# Dementia Insights: A Context-Based MultiModal Approach

Sahar Sinene Mehdoui[1] ⓘ, Abdelhamid Bouzid[1] ⓘ, Daniel Sierra-Sosa[2] ⓘ, Adel Elmaghraby[1] ⓘ

[1] Department of Computer Science and Engineering, University of Louisville
adel.elmaghraby@louisville.edu

[2] Department of Electrical Engineering and Computer Science, The Catholic University of America
sierrasosa@cua.edu

## ABSTRACT

Dementia, a progressive neurodegenerative disorder, affects memory, reasoning, and daily functioning, creating challenges for individuals and healthcare systems. Early detection is crucial for timely interventions that may slow disease progression. Large pre-trained models (LPMs) for text and audio, such as Generative Pre-trained Transformer (GPT), Bidirectional Encoder Representations from Transformers (BERT), and Contrastive Language-Audio Pretraining (CLAP), have shown promise in identifying cognitive impairments. However, existing studies generally rely heavily on expert-annotated datasets and unimodal approaches, limiting robustness and scalability. This study proposes a context-based multimodal method, integrating both text and audio data using the best-performing LPMs in each modality. By incorporating contextual embeddings, our method improves dementia detection performance. Additionally, motivated by the effectiveness of contextual embeddings, we further experimented with a context-based In-Context Learning (ICL) as a complementary technique. Results show that GPT-based embeddings, particularly when fused with CLAP audio features, achieve an F1-score of $83.33\%$, surpassing state-of-the-art dementia detection models. Furthermore, raw text data outperforms expert-annotated datasets, demonstrating that LPMs can extract meaningful linguistic and acoustic patterns without extensive manual labeling. These findings highlight the potential for scalable, non-invasive diagnostic tools that reduce reliance on costly annotations while maintaining high accuracy. By integrating multimodal learning with contextual embeddings, this work lays the foundation for future advancements in personalized dementia detection and cognitive health research.

## 1 Introduction

Dementia is an increasing global health challenge, affecting millions of individuals and imposing significant burdens on healthcare systems World Health Organization [2023]. Early diagnosis is critical for facilitating timely interventions that can slow disease progression and improve patient outcomes. Although effective, traditional diagnostic methods such as neuroimaging and biomarker analysis are often costly, invasive, and inaccessible in many settings El Abiad et al. [2024], Chávez-Fumagalli et al. [2021], Zhao et al. [2023], Lee et al. [2017]. Consequently, research has increasingly focused on non-invasive, scalable diagnostic methods, particularly those utilizing speech and language analysis to detect cognitive impairments through linguistic and acoustic markers.

Advancements in artificial intelligence (AI) and computational methods have significantly transformed dementia research. Early machine learning models, such as logistic regression and Random Forests, effectively identified cognitive decline using linguistic indicators such as coherence, hesitations, and repetitions [Shah et al., 2021, Jahan et al., 2024]. The emergence of deep learning, particularly transformer-based architectures such as BERT Devlin et al. [2019] and GPT Radford et al. [2019], Brown et al. [2020], Adler et al. [2024], has further improved the detection of subtle linguistic patterns associated with cognitive impairments. Specialized adaptations, like AD-BERT [Mao et al., 2023], have improved early dementia detection through pre-training on Alzheimer's disease-related datasets. Large language models (LLMs), such as GPT-4 Adler et al. [2024], excel in few-shot learning and effectively capture nuanced linguistic expressions linked to cognitive decline Du et al. [2024]

Beyond text-based analyses, audio data provide critical insights into dementia-related speech characteristics by examining features such as pitch, pauses, articulation changes, and prosody. Advanced signal processing techniques, including log-Mel spectrograms and delta coefficients, have proven to be effective in identifying cognitive impairments Meghanani et al. [2021]. Vision Transformers (ViTs) Dosovitskiy et al. [2021] have outperformed traditional CNN-

based models in capturing global dependencies from speech signals [Ilias et al., 2023], further advancing audio-based dementia detection.

Integrating text and audio data through multimodal approaches holds significant potential for improving classification accuracy. By combining the richness of textual features with complementary speech-based cognitive markers [Han et al., 2022], multimodal frameworks such as GP-Net, which Liu et al. [2022] proposed, have demonstrated superior performance over conventional transformer-based models in dementia detection [Liu et al., 2022]. These integrated approaches underscore the value of harnessing symbolic annotations, linguistic structures, and acoustic properties in providing a more comprehensive assessment of cognitive health.

Despite these advancements, several challenges persist. Many studies rely on structured clinical notes or specialized datasets, thereby limiting their applicability to real-world settings. Moreover, while transcribed speech data include valuable linguistic annotations, such as self-corrections and repetitions, their practical impact on machine learning performance remains underexplored. There is a critical need for dementia detection models that are accurate, scalable, generalizable, and interpretable.

In this paper, we introduce a context-driven multimodal model as our primary contribution, integrating both text and audio features for dementia detection using large pre-trained models. Unlike prior studies, our method explicitly incorporates contextual information across both modalities, enhancing interpretability and robustness. Additionally, recognizing the importance of contextual learning in multimodal analysis, we explore an In-Context Learning (ICL) strategy to evaluate its potential as a complementary technique. Our experiments demonstrate that while ICL shows promise, our proposed framework significantly outperforms alternative methods, making it the core innovation of this study. Previous works, such as Pan et al. [2025], have leveraged cross-attention mechanisms with self-supervised learning (SSL) models like wav2vec2.0 and transformer-based architectures such as BERT, but these primarily relied on structured datasets and emphasized audio features. This study addresses those limitations by integrating multimodal data with large language models and evaluating both expert-annotated and raw transcriptions to improve scalability and robustness. We leverage the best available pre-trained models in each modality, including CLIP OpenAI [2021], GPT, and BERT for text embeddings, and CLAP Elizalde et al. [2023] for audio. Using the Pitt Corpus MacWhinney et al. [2011], we conduct a broad comparative analysis, assessing the efficacy of expert-annotated versus raw datasets. Our findings offer valuable insights into how various text and audio features contribute to dementia detection, underscoring the potential of scalable, non-invasive diagnostic tools.

The structure of this paper is as follows: section 2 provides an overview of existing dementia detection methodologies, including text-based, audio-based, and multimodal approaches. Section 3 describes the *Cookie Theft Test* Cummings [2019] and the Codes for the Human Analysis of Transcripts (CHAT) MacWhinney [2024], both of which form the basis of our experimental framework. Section 4 describes the dataset utilized in this study. Section 5 presents the methodologies employed in this study. Section 6 reports the experimental results and their analysis. Finally, section 7 discusses the implications of our findings and proposes avenues for future research.

## 2 Related Work

Dementia detection has advanced substantially with computational methods, offering innovative techniques for identifying cognitive decline Du et al. [2024], Mao et al. [2023], Han et al. [2022], Bouazizi et al. [2023], Ilias et al. [2023], Balagopalan et al. [2020], Zhu et al. [2021], Pan et al. [2025], Baevski et al. [2020], Haulcy and Glass [2021], BT and Chen [2024]. While traditional diagnostic tools such as neuroimaging and biomarkers remain effective, they are often invasive, expensive, and inaccessible in many settings. As a non-invasive, scalable alternative, speech and language analysis has emerged as a promising method for detecting cognitive impairments by examining text and audio changes indicative of dementia. This section reviews related work, categorized by data modality: text-based, audio-based, and multimodal approaches.

