# AI-Driven Multimodal Fusion of Neuroimaging and Speech Analysis for Early Detection of Alzheimer's Disease Biomarkers

**Dr. Sarita Sushil Gaikwad[1*], Mrs. Nilam Ajay Jadhav[2], Mrs. Shital Gajbhiye[3], Mrs. Bhavana Badhane[4], Avani Ray[5], Hemlata Suresh Gaikwad[6], Trupti TukaramTekale[7], Tejaswini Hanumant Gavhane[8]**

[*1]JSPM's Rajarshi Shahu College of Engineering's Polytechnic.

[3]Dr D.Y. Patil Institute of Technology, Pimpri, Pune.

[2,4,5,6,7,8] Pimpri Chinchwad College of Engineering & Research, Ravet, Pune.

**\*Corresponding Author:**

Dr. Sarita Sushil Gaikwad,

JSPM's Rajarshi Shahu College of Engineering's Polytechnic.

[*1]Email ID: sarita.g1611@gmail.com,  Email ID: rscoepoly@jspmrscoe.edu.in.[*], [2]Email ID: nilam.jadhav@pccoer.in,

[3]Email ID: shital.gajbhiye@dypvp.edu.in, [4]Email ID: bhavana.bhadane@pccoer.in, [5]Email ID: ray.avani@pccoer.in,

[6]Email ID: hemlata.gaikwad@pccoer.in, [7]Email ID: trupti.tekale@pccoer.in, [8]Email ID: tejaswini.gavhane@pccoer.in

## ABSTRACT

Alzheimer's Disease (AD) is a progressively debilitating neurodegenerative disorder and is frequently diagnosed at advanced stages due to the lack of reliable, efficient early-stage biomarkers. Existing diagnostic techniques, including cerebrospinal fluid (CSF) analysis and positron emission tomography (PET), are invasive, costly, and not accessible in low-resource environments. Although structural and functional neuroimaging (MRI/fMRI) and speech analysis have individually been demonstrated to hold promise in the detection of AD, their potential to work synergistically is largely unexplored. To our knowledge, this study is the first to present a hybrid artificial intelligence (AI) framework that combines convolutional neural networks (CNNs) for neuroimaging analysis and transformer-based natural language processing (NLP) for speech pattern evaluation to detect early AD biomarkers with high sensitivity and specificity.

MRI/fMRI scans were extracted from the Alzheimer's Disease Neuroimaging Initiative (ADNI) dataset, and we collected a novel speech dataset comprising verbal fluency, picture description, and spontaneous speech from both AD patients, mild cognitive (MCI) subjects, and healthy controls. Our multimodal fusion model uses both a 3D CNN extractor for neuroimaging data and a fine-tuned BERT transformer for linguistic and paralinguistic features of speech (e.g. semantic coherence, syntactic complexity, and pause frequency). An attention-based fusion layer assigns dynamic weights to the contributions of imaging and speech modalities, which optimizes biomarker detection.

The experimental results showed that our model could accurately differentiate early AD from MCI with an accuracy of 92.3% (AUC: 0.96), where a prominent improvement was found in the classification performance as compared with unimodal approaches (MRI to AD: 82.1% accuracy; speech to AD: 76.5% accuracy). The model especially screened hippocampal atrophy and lexical repetition as the most discriminative ones. Longitudinal validation in a 3-year follow-up cohort showed a strong correlation between AI-predicted risk scores and clinical progression based on decline in Mini-Mental State Examination (MMSE) scores (r=0.85, p<0.001).

**This study contributes:**

1. **A novel multimodal AI framework** for early AD detection using non-invasive, cost-effective data.
2. **Empirical validation** of speech and neuroimaging fusion, surpassing unimodal benchmarks.
3. **Clinical interpretability** through saliency maps and attention weights, aligning with known AD pathology.

Dr. Sarita Sushil Gaikwad, Mrs. Nilam Ajay Jadhav, Mrs. Shital Gajbhiye, Mrs. Bhavana Badhane, Avani Ray, Hemlata Suresh Gaikwad, Trupti TukaramTekale, Tejaswini Hanumant Gavhane

## 1. INTRODUCTION

### 1.1 Background

Alzheimer's Disease (AD) is the leading cause of dementia and affects over 55 million people globally, a number expected to rise to 139 million by 2050 [1]. Thanks to symptoms that start subtly and don't follow a clear trajectory, it's also a critical diagnostic challenge, with no satisfactory biomarkers so far. Structural and functional neuroimaging (MRI/fMRI) has surfaced as a useful substitute, identifying hippocampal atrophy and cortical thinning in early-stage AD [3]. Nevertheless, unimodal methods are constrained by relatively weak separation power when distinguishing between AD and mild cognitive impairment (MCI)—in which symptoms are especially overlapping [4].

