

# Early Detection of Alzheimer’s Disease from Speech with Deep Learning

## A Research Proposal

Raftopoulos Michail

National Technical University of Athens, Greece

### Abstract

Early detection of Alzheimer’s Disease (AD) and Mild Cognitive Impairment (MCI) is critical for timely intervention and effective disease management. Speech analysis has emerged as a promising, non-invasive, and cost-effective approach for detecting the cognitive decline associated with AD and MCI. This research proposal presents a comprehensive literature review of existing methodologies and identifies critical gaps within the field. Despite significant advancements, current models often lack the generalizability and interpretability necessary for clinical deployment, while data scarcity remains a persistent challenge. Furthermore, existing methodologies frequently overlook the conversational dynamics inherent in the data, potentially leading to biased evaluations. To address these limitations, we propose a suite of methodologies including conversational modeling, longitudinal patient analysis, self-supervised learning on unlabeled datasets, and advanced data augmentation techniques. Collectively, these approaches aim to enhance the clinical relevance of the models and maximize the utility of available data resources.

## 1 Introduction

Dementia is a progressive neurodegenerative disorder that causes irreversible brain damage, significantly impairing activities of daily living. It is estimated to affect over 57.4 million people, a figure that is expected to rise dramatically due to the aging population [1]. Alzheimer’s Disease (AD) is the most prevalent form of dementia, accounting for 60-70% of all cases [2].

While there is no known cure for AD, early detection can greatly help patients manage their symptoms and improve the quality of life for them and their families. In particular, researchers highlight the need for detecting a cognitive decline stage known as Mild Cognitive Impairment (MCI) which may or may not progress to dementia. Speech analysis has emerged as a promising approach for detecting MCI, offering a non-invasive, cost-effective, and accessible means of assessment [3], [4].

The field has witnessed significant academic interest in recent years, characterized by frequent competitions and a growing body of literature. However, it continues to face substantial challenges that hinder clinical translation. Although many models report high evaluation metrics, performance levels remain insufficient for reliable clinical application. Moreover, these models are frequently trained and evaluated on small or outdated datasets, raising concerns regarding their generalizability. In this research proposal, we perform a comprehensive literature review of existing methodologies and the

evaluation metrics achieved. We then identify gaps in the literature and missed opportunities, and finally propose a set of promising methodologies to address these limitations.

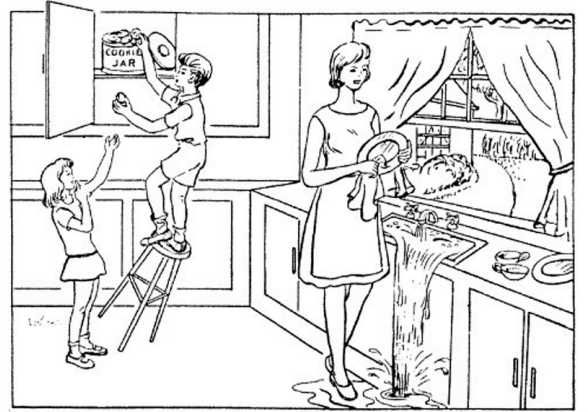


Figure 1: The Cookie Theft Picture, used in the picture description task of the Pitt corpus and other datasets from DementiaBank.

## 2 Literature Review

### 2.1 Datasets and Evaluation Challenges

**DementiaBank:** Most datasets currently available are provided by the DementiaBank database. Dementia-

Bank is a large collection of audio recordings, consisting of various cognitive assessment tasks, administered by an interviewer. Most of these recordings are manually transcribed in the CHAT format, a transcription method that incorporates various linguistic observations such as pauses and stutters, speaker roles, and timestamps along with the spoken words. The most widely used of the provided datasets is the Pitt corpus, which includes mostly the "Cookie Theft Picture Description" task, where participants are asked to describe the image shown in Figure 1. There has also been work toward the creation and expansion of the Delaware corpus, a new dataset that focuses on the binary MCI vs HC distinction. It provides a large variety of cognitive tasks and includes subjects of a wide ethnic and cultural diversity [4].

