

High frequency post-pause word choices and task-dependent speech behavior characterize connected speech in individuals with mild cognitive impairment

Michael J. Kleiman, PhD^{1*} and James E. Galvin, MD, MPH¹

¹ Comprehensive Center for Brain Health, Department of Neurology, University of Miami Miller School of Medicine, Boca Raton, FL 33433

Running Title: Post-pause speech behavior in MCI (34 characters)

Title: 18 words, 149 characters (with spaces)

Abstract: 249 words (250 max)

Text: 5772 words

Tables: 6

Figures: 4

References: 68

Total Pages (including references, figure, figure legends, and tables): 38

***Corresponding Author:**

Michael J. Kleiman, PhD

Comprehensive Center for Brain Health

University of Miami Miller School of Medicine

7700 W Camino Real, Suite 200

Boca Raton, FL 33433

E-mail: mjkleiman@miami.edu

ABSTRACT

Background: Alzheimer’s disease (AD) is characterized by progressive cognitive decline, including impairments in speech production and fluency. Mild cognitive impairment (MCI), a prodrome of AD, has also been linked with changes in speech behavior but to a more subtle degree.

Objective: This study aimed to investigate whether speech behavior immediately following both filled and unfilled pauses (post-pause speech behavior) differs between individuals with MCI and healthy controls (HCs), and how these differences are influenced by the cognitive demands of various speech tasks.

Methods: Transcribed speech samples were analyzed from both groups across different tasks, including immediate and delayed narrative recall, picture descriptions, and free responses. Key metrics including lexical and syntactic complexity, lexical frequency and diversity, and part of speech usage, both overall and post-pause, were examined.

Results: Significant differences in pause usage were observed between groups, with a higher incidence and longer latencies following these pauses in the MCI group. Lexical frequency following filled pauses was higher among MCI participants in the free response task but not in other tasks, potentially due to the relative cognitive load of the tasks. The immediate recall task was most useful at differentiating between groups. Predictive analyses utilizing random forest classifiers demonstrated high specificity in using speech behavior metrics to differentiate between MCI and HCs.

Conclusions: Speech behavior following pauses differs between MCI participants and healthy controls, with these differences being influenced by the cognitive demands of the speech tasks. These post-pause speech metrics can be easily integrated into existing speech analysis paradigms.

Keywords: Mild cognitive impairment, speech, verbal behavior, neuropsychological assessment, machine learning

1 INTRODUCTION

2 Alzheimer’s disease (AD) is characterized by progressive degradation of cognitive abilities that
3 interfere with everyday functioning. While episodic memory is often impacted early, other
4 domains also may show early deficits. The regions associated with language and speech production
5 are often among the earliest affected [1,2]. Recent studies have provided insight into how AD
6 affects speech production, including the increased use of more common words [3,4], a potentially
7 compensatory mechanism of declining executive functioning, and less content-filled language [5]
8 which may reflect degradation of semantic memory. Individuals with AD have also been found to
9 exhibit increased pauses during speech [5–7], potentially indicative of difficulties in planning and
10 increased lexical search.

11 Mild cognitive impairment (MCI), a prodromal stage of AD [8], has been associated with the
12 degradation of more subtle measures of speech behavior in its earliest stages. Measures of speech
13 fluency, including rate of speech and syllable count, have been shown to significantly degrade in
14 MCI [9,10]. These are generally measures of connected speech performance; a form of speech in
15 which each individual word is “connected” lexically to the next, as in a passage or narrative. One
16 of the most common methods of obtaining a sample of connected speech is directing the participant
17 to describe a scene or image, such as the Cookie Theft image [11], which enables participants to
18 produce a descriptive sample of speech that is also restricted to a particular topic to promote
19 scoring ability. A more open-ended approach is to simply ask the participant to freely answer a
20 given question, such as “what did you do yesterday?”, which prompts the production of more
21 natural and conversational speech behavior. When scoring episodic memory is more important
22 than capturing naturalistic speech, a paragraph recall task may also be used, as in the Wechler
23 Logical Memory or Craft Story tests [12,13]. In these tasks, a story is read or shown to a participant

24 who is then asked to recall as many details about the story as possible both immediately following
25 the story, as well as after a delay.