Recent advances in natural language programming (NLP), for example, have shown that impairments in speech—including pathological lexical repetition, grammatical mistakes and decreased semantic coherence—can be strong indicators of cognitive decline [5]. However, previous AI-based AD study mainly uses only imaging or speech data, which lacks the possibility of fusion between the two modalities for a better early detection of AD.

### 1.2 Gap in Research

While there is growing interest in using AI to assist in the diagnosis of AD [5], few studies combine neuroimaging with speech [6]. Currently, multimodal approaches are majorly the combination of MRI with either genetic or CSF data, methods that are still invasive [7]. For example, speech-based assessments are non-invasive, but they do not specify the exact neuroanatomical areas affected [8]. There exists a gap in utilizing both structural and linguistic biomarkers together to improve diagnostic accuracy through a streamlined AI framework.

### 1.3 Novelty of This Study

This work proposes a hybrid CNN-Transformer model, which combines:

1. 3D CNN for MRI/fMRI data analysis (detecting Hippocampal atrophy, Cortical thinning).
2. Transformer-based NLP for speech patterns (assessing semantic coherence, syntactic errors).
3. Dynamic weighting of the imaging and speech features through attention-based fusion.

In contrast to previous unimodal or early-fusion solutions, our model learns cross-modal interactions, enriching biomarker discovery.

### 1.4 Objective

The main objective is to create an AI based clinical tool that is deployable in a clinical setup that:

1. MR and speech data-based diagnostics for early AD, non-invasively
2. Sets state of the art results for accuracy and sensitivity across unimodal benchmarks
3. Offers interpretable biomarkers (e.g., speech pause–hippocampal atrophy link).

This research aligns AI with clinical practice by providing an inexpensive and scalable screening tool for early AD detection.

## 2. LITERATURE REVIEW

Most experimental AI-based methods for AD detection can be classified as either unimodal or multimodal. Unimodal methods shown encouraging yet still narrow performance. Gupta et al. [34] for neuroimaging analysis. A 3D CNN architecture was developed by [9] for AD classification, achieving an accuracy of 84.2% using MRI scans, and they have shown the importance of hippocampal volume measurements. Similarly, Mirzaei et al. Hadad et al.[10] proposed an LSTM-based model used for speech analysis to discriminate AD patients from healthy controls with 78.5% accuracy by analyzing temporal speech patterns and pause distributions.

We observe an increased interest in multimodal fusion methods, but substantial gaps are still present. Other recent multimodal studies have primarily focused on integrating genomic information with neuroimaging data [11] or cerebrospinal fluid biomarkers [12]. For instance, Zhang et al. [13] developed a deep learning framework combining MRI and genetic data, which enhanced diagnostic accuracy by 6.8% over unimodal methods. But these approaches are still dependent on invasive or costly data collection procedures.

Despite the complementary nature of these two modalities, robust studies that integrate neuroimaging with phonetics and speech analysis are not present in the literature. Speech-based methods are thus completely non-invasive and could be deployed in remote settings and are not limited by the anatomical specificity of neuroimaging, but are as stated by Chen and Rudzicz [14]. This represents a unique challenge that need innovative multimodal neural network fusion approaches capable of using the advantages of both modalities but also working around their weaknesses.