### 2.1 Text-Based Methods

Text-based methods have played a pivotal role in dementia detection, primarily focusing on the analysis of clinical notes, structured speech tasks, and spontaneous transcriptions Gauder et al. [2021], He et al. [2016], Haider et al. [2019], Hershey et al. [2017], Pan et al. [2025], Du et al. [2024]. Early studies extracted linguistic features such as word frequency, syntactic complexity, and hesitation markers, applying traditional machine learning algorithms such as Random Forests and logistic regression for classification Shah et al. [2021].

The introduction of deep learning into dementia research has significantly improved detection models. Transformer-based architectures, such as BERT, BioBERT Lee et al. [2020], and ClinicalBERT AI [2025], have transformed the analysis of unstructured clinical data. Specialized models such as AD-BERT, pre-trained on Alzheimer's disease-specific

datasets, demonstrated competitive performance in forecasting disease progression Mao et al. [2023]. Large Language Models (LLMs), including GPT and GPT-4, have further advanced dementia detection by leveraging few-shot learning capabilities, enabling robust performance with limited labeled data Du et al. [2024], BT and Chen [2024], Chen et al. [2024].

Despite these advancements, many models treat transcribed speech as standard text, overlooking critical linguistic annotations that provide deeper insights. For instance, Jahan et al. [2024] used the Pitt Corpus, extracting linguistic features such as hesitations, repetitions, and grammatical errors for traditional machine learning models, with Random Forests achieving the best performance. However, their methodology did not fully leverage the structural nuances of CHAT files, potentially missing critical indicators of cognitive decline.

This study addresses these limitations by incorporating both raw text and annotated transcripts from the Pitt Corpus. We focus on symbolic annotations to capture subtle linguistic cues, enhancing model interpretability and predictive performance. Additionally, we integrate transfer learning from pre-trained models to improve generalizability across diverse datasets.

## 2.2 Audio-Based Methods

Audio-based techniques analyze acoustic features such as pitch, speech rate, articulation patterns, and pauses to detect cognitive impairments. Early studies relied on handcrafted features combined with machine learning classifiers such as support vector machines and Random Forests Haulcy and Glass [2021].

Recent advances in deep learning have led to the adoption of more sophisticated architectures for analyzing audio data. Vision Transformers (ViTs) have outperformed traditional CNN-based models by capturing global dependencies in speech signals, thereby improving the identification of dementia-related patterns Ilias et al. [2023]. The use of audio-based techniques is particularly advantageous for capturing prosodic and paralinguistic cues that are often lost in text transcriptions.

However, these models face challenges related to dataset variability, noise interference, and differences in recording quality, which can limit their generalizability. Our research addresses these issues by incorporating robust pre-trained audio models such as CLAP, allowing for effective feature extraction from noisy, real-world audio samples.

## 2.3 MultiModal Approaches

Combining text, audio, and sometimes visual data has demonstrated significant potential for improving dementia detection accuracy Bouazizi et al. [2023], Zhu et al. [2023], Han et al. [2022], Lin and Washington [2024]. By leveraging complementary strengths across different modalities, these models provide a more holistic assessment of cognitive markers Han et al. [2022].

Structured tasks such as the Cookie Theft Picture Description Cummings [2019] are widely used to collect multimodal data, enabling joint analysis of linguistic coherence and audio features. Recent studies have demonstrated the effectiveness of integrating BERT-based text embeddings with spectrogram-derived audio features, establishing improved benchmarks for dementia classification Han et al. [2022]. Additionally, pre-trained multimodal models such as CLIP Radford et al. [2021] and BLIP-2 Li et al. [2023] have been explored for their ability to align textual and visual data, enabling zero-shot analysis of cognitive impairments Zhu et al. [2023].

Despite these advancements, the integration of text and audio data remains underexplored. Many studies focus on a single modality, missing opportunities to leverage the synergy between linguistic and acoustic cues. Moreover, existing multimodal models often require large, labeled datasets, which can constrain their scalability.

Our research addresses these limitations by employing state-of-the-art language models (GPT, BERT) and advanced audio processing techniques (CLAP). We assess the impact of symbolic linguistic annotations and raw transcriptions, evaluating their contributions to dementia detection. Through multimodal embeddings, our framework enhances diagnostic accuracy, robustness, and scalability.

Ultimately, this study highlights the importance of leveraging advanced multimodal AI techniques to improve early dementia detection, paving the way for more accessible, non-invasive diagnostic tools suitable for diverse clinical and community settings.

# 3 Preliminaries

## 3.1 The Cookie Theft Picture Description Task

The Cookie Theft picture description task, as shown in Figure 1, is a well-established cognitive assessment tool used to evaluate multiple cognitive functions. In this task, participants are presented with a standardized image depicting a detailed scene, known as the "Cookie Theft" picture, and are asked to verbally describe their observations in detail.
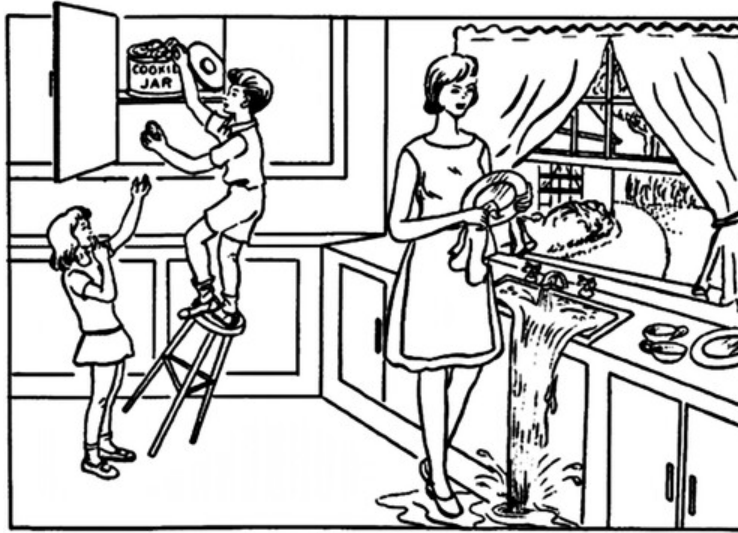


Figure 1: The Standardized 'Cookie Theft' picture for assessing cognitive abilities through verbal description. Source: Boston Diagnostic Aphasia Examination, available under Creative Commons Attribution 4.0 International (CC BY 4.0). Original image retrieved from Examination [2025]. No changes were made to the original image.

This task provides insights into several cognitive domains often impaired in dementia. It assesses attention and perception by examining the participant's capacity to explore the image visually, identify key details, and interpret visual elements. Language skills are examined by analyzing clarity, coherence, grammatical accuracy, vocabulary, and fluency in the participant's verbal description. Additionally, executive functioning is assessed through an analysis of the participant's ability to organize their thoughts, structure a coherent narrative, and logically describe the scene. Memory, particularly working memory, is evaluated based on the participant's capacity to retain and recall previously mentioned details while tracking multiple elements within the image.

