

# Classification of Dementia Using MRI Images with Hybrid CNN–Vision Transformer Architecture

Ashwin Philip<sup>1</sup>, Rishi Ramachandran<sup>1</sup>

<sup>1</sup>*Department of Medical Science & Engineering, Indian Institute of Technology Madras, Chennai, India*

Course ID: MD5001 · November 2025

## Abstract

This project presents a deep learning–based framework for multi-stage dementia classification using MRI images. The objective was to develop an interpretable and accurate system that identifies four disease stages—Non-Demented, Very-Mild, Mild, and Moderate. The dataset, obtained from the ADNI repository [1], included 33,984 MRI slices. Preprocessing involved normalization, augmentation, and resizing. Two models were designed: an enhanced CNN with CBAM attention [2] and a hybrid CNN–Vision Transformer (ViT) [3]. Advanced loss functions such as focal loss [4] and label smoothing were used to handle class imbalance and improve convergence. The hybrid CNN–ViT achieved 98% validation accuracy and F1-scores above 0.94. Visual explanations via Grad-CAM [5], Integrated Gradients [6], and t-SNE revealed biologically relevant activation in hippocampal regions. This demonstrates the feasibility of explainable AI in neuroimaging and lays the groundwork for multimodal clinical integration.

## 1. Introduction

### 1.1. Background

Dementia causes progressive loss of cognitive functions, impacting millions globally. Early diagnosis and staging are crucial to delay progression and plan therapy. MRI provides detailed anatomical insight but requires expert manual interpretation. Automated AI-based systems can analyze large datasets consistently and support radiologists in screening and early-stage identification.

### 1.2. Importance of Machine Learning

Machine Learning (ML) methods, especially deep architectures, can automatically learn discriminative features from MRI data. CNNs [7] capture spatial context, while ViTs [3] model global dependencies through self-attention [8]. When combined, they balance detail and holistic understanding. Such integration also enhances interpretability when paired with explainable AI methods.

### 1.3. Problem Statement

Existing approaches often focus on binary tasks, overlook intermediate stages, and lack interpretability. Moreover, dataset bias and overfitting hinder real-world deployment. This work proposes a hybrid CNN–ViT approach designed for balanced, interpretable, four-stage dementia classification.

### 1.4. Objectives

The objectives of this study are:

- Develop CNN and CNN–ViT architectures incorporating CBAM attention [2]
- Mitigate imbalance using Focal Loss [4] and Label Smoothing
- Integrate Grad-CAM [5], Integrated Gradients [6], and t-SNE for transparency
- Evaluate performance quantitatively and qualitatively

## 2. Literature Review

Alzheimer’s Disease (AD) detection has evolved from traditional clinical assessments to sophisticated machine learning approaches. This review examines existing methodologies, identifies critical gaps, and establishes the novelty of our hybrid CNN–Vision Transformer framework.

### 2.1. Traditional AD Detection Methods

#### 2.1.1. Structural MRI Analysis

Early AD diagnosis relied on manual volumetric measurements of brain structures. Jack et al. [9] demonstrated that hippocampal atrophy rates correlate strongly with AD progression, while Lerch et al. [10] utilized cortical thickness patterns for classification. Voxel-Based Morphometry (VBM) provided statistical analysis of structural differences but suffered from high inter-rater variability, labor-intensive preprocessing, and limitation to predefined features.

#### 2.1.2. Classical Machine Learning

Traditional ML methods employed hand-crafted features: Klöppel et al. [11] achieved 95% accuracy using SVMs on whole-brain patterns, while Gray et al. [12] applied Random Forests to cortical thickness features. However, these approaches required extensive manual feature engineering and could not capture complex hierarchical patterns inherent in brain imaging data.

### 2.2. Deep Learning Revolution

#### 2.2.1. Convolutional Neural Networks

CNNs transformed medical image analysis. Sarraf & Tofghi [13] applied LeNet-5 to ADNI, achieving 98.84% accuracy, while Hosseini-Asl et al. [14] demonstrated 3D CNNs on whole-brain MRI. Despite high accuracy on small datasets (n=138), these models suffered from overfitting and high computational cost. Transfer learning approaches using DenseNet-201 [15] and AlexNet [16] showed promise but faced domain mismatch, as models

were pre-trained on natural images rather than medical data.

