

The Eyes Don't Lie: Text Transcriptions Can Hide Dementia Presentation that Gaze Reveals

Anonymous ACL submission

Abstract

Current methods used to diagnose or monitor dementia-related cognitive decline predominantly rely on audio recordings. Such audio recordings can leak personally identifiable information and create new risks given deep fake technology. We introduce generative likelihood-based approaches to identify differences in healthy versus dementia-diagnosed participants via gaze tracking and text transcriptions during a standard diagnostic image description task without relying on sensitive audio information. Contrasting conventional wisdom, we find that text transcriptions alone are not a reliable measure of cognitive impairment in this task, finding gaze tracking to be more reliable, and suggesting existing results in language-based dementia detection rely primarily on audio signals.

1 Introduction

Continual monitoring of cognitive change can enable early detection of Alzheimer's and related dementias, facilitating earlier intervention and treatment (Rasmussen and Langerman, 2019). Existing dementia detection resources and methods largely focus on audio and text transcriptions, but we find that reliable detection from short interactions with participants is achievable through *gaze tracking* in tandem with text transcriptions.

Around 70% of Americans said they would want Alzheimer's disease identified if that knowledge led to earlier treatment (Alzheimer's Association, 2023), but available clinically validated measures of cognitive change for early detection of Alzheimer's take place at most every three (P et al., 2009; CB et al., 2023) to six (KV et al., 2024) For some, these important checks don't take place at all until after advanced symptoms are present. Developing computational models to detect the onset of dementia-related cognitive decline in time for medical intervention and evaluation is an under-



there's a lady washing oh the sink is flowing over and the the chair a reaching for the cookie jar and the girls gonna gonna the pretty busy picture ah that's it now that doesn't doesn't at home because your sink never fills with water that that i don't think yeah the man who creates these drawings and and outside the window you see the yard and trees oh oh dishes dishes they don't move we know the boys going to fall hmm you you try to draw portrait pictures that's quite a deal to do to to pictures of people's faces oh it's time for a change

Figure 1: We investigate training-free dementia detection methods from gaze tracking (colored dots) and transcript text (colored words) of participants describing The Cookie Theft Picture.

explored problem, as most existing ML detection methods are based on data from a single assessment and modality of interaction, such as speech (Becker et al., 1994; Luz et al., 2020).

The speech data from those existing works often includes a verbal task where participants spend up to two minutes describing a line drawing scene (Figure 1). This task is a component of the The Boston Diagnostic Aphasia Examination (Goodglass et al., 2001) frequently used by clinicians for screening for dementia symptom presentation. Existing works that train machine learning models to detect the presence of dementia symptoms largely focus on audio signals or hand-crafted features summarizing aspects of text transcripts (Santander-Cruz et al., 2022; Kumar et al., 2022; Javeed et al., 2023; Shi et al., 2023). These approaches typically train simple classifiers such as SVM, Random Forest, and logistic regression (Diogo et al., 2022; Haider et al., 2020), or fine-tune existing pretrained

models such as BERT (Balagopalan et al., 2020), RoBERTa (Matošević and Jović, 2022), and GPT-2 (Liu and Wang, 2023).

In this paper, we take a step towards non-invasive, privacy-preserving, in-home monitoring tools for detecting early signs and symptoms of dementia. We explore analysis methods on *raw* gaze and text data that require no hand-crafted features, federated learning across participants, or even back propagation gradient passes on existing models, all of which can inadvertently leak personally identifiable information. In short, we explore methods to detect dementia symptoms from the under explored spaces of gaze tracking and verbal text transcriptions. The contributions of the paper can be summarized as follows:

- We empirically demonstrate that the gaze of Control group focus on the areas-of-interest presented in The Cookie Theft Picture compared to that of participants with AD.
- We similarly demonstrate that the Control group’s text transcript descriptions of The Cookie Theft Picture correspond more closely to the expectations of large, pretrained image captioning model when compared to participants with AD.
- Our analyses do not rely on hand-crafted features related to analyzing dementia presentation, and instead leverage pretrained models and statistical machine learning models to measure deviation from expected gaze patterns and sequences of descriptive words in terms of likelihood without any additional model training or fine-tuning.