Despite being based on a single standardized image, the task offers a valuable dataset for evaluating a wide range of cognitive functions. This task evaluates not only what participants describe but also how they structure and articulate their responses, providing key insights into dementia-related cognitive impairments. This task was chosen due to its extensive use and demonstrated effectiveness in cognitive assessment.

## 3.2 Codes for the Human Analysis of Transcripts (CHAT)

The Codes for the Human Analysis of Transcripts (CHAT) system, developed for the TalkBank project, provides a structured framework for transcribing and analyzing spoken language MacWhinney [2024]. It uses standardized codes to annotate various linguistic features, including hesitations, grammatical errors, repetitions, self-corrections, and disfluencies. Annotations such as self-corrections, repetitions, aborted speech, and grammatical errors enable a detailed analysis of language patterns that may indicate cognitive decline. These annotations capture nuanced linguistic disruptions often associated with dementia, providing rich data for machine learning models. Table 1 highlights the linguistic differences between the Control and Dementia groups.

While CHAT annotations offer in-depth linguistic insights, they present notable challenges. The manual annotation process is time-intensive and requires expertise from linguists and speech therapists, making it resource-demanding and less scalable for large datasets or real-world clinical applications. This reliance on specialized annotators limits its feasibility for widespread use, particularly in settings where rapid, automated assessments are needed. Additionally, the structured nature of CHAT may not fully capture the variability and complexity of natural speech, potentially restricting the generalizability of models trained solely on annotated data.

Table 1: Summary of Linguistic Differences Between Control and Dementia Groups Based on Visualized CHAT Transcribed Text Data Annotations

| Aspect | Control Group | Dementia Group |
|---|---|---|
| Coherence | Clear and structured | Fragmented and incoherent |
| Disfluencies | Few | Frequent |
| Grammatical Accuracy | Minimal errors | Frequent errors |
| Repetition | Minimal | Frequent repetition |
| Hesitation Sounds | Limited | Frequent |
| Content Relevance | Focused and relevant | Unrelated or nonsensical |

In this study, we integrate CHAT-annotated transcripts with raw, unannotated text data to evaluate their effectiveness in dementia detection. By comparing models trained on both data types, we aim to determine whether the granular linguistic markers from CHAT significantly enhance model performance compared to raw text processed through advanced language models such as GPT. Furthermore, we combine these linguistic features with audio data to enable multimodal analysis, correlating text-based markers with acoustic cues such as prosody, pitch, and speech rhythm. This approach balances detailed linguistic analysis with scalability, enhancing the potential for practical, non-invasive dementia screening tools.

# 4    Dataset Description

## 4.1    Overview

The dataset used in this study is obtained from the Pitt Corpus. This dataset is part of the Alzheimer's and Related Dementias Study, collected at the University of Pittsburgh School of Medicine Becker et al. [1994]. The dataset consists of transcribed text and audio recordings from individuals categorized into two groups: a control group (cognitively healthy individuals) and a dementia group (individuals diagnosed with dementia). For this study, the dataset contains 243 tests from 99 control participants and 309 tests from 194 dementia participants. Each test comprises a text file containing the transcript of the verbal description and a corresponding *MP3* audio file of the recorded speech.

All participants completed the Cookie Theft Picture Description (Cookie) test. While the other three tests: Verbal Fluency (Fluency), Delayed Story Recall (Recall), and Sentence Repetition (Sentence) have significantly fewer participants in the Dementia group and almost no participants in the Control group. Moreover, only text files are available for these tests, with no corresponding audio recordings. This imbalance in participation and the lack of multimodal data for these tests limited their inclusion in certain analyses, reinforcing the primary role of the Cookie Theft test in this study.

This dataset includes a variable number of participants across different test categories. Table 2 summarizes the distribution of unique participants for each test, separated by control and dementia groups.

Table 2: Number of Unique Participants in Each Test Category Across Four Speech-Based Cognitive Assessments for Dementia Detection in the Pitt Corpus data

| Test | Control Group | Dementia Group |
|---|---|---|
| Cookie Theft (Cookie) | 99 | 194 |
| Verbal Fluency (Fluency) | 2 | 163 |
| Delayed Story Recall (Recall) | 1 | 178 |
| Sentence Repetition (Sentence) | 1 | 161 |

# 5    Methodologies

## 5.1    Context-Based In-Context Learning

In-Context Learning (ICL) is a transformative approach within large pretrained language models, enabling dynamic adaptation to various tasks based solely on contextual information provided during inference. This process eliminates the need for traditional fine-tuning or gradient updates. Instead, ICL utilizes the extensive pre-trained knowledge embedded within these models, using input context such as task instructions, examples, or specific queries to infer

task requirements and generate task specific outputs. The core strength of ICL lies in the model's ability to recognize patterns, relationships, and task structures directly from input data.

In this study, we use an ICL-based framework, described in Figure 2, and detailed in Appendix A, to classify speech transcriptions derived from the Cookie Theft Test, categorizing them as belonging to either the dementia or control group. The transcriptions exclusively contain patient speech content, omitting any extraneous metadata to maintain focus on linguistic features relevant to cognitive assessment.
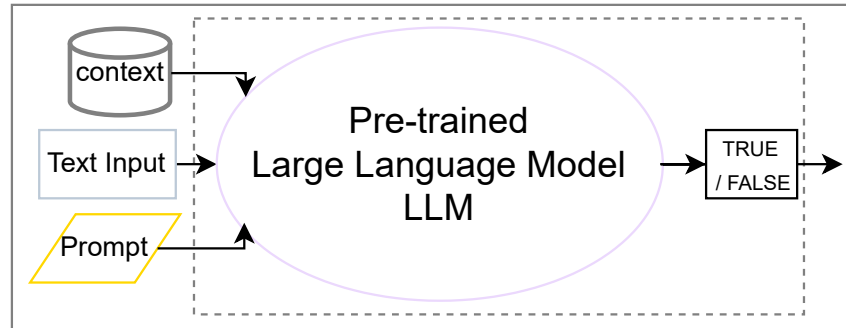


Figure 2: In-Context Learning (ICL) Model Architecture for Dementia Classification Using Structured Prompts with Large Language Models.

The ICL model architecture integrates a structured prompt designed to effectively guide the model in performing its classification task. This prompt comprises three key components:

- Task Definition: The model is instructed to act as an expert in the analysis of conversational data for dementia detection, with specific emphasis on linguistic patterns present in Cookie Theft Test descriptions.

- Classification Criterion: The decision-making process is grounded in the identification of linguistic markers associated with dementia, such as disfluencies, topic drift, reduced syntactic complexity, and other cognitive-linguistic anomalies.

- Response Format: The model is restricted to generating binary outputs, specifically 'TRUE' to indicate the presence of dementia-related markers and 'FALSE' otherwise.

Since LLMs are trained to perform well on general tasks, they struggle to achieve reliable performance on domain-specific data. Thus, providing explicit instructions is essential. An effective way of aiding LLMs in making more accurate decisions is to condition them with a context that includes several examples of data points. Upon receiving this prompt, the pretrained language model processes the transcription, leveraging its in-context reasoning capabilities to generate a classification output.

