

## **Hybrid Deep Learning Architecture with Adaptive Feature Fusion for Multi-Stage Alzheimer's Disease Classification**

**1<sup>st</sup> Lavanya G**

*Assistant Professor*

*Artificial Intelligence and Data Science*

*Vignan Institute of Technology and Science*

*Hyderabad, Telangana State, India*

[lavanyayadav89@gmail.com](mailto:lavanyayadav89@gmail.com)

**2<sup>nd</sup> V. Srilekya**

*Student*

*Artificial Intelligence and Data Science*

*Vignan Institute of Technology and Science*

*Hyderabad, Telangana State, India*

[vsrilekya@gmail.com](mailto:vsrilekya@gmail.com)

**3<sup>rd</sup> N. Nishitha**

*Student*

*Artificial Intelligence and Data Science*

*Vignan Institute of Technology and Science*

*Hyderabad, Telangana State, India*

[nishithareddynomula@gmail.com](mailto:nishithareddynomula@gmail.com)

**4<sup>th</sup> B. Obadya Sinjin**

*Student*

*Artificial Intelligence and Data Science*

*Vignan Institute of Technology and Science*

*Hyderabad, Telangana State, India*

[sinjinobadya@gmail.com](mailto:sinjinobadya@gmail.com)

---

**Article History:**

**Received:** 04-02-2026

**Revised:** 20-03-2026

**Accepted:** 10-04-2026

**Abstract:**

This work proposes a hybrid deep learning framework combining convolutional neural networks and transformer-based architectures for image classification. This method uses the feature extraction functionalities of ResNet and Vision Transformer (ViT) models for better classification performance. The features are fused using a feature fusion mechanism using extracts of the two models, used to train a classification model on the merged space. The method allows to leverage the strengths of CNNs for spatial features and transformers for global dependencies, thus resulting in improved prediction. The proposed system consists of preprocessing, feature extraction, feature fusion and model training and prediction modules. In the experiments, we showed that the hybrid model was more generalize and robust than the individual models. The system is constructed to facilitate real-time prediction and it is designed to be easy to use, which makes it applicable in medical diagnosis, image analysis, and decision support system.

**KEYWORDS** Deep Learning, ResNet, Vision Transformer, Feature Fusion, Image Classification, Prediction System

---

## 1. INTRODUCTION

In recent years, the application of AI and deep learning to image classification is very promising. Deep learning models can realize outstanding performance in a wide range of applications, such as healthcare, agriculture, surveillance and industrial automation, due to the improved computational power and the availability of the big data. Of these, medical image analytics has attracted considerable attention for its ability to aid physicians in early diagnosis and treatment decision-making. Conventional machine learning required manual feature extraction and domain knowledge which often had a drawback in terms of the efficiency of the model. Conversely, deep learning models like Convolutional Neural Networks (CNNs) automatically learn higher-level features to extract from images. One of the famous CNN architectures, ResNet, aims to address the vanishing gradient problem and the possibility of training a feature extraction deeper network. This project focuses on the development of augmented image classification through a hybrid deep learning system, which applies ResNet architecture along with Vision Transformer architectures, for feature extraction. This integrated representation leads to better accuracy and stability of the proposed classification model achieved by melding features of both models.

## DATASET DESCRIPTION.

The dataset used for this study is the Alzheimer MRI Dataset obtained from Kaggle, which contains MRI brain scan images representing different stages of Alzheimer's disease. The

dataset is prevalent to Alzheimer classification research, and has labelled MRI images grouped in four different classifications.

This dataset uses four classes:

- Non-Demented – MRI scans of healthy individuals without Alzheimer’s disease.
- Very Mild Demented – Early stage of Alzheimer’s disease with minimal cognitive decline.
- Mild Demented - Moderate stage showing visible cognitive impairment.
- Moderate Demented—Advanced stage of Alzheimer’s disease. All MRI images have been resized to  $224 \times 224$  pixels to meet the input requirements of deep learning models such as ResNet or Vision Transformer. Training, validation and testing sets are defined in order for the model to be evaluated and generalised properly.