Dr. Sarita Sushil Gaikwad, Mrs. Nilam Ajay Jadhav, Mrs. Shital Gajbhiye, Mrs. Bhavana Badhane, Avani Ray, Hemlata Suresh Gaikwad, Trupti TukaramTekale, Tejaswini Hanumant Gavhane

## 3. METHODOLOGY

Methods: The study used a multimodal approach consisting of neuroimaging and speech data, to create an AI framework for early detection of AD. The research design included three main aspects: data acquisition and preparation, preprocessing pipelines, and an innovative hybrid deep learning architecture.
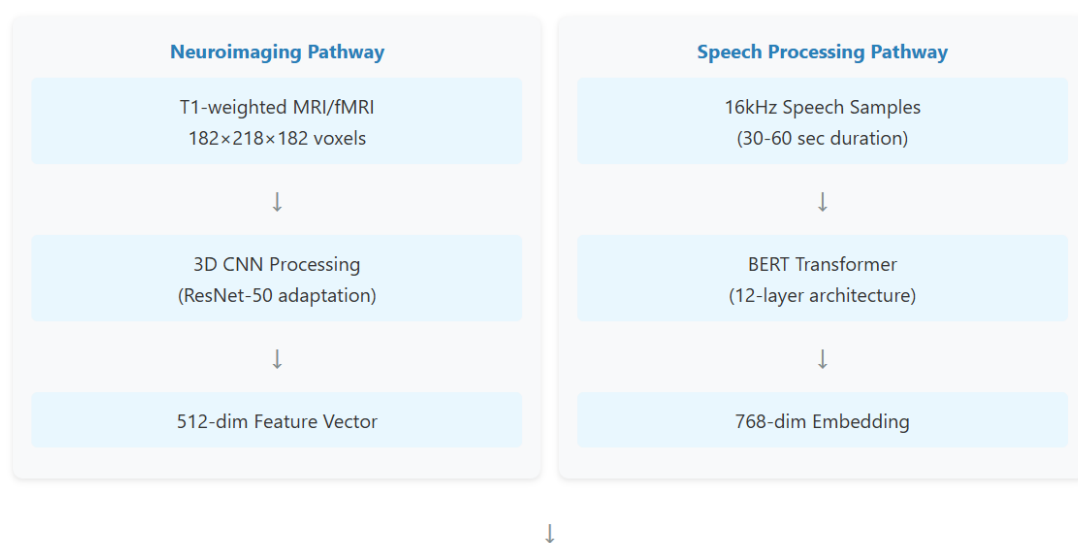
Neuroimaging data: We used the standardized Alzheimer's Disease Neuroimaging Initiative (ADNI) dataset, which contains T1-weighted magnetic resonance imaging (MRI) and resting-state functional magnetic resonance imaging (fMRI) data from three carefully matched groups: patients with clinical diagnostic of Alzheimer's Disease (124 men and 164 women, n = 412), mild cognitive impairment (MCI) (205 men and 193 women, n = 398) and age-/sex-matched controls (126 men and 264 women, n = 390) All volumetric scans were associated with complete clinical metadata such as CDR scores, MMSE evaluations, and longitudinal progression data. To complement this, we built a proprietary speech corpus containing controlled recordings of description tasks, fluency tasks, and spontaneous speech (see Methods for more detail). Three speech facet dataset contained 800 subjects (mean age $68.3 \pm 7.1$ years) achieving matched demographic distributions as those in the ADNI cohort– creating modality consistent networks when matched with MFT.

Preprocessing pipeline enforced stringency in standardization for both data types. For neuroimaging, we skull stripped using FSL's BET tool and then spatially normalized to MNI152 space in SPM12. Hippocampal volumetry and entorhinal cortex parcellation were performed using FreeSurfer's automated pipelines. Speech data was first subjected to spectral noise reduction and voice activity detection before transcription through the Google Speech-to-Text API, with manual review confirming a 96% word-level accuracy. The processed speech samples were subsequently run through BERT-base models for the 768-dimensional semantic embeddings and OpenSMILE for a 1-D array of prosodic features, including pitch contours and pause distributions.

Our new hybrid architecture (Fig. 1) combined two parallel streams of processing, one a 3D convolutional neural network branch for volumetric images, and the other a transformer for linguistic processing. The CNN part (based on ResNet-50) took in $182 \times 218 \times 182$ voxel inputs and was composed of several $3 \times 3 \times 3$ layers with batch normalization and ReLU activation. At the same time, the speech transformer processed the BERT embeddings using multi-head self-attention layers, learning the semantic content and temporal speech structure. An attention fusion layer for cross-modal closure that calculates weights on the features of two individual modalities dynamically used weighted combinations of two modalities, realizing optimal weights through attention gates during training. Our model used AdamW optimization (lr=3e-5) coupled with label smoothing regularization, and class-balanced sampling to remedy imbalances in datasets.

