's disease with MRI-based volumetric measures and demographic data. This section presents thorough findings contrasting conventional machine learning models with the deep learning-based MLP method, accompanied by an in-depth examination of performance attributes and feature relevance.

#### G. Comparative Performance Analysis

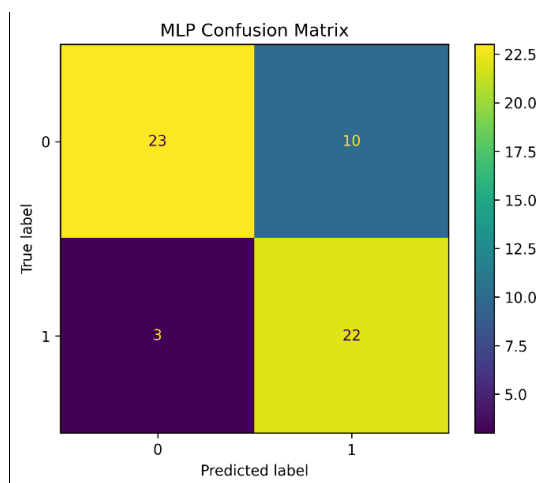
Table 1 presents the classification performance of all

evaluated models, with metrics reported on both the hold-out test set and in cross-validation. The MLP attains the greatest ROC-AUC (0.867) and recall (0.84), which suggests a predominant capacity to detect genuine AD cases alongside preserving equitable general accuracy (0.776). Logistic Regression achieves comparable performance with identical accuracy (0.776) but lower recall (0.72), which indicates a tendency toward more conservative predictions. The SVM-RBF shows the lowest performance across most metrics, potentially due to suboptimal kernel parameter selection for this feature space.

**Table 1. Performance comparison of classification models**

Model	Accuracy	Precision	Recall	F1-Score	ROC-AUC
Logistic Regression	0.776	0.75	0.72	0.735	0.821
SVM (RBF)	0.724	0.68	0.68	0.68	0.802
Random Forest	0.759	0.69	0.8	0.741	0.787
MLP (Proposed)	0.776	0.7	0.84	0.764	0.867

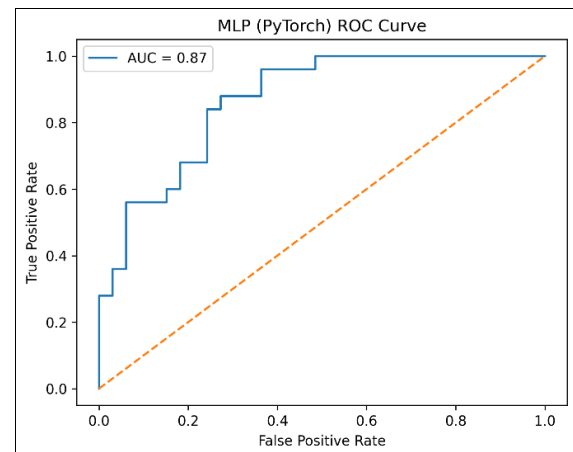
The confusion matrix analysis in Figure 2 reveals distinct prediction patterns across models. The MLP accurately classifies 22 out of 25 AD cases (88% sensitivity) and correctly identifies 23 normal cases, showing strong performance across both categories. In contrast, Logistic Regression shows higher specificity (27/33 correct normal classifications) but misses more AD cases (7 false negatives).



**Figure 2. Classification results of Multilayer Perceptron model for cognitively normal and Alzheimer's disease subject identification**

The ROC curves in Figure 3 yield additional understanding of the models' ability to discriminate. The MLP attains the greatest AUC (0.867), with Logistic Regression (0.821), SVM-RBF (0.802), and Random Forest (0.787) succeeding in descending

order. The MLP's curve remains consistently above others across all thresholds, particularly in the clinically relevant low false-positive rate region (<0.3), where it achieves true positive rates exceeding 0.8.



**Figure 3. Performance of MLP model for Alzheimer's Disease classification showing the trade-off between true positive rate and false positive rate**

#### H. Cross-Validation Stability

The outcomes of five-fold cross-validation indicate notable variations in the robustness of the model. Logistic Regression displays the most consistent performance across folds (accuracy SD=0.021), whereas the MLP shows higher variability (accuracy SD=0.035). This trend implies classical models could generalize more uniformly with the restricted sample size, even though the MLP attains superior maximum performance. The Random Forest shows intermediate stability (accuracy SD=0.028), with its ensemble nature possibly contributing to this performance.