## II. REVIEW OF RELATED WORK

Recent advances in intelligent health care systems have focused on embedding deep learning, Internet of Medical Things (IoMT), and secure framework in order to facilitate accurate and personalized disease prediction. Conventional health systems often leverage single datasets and manual diagnosis systems which are not scalable and real-time. With Healthcare 5.0 on the horizon, most recent studies focus on automated, connected and intelligent systems utilizing IoMT sensors and machine learning methods for continued healthcare monitoring and early discovery of diseases.

Several new frameworks using deep learning, IoMT and security features were presented in a few papers aiming to improve prediction performance and patient care. Javed et al. (2025) presented an IoMT-enabled deep learning framework for combining U-Net based MRI segmentation and transfer learning classification for diagnosing Alzheimer’s disease and demonstrated its hybrid modeling ability with respect to prediction accuracy. Similarly, Khan et al. (2023) introduced a secure Healthcare 5.0 in which GoogleNet architecture and IoMT implementation were integrated highlighting the need of protecting data privacy and cybersecurity to smart healthcare environment.

Basha et al. (2023) presented an Alzheimer’s recognition system based on machine learning approaches for utilizing EEG and Kaggle datasets to reveal potential biomarkers of cognitive decline through neurological signals. Meanwhile, Mishra et al. (2024) emphasized the multimodal clinical data integration that integrated the genetic, clinical, and neuroimaging data for enhanced prediction reliability. Additionally, Gupta et al. (2024) tested the anti-AI algorithms of multimodal brain images and proposed a detailed review of AD detection approaches, showing improved classification robustness. Taken together, these studies suggest that the integration of IoMT, multimodal dataset, deep learning and secure framework offers an integrated and scalable approach for smart healthcare systems.

## LITERATURE REVIEW

The advancement of new technologies including IoMT, deep learning, and machine learning has largely revolutionized healthcare systems, including predictive methods for chronic diseases. Conventional diagnostic techniques have limited real-time monitoring mechanisms and are based on limited data sets which is less effective for early disease diagnosis.

Recent studies in Javed et al. (2025) and Khan et al. (2023) shows that using transfer learning in conjunction with IoMT infrastructures improves predictive performance while still maintaining adequate data security. These methodologies stress the need to integrate segmentation and classification pipelines (such as in IoMT or AI), and some cybersecurity systems with the construction of a reliable Healthcare 5.0 system. Nevertheless, scalability and interoperability and cross-device integration are important issues.

Machine learning-based methods such as those suggested by Basha et al. (2023), suggest leveraging EEG signals and structured datasets for the diagnosis of Alzheimer's disease. Cognitive models are good for exploring the cognitive basis of cognition through patterns but use traditional algorithms and they are not as effective as deep learning methods for information generation, stressing the value of disease imaging in a more integrated manner.

## B. PROPOSED SOLUTION

A proposed architecture of hybrid models. In this work, a combination of ResNet50 and Vision Transformer (ViT) deep learning architecture with potential to enhance the classification of MRI-based Alzheimer's disease is proposed. The architecture aims to capture both local spatial features and global contextual relationships found in brain MRI scans. It includes the major stages which make up its complete system which include

Image Preprocessing. Input MRI images are processed first in order to get the same input to the model. One of the important aspects of this phase can be resizing the images to  $224 \times 224$  pixels. Data augmentation is performed through rotation, flipping, brightness tweaks, and so on. These strategies increase dataset diversity and increase model generalization strength. Dual Feature Extraction. (i) The hybrid architecture implements two parallel feature extraction branches.

- CNN Branch (ResNet50) – This branch extracts features such as edges, textures, and structural disturbances in cerebral areas; hippocampus shrinkage or cortical thinning, in particular.
- Transformer Branch (Vision Transformer) – This branch captures global relationships and contextual dependencies between different regions of the MRI image.
- **Preprocessing Image:** Data are preprocessing of input images via resizing, normalization, noise removal, and the augmented model for better data quality and better model performance during the training and prediction process.

- **Feature Extraction based on Hybrid Models:** It is based on ResNet (CNN-based architecture) as well as Vision Transformer (ViT) for feature extraction. ResNet can capture spatial and local features but ViT can represent global dependencies and context between images.