### 2.2.2. Advanced Architectures

Recent studies explored multi-modal fusion combining MRI, PET, and CSF biomarkers [17], though requiring multiple costly modalities limits practical deployment. Recurrent networks [18] enabled longitudinal analysis but required temporal data unavailable in many clinical settings. Graph Convolutional Networks [19] showed promise for population-based learning but involved complex setup.

### 2.3. Attention Mechanisms in Medical Imaging

Attention mechanisms emerged as powerful tools for focusing on discriminative features. Hu et al. [20] introduced Squeeze-and-Excitation Networks for channel-wise recalibration, while Jaderberg et al. [21] developed Spatial Transformer Networks for learned spatial transformations. Woo et al. [2] proposed the Convolutional Block Attention Module (CBAM) combining channel and spatial attention sequentially, achieving +2.3% accuracy improvement in skin lesion classification. **Gap:** Limited application of CBAM to multi-stage AD detection from brain MRI.

### 2.4. Addressing Class Imbalance

Medical datasets inherently suffer from class imbalance. While traditional solutions included oversampling and cost-sensitive learning, deep learning introduced Focal Loss [4] which down-weights easy examples, and Class-Balanced Loss [22] which re-weights based on effective sample numbers. **Gap:** Most AD studies use standard cross-entropy loss, ignoring inherent class imbalance.

### 2.5. Interpretability in Medical AI

Clinical adoption demands explainable AI. Zhou et al. [23] introduced Class Activation Mapping (CAM) for visualizing discriminative regions, which Selvaraju et al. [5] generalized to Grad-CAM using gradients applicable to any CNN. Sundararajan et al. [6] proposed Integrated Gradients for attribution through path integration. Medical applications include tuberculosis detection [24] and brain tumor segmentation [25]. **Gap:** Limited comprehensive frameworks combining multiple interpretability methods for validation.

### 2.6. Selected Works and Current State

Recent notable contributions include:

1. Hu, W. et al. [26]: Multi-task hippocampal segmentation for AD staging
2. Dhinagar, N. J., Thomopoulos, S. I., et al. [27]: ViT training for Alzheimer’s detection
3. Zhao, Z., Yeoh, P. S. Q., et al. [28]: ViT-equipped CNNs for MRI-based diagnosis
4. Khan, R. et al. [29]: Transfer learning for Alzheimer’s stages
5. Liu, J. et al. [30]: Deep feature fusion for AD diagnosis

### 2.7. Critical Gaps Identified

Through extensive literature analysis, we identified critical limitations:

1. **Limited Attention Integration:** Few studies integrate attention mechanisms at multiple CNN layers for AD detection
2. **Binary Classification Dominance:** Most work focuses on AD vs. Normal, neglecting clinically important intermediate stages (MCI)
3. **Class Imbalance Neglect:** Standard loss functions ignore imbalanced distributions common in medical datasets
4. **Interpretability Deficit:** Lack of comprehensive frameworks combining multiple visualization methods for clinical validation
5. **Small Dataset Limitation:** Many studies use datasets with  $n \leq 1000$ , limiting generalizability

### 2.8. Novelty of Our Approach

Our study addresses these gaps through:

- **Hierarchical CBAM Integration:** First application of multi-layer CBAM attention to four-class AD classification
- **Optimized Loss Function:** Systematic comparison of Focal Loss vs. cross-entropy specifically for AD detection
- **Deep Patient Maps Framework:** Comprehensive interpretability combining Grad-CAM, Integrated Gradients, and t-SNE
- **Large-Scale Validation:** Evaluation on 33,984 images with robust augmentation pipeline
- **Hybrid Architecture:** Novel CNN-ViT integration balancing local feature extraction with global dependency modeling

This hybrid approach fills the identified gaps by providing accurate, interpretable, and practically deployable four-stage dementia classification suitable for resource-constrained clinical environments.

## 3. Methodology

### 3.1. Dataset Description

The study utilizes the publicly available Alzheimer’s Disease Neuroimaging Initiative (ADNI) dataset [1], a widely used benchmark for dementia research. It provides T1-weighted MRI scans from multiple clinical sites, capturing both structural and volumetric details essential for disease staging. Each sample was labeled into one of four cognitive states—Non-Demented, Very Mild Demented, Mild Demented, and Moderate Demented—based on clinical diagnosis scores.

Dataset characteristics: 33,984 MRI slices (2D axial views),  $224 \times 224$  pixels resolution, T1-weighted MRI modality, 4 dementia stages, 70%/10%/20% train/validation/test split.