**WLS:** The WLS is a large-scale, extended longitudinal study of a random sample of 10,317 men and women who graduated from Wisconsin high schools in 1957. The WLS participants were interviewed up to 6 times between 1957 and 2011. DementiaBank provides access to a subset of audio recordings from the 2003 and 2011 interview rounds, amounting to about 1300 different speakers. Although this is a dataset of substantial size, it lacks clinical labels for cognitive impairment, limiting its direct applicability to supervised learning tasks.

**ADReSS, ADReSSo and ADReSS-M:** The ADReSS Challenge was introduced in the Interspeech 2020 conference. It provides a balanced subset of the Pitt corpus, with respect to age and gender [5]. The ADReSSo Challenge followed in 2021, introducing a more difficult task of AD detection using only speech samples, without manually created transcriptions. It also utilized a subset of the Pitt corpus [6]. Lastly, the ADReSS-M challenge, introduced in the ICASSP 2023 conference, focused on the binary classification task of AD detection in the Greek language. It provided a subset of the Pitt corpus, as well as a new dataset with Greek speech samples [7]. Even after the end of these challenges, the datasets remain available and have been widely used as benchmarks for new approaches.

**MultiConAD:** The MultiConAD dataset [8] is a recent effort to tackle the problems of 3-class classification and multilingual generalization. It combines multiple existing datasets, mostly the ones provided by DementiaBank, to create a large and diverse multilingual dataset, with a variety of cognitive assessment tasks. It includes audio and transcription samples in English, Spanish, Chinese, and Greek. Additionally, it provides an important set of baseline models that act as a starting point for future research.

**PROCESS:** The ICASSP 2025 PROCESS Grand Challenge [9] introduced a modern dataset to serve as a benchmark for the 3-class detection problem. It provides audio recordings and manual transcripts (with only audio provided for the test set) from three cognitive tasks: Semantic Fluency, Phonemic Fluency, and Picture Descrip-

tion. The contestants were also tasked with performing regression on the MMSE score, a widely used clinical assessment for cognitive impairment.

## 2.2 AD and MCI Detection from Speech

In this section, we summarize the general trends in methodologies tackling AD and MCI detection from speech. We organize our review by general methodologies; however, these subsections loosely resemble the chronological development of approaches.

### 2.2.1 Feature Engineering

Initial work emphasized the extraction of information-rich features, used in conjunction with traditional ML methods. These approaches proved very effective in the challenge settings as well, where the applications were limited to the small challenge datasets. These methods extract mostly hand-crafted acoustic (e.g., pauses, low-level descriptors, etc.) and linguistic (e.g., verbal richness, filler words, etc.) features and/or semantic representations, mostly from the BERT [10] language model.

Such approaches have led to winning submissions in the ADReSS and ADReSSo Challenges [11], [12]. In the more recent PROCESS Challenge, the overall winner leveraged linguistic features, extracted from ASR transcriptions, and an ensemble of traditional machine learning models, to achieve an F1 score of 0.649. These linguistic features include cognitive-task-specific indicators (count of correct words in Verbal Fluency Tasks) and algorithmically extracted speech pause descriptors [13]. The best performing submission specific to the classification task achieved an F1 score of 0.696, using a self-developed Digital Linguistic Biomarker (DLB) extractor [14].

In a notable paradigm of this methodology, the authors of ADscreen [15] extract a broad range of features, including metrics related to phonetic motor planning, semantic disfluency (e.g., word repetition and pausing), lexical diversity, and syntactic structure. They further incorporate BERT embeddings to capture verbal disfluencies, as well as psycholinguistic features such as LIWC and GeMAPS. Similarly, Haider et al. focus on the emotional dimension of speech by introducing the Affective Behaviour Representation (ABR), which employs a machine learning model to label each speech segment with an emotion and summarize the emotional content of an entire recording into a single vector [16].

### 2.2.2 Deep Learning Architectures and Multimodal Fusion

While feature engineering relies on extracting predefined markers, deep learning approaches aim to learn latent representations directly from the data. Given the success

of deep learning in the more general field of paralinguistics, researchers naturally sought to implement it in the AD detection from speech domain. Deep learning architectures were used to produce meaningful latent representations and fuse acoustic and textual cues.