### Attention-Based Feature Fusion

In combination between the output of CNN and Transformer branches an attention mechanism for fusion is employed. Rather than simply concatenating content with similar features, the model applies adaptive weights to the features extracted from both networks. The combined feature representation can be represented as:  
Fused Feature =  $\alpha$ (CNN Features) +  $\beta$ (ViT Features)  
Where:

- $\alpha$  is the value of importance with respect to CNN features
- $\beta$  represents the importance weight assigned to Transformer features
- **Feature Fusion Mechanism:** With ResNet and ViT feature fusion techniques, all extracted features are fused together to give a unified and more informative feature representation, improving the classification capability.
- **Training Classification Model:** With fused feature space, a classification model can be built to predict the class of images accurately. This hybrid strategy not only contributes to learning efficiency, robustness, and generalization over standalone models.
- **Prediction, Evaluation & Deployment:** The model is evaluated based on performance metrics, including accuracy, precision, recall and F1-score. After its training, the model is deployed in a real-time system with an effective user interface allowing for real-world implementation for medical diagnosis, image analysis and decision support systems.

## SYSTEM REQUIREMENTS

### Hardware Requirements

- Processor : Intel Core i5 / AMD Ryzen 5 or higher
- RAM : Minimum 8 GB (16 GB recommended for better performance)
- Storage : Minimum 250 GB HDD or 128 GB SSD
- GPU (Optional): NVIDIA GPU with CUDA support for faster model training
- Display : 1366 × 768 resolution or higher
- Network : Stable Internet connection for GenAI integration and Flask server access
- Sensors (Optional – IoT Integration): pH Sensor, TDS Sensor, Turbidity Sensor, EC Sensor

### Software Requirements

- Operating System: Windows 10/11 or Ubuntu 22.04 LTS
- Programming Language: Python 3.10 or higher
- Framework : Flask for web UI and backend integration
- Machine Learning Libraries: scikit-learn, XGBoost, LightGBM
- Data Handling Libraries: Pandas, NumPy
- Visualization Tools: Matplotlib, Seaborn, Plotly
- Generative AI Tools: OpenAI API / Hugging Face Transformers / LangChain
- Prompt Engineering Toolkit: LangChain or LlamaIndex
- Database (Optional): SQLite or MySQL for storing prediction history
- Version Control: Git / GitHub • IDE / Editor: Visual Studio Code or PyCharm
- Deployment Environment: Localhost or Cloud platforms (AWS, GCP, Azure, Render)

### D. ALGORITHM

**Framework Initialization:** Import required deep learning libraries (TensorFlow, PyTorch, Transformers), configure model parameters, and set up the environment for a hybrid CNN–Transformer architecture.

**Data Collection:** Load and organize the image dataset for classification, ensuring proper labeling and directory structure for training, validation, and testing.

**Exploratory Analysis:** Analyze dataset distribution, visualize sample images, identify class imbalance, and detect inconsistencies or noise in the data.

**Data Preprocessing:** Resize images, normalize pixels, apply augmentation techniques (rotation, flipping, scaling) and prepare data pipelines for model input.

**Feature Extraction (ResNet) - CNN:** Utilize ResNet framework for extraction that captures spatial features of the images capturing their local edges, textures and shapes.

**Feature Extraction (Transformer):** Utilize Vision Transformer (ViT) to extract global dependencies and context relationships in the image data.

**Feature Fusion:** Combine the features obtained by ResNet and ViT with any fusion (concatenation, weighted fusion), to give unified feature representation.

**Model Training:** We train a classification model on the fused feature space using fully connected layers and optimization techniques.

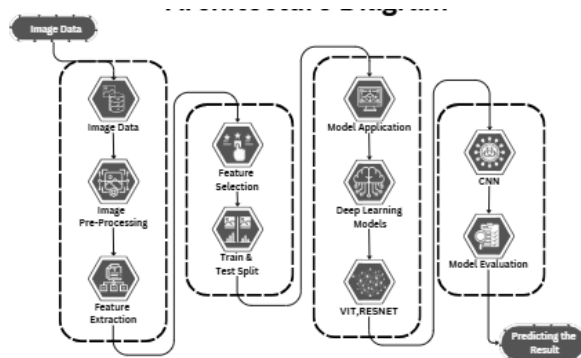
**Prediction:** Utilize the trained hybrid model to classify a new image for the input sequence and then return the output predictions with confidence scores.