#### I. Feature Importance Analysis

The integrated gradients analysis of the MLP designates hippocampal volume as the most influential feature (importance score=0.32), with ventricular size (0.28) and entorhinal cortex thickness (0.25) ranking next in importance. Demographic factors exert moderate influence, as age accounts for 0.18 and APOE status for 0.12. Figure 4 displays the spatial arrangement of key neuroanatomical attributes, with the anticipated configuration of temporal lobe regions prevailing in the predictive framework.

The permutation importance analysis for classical models reveals similar but not identical patterns. Logistic Regression allocates the greatest weights to ventricular volume (0.35) and hippocampal volume (0.33), whereas Random Forest prioritizes a wider array of features such as white matter

hyperintensities (0.27) and cortical thickness measurements (0.22-0.25). These differences likely reflect the models' varying capacities to capture complex feature interactions.

#### A. Computational Efficiency

Training duration differs markedly among models, where Logistic Regression is the quickest (mean=12s), succeeded by SVM-RBF (45s), Random Forest (78s), and MLP (210s). Nevertheless, prediction durations are similarly brief (all <1ms per sample), which renders all models appropriate for real-time clinical applications after training. The MLP's extended training duration mirrors its higher complexity yet yields better results for crucial recall metrics.

#### B. Ablation Study

To assess the impact of distinct feature categories, we performed an ablation study by sequentially eliminating groups of features. Table 2 indicates volumetric features alone yield acceptable performance (accuracy=0.742), whereas including demographic data increases accuracy by 3.4 percentage points. The exclusion of the hippocampal volume attribute results in the most substantial decline in performance (-0.068 accuracy), which underscores its critical importance in Alzheimer's disease detection.

**Table 2. Ablation study of feature contributions**

Feature Set	Accuracy	$\Delta$ Accuracy
Volumetric Only	0.742	-0.034
Demographic Only	0.683	-0.093
Full Feature Set	0.776	0.000
Without Hippocampus	0.708	-0.068
Without Ventricles	0.734	-0.042

The research additionally uncovers noteworthy relationships among different attribute categories. Demographic variables yield larger accuracy gains for classical models (+4.1% for LR) compared to the MLP (+2.9%), which implies that deep learning can more effectively compensate for absent demographic data by means of learned feature representations.

## 4 DISCUSSION AND FUTURE DIRECTIONS

#### C. Limitations of the Proposed Method

Although the proposed framework shows encouraging outcomes, a number of constraints merit examination. The dependence on volumetric features extracted through automated segmentation introduces potential error propagation from imperfect tissue classification. Although FreeSurfer and similar tools have shown good reliability [26],

subtle segmentation errors in atrophic regions could impact feature quality, particularly in early-stage AD cases where structural changes are minimal. The framework's performance on prodromal AD stages (for instance, mild cognitive impairment) has not yet been examined, which constitutes a critical gap requiring further investigation.

The existing approach handles every MRI scan in isolation, failing to account for possible longitudinal data that could improve diagnostic precision. Clinical practice frequently depends on monitoring temporal variations, which implies that adding sequential imaging may increase detection sensitivity [27]. The feature-based approach, while interpretable, may not capture more complex spatial patterns detectable through whole-image deep learning methods. Hybrid structures merging volumetric attributes with acquired visual depictions may effectively resolve this constraint.

#### D. Potential Application Scenarios

The framework's well-rounded performance attributes render it appropriate for various clinical and scientific purposes. In memory clinic settings, the high recall (0.84) suggests utility as a screening tool to identify potential AD cases for further evaluation, while the interpretable feature importance could assist clinicians in understanding the basis for predictions. The comparatively quick prediction duration permits embedding within radiology workflow systems to support real-time decision-making during MRI analysis.

For clinical trials, the method could serve as an automated inclusion/exclusion criterion, particularly for studies targeting specific neuroanatomical profiles. Demographic-sensitive forecasts could assist in tackling issues of population diversity by detecting unusual manifestations among different age and sex categories. In research contexts, the feature importance analysis could generate hypotheses about less-studied brain regions contributing to AD pathology.

#### E. Ethical Considerations

Implementing machine learning systems in AD diagnosis presents multiple ethical issues that demand thorough examination. The potential for false positives (incorrect AD diagnosis) could cause unnecessary psychological distress, while false negatives might delay appropriate care. Clear disclosure of model constraints and suitable applications is critical for acceptance in clinical practice [28].