Within early works, Liu et al. [17] feed MFCCs into a DNN encoder to produce low-dimension bottleneck features of 40ms time frames, followed by CNN and BiLSTM layers. Evaluating their model with 10-fold CV on the Pitt corpus, they achieved an accuracy of 82.59% and an F1 of 82.94%.

In the multimodal setting, a wide variety of methods and architectures emerged. These approaches experiment with various combinations of features, such as BERT embeddings, traditional acoustic features, and other neural-network extracted features [18]–[22]. In an interesting approach, Lee et al. achieved state-of-the-art results in the ADReSSo Dataset by processing the Cookie Theft Picture description task image with VLMs and comparing it with the patient’s text modality to check the validity of the response. They additionally utilized the Shapley value from game theory to introduce a new auxiliary loss function, which includes information about each modality’s contribution [23].

### 2.2.3 Fine-tuning and Foundational Models

Foundational models have been influencing the domain of AD detection from speech since the early days of the ADReSS Challenge. As mentioned in the preceding paragraphs, BERT was frequently used to extract semantic embeddings that were later used as features for classification. Today, Large Language Models are still widely used [24]–[26], utilizing ASR transcriptions from state-of-the-art models and leading to promising results. An extensive study that compared various text-based ML methods showed that the fine-tuning of pretrained models is the most performant approach [27], at least in the text-only setting.

However, limiting the input to text only discards very crucial biomarkers embedded in the patient’s voice. While some have resorted to encoding audio information into the transcripts [24], others have fine-tuned the ASR models themselves to allow the processing of purely acoustic inputs [28]–[30].

Most recent works merge acoustic and linguistic information by incorporating the use of Large Audio-Language Models [31], [32]. Zolnour et al. use a late fusion architecture, combining the predictions of a classifier that is fed pretrained encoder features, and a classifier that is fed handcrafted linguistic features. They also perform extensive experimentations with unimodal and multimodal LLMs, fine-tuned on classification from text and audio. Lastly, they generate synthetic text by prompting foundational LLMs. Overall, the late fusion scheme outperformed the fine-tuned LLMs [31]. Shanin

et al. bypass fine-tuning entirely by simply prompting an Audio-Language model. Although the results were not optimal, they were comparable with supervised methods [32].

Lastly, regarding intermediate pretraining, Zhu et al. [33] used language model perplexity metrics to select large datasets from the GLUE benchmark that are maximally similar to the ADReSS Challenge dataset. Then, they pre-trained a language model on the chosen datasets, followed by fine-tuning on the AD detection task. They additionally utilized their used perplexity-based metrics to invent a new sample-level pretraining technique, where samples that don’t reduce perplexity are discarded.

### 2.2.4 Data Augmentation

Despite the known issue of data scarcity in the field, the applications of data augmentation are limited. To our knowledge, only two publications in the field study explicitly the creation of synthetic data. The authors of CDA [34] propose a contrastive data augmentation technique that simulates cognitive decline by removing random words from sentences, and generates positive samples with multiple passes in conjunction with dropout. Hlédiková et al. [35] perform extensive experimentation with various data-space augmentation techniques. The methods included classical acoustic and verbal perturbations and deep learning-based ones, namely voice conversion with the FragmentVC model, lexical paraphrasing with the Pegasus model and text generation using GPT-2. The authors conclude that their tested neural-based methods perform similarly to traditional ones, still achieving significantly high results.

Some other works have incorporated data augmentation into their methodologies. Specifically, Liu et al. [17] utilize SpecAugment [36], Runde et al. [26] apply Synthetic Minority Over-sampling to balance out the datasets, and Lin and Washington perform Synonym Replacement for their text-based model. Additionally, the authors of LLMCARE [31] experiment with various LLMs to generate synthetic transcripts.

### 2.2.5 State Of The Art

In Table 1, we summarize the datasets that serve as benchmarks in the field, detailing their sample sizes and number of classes. Alongside each dataset, we report the most successful model we could locate in the literature and its corresponding performance metrics. However, we emphasize that the literature domain is currently quite fragmented. The standardized contest datasets are limited in size, and larger alternatives often lack determined train-test splits or standardized evaluation methods (e.g., simple inference vs. cross-validation). This lack of standardization makes objective comparisons and the identification of a true SOTA challenging. Despite these

limitations, this overview highlights the primary data resources available and the current state of performance.