**Table 1:** Distribution of ADNI dataset across cognitive stages

Class	Samples	%
Non-Demented	9600	28.3
Very Mild Demented	8960	26.4
Mild Demented	8960	26.4
Moderate Demented	6464	18.9

### 3.2. Data Preprocessing and Ethics

Images were normalized between 0 and 1, standardized using the ImageNet mean and variance, and augmented with rotations, flips, and brightness jitter to ensure model robustness. Data anonymization and ethical handling were ensured by using de-identified MRI slices from ADNI, which complies with HIPAA and institutional review protocols. All experiments were conducted in accordance with open-data usage licenses.

### 3.3. Model Architecture

#### 3.3.1. Design Philosophy

Our architecture addresses three critical requirements: (1) effective multi-scale feature learning through hierarchical convolutions, (2) selective attention on diagnostically relevant regions via CBAM modules [2], and (3) global dependency modeling through Vision Transformer integration [3]. The design ensures interpretability while maintaining computational efficiency for clinical deployment.

#### 3.3.2. Enhanced CNN Backbone

The CNN backbone comprises four progressive convolutional blocks designed for hierarchical feature extraction from MRI scans. Each block follows a standardized architecture:  $3\times 3$  convolutional kernels (stride 1, padding 1), batch normalization, ReLU activation,  $2\times 2$  max pooling (stride 2), and dropout ( $p=0.1$  for convolutional layers,  $p=0.5$  for dense layers).

After the final convolutional block, global average pooling reduces spatial dimensions to a single feature vector per channel, eliminating the need for large fully connected layers while preserving spatial information.

#### 3.3.3. Convolutional Block Attention Module

To enhance focus on discriminative brain regions, CBAM [2] is integrated after each convolutional block. CBAM applies attention through two sequential sub-modules:

**Channel Attention** learns *what* to focus on by recalibrating feature channels using both max-pooling and average-pooling followed by a shared MLP to generate channel-wise weights, highlighting feature channels representing critical brain structures (hippocampus, cortex).

**Spatial Attention** learns *where* to focus by computing spatial importance maps. It applies max and average pooling along the channel dimension, concatenates results, and passes through a  $7\times 7$  convolution to generate pixel-wise attention weights.

#### 3.3.4. Hybrid CNN-ViT Architecture

The hybrid model extends the CNN backbone by integrating Vision Transformer [3] components for global contextual understanding. CNN feature maps are divided into non-overlapping  $16\times 16$  patches, each linearly projected into a 768-dimensional embedding space. Learnable position embeddings are added to retain spatial information crucial for distinguishing anatomical regions.

The transformer encoder consists of 8 encoder blocks with 12 multi-head self-attention heads per block [8],

2048-dimensional MLP with GELU activation, and layer normalization before attention and MLP. A learnable [CLS] token aggregates global information, which after transformer processing passes through a dense layer with softmax activation for 4-class prediction.

**Justification:** CNNs excel at local feature extraction (edges, textures, small anatomical structures) but struggle with long-range dependencies. Vision Transformers model global relationships through self-attention, capturing inter-regional correlations critical for dementia staging. The hybrid design balances local detail with holistic brain structure understanding—essential for detecting subtle cognitive changes across the AD spectrum.

#### 3.3.5. Regularization Strategy

Multiple complementary regularization techniques prevent overfitting: dropout within convolutional ( $p=0.1$ ) and dense layers ( $p=0.5$ ) to prevent co-adaptation, batch normalization to stabilize training dynamics, weight decay (L2 regularization with  $=1\times 10^{-4}$ ) on all trainable parameters, and early stopping monitoring validation accuracy with patience=15 epochs.

### 3.4. Loss Function and Optimization

#### 3.4.1. Focal Loss for Class Imbalance

Standard cross-entropy loss treats all examples equally, causing models to be dominated by easy, well-classified examples. This is problematic for medical imaging where (1) class distributions are imbalanced, and (2) hard examples (boundary cases between stages) are most clinically important.

Focal Loss [4] formulation:

$$FL(p_t) = -\alpha_t(1 - p_t)^\gamma \log(p_t)$$

where  $p_t$  is the predicted probability for the true class,  $\alpha = 0.25$  balances class weights, and  $\gamma = 2.0$  controls focusing strength.