To assess the efficacy of the ICL approach, we conducted experiments using multiple state-of-the-art large language models. Additionally, recognizing the impact of context-based ICL approach, we explored a context-based multimodal model, which is the primary contribution of this work and the most effective solution.

## 5.2 Context-Based MultiModal Model

The primary contribution of this work is a context-based multimodal model architecture, which improves dementia detection by integrating text and audio data while incorporating contextual information during processing. Our approach leverages advanced pre-trained models, ensuring that both text and speech signals contribute meaningfully to classification. This architecture leverages the complementary strengths of pre-trained Large Language Models (LLMs) and Large Audio Models to capture nuanced cognitive markers. The components of this architecture are illustrated in Figure 3 and Appendix B.

The model processes two primary data modalities:

- Audio Data Processing: The audio recordings are first processed by a pre-trained Large Audio Model, specifically CLAP, which is optimized for capturing intricate acoustic patterns. This model extracts high-

dimensional audio embeddings representing features such as prosody, pitch variations, speech rate, and pauses, which serve as critical markers indicative of cognitive decline.

- Text Data Processing: Simultaneously, the corresponding speech transcripts are processed by a pre-trained LLM, such as GPT, BERT or CLIP, to generate dense text embeddings. These embeddings encapsulate semantic richness, syntactic complexity, and other linguistic markers relevant to dementia detection.

The audio and text embeddings are then fused through element-wise addition, creating a unified multimodal representation. This fusion strategy ensures that both modalities contribute equally to the downstream classification task, allowing the model to learn joint representations that capture the synergy between text and audio cues. The fused embeddings are then fed into a sophisticated classification model, detailed in Figure 4, which employs a cross-attention mechanism. This mechanism enables dynamic interaction between the multimodal embeddings and additional contextual data, enhancing classification performance.

The context consists of balanced samples from both the control and dementia groups, consisting of five samples per group, which are dynamically selected during training to:

- Expose the model to a diverse range of cognitive patterns, thereby enhancing generalization.
- Enable context-aware conditioning, allowing curated contextual inputs during inference to enhance performance without requiring re-training.
- Reduce over-reliance on fixed contexts, promoting adaptability to varying input conditions.

To reduce information loss often associated with attention mechanisms, we incorporate a residual connection. This reintroduces the original embedded inputs into the cross-attention outputs, ensuring the retention of critical features and enhancing model robustness. The output of the cross-attention layer flows through a fully connected classification head, designed for binary classification (dementia vs. control). This head consists of dense layers with non-linear activations, culminating in a softmax output that generates probabilistic predictions.
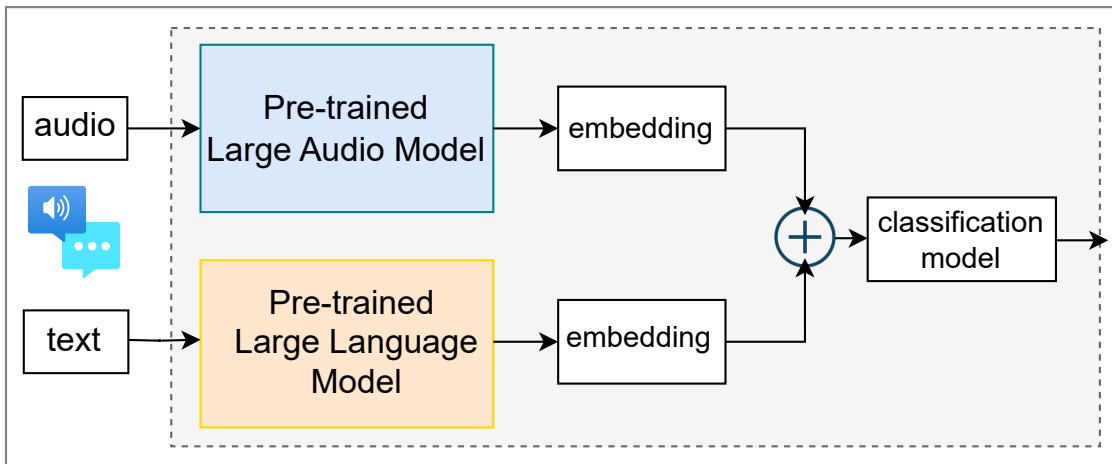


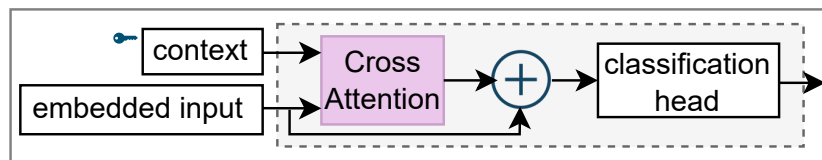Figure 3: Schematic Representation of the Multimodal Model Integrating Text and Audio Features.



Figure 4: Illustration of Cross-Attention Mechanism within the Multimodal Classification Framework.

# 6 Experiments

## 6.1 Experiments Set-up

Informed by the dataset description in Section 4 and the details of the Cookie Theft test outlined in Section 3.1, we designated the Cookie Theft test as the primary task for developing and evaluating methods for early dementia diagnosis. The dataset comprises 99 control participants who collectively contributed 243 tests, and 194 dementia participants with 309 tests. Each test includes both a CHAT-transcribed text file and an accompanying *MP3* audio recording, providing a rich multimodal dataset.

To ensure robust and unbiased evaluation, we employed a repeated stratified sampling approach similar to 10-fold cross-validation but with a fixed test set proportion. Specifically, for each iteration, we randomly selected 30% of the positive class (dementia group) for the test set and an equal number of samples from the negative class (control group) to maintain balance. The remaining data was used for training. This process was repeated 10 times, ensuring that different subsets of the data were evaluated while maintaining a consistent 30% test set size. We opted for this approach instead of strict 10-fold cross-validation to ensure a sufficiently large test set, allowing for a more stable evaluation of model performance. To prevent data leakage and ensure fair evaluation, we enforced a strict constraint to ensure that no participant appeared in both the training and testing sets within the same fold. Additionally, each test fold was balanced to contain an equal representation of control and dementia samples, thereby enhancing the generalizability of the model across unseen individuals. More details about the experimental setup, such as the configurations of the neural networks or specifics about the data preprocessing steps can be found in the appendix section.

For both the multimodal architecture and the In-Context Learning (ICL) approach, we maintained consistency in the context size. Each model was provided with a context window comprising 10 samples, with 5 from control participants and 5 from dementia participants. This balanced context selection was instrumental in conditioning the models to learn from diverse cognitive patterns.

- Multimodal Architecture: The architecture incorporates two standard cross-attention blocks. Information flows through these blocks using residual connections (skip connections) to mitigate vanishing gradient issues and preserve essential feature representations.
- Classification Head: The final classification head consists of two fully connected linear layers. These layers have dimensions of 32 and 16, optimized to handle the fused multimodal embeddings and output binary classification results.