26 In connected speech tasks, some studies find that MCI is associated with more common, less
27 complex language [14–16], which may be related to a delay in accessing lexical information [17]
28 resulting in a tendency to use high frequency (more common) words, however some other studies
29 do not observe this increase in lexical frequency within MCI participants [18]. This trend towards
30 common language is most visible in more severe forms of dementia, including primary progressive
31 aphasia [19,20] and AD [3,4], suggesting that while the limited findings of increased lexical
32 frequency due to MCI may be reflective of markers of future decline, the heterogeneity of the
33 disorder may make it more difficult for these markers to be detected.

34 Pauses, or periods between clauses where a speaker hesitates or stops speaking, have also been
35 shown to be affected in MCI. Pauses are characterized as either filled or unfilled. While unfilled
36 pauses consist of extended pauses between words without any vocalization, often defined at greater
37 than 250ms [21] but occasionally at a lower threshold or with no minimum threshold [22], filled
38 pauses include utterances such as “uh”, “um”, and “er”, and are often used when a speaker is
39 attempting to think of the next word or phrase, marking a region in which the speaker is engaging
40 in word retrieval or language planning [23]. The usage of different “forms” of filled pauses have
41 also been shown to vary in their usage; “uh” is typically used to indicate a short pause, while “um”
42 is more often followed by longer pauses [24]. The frequency of these pause forms differ based on
43 the task administered [25] which matches findings showing that pause usage in MCI changes based
44 on the type of assessment; MCI is associated with increases in only filled pauses in free-response
45 tasks [26,27] and either only unfilled pauses [27] or both filled and unfilled pauses [28] in narrative
46 recall tasks.

47 Given that pauses are often produced during word searching behavior [23] and impaired word
48 searching behavior has been found to lead to increased lexical frequency in MCI [14,16,17], we
49 hypothesize that words immediately following pauses (“post-pause” behavior) are more likely to
50 be higher frequency and syntactically simple in MCI participants compared to post-pause word
51 choices in cognitively-intact healthy controls (HCs). Further, previous studies have shown that
52 speech behavior is variably affected by task [16,27,29], as well as by the cognitive demands
53 involved in performing that task [10]. As a result, we expect that differences in pauses and post-
54 pause word choices will differ based on the degree of task demands for each of the three tasks that
55 we will examine: narrative recall (immediate and delayed), picture descriptions, and free response.
56 In particular, both the immediate and delayed narrative recall tasks which involve participants
57 engaging in episodic memory recall are more difficult and require more attention and focus
58 compared to tasks with lesser demands such as describing a picture or simply answering a question.
59 We thus secondarily hypothesize that not all pauses are lexically-driven; more difficult tasks,
60 especially the delayed narrative recall task, may increase pauses but without affecting post-pause
61 lexical frequency, and vice versa.

62 **METHODS**

63 **Participants**

64 53 participants (39 HCs and 14 participants with MCI) were selected from the Healthy Brain
65 Initiative, a longitudinal study of brain health and cognition at the University of Miami
66 Comprehensive Center for Brain Health [30] based on their consensus clinical diagnosis (no
67 cognitive impairment or mild cognitive impairment) and their age (greater than 60 years old). A
68 full description of the study protocol is described elsewhere [30], but in-brief every participant
69 undergoes identical annual comprehensive clinical, cognitive, functional, and behavioral

70 assessments, including the Clinical Dementia Rating (CDR) [31], a complete physical and
71 neurological examination, and neuropsychological test battery from the Uniform Data Set [32],
72 supplemented with additional tests. Healthy controls are defined as individuals with a global CDR
73 of 0, unimpaired activities of daily living, and normal age- and education-normed
74 neuropsychological test performance. Individuals with MCI were defined as a global CDR 0.5,
75 unimpaired activities of daily living, and neuropsychological test performance in at least 1 domain
76 greater than 1.5 standard deviations below the norm. Imaging and plasma biomarker collection is
77 ongoing so only clinical diagnoses were considered here. Exclusion criteria for this study included
78 presence of clinically diagnosed aphasia, schizophrenia, generalized anxiety disorder or major
79 depressive disorder, non-Alzheimer’s dementias, other neurological disorders including
80 Parkinson’s disease or stroke, and cancer diagnosed within the past five years has not been
81 determined to be in remission by a physician. All participants identified English as their primary
82 language. All protocols were approved by the Institutional Review Board at the University of
83 Miami Miller School of Medicine.