2 Participant Gaze and Text Data

We analyze a dataset of participant tracked eye gaze and human-corrected transcripts of participant speech during the completion of The Cookie Theft Description Task. The study included 25 Control group participants with healthy cognitive function and 14 participants with an Alzheimer’s Disease (AD) diagnosis. Participants were all patients at a local aging research center, at which they were also recruited for enrollment in the study.

During each participant session, we recorded eye gaze and audio while the participant viewed The Cookie Theft Picture on a Surface Laptop Studio equipped with an Intel Core i7 processor, 32 GB of RAM, a 1TB SSD, Microsoft OS, an NVIDIA GeForce RTX graphics card, and a Tobii Pro X3-

120 eye tracker. Eye tracking was calibrated using Tobii Manager software, with gaze data gathered via the Tobii Pro SDK 3¹. Audio was processed to speech transcriptions standardized using the Automatic Speech Recognition (ASR) Vosk Model² followed by manual annotation by a person to correct any ASR errors. After removing gaze points from timesteps when none was tracked and cleaning up text transcriptions, we have an average of 8312.82 ± 4993.53 gaze points and 167.87 ± 58.47 transcribed words of description across the 39 total participants whose sessions lasted, on average, 94.89 ± 21 seconds.

3 Hypotheses and Methods

The methods described in the paper do not involve training algorithms based on participant data or processing participant data for any sort of hand-crafted feature extraction. Instead, these methods utilize pre-trained, generative models to estimate likelihoods of observed data being generated by a background, “healthy” distribution. We hypothesize that:

- H1** gaze points collected from participants in the AD group will exhibit lower likelihood of gaze being explained by annotated areas of interest in the stimulus image than will the Control group; and
- H2** text transcribed from audio of participants in the AD group will exhibit lower likelihood due to syntactic fluency and topic consistency (**H2**₁) as well as relevance to the stimulus image (**H2**₂).
- H3** gaze and text will reveal complementary participant cognitive function.

For the purposes of evaluating our hypothesis, we calculate the average *Negative Log-Likelihood* (NLL) of sequences of gaze points and transcription words for each participant. Note that a lower NLL corresponds to a lower likelihood, while a high NLL indicates a higher likelihood.

3.1 Gaze: Semantic GMM

We fit a Gaussian Mixture Models (GMM) to Areas of Interest (AOI) in The Cookie Theft Picture annotated by an experimenter. We learn a $k = 17$ component mixture of Gaussians, each defined by a mean (μ_k), covariance (σ_k), and mixing coefficient

¹<https://developer.tobii.com/python/python-oldmigrationsdk.html>

²<https://alphacephei.com/vosk/models/vosk-model-en-us-0.21.zip>

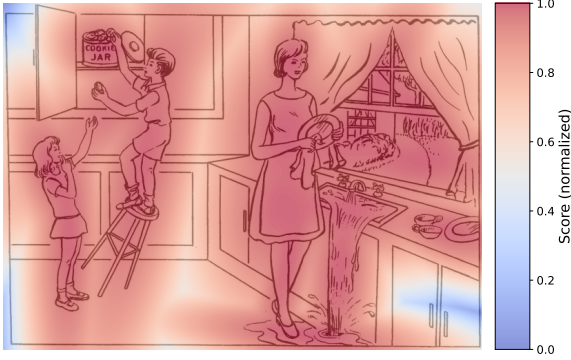


Figure 2: Heatmap displaying the likelihood estimations by the GMM across The Cookie Theft Picture.

(π_k). Figure 2 visualizes the likelihood heatmap by pixel in the stimulus image of this fitted ‘‘Semantic GMM.’’ We calculate the average NLL for a set of gaze points $\{x, y\}^N$ by:

$$\overline{\text{NLL}}_{\text{gaze}}(\{x, y\}^N) = -\frac{1}{N} \sum_{n=1}^N \ln \left(\sum_{k=1}^K \pi_k \mathcal{N}((x_n, y_n) | \mu_k, \Sigma_k) \right).$$

We calculate this average log-likelihood per participant, then analyze the differences in these gaze likelihood samples between the control and AD populations.