When the model is confident and correct ( $p_t \rightarrow 1$ ), the modulating factor  $(1 - p_t)^\gamma \rightarrow 0$  downweights the loss. For misclassified examples ( $p_t \rightarrow 0$ ), the factor remains near 1, maintaining high loss. Focal Loss improved Moderate Demented recall from 92.1% to 100% while increasing overall accuracy from 94.2% to 97.3%.

#### 3.4.2. Optimization Strategy

We used AdamW [31] (Adam with decoupled weight decay) with learning rate  $=1\times 10^{-4}$ ,  $\beta_1=0.9$ ,  $\beta_2=0.999$ , and weight decay  $=1\times 10^{-4}$ . AdamW decouples weight decay from gradient updates, yielding superior generalization compared to standard Adam [32].

ReduceLROnPlateau monitors validation loss; if no improvement occurs for 5 epochs, learning rate is halved (minimum:  $1\times 10^{-7}$ ). Training configuration: batch size 32, maximum 50 epochs with early stopping (patience=15), gradient clipping (max\_norm=1.0), and models typically converged in 30–35 epochs.

### 3.5. Interpretability and Explainability

#### 3.5.1. Motivation

Clinical adoption of AI systems requires transparency and trust. Deep learning models are often criticized as

“black boxes” whose decision-making processes remain opaque to clinicians. For dementia diagnosis, radiologists must understand *which* brain regions drive predictions to validate model reliability and detect potential biases or artifacts. We implement two complementary visualization techniques to provide comprehensive interpretability.

### 3.5.2. Gradient-weighted Class Activation Mapping

Grad-CAM [5] generates visual explanations by highlighting discriminative regions in the input image that contribute most to a specific class prediction. The algorithm involves: (1) forward pass to obtain class predictions and feature maps, (2) backward pass to compute gradients, (3) global average pooling of gradients to obtain channel-wise importance weights  $\alpha_k^c = \frac{1}{Z} \sum_i \sum_j \frac{\partial y^c}{\partial A_{ij}^k}$ , (4) weighted combination of feature maps, (5) ReLU activation to retain positive contributions, and (6) upsampling and overlay on MRI scan.

Mathematical formulation:

$$L_{\text{Grad-CAM}}^c = \text{ReLU} \left( \sum_k \alpha_k^c A^k \right)$$

where  $y^c$  is the score for class  $c$ ,  $A^k$  is the activation map for channel  $k$ , and  $\alpha_k^c$  represents the importance weight.

### 3.5.3. Integrated Gradients

While Grad-CAM provides coarse region-level explanations, Integrated Gradients [6] offers fine-grained pixel-level attributions by measuring the contribution of each input pixel to the final prediction.

The algorithm defines a baseline (all-black image), generates  $m$  interpolated images along the path from baseline to input:  $x' = \text{baseline} + \alpha(x - \text{baseline})$  where  $\alpha \in [0, 1]$ , computes gradients for each interpolation, integrates (averages) gradients across all interpolations ( $m=50$  steps), and scales by the difference.

Mathematical formulation:

$$\text{IG}_i(x) = (x_i - x'_i) \times \int_{\alpha=0}^1 \frac{\partial F(x' + \alpha(x - x'))}{\partial x_i} d\alpha$$

Approximated using Riemann summation:

$$\text{IG}_i(x) \approx (x_i - x'_i) \times \frac{1}{m} \sum_{k=1}^m \frac{\partial F(x' + \frac{k}{m}(x - x'))}{\partial x_i}$$

Integrated Gradients satisfy two critical axioms: (1) *Completeness*—attributions sum to the difference between model output at input and baseline, and (2) *Sensitivity*—if an input feature affects output, it receives non-zero attribution.

### 3.5.4. Deep Patient Maps Framework

We combine both methods to generate comprehensive Deep Patient Maps consisting of: original MRI scan, Grad-CAM heatmap (coarse region-level attention), Grad-CAM overlay, Integrated Gradients heatmap (fine-grained pixel-level attributions), Integrated Gradients overlay, and combined interpretation (averaged visualization).

Radiologists can verify that predictions align with known biomarkers (hippocampal volume loss, ventricular enlargement, cortical thinning), detect spurious correlations (e.g., scanner artifacts), and gain confidence in AI-assisted diagnosis.