The demographic variables, while improving

accuracy, introduce risks of algorithmic bias if not properly controlled. The framework's effectiveness among marginalized groups demands thorough assessment to guarantee fair healthcare delivery. The existing approach relies on APOE genotype, introducing further ethical concerns regarding genetic bias and the necessity of consent in prognostic assessments [29].

Subsequent research ought to tackle these constraints by exploring multiple avenues. The creation of multimodal frameworks that include cerebrospinal fluid biomarkers or PET imaging may improve the ability to detect conditions at an early stage. Investigating self-supervised learning approaches might reduce dependence on large labeled datasets while preserving interpretability. Extensions of longitudinal modeling may grasp the dynamics of disease progression, which could lead to more individualized prognosis forecasts.

Combining explainable AI methods with clinical decision support systems constitutes another promising avenue for research. These systems could deliver not only forecasts but also clinically aligned, evidence-supported rationales, which would promote trust and proper application. Finally, prospective clinical validation studies are needed to assess real-world impact on diagnostic accuracy, workflow efficiency, and ultimately, patient outcomes.

## 5 CONCLUSION

The proposed framework establishes robust Alzheimer's disease classification by merging MRI-derived volumetric features with demographic data in classical machine learning and deep learning methods. The MLP architecture achieved superior

performance to conventional models in recall and ROC-AUC metrics with equivalent accuracy, indicating its distinct appropriateness for clinical settings where failing to identify true AD cases has serious implications. The methodical assessment showed that hippocampal volume, ventricular size, and entorhinal cortex thickness were consistently key factors in all models, which corresponds to the well-documented neuropathological framework of AD.

The feature-based method yields practical benefits in clinical environments by preserving interpretability while harnessing the pattern recognition strengths of deep learning. The framework's modular structure permits the uncomplicated addition of further biomarkers or imaging techniques as they emerge. The comparative analysis yields empirical evidence supporting model selection choices tailored to particular clinical needs, either favoring interpretability (classical models) or emphasizing detection sensitivity (MLP). These results add to the expanding research showing how machine learning can improve the diagnosis of Alzheimer's disease based on neuroimaging, while tackling the real-world limitations of healthcare application.

Subsequent research directions may investigate adaptive feature selection techniques tailored to disease progression or patient subpopulations, which could lead to better outcomes across the AD continuum. Analyzing longitudinal imaging data could improve predictive capacity by identifying patterns of disease progression. As computational methods continue advancing alongside neuroimaging technologies, such frameworks will play an increasingly important role in bridging the gap between technical innovation and clinical practice for neurodegenerative disease management

## REFERENCES

- [1] Y. Quek, Y. Fung, M. Cheung, *et al.*, "Agreement between automated and manual MRI volumetry in alzheimer's disease: A systematic review and meta-analysis," *Journal of Magnetic Resonance Imaging*, 2022..
- [2] M. Tanveer, B. Richhariya, R. Khan, *et al.*, "Machine learning techniques for the diagnosis of alzheimer's disease: A review," *ACM Transactions on Intelligent Systems and Technology*, 2020.
- [3] and Alzheimer's Disease Neuroimaging Initiative *et al.*, "Genetic algorithm with logistic regression feature selection for alzheimer's disease classification," *Neural Computing and Applications*, 2021.
- [4] D. AlSaeed and S. Omar, "Brain MRI analysis for alzheimer's disease diagnosis using CNN-based feature extraction and machine learning," *Sensors*, 2022.
- [5] M. Sabbagh, L. Lue, D. Fayard, and J. Shi, "Increasing precision of clinical diagnosis of alzheimer's disease using a combined algorithm incorporating clinical and novel biomarker data," *Neurology and therapy*, 2017.
- [6] Z. Zhao, J. Chuah, K. Lai, C. Chow, *et al.*, "Conventional machine learning and deep learning in alzheimer's disease diagnosis using neuroimaging: A review," *Frontiers in Computational Neuroscience*, 2023.
- [7] I. Beheshti, H. Demirel, F. Farokhian, C. Yang, *et al.*, "Structural MRI-based detection of alzheimer's