All mentioned dataset, besides the Pitt Corpus and MultiConAD were introduced as evaluation challenges. ADReSS and ADReSSo are subsets of the Pitt Corpus. MultiConAD is an aggregated dataset combining multiple DementiaBank corpora. Although it is multilingual, we focus only on the English subset. Since MultiConAD is a very recent dataset, the only model for comparison is the baseline set by the authors.

## 2.3 Limitations and Opportunities

### 2.3.1 Generalizability Concerns

While the scientific community has achieved remarkable evaluation metrics and made impressive strides in the field, these advancements are currently limited by the small, often outdated and noisy available datasets—often containing only a few hundred samples. Consequently, the applicability of these models in a clinical setting must be questioned.

In a study by Runde et al. [26], the authors managed to achieve 0.99 accuracy and F1 on the Pitt Corpus - a superset of the ADReSS and ADReSSo datasets - by utilizing Wav2Vec transcripts and ada-002 text embeddings. However, when tested on a 10-fold cross validation scheme, these scores dropped to 0.79. A more concerning example involves researchers achieving nearly 100% accuracy on the Pitt Corpus, using solely the silent segments of the audio recordings. This demonstrates the presence of a "Clever Hans" effect, where models achieve high accuracy not by actually learning the underlying mechanisms of the problem, but through spurious correlations in the training data [41]. These findings strongly emphasize the need for a larger, more diverse dataset, a gap promised to be filled by MultiConAD.

### 2.3.2 Interpretability

Clinical adoption necessitates model interpretability, as healthcare professionals must understand which speech characteristics drive diagnostic predictions in order to integrate findings with other clinical indicators and establish professional trust. However, most current approaches utilize complex deep learning architectures and pre-trained language models that function as "black boxes," thereby obscuring the decision-making process. While the domain may not yet have prioritized interpretability as an immediate prerequisite, future methodologies must be designed with the understanding that clinical viability ultimately depends on the transparency and interpretability of the model.

### 2.3.3 Untreated Conversational Nature of Data

With the exception of the 100 samples of the VAS (Voice Assistant System) corpus, all of the DementiaBank data consist of recorded dialogues, where an interviewer administers a cognitive test with the patient. To our knowledge, no recent study has accounted for this conversational structure, instead treating the recordings as homogeneous samples. Consequently, models analyze semantic and paralinguistic information from both the patient and the interviewer indiscriminately. Previous studies have highlighted this limitation, demonstrating that the interviewer's speech significantly impacts observed linguistic features, thereby potentially introducing bias and confounding model predictions [42].

### 2.3.4 Addressing Data Limitations

The most significant constraint within the studied domain is the limited size of available datasets, which restricts methodological innovation and hampers the generalization capabilities of proposed models. Furthermore, the acquisition and dissemination of such data are inherently slow processes due to severe data privacy concerns. Consequently, it is expected that data availability in this field will consistently lag behind industry standards established for more general tasks. While advanced techniques for mitigating data scarcity—such as data augmentation, synthetic data generation, and advanced fine-tuning—should be central to research in this area, they appear to be underrepresented in the current literature. Moreover, the WLS dataset remains significantly underutilized. Despite being unlabeled, this dataset offers substantial potential for pre-training and silver labeling strategies prior to fine-tuning on labeled data. To the best of our knowledge, the only studies that have leveraged WLS in this manner are [8] and [43], which manually determine potential labels based on task performance criteria.

## 3 Methodology

In this section, we propose a set of promising methodologies and demonstrate their potential to address gaps in the literature. Each of these concepts may evolve into a standalone study, capable of being pursued independently or, in some cases, in conjunction with one another.

### 3.1 Conversational Modeling

As mentioned in 2.3.3, most studies disregard the conversational nature of the available data, leading to the undetermined involvement of the interviewer's acoustic and linguistic features. A promising solution is to model the dialogue in a manner that distinguishes between



Table 1: State-of-the-Art (SOTA) Models for Alzheimer’s Disease Detection from Speech.