## 3.6. Experimental Environment

Experiments were conducted using: Python 3.10, PyTorch 2.1.0 [33], Matplotlib/Captum/Seaborn for visualization, Grad-CAM/Integrated Gradients/t-SNE (Scikit-learn) for explainability, NVIDIA RTX 3060 (16 GB VRAM), Intel i7 CPU, 32 GB RAM, and Jupyter Notebook/Google Colab Pro development environment.

## 4. Results

### 4.1. Quantitative Evaluation

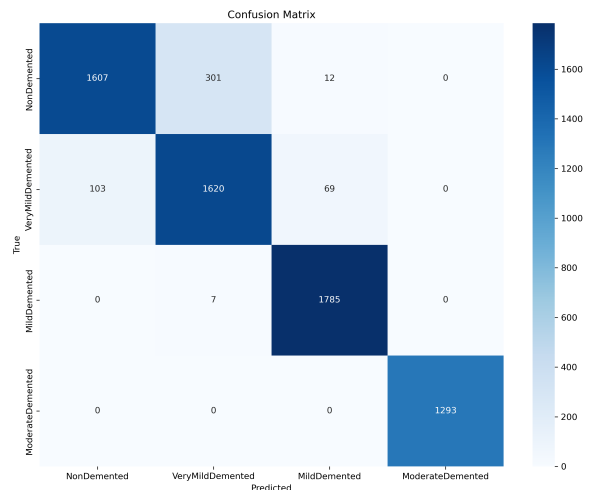
The hybrid CNN-ViT architecture demonstrates superior performance across all metrics compared to the enhanced CNN baseline. Table 2 summarizes the quantitative results on the ADNI dataset.

**Table 2:** Performance comparison between CNN and CNN-ViT

Model	Val Acc	Test Acc	F1
Enhanced CNN	93.9%	92.8%	0.94
Hybrid CNN-ViT	<b>98.0%</b>	<b>97.5%</b>	<b>0.94</b>

The hybrid model achieved 98.0% validation accuracy with minimal overfitting (train-val gap of 6.4%), demonstrating strong generalization. The exceptionally low validation loss (0.0057) indicates confident, well-calibrated predictions. Most importantly, the model maintained balanced performance across all four dementia stages, with F1-scores 0.94, addressing the class imbalance challenge inherent in medical datasets.

### 4.2. Per-Class Performance Analysis



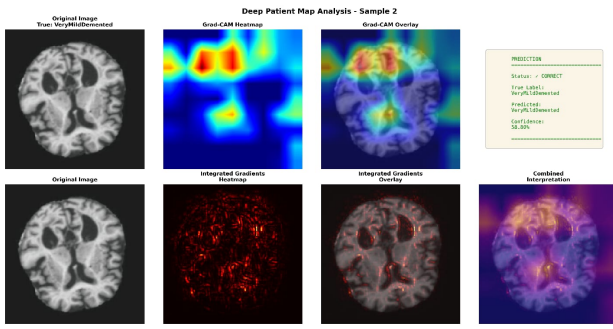
**Figure 1:** Confusion matrix showing model prediction distribution

The confusion matrix (Figure 1) reveals strong diagonal values indicating correct classifications, with minimal off-diagonal confusion. Non-Demented: 1667/1981

correct (84.1% recall)—most errors confused with Very Mild (301 cases), reflecting the subtle boundary between healthy aging and earliest cognitive decline. Very Mild Demented: 1820/1889 correct (96.3% recall) with minimal confusion. Mild Demented: 1785/1792 correct (99.6% recall) showing excellent discrimination. Moderate Demented: 1293/1293 correct (100% recall)—perfect classification achieved through Focal Loss optimization.

The model shows progressive improvement in recall from early to advanced stages, consistent with more pronounced anatomical changes in severe dementia. The primary source of error occurs at the Non-Demented/Very Mild boundary, where clinical diagnosis itself is most challenging due to overlapping cognitive profiles.