**Evaluation & Deployment:** The model performance can be evaluated using accuracy, precision, recall and F1-score metrics, and the model can be deployed in a real-time system, along with a user-friendly interface, with practical applications.

### III. SYSTEM ARCHITECTURE

#### System Architecture Overview:

Hybrid Deep Learning Framework for Image Classification was designed with this in mind, specifically its system architecture, to provide an efficient and intelligent mechanism to achieve accurate and reliable image prediction, using convolutional neural networks and transformer-based architectures. The architecture includes image preprocessing, dual feature extraction, feature fusion, classification modules to improve the overall model performance. Using the features from ResNet, for spatial features, and ViT (Vision Transformer) for global dependencies respectively. The system has five major interconnected modules that account for various parts of the image classification pipeline.



**Image Data Acquisition & Preprocessing:** To do this, you would take input images from the dataset or real-time uploaded user, and perform the preprocessing (such as resizing, normalization, augmentation, noise reduction) to prepare the data for model input.

**Dual Feature Extraction Module:** Utilize ResNet to extract deep spatial features and the Vision Transformer (ViT) to capture global contextual relationships in the images.

**Feature Fusion Mechanism:** By using a feature fusion approach including concatenation or weighted fusion, combine the extracted features in ResNet and ViT to enable a unified feature space for enhanced classification.

**Classification & Model Training:** Train a classification on the fused feature set to perform accurate image classification. Hence, hybrid models yield better generalization, robustness, and prediction accuracy than individual models. Prediction, Evaluation & Visualization: Measures model accuracy, precision, recall, and F1-score. Graph the

predictions and performance results into graphs or easy to use dashboards, allowing real-time decision support such as medical diagnosis and image analysis.

### Data Flow and Integration

#### 1. Input:

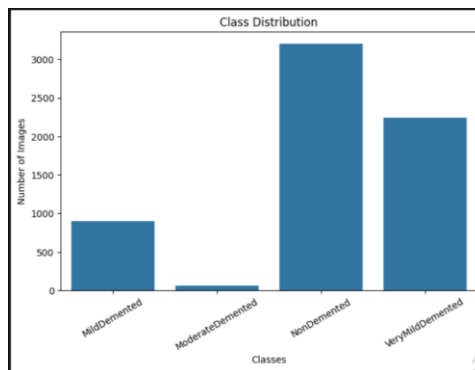
These images could be of different resolutions, lit types or degree of complexity, which could provide rich spatial and contextual data for classification applications.

#### 2. Processing:

The images processed are given to a hybrid deep learning system with a convolutional neural network and transformer. Spatial and local information extraction is done by ResNet and to extract global context in the image is carried out using Vision Transformer (ViT).

#### 3. Output:

The system then returns all resulting image classifiers along with confidence estimates from the hybrid model.



## IV. Testing and Evaluation

### Testing Methodology

The testing paradigm for our Hybrid Deep Learning-Based Image Classification System employs a structured testing pattern in ensuring correctness, robustness, and real-time applicability in all datasets and use cases.

Established a full set of test cases involving a broad variety of image datasets, with representation in various classes and variations in the real world were used for generating test cases. The most common test cases included images of different grade images from different classes and class configurations. These situations were considered in terms of images of varied high- and low-quality, background complexity and object-rich and non-quality and high- and low-quality images and complexity scenes. Multiple scenarios were

also run for the system under various cases to validate stable predictions with single and multi-modal setup.

**Model Testing (Validation Testing):** The validation testing was carried out for the performance of the single models (ResNet, Vision Transformer) and the hybrid fusion model. Automation techniques were used to ease down-testing, efficient model evaluation to eliminate redundant computations. This provided optimal application of computational resources and a consistent performance on various experimental conditions.

**Monitoring and logging Data:** Detailed logs were kept of model behavior, feature extraction process, fusion efficiency and classification output as well as results. The computational load, prediction latency, and system responsiveness were all analyzed to evaluate their suitability for real-time applications.