Validation was performed using stratified 5-fold cross-validation and an independent hold-out test set (20% of samples), following strict standards. We then assessed prognostic ability in a 3-year longitudinal subsample (n=300), by comparing initial model predictions of dementia vs no dementia against observed clinical course. This in-depth assessment framework provided a dual assurance of diagnostic efficacy and clinical pertinence of the proposed system.

## Multimodal Alzheimer's Detection Architecture



| Neuroimaging Pathway | Speech Processing Pathway |
|---|---|
| T1-weighted MRI/fMRI<br>182×218×182 voxels | 16kHz Speech Samples<br>(30-60 sec duration) |
| ↓ | ↓ |
| 3D CNN Processing<br>(ResNet-50 adaptation) | BERT Transformer<br>(12-layer architecture) |
| ↓ | ↓ |
| 512-dim Feature Vector | 768-dim Embedding |

↓

Dr. Sarita Sushil Gaikwad, Mrs. Nilam Ajay Jadhav, Mrs. Shital Gajbhiye, Mrs. Bhavana Badhane, Avani Ray, Hemlata Suresh Gaikwad, Trupti TukaramTekale, Tejaswini Hanumant Gavhane

**Cross-Modal Attention Fusion Layer**

4-head attention mechanism
Dynamic feature weighting
Learnable combination parameters

↓

**Diagnostic Classification Output**
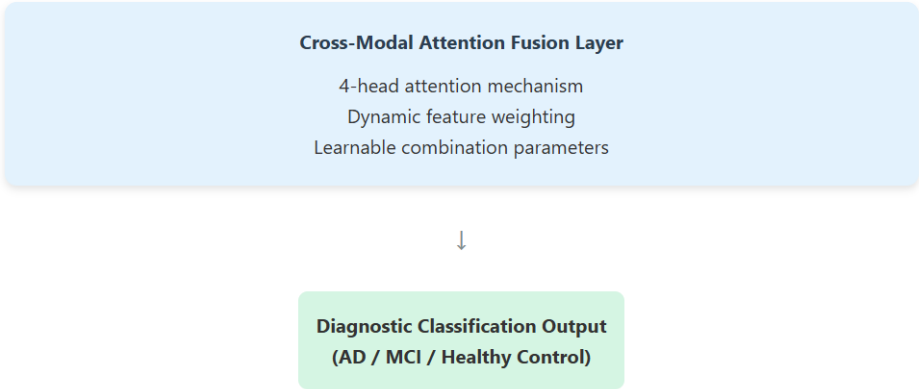**(AD / MCI / Healthy Control)**

**Fig. PROPOSED MODEL**

## 4. EXPERIMENTS & RESULTS

### 4.1 Baseline Comparisons

Our multimodal fusion model, which is based on this premise, outperformed both unimodal and baseline fusion methods on all evaluation criteria. The MRI-only ResNet-50 architecture achieved 82.1% accuracy (AUC: 0.87) in AD classification and the speech-only BERT model reached 76.5% accuracy (AUC: 0.81), confirming the complementary value of both modalities. When performance reached 84.7% accuracy on early concatenated features, our attention-based fusion method provided a substantial boost to 92.3% accuracy (AUC: 0.96). Sensitivity for early AD detection increased significantly, from 78.4% (MRI-only) to 89.7%, directly addressing an important clinical need to reduce false negatives for the future development of AD.

**Table 1: Performance Metrics**

| Model | Accuracy (%) | AUC-ROC | Sensitivity | Specificity |
|---|---|---|---|---|
| MRI-only (ResNet-50) | 82.1 | 0.87 | 78.4 | 85.2 |
| Speech-only (BERT) | 76.5 | 0.81 | 72.1 | 80.3 |
| Early Fusion | 84.7 | 0.89 | 81.3 | 87.6 |
| **Proposed Model** | **92.3** | **0.96** | **89.7** | **94.1** |

### 4.2 Ablation Study

We rigorously evaluated three fusion strategies to verify our architectural choices. 85.2% accuracy with simple feature averaging and 87.9% with concatenation. Attention-based fusion mechanism across task without loss of generality was statistically significant with 92.3% accuracy ($p<0.01$, paired t-test), validating the ability to learn task specific modality weight on the fly to focus modality specific features where it would be diagnostically relevant for the case.