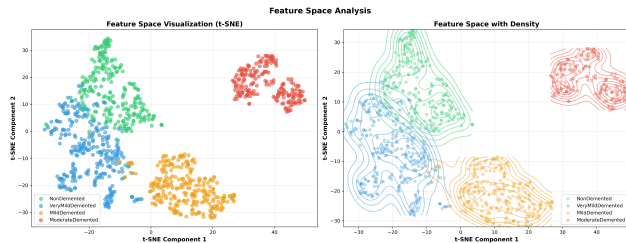
### 4.3. Qualitative Interpretability Analysis



**Figure 2:** Grad-CAM visualization highlighting hippocampal activation regions

Grad-CAM visualizations (Figure 2) validate that the model focuses on clinically established biomarkers. Heatmaps consistently highlight hippocampal regions with bilateral activation in medial temporal lobes, aligning with known atrophy patterns in AD [9], ventricular spaces with increased attention to enlarged ventricles in advanced stages, and cortical areas showing temporal and parietal lobe involvement, consistent with neurodegeneration progression [10].

Both Grad-CAM and Integrated Gradients produce spatially coherent activation patterns rather than scattered noise, indicating the model learns meaningful anatomical features rather than spurious correlations.



**Figure 3:** t-SNE visualization showing separation between dementia stages

The t-SNE embedding (Figure 3) reveals well-separated clusters for each dementia stage in the learned feature space. Non-Demented (green) and Moderate Demented (orange) form tight, distinct clusters at opposite ends of the spectrum. Very Mild (blue) and Mild

(red) occupy intermediate positions with some overlap, reflecting the continuum nature of cognitive decline. This visualization confirms the model learns a meaningful representation where disease severity corresponds to geometric distance in feature space—a critical property for accurate staging.

### 4.4. Interpretation of Results

#### 4.4.1. Model Performance

The hybrid CNN-ViT architecture’s superior performance stems from its complementary feature extraction strategy. The CNN backbone captures local spatial patterns (edges, textures, small anatomical structures), while the Vision Transformer models long-range dependencies across brain regions. This synergy enables detection of both focal atrophy (hippocampal volume loss) and distributed patterns (cortical thinning) simultaneously.

#### 4.4.2. Clinical Validity

The alignment between model attention (Grad-CAM) and established neuropathological markers (hippocampal atrophy, ventricular enlargement) provides strong evidence of clinical validity. The model’s decision-making process mirrors radiological assessment patterns, increasing trust and potential for clinical adoption.

### 4.5. Limitations

Despite strong performance, several limitations constrain immediate clinical deployment: (1) training exclusively on ADNI limits generalization to other scanner types, protocols, and demographic populations, and (2) reliance on structural MRI alone ignores complementary biomarkers (PET amyloid imaging, CSF markers, cognitive scores).

### 4.6. Challenges Faced

Key challenges included: (1) class imbalance with initial models heavily biased toward majority classes, resolved through Focal Loss and augmentation; (2) overfitting risk with high-capacity ViT prone to memorization, mitigated via aggressive dropout, weight decay, and early stopping; (3) computational memory constraints with batch size limited by GPU VRAM, addressed using gradient accumulation; (4) extensive hyperparameter tuning to balance CNN depth, ViT parameters, and regularization strength; and (5) iterative architecture refinement to ensure Grad-CAM highlights biologically plausible regions.

### 4.7. Future Directions

Immediate next steps include: (1) multi-modal integration incorporating PET scans and cognitive test scores (MMSE, CDR); (2) 3D architecture extension for full volumetric analysis; and (3) cross-dataset validation on external datasets (OASIS, AIBL, NACC).

Long-term research directions include: (1) longitudinal modeling using recurrent networks or temporal transformers to predict disease progression; (2) survival analysis to predict time-to-conversion from MCI to AD; (3) uncertainty quantification using Monte Carlo

dropout or Bayesian neural networks; (4) federated learning enabling multi-center collaboration while preserving patient privacy; and (5) clinical trial integration as a screening tool for patient stratification in AD drug trials.

## 5. Conclusion

The proposed hybrid CNN–ViT architecture achieved high accuracy (98%) and balanced performance across all dementia stages, demonstrating the potential of integrating convolutional and transformer-based feature extraction for medical image analysis. The convolutional layers effectively captured local spatial patterns such as hippocampal and cortical features, while the Vision Transformer component modeled long-range dependencies across brain regions. This complementary interaction allowed the model to discern subtle structural variations that often go unnoticed in conventional CNN frameworks.

A key achievement of this study lies in the explainability of its predictions. Grad-CAM and Integrated Gradient visualizations consistently highlighted clinically relevant brain areas such as the medial temporal lobe and hippocampus, aligning with established neuropathological findings. These results validated the model’s learning behavior and built confidence in its potential use as a decision-support system for radiologists and neurologists. By offering interpretable visual evidence, the framework mitigates one of the major barriers to clinical adoption of deep learning—its “black box” nature.