84 **Equipment and Materials**

85 Participants were seated in a quiet area with minimal visual and auditory distractions in front of a
86 24” computer monitor attached to a Windows 10 machine. Participants interacted with the machine
87 using custom-developed software written in Python and C, as well as a backlit keypad for
88 indicating intention (finished with response, need help, etc.) and a volume knob. Participants wore
89 Sony WH1000XM4 noise-cancelling headphones to further reduce auditory distractions. The
90 software displayed visual and auditory stimuli for all tasks as well as written and audio instructions.
91 Participants were given the opportunity to indicate whether they did not understand instructions
92 and needed further explanation, after which an experimenter would assist them.

93 **Procedure**

94 *Speech Tasks*

95 Three speech-focused tasks were administered to participants: narrative recall with delay, picture
96 description, and a spontaneous free response. Participant audio was recorded using a hyper-
97 cardioid microphone to minimize background noise.

98 The **narrative recall (NR)** task uses a custom-developed story (“Maddy the dog”) structured to
99 be similar to the Craft Story 21 and Weschler Logical Memory narratives, with 25 content units of
100 varying complexity [33] and developed to be used alongside the Craft Story 21 [13] which is
101 administered to all participants in a separate visit as part of their longitudinal evaluations. This
102 story differs from traditional narrative assessment in that it is both verbally read and visually
103 presented to participants, promoting stronger multimodal encoding than simple auditory
104 presentation. Additionally, the delayed portion was collected between fifteen and twenty minutes
105 after initial presentation of the story, in contrast to the traditional 30 minutes, to both streamline
106 data collection and explore differences between the already collected Craft Story delay. During
107 this delay, additional speech tasks were administered. Participants were given two minutes for both
108 the immediate and delayed responses.

109 **Picture descriptions (PD)** for a black and white version of the Modern Cookie Theft image [34]
110 were collected from each participant, chosen as it is a highly validated and commonly-used
111 stimulus for eliciting descriptive spontaneous speech. The updated image, as opposed to the
112 original line-drawing version [11], is a more modern depiction of a kitchen scene which avoids
113 stereotypical gender roles, adds additional content units, and incorporates more robust shading and
114 coloring as opposed to the flat line drawings of the original which may aid participants in visually

115 exploring the scene [35]. This updated version has been shown to elicit more detailed descriptions
116 even under time constraints [36]. We further transformed this version by desaturating the image
117 while still maintaining high contrast, as preliminary study revealed that many participants overly
118 relied on color-based descriptions and less on describing action. Participants are given ninety
119 seconds to describe each image, with a visual indicator at the 80-second mark. Additionally, if they
120 choose to end their description before the time limit, they are asked “is there anything else?”.

121 **Spontaneous free responses (FR)** to the prompt “Describe your typical morning routine” were
122 also collected. Participants were given no additional prompting, with their responses recorded until
123 they indicated completion or until 90 seconds had passed.

124 *Transcripts*

125 Audio recordings from each task were processed using *OpenAI Whisper*’s Large-v2 model [37], a
126 speech-to-text model that generates transcripts at higher accuracy than many other free and paid
127 options and is capable of being run on in-house hardware. Each automatically generated transcript
128 was then manually corrected as needed by trained research staff. The majority of corrections
129 included the removal of experimenter instructions, hallucinations including “The end” and
130 “Thanks for watching”, and obscured words due to high ambient noise, however the rate of
131 corrections including the word error rate was not tracked. Generated transcripts included start and
132 stop times in seconds for each word as well as confidence for each predicted word which helped
133 highlight areas of potentially necessary correction.