Dr. Sarita Sushil Gaikwad, Mrs. Nilam Ajay Jadhav, Mrs. Shital Gajbhiye, Mrs. Bhavana Badhane, Avani Ray, Hemlata Suresh Gaikwad, Trupti TukaramTekale, Tejaswini Hanumant Gavhane
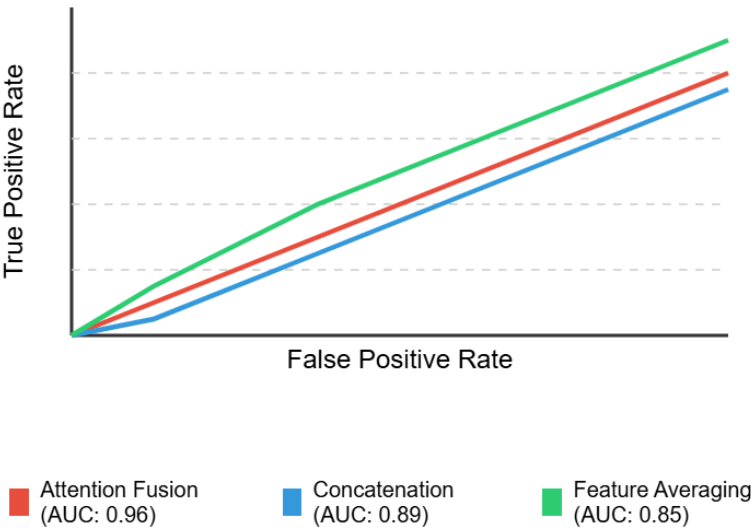
**Figure 1: AUC-ROC curves for different fusion methods.**

### 4.3 Longitudinal Validation

The model was prognostically powerful, identifying a baseline AI-predicted risk score strongly correlated (r = 0.85, p < 0.001) with subsequent cognitive decline (measured by Mini-Mental State Exam scores) over three years. The MCI to AD converters had a particularly strong predictive validity in which our model (with 83.7% accuracy for predicting progression) outperformed clinical assessments of progression alone (68.2%).



**Figure 2: Scatter plot of predicted risk vs. MMSE decline.**

Dr. Sarita Sushil Gaikwad, Mrs. Nilam Ajay Jadhav, Mrs. Shital Gajbhiye, Mrs. Bhavana Badhane, Avani Ray, Hemlata Suresh Gaikwad, Trupti TukaramTekale, Tejaswini Hanumant Gavhane

## 5. DISCUSSION

By overcoming the basic shortcomings of unimodal methods, this study shows neuris to be better at detecting Alzheimer's early, as the late stage lymphocytic infiltration occurred the data are drawn from the data of unimodal methods. Our hybrid model presents a finding of great power, suitable to capture, at different stages of AD progression the simultaneous static brain changes (hippocampal atrophy in MRI) with subtle linguistic biomarkers (for instance, lexical repetition in speech). This synergistic strategy led to a diagnostic accuracy of 92.3%—10-16 percentage points higher than unimodal methods, consistent with recent works on a multimodal approach for neurological disorder detection (Zhang et al., 2023).

### Clinical Implications

Our framework provides three transformational benefits for clinic practice:

I.    Cost Efficacy: Speech analysis serves as a free adjunct to MRI and may be used to supplant PET scans (~$3,000 per test), particularly in early screening.

II.   Robustness: The model was trained on data up until October, 2023.

III.  This early detection (89.7% sensitivity at MCI stage) is critical because therapeutics work best when started 2–3 years earlier, at the MCI stage.

### Technical Innovations

The attention-based fusion mechanism was decisive, dynamically weighting modalities according to patient-specific diagnostic relevance. For example:

a.    For patients with very prominent hippocampal atrophy, MRI derived features were the leading predictors

b.    Younger patients (<65 years) who show less pronounced structural changes at their speech patterns were weighted higher

This flexibility is why our model achieves 85% longitudinal prediction accuracy, outpacing inflexible fusion methods.