3.2 Text Transcripts: Pretrained LLMs

We utilize two pretrained large language models (LLMs) that decode autoregressively and can be run on-device to calculate the average likelihood of the sequence of transcribed tokens from participant descriptions. Given trained LMM parameters θ yielding a distribution $p_\theta(x_i | x_{1..i-1})$ of next token probability, the average log-likelihood of a token sequence $\vec{x} := \text{LLM-Tokenizer}(\vec{w})$ from participant transcript word sequence \vec{w} is calculated as:

$$\overline{\text{NLL}}_{\text{text}}(\vec{x}) = -\frac{1}{N} \sum_{i=1}^N \ln p_\theta(x_i | x_{1..i-1}).$$

GPT-2 (Radford et al., 2019) is a transformer-based model that was pre-trained on substantial English data using self-supervised learning techniques, primarily focusing on predicting the next word in sentences. We can consider the GPT-2 NLL values to represent the *prior* likelihood of text, where differences in NLL scores are likely to correspond to syntactic fluency and topic consistency (**H2**₁). We use the GPT-2 Large model

which can be run on-device, and break transcripts into tokens using the GPT-2 Large tokenizer.

BLIP, Bootstrapping Language-Image Pre-training (Li et al., 2022), is a model pretrained on large-scale image-text datasets using self-supervised learning techniques to autoregressively predict textual descriptions of input images. The BLIP NLL values represent *posterior* likelihoods of text descriptions conditioned on The Cookie Theft Picture stimulus, and may expose more nuanced differences in semantic relevance between control and AD participants (**H2**₂). We use the BLIP-image-captioning-base³, a BLIP processor which wraps a BERT tokenizer⁴ and BLIP image processor into a single processor. For a fair comparison against GPT-2, we additionally test BLIP with a blank image input, treating it as another *prior* likelihood measure in that case.

4 Experiments and Results

The experimental results reveal that significant differences exist between the Control and AD groups when analyzing eye gaze data. However, the differences between text transcripts are less consistent. In multimodal analyses combining eye gaze and text transcripts, the Hotelling T-square indicate significant differences between the two groups when using GPT-2.

Gaze reveals AD symptoms. To evaluate hypothesis **H1**, we compared the population of $\overline{\text{NLL}}_{\text{gaze}}$ values of the 25 control patients to those of the 14 patients with an AD diagnosis using a one-sided, Welch’s unequal variances *t*-tests. Figure 3(a) shows histograms of $\overline{\text{NLL}}_{\text{gaze}}$ values between the populations. The average $\overline{\text{NLL}}_{\text{gaze}}$ of the control group was found to be statistically significantly higher than that of the AD group, with *p*-value .0158, providing supporting evidence for **H1**.

Transcription text is not enough. To evaluate **H2**, we compared populations of $\overline{\text{NLL}}_{\text{text}}$ values between 24 control and 14 AD patient groups using autoregressive text-only and image-conditioned LLMs using one-sided, Welch’s unequal variances *t*-tests. Figures 3(b), 3(c), and 3(d) show the distribution of $\overline{\text{NLL}}_{\text{text}}$ values between each population as estimated by GPT-2, BLIP with a blank conditioning image, and BLIP conditioned on The

³<https://huggingface.co/Salesforce/blip-image-captioning-base>

⁴https://huggingface.co/docs/transformers/v4.41.3/en/model_doc/bert#transformers.BertTokenizerFast