Dataset	Train/Test Split	Classes	Modalities	Performance	Source
Pitt Corpus	552 total (10-fold CV)	2 (AD, HC)	Audio, CHAT Transcripts	Acc: 0.95 F1: 0.95 AUR: 0.93	[37]
ADReSS	Train: 108 Test: 48	2 (AD, HC)	Audio, CHAT Transcripts	Acc: 0.94	[38]
ADReSSo	Train: 166 Test: 71	2 (AD, HC)	Audio	Acc: 0.96 F1: 0.96	[39]
ADReSS-M	Train: 237 Test: 46	2 (AD, HC)	Audio	Acc: 0.87 RMSE: 3.73	[40]
TAUKADIAL	Train: 387 Test: 120	2 (MCI, HC)	Audio	UAR: 0.86	11249319
PROCESS	Train: 157 Test: Hidden	3 (AD, MCI, HC)	Audio, Transcripts	Macro-F1: 0.696	[14]
MultiConAD (English subset)	Train: 2201 Test: 210	3 (AD, MCI, HC)	Audio, Transcripts	Bin Acc: 0.90 Ter Acc: 0.65	[8]

the two speakers, drawing inspiration from recent advancements in paralinguistics and Emotion Recognition in Conversation (ERC). Various methods exist for modeling conversation, such as context modeling—where each turn is isolated and enriched with speaker-level information [44]—or graph modeling, where utterances serve as nodes connected across speakers and processed via GNNs[45]. Such an approach would align more closely with the inherent nature of the DementiaBank data, allowing the modeling and control of interviewer involvement and potentially yielding more accurate predictions and enhancing model explainability.

### 3.2 Longitudinal Speaker Analysis

As a neurodegenerative disorder, Alzheimer’s disease is inherently progressive. Consequently, there is significant value in analyzing the temporal progression of individual patients by benchmarking speech markers against their own historical baselines. Currently, the only longitudinal data available through DementiaBank comprises the unlabeled WLS samples. However, a recent study conducted an additional sampling round, assessing the cognitive status of 5,414 WLS participants and providing reliable labels for cognitive impairment [46]. This development paves the way for longitudinal patient modeling, provided that access to the full WLS dataset can be secured.

### 3.3 Self-Supervised Learning

Data scarcity is one of the most prevalent challenges in the field, with labeled data amounting to fewer than 3,000 total recordings. However, a significantly larger dataset containing highly similar cognitive tests exists: the Wisconsin Longitudinal Study (WLS). From its most recent sampling round in 2011, DementiaBank provides recordings from over 1300 speakers, yielding approximately 6500 recordings, with similar counts observed in the 2003 round. Consequently, utilizing this significant amount of unlabeled data via self-supervised learning techniques is a logical step. There is a rich body of literature supporting this domain adaptation approach, exemplified by the work of Paraskevopoulos et al. [47], where self-supervision was applied to perform ASR in the resource-limited greek language.

### 3.4 Data Augmentation

As discussed in the literature review, data scarcity remains a prevalent challenge. Although the MultiConAD dataset represents an improvement in this regard, it remains constrained in size for deep learning applications and exhibits significant class imbalance (575 AD, 299 MCI, 1537 HC). Consequently, data augmentation remains crucial, necessitating the exploration of modern methodologies. Broadly, these techniques operate in either the data space or the feature space. In the data space, the objective is to generate synthetic speech samples that mimic full interviews—for instance, via Voice Conversion [48]. However, generating realistic, coherent interviews is a non-trivial task. Conversely, one may oper-

ate in the feature space, like in [49] where spectrograms are generated with GANs to address dataset imbalance. Amidst these complex generative approaches, it is crucial not to overlook simpler, proven interpolation methods like Mixup, which remain highly effective and reliable. Nevertheless, generative methods face significant hurdles, as they must account for the multimodal nature of the task, specifically the challenge of accurately aligning generated text with acoustic modalities.