Given the binary nature of the dementia classification task, we assessed model performance using precision, recall, and F1-score. These metrics provide a comprehensive assessment of the models' capability to accurately detect dementia cases (true positives) while minimizing false positives and false negatives. This evaluation framework ensures a balanced analysis of both sensitivity and specificity, which is critical for diagnostic applications.

## 6.2 Context-Based In-Context Learning Experiments

In this section, we explore the effectiveness of In-Context Learning (ICL) paradigms by employing state-of-the-art large language models (LLMs), including GPT-4o OpenAI [2024], Gemini Pro Anil et al. [2023], Gemini 1.5 Pro Team [2024], Claude 3 Anthropic [2024a], Claude 3.5 Sonnet Anthropic [2024b], in the context of early dementia detection. The models were provided with contextual information to contextualize dementia-related linguistic patterns, based on expert annotations highlighting specific markers such as language disfluencies, topic drift, and reduced syntactic complexity. The goal was to determine whether these models could leverage such in-context cues to identify dementia-related language patterns without the need for extensive fine-tuning.

The experimental setup consisted of prompting each LLM with raw text data, accompanied by a carefully crafted context that emphasized key dementia indicators. This approach was designed to simulate a zero-shot or few-shot learning scenario in which model performance depends primarily on its capacity to generalize from the provided context. The results, summarized in Table 3, reveal nuanced differences in performance across the models.

GPT-4o achieved an F1 score of 64.98%, exhibiting strong performance, albeit trailing the top-performing models. Gemini Pro and Gemini 1.5 Pro showed moderate gains, with Gemini 1.5 Pro achieving an F1 score of 68.16%, suggesting that improvements in model architecture and training methodologies enhance ICL performance. Claude 3 outperforms Claude 3.5 with an F1 score of 67.82% compared with 66.83%, indicating strong consistency in detecting dementia-related linguistic features.

Notably, despite the advanced capabilities of these LLMs, their performance in the ICL setting did not surpass the multimodal models that integrated both text and audio data. For instance, the GPT+CLAP multimodal configuration

Table 3: Performance metrics for In Context Learning.

| Text data type | Large Language Model | Precision | Recall | F1 Score | Accuracy |
|---|---|---|---|---|---|
| Raw Text | GPT-4o | 67.86% | 62.32% | 64.98% | 64.34% |
| Raw Text | Gemini Pro | 64.37% | 56.16% | 59.99% | 60.00% |
| Raw Text | Gemini 1.5 Pro | **68.38%** | **67.93%** | **68.16%** | **68.16%** |
| Raw Text | Claude 3 | 72.18% | 63.95% | 67.82% | 68.10% |
| Raw Text | Claude 3.5 sonnet | 70.40% | 63.59% | 66.83% | 67.12% |

achieved an F1 score of 83.33% when applied to raw text data, significantly outperforming the best ICL models. This disparity highlights the benefits of multimodal inputs, where audio features complement textual information to provide a richer, more comprehensive representation of cognitive markers.

In text-only configurations, GPT's fine-tuned model achieved an F1 score of 81.96%, outperforming its ICL counterpart 68.16%. This gap suggests that while LLMs are adept at leveraging in-context information, explicit fine-tuning still offers superior performance in specialized tasks like dementia detection. Additionally, the decline in recall across ICL models, particularly for Gemini Pro 56.16%, highlights difficulties in reliably identifying all dementia-related instances, likely stemming from the subtlety and variability of linguistic symptoms.

Overall, these findings underscore the potential of ICL as a flexible, resource-efficient approach for dementia detection. However, they further highlight the importance of multimodal data and fine-tuning in achieving the highest levels of diagnostic accuracy. The complementary nature of these methodologies suggests a promising avenue for future research, in which the strengths of ICL and multimodal learning may be effectively integrated to advance cognitive health assessments.

## 6.3 Context-Based MultiModal Model Experiments

In this section, we analyze the performance of our proposed dementia detection models using two reported tables: Table 4 reports results on an expert-annotated version of the CHAT-transcribed dataset, and Table 5 reports results on the same dataset without expert annotations (i.e., using raw, unannotated transcriptions). Overall, these experiments evaluate various combinations of pre-trained language models (CLIP, BERT and GPT) alongside an audio model (CLAP) under both single-modal (text-only or audio-only) and multimodal (text+audio) settings.

Table 4: Performance of Pre-trained Language and Audio Models on Expert-Annotated Chat Data for Dementia Detection.

| Input Data | Pretrained LLM | Pretrained Audio Model | Precision | Recall | F1 Score | Accuracy |
|---|---|---|---|---|---|---|
| Text Only | CLIP | — | 71.97% | 72.02% | 72.00% | 72.05% |
| Text Only | BERT | — | 77.25% | 76.19% | 76.72% | 76.72% |
| Text Only | GPT | — | 79.15% | 78.57% | 78.86% | 78.86% |
| Audio Only | — | CLAP | 71.36% | 70.83% | 71.09% | 71.09% |
| Text and Audio | CLIP | CLAP | 73.17% | 73.21% | 73.19% | 73.22% |
| Text and Audio | BERT | CLAP | 79.92% | 78.57% | 79.24% | 79.20% |
| Text and Audio | GPT | CLAP | **80.32%** | **80.36%** | **80.34%** | **80.36%** |

Across text-only configurations, GPT consistently outperforms BERT and CLIP in both expert-annotated and raw-data scenarios. In the expert-annotated dataset, GPT achieves an F1 score of 78.86%, while BERT attains 76.72% and CLIP attains 72.00%. For the raw dataset, GPT's F1 score increases to 81.96%, surpassing BERT's 79.54% and CLIP's 75.56%. This finding indicates GPT's robustness in capturing linguistic cues indicative of dementia, highlighting its capacity to adapt to diverse conversational contexts.

When examining audio-only input using CLAP, performance remains the lowest among all tested configurations, with an F1 score of approximately 71.09%. While audio features provide useful information, relying exclusively on audio is insufficient relative to text-based or multimodal approaches. Nonetheless, combining the audio model with text models substantially improves overall performance. The fusion of text (whether CLIP, BERT, or GPT) with CLAP

Table 5: Performance of Pre-trained Language and Audio Models on Raw (Non-Expert-Annotated) Chat Data for Dementia Detection.

| Input Data | Pretrained LLM | Pretrained Audio Model | Precision | Recall | F1 Score | Accuracy |
|---|---|---|---|---|---|---|
| Text Only | CLIP | — | 76.12% | 75.01% | 75.56% | 75.60% |
| Text Only | BERT | — | 81.37% | 77.79% | 79.54% | 80.50% |
| Text Only | GPT | — | 83.01% | 80.95% | 81.96% | 81.93% |
| Audio Only | — | CLAP | 71.36% | 70.83% | 71.09% | 71.09% |
| Text and Audio | CLIP | CLAP | 78.96% | 78.57% | 78.76% | 78.66% |
| Text and Audio | BERT | CLAP | 81.83% | 81.55% | 81.68% | 81.62% |
| Text and Audio | GPT | CLAP | **83.34%** | **83.33%** | **83.33%** | **83.33%** |

yields higher F1 scores than either text-only or audio-only pipelines. Notably, GPT+CLAP obtains the strongest results, achieving an F1 score of 80.34% on the expert-annotated set and 83.33% on the raw dataset. These trends confirm the complementarity between text and audio features, which together capture a broader spectrum of dementia-related signals.