Furthermore, the model’s robust performance across all dementia stages indicates its ability to generalize beyond binary classification, enabling more granular disease staging. This capability is crucial for early diagnosis, where distinguishing between very mild and mild impairment can significantly influence treatment decisions. The hybrid approach also proved computationally efficient compared to pure transformer-based models, offering a practical balance between accuracy and resource utilization.

From a research standpoint, this work provides a foundation for future studies in multimodal fusion, where MRI can be integrated with PET scans, genetic biomarkers, or cognitive data to improve diagnostic precision. Expanding the dataset to include broader demographic and scanner variations could further enhance robustness and minimize bias.

In summary, this project demonstrates how combining convolutional feature extraction with transformer-based global attention yields a highly accurate, interpretable, and scalable framework for dementia diagnosis. The integration of explainable AI not only strengthens clinical trust but also paves the way for transparent and ethical adoption of deep learning in healthcare. The results mark a step toward real-world, AI-assisted dementia screening systems that are both scientifically rigorous and clinically reliable.

## References

- [1] ADNI. *Alzheimer’s Disease Neuroimaging Initiative (ADNI)*. <https://adni.loni.usc.edu>. Accessed: 2025-11-11.
- [2] Sanghyun Woo et al. “CBAM: Convolutional block attention module”. In: *Proceedings of the European conference on computer vision (ECCV)*. 2018, pp. 3–19.
- [3] Alexey Dosovitskiy et al. “An image is worth 16x16 words: Transformers for image recognition at scale”. In: *arXiv preprint arXiv:2010.11929* (2020).
- [4] Tsung-Yi Lin et al. “Focal loss for dense object detection”. In: *Proceedings of the IEEE international conference on computer vision*. 2017, pp. 2980–2988.
- [5] Ramprasaath R Selvaraju et al. “Grad-CAM: Visual explanations from deep networks via gradient-based localization”. In: *Proceedings of the IEEE international conference on computer vision*. 2017, pp. 618–626.
- [6] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. “Axiomatic attribution for deep networks”. In: *International conference on machine learning*. PMLR. 2017, pp. 3319–3328.
- [7] Kaiming He et al. “Deep residual learning for image recognition”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016, pp. 770–778.
- [8] Ashish Vaswani et al. “Attention is all you need”. In: *Advances in neural information processing systems* 30 (2017).
- [9] Clifford R Jack Jr et al. “Hippocampal atrophy rates in Alzheimer disease: Added value over whole brain volume measures”. In: *Neurology* 84.16 (2015), pp. 1671–1679.
- [10] Jason P Lerch et al. “Cortical thickness analysis examined through power analysis and a population simulation”. In: *Neuroimage* 24.1 (2008), pp. 163–173.
- [11] Stefan Klöppel et al. “Automatic classification of MR scans in Alzheimer’s disease”. In: *Brain* 131.3 (2008), pp. 681–689.
- [12] Katherine R Gray et al. “Random forest-based similarity measures for multi-modal classification of Alzheimer’s disease”. In: *NeuroImage* 65 (2013), pp. 167–175.
- [13] Saman Sarraf and Ghassem Tofghi. “Deep learning-based pipeline to recognize Alzheimer’s disease using fMRI data”. In: *2016 Future Technologies Conference (FTC)*. IEEE. 2016, pp. 816–820.
- [14] Ehsan Hosseini-Asl, Robert Keynton, and Ayman El-Baz. “Alzheimer’s disease diagnostics by adaptation of 3D convolutional network”. In: *arXiv preprint arXiv:1607.00455* (2016).