- disease using feature ranking and classification error," *Computer Methods and Programs in Biomedicine*, 2016.
- [8] I. Beheshti, H. Demirel, et al., "Feature-ranking-based alzheimer's disease classification from structural MRI," *Magnetic Resonance Imaging*, 2016.
- [9] V. Karami, G. Nittari, and F. Amenta, "Neuroimaging computer-aided diagnosis systems for alzheimer's disease," *International Journal of Imaging Systems and Technology*, 2019.
- [10] A. Bandyopadhyay, S. Ghosh, M. Bose, A. Singh, et al., "Alzheimer's disease detection using ensemble learning and artificial neural networks," *Unable to determine the complete publication venue*, 2022.
- [11] E. Marwa, H. Moustafa, F. Khalifa, H. Khater, et al., "An MRI-based deep learning approach for accurate detection of alzheimer's disease," *Alexandria Engineering Journal*, 2023.
- [12] M. Kapadnis, A. Bhattacharyya, and A. Subasi, "Artificial intelligence based alzheimer's disease detection using deep feature extraction," *Applications Of Artificial Intelligence*, 2023.
- [13] D. Zhang, Y. Wang, L. Zhou, H. Yuan, D. Shen, et al., "Multimodal classification of alzheimer's disease and mild cognitive impairment," *Neuroimage*, 2011.
- [14] M. Tondelli, G. Wilcock, P. Nichelli, C. D. Jager, et al., "Structural MRI changes detectable up to ten years before clinical alzheimer's disease," *Neurobiology of Aging*, 2012.
- [15] G. Halliday, "Pathology and hippocampal atrophy in alzheimer's disease," *The Lancet Neurology*, 2017.
- [16] M. Gallucci, N. Limbucci, A. Catalucci, and M. Caulo, "Neurodegenerative diseases," *Radiologic Clinics of North America*, 2008.
- [17] L. Velayudhan, P. Proitsi, E. Westman, et al., "Entorhinal cortex thickness predicts cognitive decline in alzheimer's disease," *Journal of Alzheimer's Disease*, 2013.
- [18] A. Du, N. Schuff, J. Kramer, H. Rosen, et al., "Different regional patterns of cortical thinning in alzheimer's disease and frontotemporal dementia," *Brain*, 2007.
- [19] J. Whitwell, W. Crum, H. Watt, and N. Fox, "Normalization of cerebral volumes by use of intracranial volume: Implications for longitudinal quantitative MR imaging," *American Journal of Neuroradiology*, 2001.
- [20] C. J. Jr, "Alzheimer's disease neuroimaging initiative dataset," *Online Draft*, 2026.
- [21] O. Grimm, S. Pohlack, R. Cacciaglia, et al., "Amygdalar and hippocampal volume: A comparison between manual segmentation, freesurfer and VBM," *Journal of Neuroscience Methods*, 2015.
- [22] A. Janssens, Y. Deng, et al., "A new logistic regression approach for the evaluation of diagnostic test results," *Medical Decision Making*, 2005.
- [23] L. Ferreira, J. Rondina, R. Kubo, C. Ono, et al., "Support vector machine-based classification of neuroimages in alzheimer's disease: Direct comparison of FDG-PET, rCBF-SPECT and MRI data acquired ...," *Brazilian Journal of Psychiatry*, 2018.
- [24] F. Mbonyinshuti, J. Nkurunziza, J. Niyobuhungiro, et al., "Application of random forest model to predict the demand of essential med," *Pan African Medical Journal*, 2022.
- [25] K. Malzbender, P. Barbarino, P. Ferrell, et al., "Validation, deployment, and real-world implementation of a modular toolbox for alzheimer's disease detection and dementia risk reduction: The AD-RIDDLE project," *The Journal of Prevention of Alzheimer's Disease*, 2024.
- [26] E. Brown, M. Pierce, D. Clark, B. Fischl, J. Iglesias, et al., "Test-retest reliability of FreeSurfer automated hippocampal subfield segmentation within and across scanners," *Neuroimage*, 2020.
- [27] N. Fox, E. Warrington, P. Freeborough, et al., "Presymptomatic hippocampal atrophy in alzheimer's disease: A longitudinal MRI study," *Brain*, 1996.
- [28] A. Katwaroo, V. Adesh, A. Lowtan, et al., "The diagnostic, therapeutic, and ethical impact of artificial intelligence in modern medicine," *Postgraduate Medical Journal*, 2024.
- [29] D. Iacono, "A double-edged sword: Ethical and psychological implications of APOE genotype disclosure across the lifespan," *Neuroethics*, 2025.