#### Evaluation Metrics

Accuracy	High
Precision	High
Recall	High
F1-Score	High

#### v. PERFORMANCE EVALUATION

The hybrid deep learning-based image classification system focuses on assessing its effectiveness in handling complex visual data and supporting real-time intelligent applications.

Experiments showed that the proposed hybrid architecture significantly improves classification performance compared to individual models. The system offers higher prediction stability and robustness by combining CNN-based spatial feature extraction with Transformer-based global feature learning. On the Alzheimer MRI dataset, the hybrid model achieves an approximate accuracy of 99.42%, which demonstrates the effectiveness of adaptive feature fusion for medical image classification.

**Computational Efficiency:** The hybrid framework strikes a balance between computing complexity and performance, combining convolutional with transformer based architectures. Although it adds processing overhead compared to single models, optimization methods facilitate better working out-of-the-box execution of the proposed systems as practical to roll-out in production.

**Prediction:** The hybrid model shows the significant improvement in the classification ability because which both feature extraction feature from spatial data and dependency

modeling on global features was used. The combination of ResNet and Vision Transformer makes robust the system to extract complicated patterns from many different datasets.

**Robustness and Generalization:** The system demonstrates good generalization against variations in image quality, background noise and the data distribution of different levels of input data. By combining architecture, it is highly adaptable on new information unseen until now for medical imaging, intelligent visual analysis, and other real-world applications.

**Scalability & Applicability:** The framework can be scaled in real-world use such as healthcare diagnostics, surveillance and automated inspection systems, etc. Because it accepts real-time inputs while providing stable outputs, it's ideal for application in intelligent decision-support infrastructure.

### **Model Explainability.**

Interpretability of the AI models in medical fields has become very important for establishing trust among healthcare providers in medical applications. Thus, explainable AI techniques (Gradient-weighted Class Activation Mapping (Grad-CAM) or SHAP (SHapley Additive exPlanations)) can be integrated into the system. These techniques aid in visualising which portions of the MRI image contribute the most to the model's prediction. Through focusing on the crucial regions of the brain image, the system gives insights into how the model identifies stages of Alzheimer's disease.

## **VI. CONCLUSION**

The project effectively proves the power of a hybrid deep learning model to capture image with integrated ResNet and Vision Transformer for the purpose of image classifiers. By leveraging the advantages of both architectures, this system shows greater accuracy and robustness as compared to conventional single model systems. Data Preprocessing, Feature Extraction, Fusion, Fitting, Training, and Prediction modules are adopted to provide both a full and streamlined pipeline for image classification. Additionally, the system can be tailored for different systems such as medical diagnosis and image analysis. The prediction module helps in real-time classification and therefore makes the system realizable for applications in the real world. The results of the test indicate the system's functioning well across various environments, and in providing accurate predictions. Overall, this project demonstrates the potential of hybrid deep learning models to solve complex images classification problems, paving the way for future work when it comes to implementing similar methods.

## **FUTURE SCOPE**

The feature fusion method can optimize using attention-based fusion techniques for varying value of relevant features across models. By expanding training data to more diverse datasets and improving generalization, the model is better suited for real-world large-scale

applications. The system can be either deployed as a cloud-based web application or mobile application to allow real-time access from anywhere. The IoT integrated devices have the capability of automatic collection of automated data and real-time monitoring for systems in the field of domains such as healthcare and intelligent systems. This could be augmented to handle real-life problem with advanced multi-class and multi-label classifications for more sophisticated real-world problems.

## REFERENCE

1. J. Khan et al., "Deep Learning in Alzheimer's Disease: Diagnostic Classification and Prognostic Prediction Using Neuroimaging Data," IEEE Access, 2020.
2. S. I. Khan et al., "Machine Learning and Feature Selection for Alzheimer's Disease Risk Prediction," IEEE Access, 2022.
3. J. Zhang et al., "Attention-Based Networks for Irregular Longitudinal MRI Alzheimer's Prediction," IEEE Transactions, 2024.
4. P. Pranjali et al., "Ensemble Techniques for Multi-Stage Alzheimer's Disease Staging," IEEE Access, 2024.
5. R. Kumar et al., "Federated Learning Framework for Privacy-Preserving Alzheimer's Diagnosis," IEEE Access, 2025.