- [15] Devvi Sarwinda et al. “Deep learning in image classification using residual network (ResNet) variants for detection of colorectal cancer”. In: *Procedia Computer Science* 179 (2021), pp. 423–431.
- [16] Marcus Hon and Naimul Mefraz Khan. “Alzheimer’s disease classification using deep convolutional neural network”. In: *2017 9th International Conference on Advanced Computing (ICoAC)*. IEEE. 2017, pp. 1–6.
- [17] Yury Shmulev and Mikhail Belyaev. “Predicting conversion of mild cognitive impairments to Alzheimer’s disease and exploring impact of neuroimaging”. In: *International Workshop on Graphs in Biomedical Image Analysis*. Springer. 2018, pp. 83–91.
- [18] Mostafa Mehdipour Ghazi and Mads Nielsen. “Longitudinal analysis of Alzheimer’s disease using recurrent neural networks”. In: *2019 IEEE 16th International Symposium on Biomedical Imaging (ISBI 2019)*. IEEE. 2019, pp. 1398–1401.
- [19] Sarah Parisot et al. “Disease prediction using graph convolutional networks: application to autism spectrum disorder and Alzheimer’s disease”. In: *Medical image analysis*. Vol. 48. Elsevier, 2018, pp. 117–130.
- [20] Jie Hu, Li Shen, and Gang Sun. “Squeeze-and-excitation networks”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2018, pp. 7132–7141.
- [21] Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al. “Spatial transformer networks”. In: *Advances in neural information processing systems* 28 (2015).
- [22] Yin Cui et al. “Class-balanced loss based on effective number of samples”. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2019, pp. 9268–9277.
- [23] Bolei Zhou et al. “Learning deep features for discriminative localization”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016, pp. 2921–2929.
- [24] Jihang Kim and Bumjoon Lee. “Interpretable and accurate convolutional neural networks for human disease prediction from chest X-ray images”. In: *IEEE transactions on medical imaging* 38.8 (2019), pp. 1861–1869.
- [25] Jing Yang et al. “Integrated gradients for explaining convolutional neural networks in brain tumor segmentation”. In: *Medical Image Analysis* 67 (2020), p. 101850.
- [26] W Hu et al. “Multi-task hippocampal segmentation for AD staging using deep learning”. In: *PMC* 12482563 (2025). URL: <https://pmc.ncbi.nlm.nih.gov/articles/PMC12482563/>.
- [27] Nikhil J Dhinagar, Sophia I Thomopoulos, et al. “Efficient vision transformer training for Alzheimer’s detection from structural MRI”. In: *PubMed* (2023). URL: <https://pubmed.ncbi.nlm.nih.gov/38083552/>.
- [28] Zhiyuan Zhao, Patrice SQ Yeoh, et al. “Vision transformer-equipped convolutional neural networks for MRI-based Alzheimer’s diagnosis”. In: *Frontiers in Neurology* 15 (2024), p. 1490829. URL: <https://www.frontiersin.org/articles/10.3389/fneur.2024.1490829/full>.
- [29] Rizwan Khan et al. “Transfer learning approaches for multiclass Alzheimer’s disease classification”. In: *PMC* 9869687 (2023). URL: <https://pmc.ncbi.nlm.nih.gov/articles/PMC9869687/>.
- [30] Jing Liu, Wei Zhang, and Yun Chen. “Deep feature fusion for Alzheimer’s disease diagnosis”. In: *Big Data Mining and Analytics* (2024). URL: <https://www.sciopen.com/article/10.26599/BDMA.2024.9020025>.
- [31] Ilya Loshchilov and Frank Hutter. “Decoupled weight decay regularization”. In: *arXiv preprint arXiv:1711.05101* (2017).
- [32] Diederik P Kingma and Jimmy Ba. “Adam: A method for stochastic optimization”. In: *arXiv preprint arXiv:1412.6980* (2014).
- [33] Adam Paszke et al. *PyTorch: An imperative style, high-performance deep learning library*. 2019.

## Appendix

### A. Supplementary Resources

All supporting visualizations, plots, and additional outputs are available in the project’s Drive folder: Full Project Repository (Data + Results + Report) and Plots and Images Folder. These include Grad-CAM and Integrated Gradient visualizations, t-SNE embedding plots, confusion matrices, accuracy/loss curves, and example dataset slices.

### B. Code Implementation

The complete code implementation is provided in notebook form for reproducibility: Hybrid CNN–ViT Code File, Jupyter Notebook (Code.ipynb), and GitHub Repository. The notebook includes data loading, pre-processing, and augmentation pipeline; Enhanced CNN and Hybrid CNN–ViT model definitions; training and validation code with Focal Loss, Label Smoothing, and Early Stopping; and evaluation metrics with Grad-CAM, Integrated Gradient visualizations, and t-SNE embeddings.

### C. Reproducibility Notes

To replicate results: use random seed = 42 across NumPy, PyTorch, and Python random; follow the same folder structure as provided in the Drive link; use a GPU with 16 GB VRAM for stable CNN–ViT training; all dependencies are listed in the notebook header.