## References

- [1] E. Nichols, J. D. Steinmetz, S. E. Vollset, *et al.*, “Estimation of the global prevalence of dementia in 2019 and forecasted prevalence in 2050: An analysis for the Global Burden of Disease Study 2019,” *The Lancet Public Health*, vol. 7, no. 2, e105–e125, Feb. 2022.
- [2] M. V. F. Silva, C. d. M. G. Loures, L. C. V. Alves, L. C. de Souza, K. B. G. Borges, and M. d. G. Carvalho, “Alzheimer’s disease: Risk factors and potentially protective measures,” *Journal of Biomedical Science*, vol. 26, p. 33, May 2019.
- [3] “2018 Alzheimer’s disease facts and figures,” *Alzheimer’s & Dementia*, vol. 14, no. 3, pp. 367–429, Mar. 2018.
- [4] A. M. Lanzi, A. K. Saylor, D. Fromm, H. Liu, B. MacWhinney, and M. L. Cohen, “DementiaBank: Theoretical Rationale, Protocol, and Illustrative Analyses,” *American Journal of Speech-Language Pathology*, vol. 32, no. 2, pp. 426–438, Mar. 2023.
- [5] S. Luz, F. Haider, S. de la Fuente, D. Fromm, and B. MacWhinney, *Alzheimer’s Dementia Recognition through Spontaneous Speech: The ADReSS Challenge*, Aug. 2020. arXiv: [2004.06833 \[eess\]](#).
- [6] S. Luz and F. Haider, “Detecting cognitive decline using speech only: The ADReSSO Challenge,”
- [7] S. Luz, F. Haider, D. Fromm, I. Lazarou, I. Kompatsiaris, and B. Macwhinney, “An Overview of the ADReSS-M Signal Processing Grand Challenge on Multilingual Alzheimer’s Dementia Recognition Through Spontaneous Speech,” *IEEE Open Journal of Signal Processing*, vol. PP, pp. 1–12, Mar. 2024.
- [8] A. Shakeri, M. Farmanbar, and K. Balog, *Mul-tiConAD: A Unified Multilingual Conversational Dataset for Early Alzheimer’s Detection*, Feb. 2025. arXiv: [2502.19208 \[cs\]](#).
- [9] F. Tao, B. Mirheidari, M. Pahar, *et al.*, *Early Dementia Detection Using Multiple Spontaneous Speech Prompts: The PROCESS Challenge*, Dec. 2024. arXiv: [2412.15230 \[cs\]](#).
- [10] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*, May 2019. arXiv: [1810.04805 \[cs\]](#).
- [11] R. Haulcy and J. Glass, “Classifying Alzheimer’s Disease Using Audio and Text-Based Representations of Speech,” *Frontiers in Psychology*, vol. 11, Jan. 2021.
- [12] Z. S. Syed, M. S. S. Syed, M. Lech, and E. Pirogova, “Tackling the ADRESSO Challenge 2021: The MUET-RMIT System for Alzheimer’s Dementia Recognition from Spontaneous Speech,” in *Interspeech 2021*, ISCA, Aug. 2021, pp. 3815–3819.
- [13] Y. Gao, L. Guo, and H. Liu, *Leveraging Multimodal Methods and Spontaneous Speech for Alzheimer’s Disease Identification*, Feb. 2025. arXiv: [2412.09928 \[cs\]](#).
- [14] S. Zhang, N. Khelif, M. Ferro, G. Gagliardi, and F. Tamburini, “Cognitive Decline Detection using DLB Extraction Pipelines,” in *ICASSP 2025 - 2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Apr. 2025, pp. 1–2.
- [15] M. Zolnoori, A. Zolnour, and M. Topaz, “AD-screen: A Speech Processing-based Screening System for Automatic Identification of Patients with Alzheimer’s Disease and Related Dementia,” *Artificial intelligence in medicine*, vol. 143, p. 102624, Sep. 2023.
- [16] F. Haider, S. de la Fuente Garcia, P. Albert, and S. Luz, *Affective Speech for Alzheimer’s Dementia Recognition*. May 2020.
- [17] Z. Liu, Z. Guo, Z. Ling, and Y. Li, “Detecting Alzheimer’s Disease from Speech Using Neural Networks with Bottleneck Features and Data Augmentation,” in *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Jun. 2021, pp. 7323–7327.
- [18] L. Ilias, D. Askounis, and J. Psarras, “A Multimodal Approach for Dementia Detection from Spontaneous Speech with Tensor Fusion Layer,” in *2022 IEEE-EMBS International Conference on Biomedical and Health Informatics (BHI)*, Sep. 2022, pp. 1–5. arXiv: [2211.04368 \[cs\]](#).
- [19] L. Ilias and D. Askounis, “Multimodal Deep Learning Models for Detecting Dementia From Speech and Transcripts,” *Frontiers in Aging Neuroscience*, vol. 14, Mar. 2022.
- [20] Y. Pan, Y. Shi, Y. Zhang, and M. Lu, *Swin-BERT: A Feature Fusion System designed for Speech-based Alzheimer’s Dementia Detection*, Oct. 2024. arXiv: [2410.07277 \[eess\]](#).