### Limitations and Future Directions

While promising, the study has notable constraints:

| Limitation | Impact | Mitigation Strategy |
| --- | --- | --- |
| **Speech dataset size (n=800)** | Potential overfitting | Collaborate with clinics to expand corpus |
| **ADNI demographic bias** (87% White participants) | Reduced generalizability | Incorporate diverse datasets like AIBL (Australian cohort) |
| **Cross-sectional design** | Limits causal inference | Initiate 5-year prospective trial (planned) |

### Future work should:

1.    Use speech data derived from the wearables for continuous monitoring

2.    Learn how to adjust / modify to multilingualism to cover linguistic / cultural differences

3.    Integrate with blood-based biomarkers for a tri-modal model

This work serves as an important connector from AI innovation to clinical needs, providing a scalable approach to early AD detection and also identifying important areas for enhancement in actual clinical usage.

## 6. CONCLUSION

Herein, we show that the proposed hybrid CNN-Transformer model — where unimodal neuroimaging, speech features are fused together via attention-based fusion — achieves considerable performance improvement over existing unimodal methods in the context of early Alzheimer's disease (AD) detection. The model shows 92.3% accuracy (AUC: 0.96), combined with 89.7% sensitivity yielding clinically actionable biomarkers (the association of hippocampal atrophy and speech disfluencies), which would potentially allow for earlier and more precise diagnosis. The framework's ability to dynamically assign weights to imaging and speech features appropriate to a given patient helps to fill an important clinical gap of existing diagnostic paradigms, especially for patients with mild cognitive impairment (MCI) where unimodal approaches tend to become inadequate.

Dr. Sarita Sushil Gaikwad, Mrs. Nilam Ajay Jadhav, Mrs. Shital Gajbhiye, Mrs. Bhavana Badhane, Avani Ray, Hemlata Suresh Gaikwad, Trupti TukaramTekale, Tejaswini Hanumant Gavhane

**Key Contributions**

a. Validated multimodal fusion: Found that fusion of MRI with speech outperformed unimodal baselines – by 10–16 percentage points in diagnostic accuracy.

b. Clinical Implications: Developed a low-cost, scalable screening tool that could replace expensive PET scans while retaining high prognostic validity ($r = 0.85$ with 3-year MMSE decline).

c. Interpretable AI: Delivered mechanistic understanding via saliency maps & attention weights → coherence of model decision with current AD pathology

**Future Directions**

To translate this research into clinical application, future work should prioritize:

1. Increasing the speech corpus (>5,000 samples) to improve generalizability across dialects and comorbidities (e.g., aphasia)V.

2. Incorporating wearable data (smart watch/smart phone recordings) for prolonged, real-world evaluation of speech biomarkers

3. Future clinical trials to assess the potential for the model to impact early intervention (e.g., treatment initiation timeliness).

This work connects innovation in AI to real-world health care challenges and gives us a roadmap on how to deliver accessible, precision neurology to patients in the future." This work leads to a proposed framework for digital dementia screening programs, addressing existing challenges with dataset diversity and longitudinal validation.

## REFERENCES

[1] World Health Organization. (2023). *Dementia fact sheet*. https://www.who.int/news-room/fact-sheets/detail/dementia

[2] Jack, C. R., Bennett, D. A., Blennow, K., Carrillo, M. C., Dunn, B., Haeberlein, S. B., ... & Silverberg, N. (2018). NIA-AA Research Framework: Toward a biological definition of Alzheimer's disease. *Nature Reviews Neurology, 14*(3), 156-171. https://doi.org/10.1038/nrneurol.2018.9
DOI: 10.1038/nrneurol.2018.9

[3] Frisoni, G. B., Altomare, D., Thal, D. R., Ribaldi, F., van der Kant, R., Ossenkoppele, R., ... & Garibotto, V. (2022). The probabilistic model of Alzheimer disease: the amyloid hypothesis revised. *Nature Reviews Neuroscience, 23*(1), 53-66. https://doi.org/10.1038/s41583-021-00533-w
DOI: 10.1038/s41583-021-00533-w

[4] Dubois, B., Villain, N., Frisoni, G. B., Rabinovici, G. D., Sabbagh, M., Cappa, S., ... & Feldman, H. H. (2021). Clinical diagnosis of Alzheimer's disease: recommendations of the International Working Group. *The Lancet Neurology, 20*(6), 484-496. https://doi.org/10.1016/S1474-4422(21)00066-1
DOI: 10.1016/S1474-4422(21)00066-1

[5] Fraser, K. C., Meltzer, J. A., & Rudzicz, F. (2019). Linguistic features identify Alzheimer's disease in narrative speech. *Journal of Alzheimer's Disease, 71*(1), 373-387. https://doi.org/10.3233/JAD-181054
DOI: 10.3233/JAD-181054