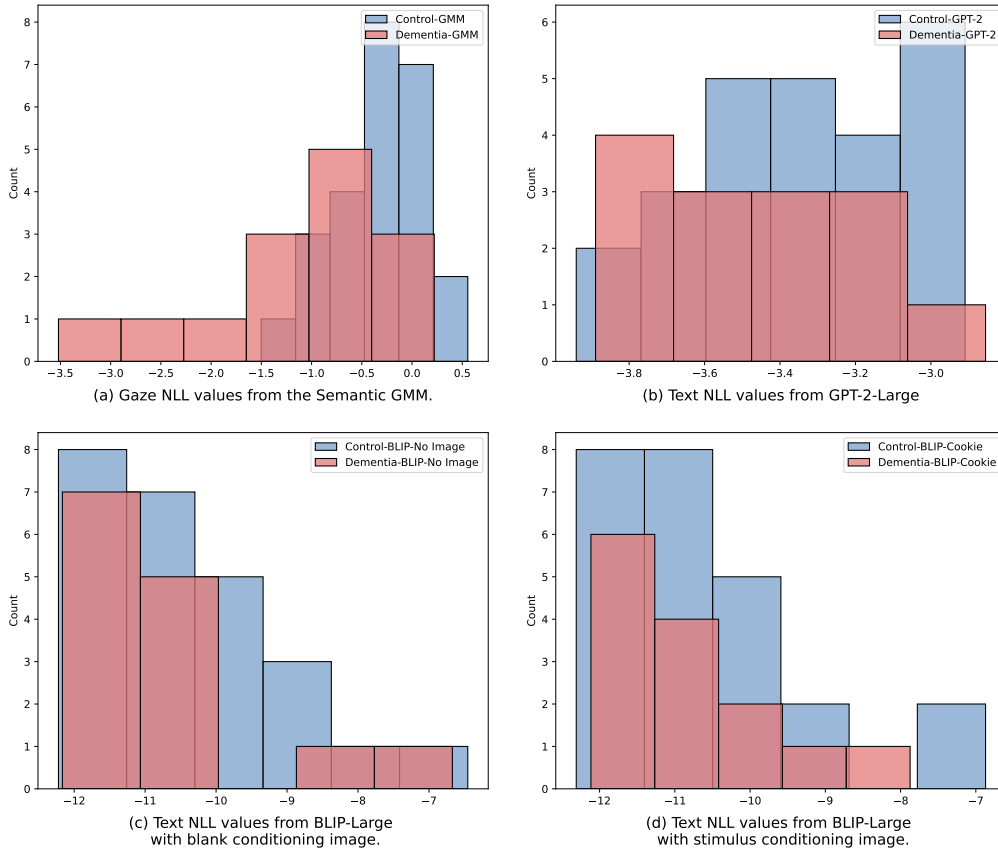


Figure 3: Average NLL values from control group gaze and text transcripts estimated via the Semantic GMM (a), GPT-2-Large (b), BLIP-Large conditioned on a blank image (c), and BLIP-Large conditioned on the stimulus image.

233 Cookie Theft Picture stimulus image. The correspond- 255
 234 ing p -values are .0795, .317, and .355, respec- 256
 235 tively. These results suggest that there may 257
 236 be support for $H2_1$, that there are measurable 258
 237 likelihood-based differences in control versus AD
 238 patient transcripts with respect to syntactic fluency
 239 and topical consistency (as measured by GPT-2;
 240 $p = 0.0795$). However, the image-conditioned
 241 BLIP model, with both a blank image and the ac-
 242 tual stimulus image, show no substantial differenti-
 243 ation in likelihood estimates of transcription tokens
 244 between the groups; we suspect this result may
 245 arise from the misalignment between BLIP’s im-
 246 age caption language pretraining data and the long
 247 form text transcription descriptions of images.

248 **Transcription May Not Complement Gaze.** We
 249 used a Multivariate Hotelling’s T-square test to
 250 compare participant \overline{NLL}_{gaze} and \overline{NLL}_{text} data si-
 251 multaneously. This multivariate population dif-
 252 ference was found statistically significant, but we
 253 repeated the test with identical 0 values substituted
 254 for \overline{NLL}_{text} for all participants also found signif-

255 icance. Our findings do not support **H3**. Partici-
 256 pant \overline{NLL}_{text} contributed no significant information
 257 about the presenc or absence of dementia symp-
 258 toms compared to \overline{NLL}_{gaze} alone.