134 *Audio Preprocessing*

135 All audio files were preprocessed prior to speech-to-text transcription. Leading and trailing silence
136 was removed using *PyDub* [38], and the end of a file was marked using a chime to encourage

137 *OpenAI Whisper* to end sentences appropriately when subjects ran out of time and the audio file
138 ended abruptly. If a chime is not used, hallucinations where *Whisper* finished sentences were
139 common, including erroneously completed sentences with predicted text or an ending message
140 (e.g., “The end.”) Background noise, including room tone and non-target extraneous speech (e.g.,
141 voices from hallway) were removed using *PySoX* [39] to further improve speech-to-text accuracy
142 as well as to facilitate audiometric analysis.

143 **Analysis**

144 *Analysis of Speech Transcripts*

145 After speech responses were transcribed and validated, the transcripts were parsed using a Python
146 script incorporating multiple NLP packages, including *SpaCy* for tokenization, lemmatization, part
147 of speech analysis, and general NLP analyses [40], *wordfreq* for examining lexical frequency using
148 its own proprietary formula [41], and *textdescriptives* for generating a wide range of statistics
149 including entropy, perplexity, sentence count and length, syllable count, and reading ease statistics
150 such as Flesch Kincaid grade level, Gunning-Fog, and the Coleman Liau Index [42,43]. Custom
151 scripts were also written that generate SUBTLEX rankings [44] for determining lexical frequency,
152 mean Yngve depth for examining passage complexity [45], and type token ratios and derivatives
153 (MATTR or “moving average type token ratio” [46] and MTLT or “measure of textual diversity”
154 [47]) for lexical diversity. Distinctions are made for which corpus (*wordfreq* or SUBTLEX) was
155 utilized when examining lexical frequency, as each produces slightly different results due to using
156 different word lists.

157 Pauses were defined as either filled or unfilled, with unfilled pauses defined as pauses lasting
158 greater than 250 milliseconds between words, based on literature examining pauses in normal

159 conversational speech [22] as well as patients with aphasia [21]. Filled pauses were defined as the
160 presence of “um”, “uh”, and “er”, chosen based on their frequency in English speech [48]. Both
161 filled and unfilled pauses were calculated automatically using Python scripts; unfilled pauses were
162 marked when durations greater than 250ms were detected between the end of the previous word
163 and the start of the current word, as determined by timestamps generated by OpenAI Whisper, and
164 filled pauses were determined by tokens containing a filler word. The four words immediately
165 following the pause were marked as “post-pause” words, and statistics were generated for both the
166 average of all four words and only the first word following the pause. These pause and post-pause
167 statistics included total counts and durations of pauses, latencies, syllable counts, SUBTLEX
168 rankings, *wordfreq* frequencies, and part of speech frequencies. Pause duration was defined as the
169 total time in milliseconds between the start of the pause (end of the previous word) and the
170 beginning of the following word. For filled pauses, latency was defined as the time from the end
171 of the utterance of the filler word to the beginning of the following word, and for unfilled pauses
172 latency was only used in the context of calculating the latencies of all pauses in which case pause
173 duration was used.

174 *Statistical and Predictive Analyses*

175 Transcripts were compared between groups for both individual and combined tasks using the
176 *pingouin v0.5.3* [49] Python package. Analyses of covariance (ANCOVA) with age as a covariate
177 were performed for all measures where assumptions of normality and homoscedasticity were met,
178 and Kruskal-Wallis Rank Sum test when assumptions were not met and age was not significantly
179 different between groups as determined by linear regression. In cases where assumptions were not
180 met and age was a significant factor, comparisons were not considered. Statistical logistic
181 regressions were also used to determine the combined effects of multiple variables, especially

182 when isolating the effects of post-pause metrics alone. Pearson's correlations were used to compare
183 speech metrics identified as significantly different by ANOVAs with cognitive and
184 neuropsychological assessments including the Montreal Cognitive Assessment (MoCA) [50], the
185 Cognivue total [51], number span backwards, the Trail Making Test B, Hopkins Verbal Learning
186 Task immediate and delayed [52], semantic verbal fluency (animals), and the Number Symbol
187 Coding Task [53]. Verbal IQ assessed by the Test of Premorbid Functioning as well as the
188 Vulnerability Index [54] and the Resilience Index [55] were also included in correlational analyses
189 with the identified speech metrics. See Besser et al. [30] for details on the administration and
190 examination of these assessments.