- [21] J. Cheng, M. Elgaar, N. Vakil, and H. Amiri, *Cog-niVoice: Multimodal and Multilingual Fusion Networks for Mild Cognitive Impairment Assessment from Spontaneous Speech*, Jul. 2024. arXiv: [2407.13660 \[cs\]](#).
- [22] K. Lin and P. Y. Washington, "Multimodal deep learning for dementia classification using text and audio," *Scientific Reports*, vol. 14, no. 1, p. 13 887, Jun. 2024.
- [23] B. Lee, H. J. Song, Y.-J. Park, and B. O. Kang, "Multimodal Alzheimer's disease recognition from image, text and audio," *Scientific Reports*, vol. 15, no. 1, p. 29 038, Aug. 2025.
- [24] X. Ke, M.-W. Mak, and H. Meng, "Optimizing Pause Context in Fine-Tuning Pre-trained Large Language Models for Dementia Detection," in *Interspeech 2025*, ISCA, Aug. 2025, pp. 1408–1412.
- [25] C. Park, A. S. G. Choi, S. Cho, and C. Kim, "Reasoning-Based Approach with Chain-of-Thought for Alzheimer's Detection Using Speech and Large Language Models,"
- [26] B. S. Runde, A. Alapati, and N. G. Bazan, "The Optimization of a Natural Language Processing Approach for the Automatic Detection of Alzheimer's Disease Using GPT Embeddings," *Brain Sciences*, vol. 14, no. 3, p. 211, Feb. 2024.
- [27] B. Ihnaini, Y. Deng, Y. He, L. Geng, and J. Xu, "Detection of Alzheimer's Disease Using Fine-Tuned Large Language Models," *Forum for Linguistic Studies*, vol. 7, no. 8, pp. 373–384, Aug. 2025.
- [28] K. Jia, J. Li, K. Li, and W.-Q. Zhang, "Whisper-Based Multilingual Alzheimer's Disease Detection and Improvements for Low-Resource Language,"
- [29] E. Akinrintoyo, N. Abdelhalim, and N. Salomons, *WhisperD: Dementia Speech Recognition and Filler Word Detection with Whisper*, May 2025. arXiv: [2505.21551 \[eess\]](#).
- [30] J. Li and W.-Q. Zhang, "Whisper-Based Transfer Learning for Alzheimer Disease Classification: Leveraging Speech Segments with Full Transcripts as Prompts," in *ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Apr. 2024, pp. 11 211–11 215.
- [31] A. Zolnour, H. Azadmaleki, Y. Haghbin, *et al.*, "LLMCARE: Early detection of cognitive impairment via transformer models enhanced by LLM-generated synthetic data," *Frontiers in Artificial Intelligence*, vol. 8, Nov. 2025.
- [32] M. Shahin, B. Ahmed, and J. Epps, *Zero-Shot Cognitive Impairment Detection from Speech Using AudioLLM*, 2025.
- [33] Y. Zhu, X. Liang, J. A. Batsis, and R. M. Roth, "Domain-aware Intermediate Pretraining for Dementia Detection with Limited Data," in *Interspeech 2022*, ISCA, Sep. 2022, pp. 2183–2187.
- [34] J. Duan, F. Wei, J. Liu, H. Li, T. Liu, and J. Wang, "CDA: A Contrastive Data Augmentation Method for Alzheimer's Disease Detection," in *Findings of the Association for Computational Linguistics: ACL 2023*, Toronto, Canada: Association for Computational Linguistics, 2023, pp. 1819–1826.
- [35] A. Hlédiková, D. Woszczyk, A. Akman, S. Demetriou, and B. Schuller, *Data Augmentation for Dementia Detection in Spoken Language*, Jul. 2022. arXiv: [2206.12879 \[cs\]](#).
- [36] D. S. Park, W. Chan, Y. Zhang, *et al.*, "SpecAugment: A Simple Data Augmentation Method for Automatic Speech Recognition," in *Interspeech 2019*, Sep. 2019, pp. 2613–2617. arXiv: [1904.08779 \[eess\]](#).
- [37] S. Latif, N. U. Islam, Z. Uddin, K. M. Cheema, S. S. Ahmed, and M. F. Khan, "Deep ensemble learning with transformer models for enhanced Alzheimer's disease detection," *Scientific Reports*, vol. 15, p. 24 720, Jul. 2025.
- [38] M. Martinc, F. Haider, S. Pollak, and S. Luz, "Temporal Integration of Text Transcripts and Acoustic Features for Alzheimer's Diagnosis Based on Spontaneous Speech," *Frontiers in Aging Neuroscience*, vol. 13, p. 642 647, Jun. 2021.
- [39] N. Ntampakis, K. Diamantaras, I. Chouvarda, M. Tsolaki, V. Argyriou, and P. Sarigiannidis, "NeuroXVocal: Detection and Explanation of Alzheimer's Disease through Non-invasive Analysis of Picture-prompted Speech," in vol. 15973, 2026, pp. 410–419. arXiv: [2502.10108 \[cs\]](#).
- [40] X. Chen, Y. Pu, J. Li, and W.-Q. Zhang, "Cross-Lingual Alzheimer's Disease Detection Based on Paralinguistic and Pre-Trained Features," in *ICASSP 2023 - 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Jun. 2023, pp. 1–2.
- [41] Y.-L. Liu, R. Feng, J.-H. Yuan, and Z.-H. Ling, *Clever Hans Effect Found in Automatic Detection of Alzheimer's Disease through Speech*, Jun. 2024. arXiv: [2406.07410 \[eess\]](#).
- [42] C. Li, Z. Sheng, T. Cohen, and S. Pakhomov, "Is There Anything Else?: Examining Administrator Influence on Linguistic Features from the Cookie Theft Picture Description Cognitive Test," Mar. 2025. arXiv: [2503.20104 \[cs\]](#).