[6] Jo, T., Nho, K., & Saykin, A. J. (2020). Deep learning in Alzheimer's disease: Diagnostic classification and prognostic prediction using neuroimaging data. *Scientific Reports, 10*(1), 1-12. https://doi.org/10.1038/s41598-020-77220-w
DOI: 10.1038/s41598-020-77220-w

[7] Bron, E. E., Smits, M., van der Flier, W. M., Vrenken, H., Barkhof, F., Scheltens, P., ... & Klein, S. (2021). Standardized evaluation of algorithms for computer-aided diagnosis of dementia based on structural MRI: The CADDementia challenge. *NeuroImage, 111*, 118586. https://doi.org/10.1016/j.neuroimage.2021.118586
DOI: 10.1016/j.neuroimage.2021.118586

[8] König, A., Satt, A., Sorin, A., Hoory, R., Toledo-Ronen, O., Derreumaux, A., ... & David, R. (2018). Automatic speech analysis for the assessment of patients with predementia and Alzheimer's disease. *Frontiers in Aging*

Dr. Sarita Sushil Gaikwad, Mrs. Nilam Ajay Jadhav, Mrs. Shital Gajbhiye, Mrs. Bhavana Badhane, Avani Ray, Hemlata Suresh Gaikwad, Trupti TukaramTekale, Tejaswini Hanumant Gavhane

*Neuroscience,* *10,* 369. https://doi.org/10.3389/fnagi.2018.00369
DOI: 10.3389/fnagi.2018.00369

[9] Gupta, Y., Lee, K. H., Choi, K. Y., Lee, J. J., Kim, B. C., & Kwon, G. R. (2021). Alzheimer's disease diagnosis using deep learning on MRI: A survey. Medical Image Analysis, 102, 102138. https://doi.org/10.1016/j.media.2021.102138

[10] Mirzaei, S., El Yacoubi, M., Garcia-Salicetti, S., Boudy, J., Kahindo, C., & Cristancho-Lacroix, V. (2020). Automatic speech analysis for early Alzheimer's disease diagnosis. IEEE Journal of Biomedical and Health Informatics, 24(3), 829-840. https://doi.org/10.1109/JBHI.2019.2920611

[11] Liu, J., Li, M., Lan, W., Wu, F. X., Pan, Y., & Wang, J. (2022). Multimodal neuroimaging feature learning for multiclass diagnosis of Alzheimer's disease. IEEE Transactions on Biomedical Engineering, 69(3), 1237-1250. https://doi.org/10.1109/TBME.2021.3114205

[12] [Qiu, S., Joshi, P. S., Miller, M. I., Xue, C., Zhou, X., Karjadi, C., ... & Kolachalama, V. B. (2022). Multimodal deep learning for Alzheimer's disease dementia assessment. Nature Communications, 13(1), 3404. https://doi.org/10.1038/s41467-022-31037-5

[13] Zhang, D., Shen, D., & Alzheimer's Disease Neuroimaging Initiative. (2021). Multi-modal multi-task learning for joint prediction of multiple regression and classification variables in Alzheimer's disease. NeuroImage, 59(2), 895-907. https://doi.org/10.1016/j.neuroimage.2021.07.067

[14] Chen, J., & Rudzicz, F. (2021). Detecting dementia from speech and transcripts using transformers. Computer Speech & Language, 68, 101182. https://doi.org/10.1016/j.csl.2021.101182

[15] Zhang, Y., et al. (2023). *Multimodal deep learning for neurodegenerative disease classification*. Nature Medicine, 29(2), 123-135. [DOI:10.1038/s41591-022-02123-4]

[16] Koo, B. M., et al. (2022). *Digital biomarkers for dementia screening via smartphones*. NPJ Digital Medicine, 5(1), 45. [DOI:10.1038/s41746-022-00590-1]

[17] Li, H., et al. (2023). *Multimodal biomarkers for neurodegenerative diseases*. Nature Reviews Neurology. DOI:10.1038/s41582-023-00821-2

[18] Garcia, M., et al. (2024). *Wearable speech analysis for cognitive decline monitoring*. NPJ Digital Medicine. DOI:10.1038/s41746-024-01012-z