191 Predictive analyses, including random forest classifiers, gradient boosted machines, logistic
192 regressions, and support vector classifiers, were performed to identify the predictive power of post-
193 pause and general speech metrics to determine binary impairment status (impaired vs. not
194 impaired). All analyses were cross-validated using a repeated stratified K-folds procedure (3-fold,
195 3-repeat), which resulted in nine combinations of train/test sets for better generalizability of model
196 results. Outputs were evaluated using sensitivity, specificity, F1 score, and area under the ROC
197 curve (AUC). The *LightGBM v4.2.0* [56] Python package's implementation was used for gradient
198 boosted machines, and all other predictive analyses and methods were performed using *scikit-learn*
199 *v1.3.2* [57].

200 Two sets of features were used in our models; one including all statistically significant features,
201 and another including only statistically significant post-pause metrics (**Table 3**). Only individual
202 task metrics were used; when significant results for were found in the mixed ANOVA containing
203 all tasks, only the task with the highest Bayes index in the post-hoc was used in an effort to
204 minimize the negative effect of multicollinearity on model results. Additionally, feature selection

205 was performed using the *BorutaSHAP v1.1* [58] selection algorithm using a random forest
206 classifier (*scikit-learn v1.3.0* [57]) as its base model. One-fourth (25%) of the stratified training
207 data was held for use in the feature selection process, to avoid leakage and subsequent overfitting
208 in the model training phase. *Optuna v3.5.0* [59] selected optimal hyperparameters for each model,
209 leveraging its implementation of define-by-run dynamic parameter search spaces and efficient
210 strategies for pruning, using the same 25% of stratified training data in the feature selection step.

211 **RESULTS**

212 **Sample Characteristics**

213 There was a significant difference in age between the two groups ($p = .003$), resulting in age being
214 included as a covariate in all following comparisons (**Table 1**). No differences between groups
215 were found for gender, years of education, race or ethnicity, vulnerability [54], or resilience [55].
216 Significant differences in cognitive tasks were observed between groups, consistent with a
217 classification of MCI.

218 **Common Findings in All Tasks**

219 When examining all tasks using a Mixed ANOVA, with task as the within-subjects variable and
220 impairment status as the between-subjects variable, differences were found between the HC and
221 MCI groups. There were significantly more filled pauses (“um”, “uh”, and “er”) in the MCI group
222 (4.21 ± 4.42) compared to the HC group (2.47 ± 2.71), $p=.006$, especially containing “uh”
223 (2.75 ± 3.53 MCI vs 1.11 ± 1.78 HC; $\chi^2(1)=13.31$, $p<.001$), with latencies between all pauses and the
224 following word increasing in the MCI group (1.22 ± 0.35 MCI vs 1.11 ± 0.33 HC; $\chi^2(1)=9.31$,
225 $p=.002$), **Figure 1**.

226 Correlational analyses with cognitive and neuropsychological assessments revealed that total word
227 count was often mildly or moderately negatively correlated with decreased cognitive functioning,
228 including the MoCA ($r=.190, p=.002$), Cognivue ($r=.168, p=.009$), trailmaking test B ($r=-.214,$
229 $p<.001$), semantic verbal fluency ($r=.198, p=.002$), and the Number Symbol Coding Task ($r=.194,$
230 $p=.002$), as well as verbal IQ as assessed by the Test of Premorbid Functioning ($r=.269, p<.001$)
231 and the Resilience Index ($r=.250, p<.001$). Despite “um” pause counts alone not being significantly
232 different between groups, moderate negative correlations were found with the MoCA ($r=-.261,$
233 $p<.001$), Cognivue ($r=-.308, p<.001$), Number Symbol Coding ($r=-.232, p<.001$), and semantic
234 verbal fluency ($r=-.196, p=.003$) tests, in addition to the Resilience Index ($r=-.251, p<.001$); a mild
235 positive correlation was also found with the Vulnerability Index ($r=.179, p=.005$).