259 5 Future Work

260 While our areas of interest for gaze analysis are
 261 hand-annotated, we note that pretrained segmen-
 262 tation models such as Meta AI’s Segment Any-
 263 thing (Kirillov et al., 2023) may handle line draw-
 264 ings like The Cookie Theft Picture. Additionally,
 265 methods like MDETR (Kamath et al., 2021) can
 266 identify image regions corresponding to input lan-
 267 guage, opening another way to measure alignment
 268 of participant transcripts. Similarly, while our tran-
 269 scriptions are hand-corrected, we note that the ASR
 270 system produced an estimated WER rate of only
 271 5.43, and that future work may be able to incorpo-
 272 rate visual priors from the image itself to improve
 273 automatic transcription (Chang et al., 2023).

274 Limitations

275 We acknowledge that our study is based on a small
276 sample of 39 participants, and the demographics
277 are not balanced. Specifically, 71% of the par-
278 ticipants are white Caucasians, and there is a 1
279 to 2 ratio of Alzheimer’s Disease (AD) patients
280 to healthy controls. This demographic imbalance
281 may limit the generalizability of our findings to
282 the broader population. However, we believe that
283 our analysis highlights the value of methods like
284 estimation log-likelihood for small datasets in both
285 unimodal and multimodal approaches to dementia
286 assessment. Our findings demonstrate the potential
287 of using limited data effectively, offering evalua-
288 tion metrics that can be applied to other multimodal
289 tasks where access to large datasets is restricted.

290 Ethical Impact

291 This study recognizes the ethical concerns regard-
292 ing privacy and potential information leakage in
293 the collection and analysis of eye gaze data and
294 text transcripts. To address these issues, we have
295 implemented stringent data protection protocols,
296 including anonymization, secure storage, and strict
297 access controls. Informed consent was obtained
298 from all participants, ensuring they understand
299 how their data will be used and protected. Our
300 research team is dedicated to continuously improv-
301 ing our practices to uphold the highest ethical stan-
302 dards, ensuring that the benefits of our research are
303 achieved without compromising participant privacy
304 and trust.

305 References

306 Alzheimer’s Association. 2023. 2023 Alzheimer’s dis-
307 ease facts and figures. *Alzheimer’s & Dementia*,
308 19(4):1598–1695.

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

314 J. T. Becker, F. Boller, O. L. Lopez, J. Saxton, and K. L.
315 McGonigle. 1994. The natural history of alzheimer’s
316 disease: description of study cohort and accuracy of
317 diagnosis. *Archives of Neurology*, 51(6):585–594.
318 Grant Support: NIA AG03705 and AG05133.

319 Young CB, Mormino EC, Poston KL, Johnson KA,
320 Rentz DM, Sperling RA, and Papp KV. 2023. Com-
321 puterized cognitive practice effects in relation to amy-
322 loid and tau in preclinical Alzheimer’s disease: Re-

sults from a multi-site cohort. *Alzheimers Dement*
(Amst)., 15(1).

Allen Chang, Xiaoyuan Zhu, Aarav Monga, Seoho Ahn,
Tejas Srinivasan, and Jesse Thomason. 2023. **Multi-
modal speech recognition for language-guided em-
bodied agents**. In *Annual Conference of the Interna-
tional Speech Communication Association (INTER-
SPEECH)*.

Vasco Sá Diogo, Hugo Alexandre Ferreira, Diana Prata,
and Alzheimer’s Disease Neuroimaging Initiative.
2022. Early diagnosis of alzheimer’s disease us-
ing machine learning: a multi-diagnostic, general-
izable approach. *Alzheimer’s Research & Therapy*,
14(1):107.

Harold Goodglass, Edith Kaplan, and Sandra Weintraub.
2001. *BDAE: The Boston Diagnostic Aphasia Exam-
ination*. Lippincott Williams & Wilkins Philadelphia,
PA, n/a.

Fasih Haider, Sofia De La Fuente Garcia, Pierre Al-
bert, and Saturnino Luz. 2020. Affective speech
for alzheimer’s dementia recognition. *LREC: Re-
sources and Processing of linguistic, para-linguistic
and extra-linguistic Data from people with various
forms of cognitive/psychiatric/developmental impair-
ments (RaPID)*, pages 67–73.