- [43] Y. Guo, C. Li, C. Roan, S. Pakhomov, and T. Cohen, "Crossing the "Cookie Theft" Corpus Chasm: Applying What BERT Learns From Outside Data to the ADReSS Challenge Dementia Detection Task," *Frontiers in Computer Science*, vol. 3, Apr. 2021.
- [44] G.-T. Lin, P. G. Shivakumar, A. Gandhe, *et al.*, *Paralinguistics-Enhanced Large Language Modeling of Spoken Dialogue*, Jan. 2024. arXiv: [2312.15316 \[cs\]](#).
- [45] C. Wu, Y. Cai, Y. Liu, *et al.*, *Multimodal Emotion Recognition in Conversations: A Survey of Methods, Trends, Challenges and Prospects*, Sep. 2025. arXiv: [2505.20511 \[cs\]](#).
- [46] V. Williams, R. Trane, K. Sicinski, *et al.*, "Dementia prevalence in the Wisconsin Longitudinal Study," *Alzheimer's & Dementia*, vol. 21, Sep. 2025.
- [47] G. Paraskevopoulos, T. Kouzelis, G. Rouvalis, A. Katsamanis, V. Katsouros, and A. Potamianos, *Sample-Efficient Unsupervised Domain Adaptation of Speech Recognition Systems A case study for Modern Greek*, Dec. 2022. arXiv: [2301.00304 \[cs\]](#).
- [48] M. Illa, B. M. Halpern, R. V. Son, L. Moro-Velazquez, and O. Scharenborg, "Pathological voice adaptation with autoencoder-based voice conversion," in *11th ISCA Speech Synthesis Workshop (SSW 11)*, ISCA, Aug. 2021, pp. 19–24.
- [49] A. Chatziagapi, G. Paraskevopoulos, D. Sgouropoulos, *et al.*, "Data Augmentation Using GANs for Speech Emotion Recognition," in *Inter-speech 2019*, ISCA, Sep. 2019, pp. 171–175.