236 For post-pause metrics, the Kruskal-Wallis non-parametric ANOVA test was required due to many
237 participants not producing pauses at all in some tasks, resulting in unequal variances between
238 groups. Lexical frequency was significantly higher following a filled pause in MCI participants,
239 supporting our hypotheses; *wordfreq* rankings following “uh” fillers ($.0087\pm.013$ MCI vs
240 $.0058\pm.012$ HC; $\chi^2(1)=7.31, p=.007$) were significantly different. Latencies following “um” fillers
241 were also significantly higher in MCI participants ($1.83 \text{ sec} \pm 0.73$ MCI vs $1.33 \text{ sec} \pm 0.56$ HC;
242 $\chi^2(1)=17.36, p<.001$). Mean latency following pause filled with “um” was negatively correlated
243 with the Cognivue ($r=-.335, p<.001$), Hopkins Verbal Learning immediate ($r=-.244, p=.005$) and
244 delayed ($r=-.259, p=.003$), semantic verbal fluency ($r=-.285, p=.002$), and the Number Symbol
245 Coding task ($r=-.323, p<.001$), in addition to a positive correlation with Trailmaking B ($r=.312,$
246 $p<.001$).

247 **Task Comparisons**

248 All differences between tasks described below exhibited a significant ANOVA or pairwise t-test
249 within-subjects score, with an alpha set at .015.

250 *Picture description*

251 In the PD task, the mean Yngve depth was shallower in MCI participants (3.36 ± 0.67 levels) than
252 for HC participants (3.98 ± 0.86 levels), $p = .012$, however no significant correlations with cognitive
253 assessments were identified. The Coleman-Liau Index, a measure of readability, also decreased in
254 MCI participants (3.65 ± 1.45 MCI vs 4.74 ± 0.88 HC; $p = .008$) (**Table 2.c**), findings supported by a
255 moderate correlation with the MoCA total ($r = .395$, $p = .001$). Total pause count was moderately
256 negatively correlated with the Cognivue ($r = -.405$, $p = .001$) in this task, and the mean *wordfreq*
257 lexical frequency was also negatively correlated with the MoCA total ($r = -.338$, $p = .007$). Other
258 significant correlations are found in **Figure 2**.

259 *Narrative recall*

260 In both narrative recall tasks, filled pauses significantly increased in MCI participants (3.23 ± 3.97
261 pauses MCI vs 1.58 ± 2.28 HC; $\chi^2(1) = 7.37$, $p = .007$), however this significance did not appear in
262 pairwise comparisons (**Table 2.a, 2.b**). Latencies following “um” fillers increased in MCI
263 participants in both NR tasks (1.82 ± 0.72 sec MCI vs 1.25 ± 0.49 sec HC; $\chi^2(1) = 8.43$, $p = .004$) and
264 were significantly correlated with measures of cognition (**Figure 2**), however post-hoc analyses
265 revealed that only the delayed NR task showed significance between groups (**Table 2.b**).

266 Mean total word counts decreased in both immediate and delayed NR tasks (74.75 ± 30.21 words
267 in MCI vs 92.88 ± 31.95 words in HC; $p = .002$), with lexical diversity measured by MTLTD found to
268 decrease in only the immediate NR task (38.79 ± 16.53 MCI vs 50.42 ± 15.64 HC; $p = .013$). Also in
269 the immediate NR task, the Coleman Liau Index decreased from 6.74 ± 1.67 in HC to 5.27 ± 2.31 in