A noteworthy observation is that the raw (non-expert-annotated) data yields performance that is equal to or even superior to the expert-annotated version. GPT+CLAP, for example, improves from an F1 score of 80.34% in the annotated dataset to 83.33% in the raw dataset. One explanation is that expert annotations may inadvertently introduce confusion by labeling behaviors that can occur in both healthy individuals and those with dementia, thus diminishing model precision. Furthermore, large pre-trained architectures such as GPT and CLAP demonstrate a strong capacity to learn text and audio patterns directly from raw data, reducing the need for detailed manual labeling.

From a practical standpoint, these results are significant because they suggest that expensive and time-consuming expert annotations are not strictly necessary for achieving robust dementia detection. By leveraging pre-trained models capable of extracting nuanced patterns, researchers and clinicians can train effective classifiers even when only raw text is available. Eliminating the requirement for expert annotation also makes the approach more scalable, enabling broader application in clinical and community-based settings. Finally, the consistent gains achieved by combining GPT with CLAP highlight the advantages of multimodal analysis, where language usage and speech properties together yield a more holistic and powerful assessment of cognitive functioning.

## 6.4 Comparison and Discussion

We conducted a comparative analysis of our proposed approach against two state-of-the-art models referenced from recent literature Pan et al. [2025] and BT and Chen [2024]. The models were selected based on their outstanding performance in dementia detection, as reported in their respective studies. To ensure a fair and rigorous comparison, we re-implemented these models by closely following the methodologies and configurations detailed in their original publications. This process included replicating the architectural frameworks, data preprocessing techniques, and training protocols as described by the authors. Before proceeding with the evaluation, we first ensured that our re-implementations replicated the reported results as closely as possible. This step was crucial to validate the correctness of our re-implementation and ensure that any observed differences in performance were due to methodological advancements rather than inconsistencies in replication. Once we confirmed the fidelity of our implementations, we evaluated all models under the same conditions.

Consistency in experimental design was paramount. Therefore, we adhered strictly to the experimental setup outlined at the beginning of Section 6, which includes the use of the Pitt Corpus and a repeated stratified sampling with a fixed 30% test set test set. This ensured that different subsets of the data were evaluated while maintaining a balanced dataset to mitigate bias. This uniform framework ensured that all models were evaluated under identical conditions, providing a robust basis for comparison.

The comparative analysis juxtaposes these state-of-the-art models with our best-performing configuration. We employed the same evaluation metrics precision, recall, F1-score, and accuracy across all models to maintain consistency and objectivity in performance assessment. Our results, shown in table 6, demonstrate that our approach significantly outperforms the selected state-of-the-art models across all metrics. Specifically, our model achieved better precision, recall, and F1-scores, reflecting its enhanced capability in accurately detecting dementia-related patterns. This performance advantage is attributed to our innovative integration of multimodal data, leveraging both text and audio

features through advanced pre-trained models, which capture a broader spectrum of cognitive markers compared to traditional approaches.

Table 6: Performance Comparison of State-of-the-Art Models and Proposed Approach.

| Model | Precision | Recall | F1 Score | Accuracy |
|---|---|---|---|---|
| Pan et al. [2025] Two-step Attention | 79.38% | 76.79% | 76.87% | 77.01% |
| BT and Chen [2024] ICL-bard (Q1+Q2) | 70.15% | 68.75% | 68.92% | 68.99% |
| Context-Based MultiModal (**ours**) | **83.34%** | **83.33%** | **83.33%** | **83.33%** |

Figure 5 presents the confusion matrix for one of the randomly selected folds, illustrating the classification performance of the best context-based multimodal configuration (GPT + CLAP). The matrix highlights the distribution of true positives, true negatives, false positives, and false negatives. Notably, the model achieved a high true positive rate, accurately identifying most dementia cases. However, the presence of false positives indicates that some control participants were misclassified as having dementia.

To better understand the causes of misclassification, we conducted an in-depth analysis of both false positives and false negatives, focusing on text, audio, and multimodal features that may have influenced the model's predictions. The goal was to identify patterns that could inform future model refinements and enhance classification robustness.
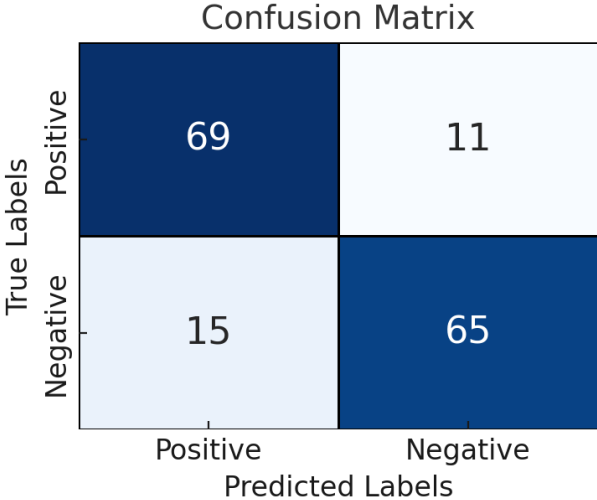


Figure 5: Confusion matrix for one of the 10 folds.

Our analysis revealed that dementia cases misclassified as control often contained structured and coherent narratives, effectively masking subtle cognitive impairments. For instance, clear and detailed descriptions, such as "The mother is drying dishes and the water is cascading onto the floor," made it difficult for the model to detect underlying dementia-related markers. Conversely, control cases misclassified as dementia exhibited frequent hesitations and verbosity, which the model erroneously interpreted as indicators of cognitive decline.

This limitation likely stems from the pre-trained model's ability to capture fine-grained linguistic details in embedding representations. Some subtle but crucial features may not have been sufficiently emphasized in the learned embeddings, leading to increased confusion between the two classes. Future work should explore enhancing feature extraction techniques to better distinguish between genuine cognitive impairment and naturally occurring speech variations.

## 7 Conclusions

In conclusion, this study demonstrates the effectiveness of our context-based multimodal approach as the primary innovation in dementia detection. By explicitly incorporating contextual information in both text and audio processing, we achieved state-of-the-art results while improving model robustness and interpretability. Additionally, inspired by the success of context in multimodal learning, we explored a context-based In-Context Learning (ICL) approach, which,

while promising, did not surpass the multimodal framework. Our findings highlight the superiority of context-aware multimodal models, paving the way for more reliable and scalable dementia detection methods. GPT-based embeddings consistently outperformed BERT and CLIP, showing superior precision, recall, and F1 scores, likely due to their ability to capture intricate linguistic nuances. The integration of audio features via the CLAP model further improved accuracy, with raw, unannotated data outperforming expert-annotated versions. This suggests that large pre-trained models can extract meaningful patterns without extensive manual annotation, enhancing scalability and real-world applicability. These findings are supported by the comparative performance results presented in Table 6, where our proposed model achieved an F1-score of 83.33% and an accuracy of 83.33%, outperforming state-of-the-art models such as the Two-step Attention model Pan et al. [2025] and the ICL-bard (Q1+Q2) model BT and Chen [2024].