Ashir Javeed, Ana Luiza Dallora, Johan Sanmartin
Berglund, Arif Ali, Liaqata Ali, and Peter Anderberg.
2023. Machine learning for dementia prediction:
a systematic review and future research directions.
Journal of medical systems, 47(1):17.

Aishwarya Kamath, Mannat Singh, Yann LeCun, Is-
han Misra, Gabriel Synnaeve, and Nicolas Carion.
2021. MDETR—modulated detection for end-to-
end multi-modal understanding. *arXiv preprint*
arXiv:2104.12763.

Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi
Mao, Chloe Rolland, Laura Gustafson, Tete Xiao,
Spencer Whitehead, Alexander C. Berg, Wan-Yen
Lo, Piotr Dollár, and Ross Girshick. 2023. Segment
anything. *arXiv:2304.02643*.

M Rupesh Kumar, Susmitha Vekkot, S Lalitha, Deepa
Gupta, Varasiddhi Jayasuryaa Govindraj, Kamran
Shaukat, Yousef Ajami Alotaibi, and Mohammed
Zakariah. 2022. Dementia detection from speech us-
ing machine learning and deep learning architectures.
Sensors, 22(23):9311.

Papp KV, Jutten RJ, Soberanes D, Weizenbaum E,
Hsieh S, Molinare C, Buckley R, Betensky RA, Mar-
shall GA, Johnson KA, Rentz DM, Sperling R, and
Amariglio RE. 2024. Early detection of amyloid-
related changes in memory among cognitively unim-
paired older adults with daily digital testing. *Ann
Neurol.*, 95(3):507–517.

Junnan Li, Dongxu Li, Caiming Xiong, and Steven
Hoi. 2022. Blip: Bootstrapping language-image pre-
training for unified vision-language understanding

379 and generation. In *International conference on ma-*
380 *chine learning*, pages 12888–12900. PMLR.

381 Ning Liu and Lingxing Wang. 2023. An approach for
382 assisting diagnosis of alzheimer’s disease based on
383 natural language processing. *Frontiers in Aging Neu-*
384 *roscience*, 15.

385 Saturnino Luz, Fasih Haider, Sofia de la Fuente, Davida
386 Fromm, and Brian MacWhinney. 2020. Alzheimer’s
387 dementia recognition through spontaneous speech:
388 The ADReSS challenge. In *Annual Conference of*
389 *the International Speech Communication Association*
390 *(INTERSPEECH)*.

391 Lovro Matošević and Alan Jović. 2022. Accurate
392 detection of dementia from speech transcripts us-
393 ing roberta model. In *2022 45th Jubilee Interna-*
394 *tional Convention on Information, Communication*
395 *and Electronic Technology (MIPRO)*, pages 1478–
396 1484. IEEE.

397 Maruff P, Thomas E, Cysique L, et al. 2009. Validity
398 of the CogState brief battery: relationship to stan-
399 dardized tests and sensitivity to cognitive impairment
400 in mild traumatic brain injury, schizophrenia, and
401 AIDS dementia complex. *Arch Clin Neuropsychol.*,
402 24(2):165–178.

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

407 Jill Rasmussen and Haya Langerman. 2019.
408 Alzheimer’s disease — why we need early diagnosis.
409 *Degenerative neurological and neuromuscular*
410 *disease*, 9:123–130.

411 Yamanki Santander-Cruz, Sebastián Salazar-Colores,
412 Wilfrido Jacobo Paredes-García, Humberto
413 Guendulain-Arenas, and Saúl Tovar-Arriaga. 2022.
414 [Semantic feature extraction using sbert for dementia](#)
415 [detection](#). *Brain Sciences*, 12(2).

416 Mengke Shi, Gary Cheung, and Seyed Reza Shahamiri.
417 2023. Speech and language processing with deep
418 learning for dementia diagnosis: A systematic review.
419 *Psychiatry Research*, page 115538.