Furthermore, this work establishes a foundational framework for future research aimed at providing more personalized insights and feedback to dementia patients. By enabling the development of diagnostic tools that go beyond detection, future studies can focus on delivering actionable cognitive health insights tailored to individual patients. This will not only support early diagnosis but also contribute to personalized intervention strategies, empowering clinicians with deeper cognitive analyses. Ultimately, this approach has the potential to improve patient outcomes, foster more precise monitoring of disease progression, and offer valuable resources for both experts and researchers to better understand the complex dynamics of dementia.

## Funding

## Ethics Statement

This study did not involve direct interaction with human participants or the creation of new patient data. The data used in this study were obtained from the publicly available Pitt Corpus dataset, which was collected and shared by the Alzheimer and Related Dementias Study at the University of Pittsburgh School of Medicine.

## Data Availability

The data used in this study are publicly available as part of the Pitt Corpus dataset. The dataset can be accessed through the DementiaBank project at `https://dementia.talkbank.org/access/English/Pitt.html`, subject to the approval of the data administrators. This study did not generate any new data.

## Conflicts of Interest

The authors declare no conflicts of interest.

## Abbreviations

The following abbreviations are used in this manuscript:

| | |
|---|---|
| BERT | Bidirectional Encoder Representations from Transformers |
| GPT | Generative Pre-trained Transformer |
| CLIP | Contrastive Language-Image Pretraining |
| CLAP | Contrastive Language-Audio Pretraining |
| CHAT | Codes for the Human Analysis of Transcripts |
| MLP | Multilayer Perceptron |
| F1 | Harmonic Mean of Precision and Recall |

## A    Context-based In-Context Learning (ICL) Setup

To ensure a controlled evaluation, all models received the same structured input.

**Note:** This appendix does not contain real patient data. Instead, we present *generic* examples to illustrate the input structure used. All transcript data in this study were anonymized and handled in accordance with ethical research standards.

The content provided to the language models is as follows:

---

**You are a neurologist specializing in speech pattern analysis for dementia detection.**
Your task is to analyze patient transcripts from the Cookie Theft picture description test and classify them as either showing signs of dementia (TRUE) or not (FALSE).
Follow the classification criteria strictly, and answer using only TRUE or FALSE.
**Task:** You will be given a transcript of a patient describing the Cookie Theft picture from the Boston Diagnostic Aphasia Examination. Your task is to classify the speech as either showing signs of dementia (TRUE) or not showing signs of dementia (FALSE) based only on the transcript.
**Key Considerations:**
✓ Normal Aging May Include:

- Occasional word-finding difficulty but coherent speech.

- Slight redundancy or mild repetition of ideas.

- Slower speech but still structured and logical.

✗ Dementia May Include:

- Disorganized or fragmented speech (sentences missing logical connections).

- Frequent grammatical errors (incorrect verb tense, missing words).

- Noticeable forgetting of prior statements (repeating an idea separately).

- Excessive hesitation sounds ("uh," "um," "you know").

- Failure to describe major elements of the picture (key omissions).

- Tangential or unrelated comments (mentioning details that do not exist in the picture).

**Expected Features of a Normal Response:**

- Mentions key elements of the picture, such as:

    – The mother washing or drying dishes.
    – The boy reaching for cookies.
    – The stool falling.
    – The girl near the boy.
    – Water overflowing from the sink.

**Response Format:** Based only on the transcript provided, answer with one word:

- TRUE → Signs of dementia are present.

- FALSE → No signs of dementia.

TRANSCRIPT: <add participant text here>

---

The following are *generic* transcripts created for illustration purposes. They are not real patient data.

Table 7: Example input-output pairs from ICL experiments.

| Input Transcript | Expected Output | Model Response |
|---|---|---|
| "A boy is on a stool reaching for cookies. The mother is washing dishes. Water is overflowing." | FALSE | FALSE |
| "The, uh, kid is, um, doing something? And, um, the lady is, uh, water? I don't know." | TRUE | TRUE |

## B  Context-Based MultiModal Model Implementation Details

### B.1  Data Preprocessing

We adopted the preprocessing pipelines of the respective pre-trained models to ensure consistency. Specifically, for the annotated text data, we ensured that only the patient's speech and the expert annotations were retained, excluding any extraneous information.

## B.2 Model Architecture Configuration

- **Hidden Dimension Alignment:** Audio embeddings are projected using a linear layer to match the hidden dimension of text embeddings.
- **Cross-Attention Mechanism:**
  - **Number of Attention Heads:** 8
  - **Hidden Dimensions:** 768 for BERT, 1536 for GPT
- **Classification Layers:**
  - Fully Connected Layers: $[32, 16]$

## B.3 Training Configuration

- **Loss Function:** Cross-Entropy Loss
- **Optimizer:** AdamW
- **Batch Size:** 4
- **Number of Epochs:** 200
- **Learning Rate:** 0.0001
- **Context Size:** 10

## References

World Health Organization. Dementia, 2023. URL `https://www.who.int/news-room/fact-sheets/detail/dementia`. Accessed: 2025-02-09.

Elsa El Abiad, Ali Al-Kuwari, Ubaida Al-Aani, Yaqoub Al Jaidah, and Ali Chaari. Navigating the alzheimer's biomarker landscape: A comprehensive analysis of fluid-based diagnostics. *Cells*, 13(22):1901, 2024.

Miguel A Chávez-Fumagalli, Pallavi Shrivastava, Jorge A Aguilar-Pineda, Rita Nieto-Montesinos, Gonzalo Davila Del-Carpio, Antero Peralta-Mestas, Claudia Caracela-Zeballos, Guillermo Valdez-Lazo, Victor Fernandez-Macedo, Alejandro Pino-Figueroa, et al. Diagnosis of alzheimer's disease in developed and developing countries: systematic review and meta-analysis of diagnostic test accuracy. *Journal of Alzheimer's Disease Reports*, 5(1):15–30, 2021.

K Zhao, P Chen, A Alexander-Bloch, Y Wei, M Dyrba, F Yang, et al. A neuroimaging biomarker for individual brain-related abnormalities in neurodegeneration (ibrain): a cross-sectional study. eclinicalmedicine 65: 102276, 2023.

Spencer AW Lee, Luciano A Sposato, Vladimir Hachinski, and Lauren E Cipriano. Cost-effectiveness of cerebrospinal biomarkers for the diagnosis of alzheimer's disease. *Alzheimer's Research & Therapy*, 9:1–14, 2017.

Zehra Shah, Jeffrey Sawalha, Mashrura Tasnim, Shi-ang Qi, Eleni Stroulia, and Russell Greiner. Learning language and acoustic models for identifying alzheimer's dementia from speech. *Frontiers in Computer Science*, 3:624659, 2021.

Zerin Jahan, Surbhi Bhatia Khan, and Mo Saraee. Early dementia detection with speech analysis and machine learning techniques. *Discover Sustainability*, 5(1):1–18, 2024.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding, 2019. URL `https://arxiv.org/abs/1810.04805`.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.

Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2, 2024.

Chengsheng Mao, Jie Xu, Luke Rasmussen, Yikuan Li, Prakash Adekkanattu, Jennifer Pacheco, Borna Bonakdarpour, Robert Vassar, Li Shen, Guoqian Jiang, et al. Ad-bert: Using pre-trained language model to predict the progression from mild cognitive impairment to alzheimer's disease. *Journal of Biomedical Informatics*, 144:104442, 2023.

Xinsong Du, John Novoa-Laurentiev, Joseph M Plasek, Ya-Wen Chuang, Liqin Wang, Gad A Marshall, Stephanie K Mueller, Frank Chang, Surabhi Datta, Hunki Paek, et al. Enhancing early detection of cognitive decline in the elderly: a comparative study utilizing large language models in clinical notes. *EBioMedicine*, 109, 2024.

Amit Meghanani, Chandran Savithri Anoop, and AG Ramakrishnan. An exploration of log-mel spectrogram and mfcc features for alzheimer's dementia recognition from spontaneous speech. In *2021 IEEE spoken language technology workshop (SLT)*, pages 670–677. IEEE, 2021.

Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale, 2021. URL `https://arxiv.org/abs/2010.11929`.

Loukas Ilias, Dimitris Askounis, and John Psarras. Detecting dementia from speech and transcripts using transformers. *Computer Speech & Language*, 79:101485, 2023.

Hee Jeong Han, Suhas BN, Ling Qiu, and Saeed Abdullah. Automatic classification of dementia using text and speech data. In *Multimodal AI in healthcare: A paradigm shift in health intelligence*, pages 399–407. Springer, 2022.

Ning Liu, Zhenming Yuan, and Qingfeng Tang. Improving alzheimer's disease detection for speech based on feature purification network. *Frontiers in Public Health*, 9:835960, 2022.

Yilin Pan, Bahman Mirheidari, Daniel Blackburn, and Heidi Christensen. A two-step attention-based feature combination cross-attention system for speech-based dementia detection. *IEEE Transactions on Audio, Speech and Language Processing*, 2025.

OpenAI. Clip: Contrastive language-image pretraining. GitHub repository, 2021. URL `https://github.com/openai/CLIP`. Accessed: 2025-02-08.

Benjamin Elizalde, Soham Deshmukh, Mahmoud Al Ismail, and Huaming Wang. Clap learning audio concepts from natural language supervision. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023.

Brian MacWhinney, Davida Fromm, Margaret Forbes, and Audrey Holland. Pitt corpus: Dementiabank, 2011. URL `https://dementia.talkbank.org/access/English/Pitt.html`. Accessed: 2025-02-02.

Louise Cummings. Describing the cookie theft picture: Sources of breakdown in alzheimer's dementia. *Pragmatics and Society*, 10(2):153–176, 2019.

Brian MacWhinney. Tools for analyzing talk – part 1: The chat transcription format. Technical report, Carnegie Mellon University, October 2024. URL `https://talkbank.org/manuals/CHAT.html`. Accessed: 2025-02-08.

Mondher Bouazizi, Chuheng Zheng, Siyuan Yang, and Tomoaki Ohtsuki. Dementia detection from speech: what if language models are not the answer? *Information*, 15(1):2, 2023.

Aparna Balagopalan, Benjamin Eyre, Frank Rudzicz, and Jekaterina Novikova. To bert or not to bert: comparing speech and language-based approaches for alzheimer's disease detection. *arXiv preprint arXiv:2008.01551*, 2020.

Youxiang Zhu, Xiaohui Liang, John A Batsis, and Robert M Roth. Exploring deep transfer learning techniques for alzheimer's dementia detection. *Frontiers in computer science*, 3:624683, 2021.

Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in neural information processing systems*, 33:12449–12460, 2020.

R'mani Haulcy and James Glass. Classifying alzheimer's disease using audio and text-based representations of speech. *Frontiers in Psychology*, 11:624137, 2021.

Balamurali BT and Jer-Ming Chen. Performance assessment of chatgpt versus bard in detecting alzheimer's dementia. *Diagnostics*, 14(8):817, 2024.

María Lara Gauder, Leonardo Daniel Pepino, Luciana Ferrer, and Pablo Riera. Alzheimer disease recognition using speech-based embeddings from pre-trained models. In *Proc. Interspeech*. International Speech Communication Association, 2021.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

Fasih Haider, Sofia De La Fuente, and Saturnino Luz. An assessment of paralinguistic acoustic features for detection of alzheimer's dementia in spontaneous speech. *IEEE Journal of Selected Topics in Signal Processing*, 14(2):272–281, 2019.

Shawn Hershey, Sourish Chaudhuri, Daniel PW Ellis, Jort F Gemmeke, Aren Jansen, R Channing Moore, Manoj Plakal, Devin Platt, Rif A Saurous, Bryan Seybold, et al. Cnn architectures for large-scale audio classification. In *2017 ieee international conference on acoustics, speech and signal processing (icassp)*, pages 131–135. IEEE, 2017.

Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. Biobert: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4):1234–1240, 2020.

Medical AI. Clinicalbert - hugging face model, 2025. URL `https://huggingface.co/medicalai/ClinicalBERT`. Accessed: 2025-02-08.

Jer-Ming Chen et al. Performance assessment of chatgpt vs bard in detecting alzheimer's dementia. *arXiv preprint arXiv:2402.01751*, 2024.

Youxiang Zhu, Nana Lin, Xiaohui Liang, John A Batsis, Robert M Roth, and Brian MacWhinney. Evaluating picture description speech for dementia detection using image-text alignment. *arXiv preprint arXiv:2308.07933*, 2023.

Kaiying Lin and Peter Y Washington. Multimodal deep learning for dementia classification using text and audio. *Scientific Reports*, 14(1):13887, 2024.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.

Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR, 2023.

Boston Diagnostic Aphasia Examination. The cookie theft picture from the boston diagnostic aphasia examination, 2025. URL `https://www.researchgate.net/figure/The-Cookie-Theft-Picture-from-the-Boston-Diagnostic-Aphasia-Examination-For-the-PD-task_fig1_349613269`. Available under Creative Commons Attribution 4.0 International (CC BY 4.0).

James T. Becker, François Boller, Oscar L. Lopez, John Saxton, and Kathleen L. McGonigle. Pitt corpus, 1994. URL `https://talkbank.org/`. Part of the Alzheimer's and Related Dementias Study, University of Pittsburgh School of Medicine.

OpenAI. Hello gpt-4o. OpenAI Website, 2024. URL `https://openai.com/index/hello-gpt-4o/`. Accessed: 2025-02-08.

Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: A family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 1, 2023.

G Gemini Team. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context (2024). *URL https://goo. gle/GeminiV1-5*, 2024.

Anthropic. Introduction to claude ai. Anthropic Documentation, 2024a. URL `https://docs.anthropic.com/en/docs/intro-to-claude`. Accessed: 2025-02-08.

Anthropic. Claude 3.5 sonnet: Advancing ai performance and efficiency. Anthropic News, 2024b. URL `https://www.anthropic.com/news/claude-3-5-sonnet`. Accessed: 2025-02-08.