270 MCI, $p=.007$. Mean word lengths (3.94 ± 0.26 characters MCI vs 4.16 ± 0.24 characters HC; $p=.003$)
271 in the immediate NR task also decreased in MCI participants, while median word length decreased
272 in the delayed NR task (3.50 ± 0.48 MCI vs 3.90 ± 0.47 HC; $p=.002$). The number of syllables per
273 word in the immediate NR task was also found to decrease for MCI participants (1.12 ± 0.06
274 syllables MCI vs 1.16 ± 0.05 syllables HC; $p=.006$). Fewer nouns were used by MCI participants
275 in the immediate NR task (9.82 ± 3.23 nouns MCI vs 12.73 ± 3.77 nouns HC; $p=.007$). Further,
276 proper nouns such as names were significantly less common for MCI participants in both the
277 immediate NR (2.71 ± 2.39 words MCI vs 5.24 ± 2.54 words HC; $p=.001$) and delayed NR tasks
278 (1.94 ± 1.69 words MCI vs 4.89 ± 2.92 words HC; $p<.001$), findings supported by significant
279 correlations for both nouns and proper nouns with most cognitive and neuropsychological
280 assessments in these tasks (**Figure 2**).

281 *Free response*

282 Word frequencies as calculated by *wordfreq* were significantly more common in MCI participants
283 after "uh" fillers in this task ($.0086\pm .0078$ MCI vs $.0025\pm .0027$ HC; $p=.005$) (**Table 2.d**), but not
284 in other tasks (**Figure 3**). Supporting this, strong correlations were found between the Number
285 Symbol Coding Task ($r=-.496$, $p=.002$) and the Vulnerability Index ($r=.483$, $p=.003$). Additionally,
286 adverb usage significantly decreased following unfilled pauses in MCI participants ($1.8\% \pm 3.9\%$
287 adverbs post-pause) compared to HCs ($16.4\% \pm 20.8\%$ adverbs post-pause), $\chi^2(1)=8.89$, $p=.003$.

288 *Controlling for individual cognitive resources*

289 To rule out the effects of varying cognitive resources within groups, the Trailmaking B time
290 (measure of attention), Number-Symbol Coding Task (measure of executive function), Hopkins
291 verbal learning task delayed (measure of episodic memory), Animal Naming (measure of

292 categorical verbal fluency), and the Test of Premorbid Functioning (measure of verbal IQ) were
293 used as covariates for all significant findings above. No results were rendered non-significant after
294 controlling for these measures.

295 **Predictive Analysis**

296 When starting with all significant features, feature selection narrowed the field to 16 of the best-
297 performing features: total pause count (immediate NR), total filled pause count (immediate NR),
298 total “uh” count (immediate and delayed NR), mean *wordfreq* frequency post-“uh” (FR), mean
299 latency post-“um” (delayed NR), word count (immediate NR), MTLN (immediate NR), mean
300 *wordfreq* frequency (immediate NR), median word length (delayed NR) Coleman-Liau Index (PD
301 & immediate NR), mean Yngve depth (PD), proper noun ratio (immediate NR), and noun count
302 (delayed NR). Using this set of features and after performing hyperparameter optimization, the
303 best performing model was the random forest classifier with a diagnostics odds ratio (DOR) of
304 13.52 (**Table 3**) and an AUC of 0.828 (**Figure 4**).

305 When using only the two most significant post-pause metrics (*wordfreq* lexical frequency
306 following “uh” in the FR task, and latencies following “um” in the delayed NR task), the best
307 performing model was again the random forest classifier with a DOR of 10.3 (**Table 3**) and an
308 AUC of 0.791 (**Figure 4**).

309 **DISCUSSION**

310 In this study, we examined whether post-pause speech is affected due to cognitive impairment, and
311 if task demands elicit differences in speech behavior between HC and MCI participants. We found
312 significant differences between groups in *wordfreq* lexical frequency following filled pauses in the
313 FR task, as well as increased latencies following “um” pauses in the NR tasks.

314 *Task-dependent pause production*

315 MCI participants tended to use more common (high frequency) language following filled “uh”
316 pauses, but only in the FR task. Our initial hypothesis that lexical frequency would decrease in
317 MCI individuals due to increased word searching behavior may have underestimated the
318 intervening effect of task demands. Just as different forms of filled pauses differ in their usage [24]
319 it is possible that each form may be produced for different reasons as task conditions vary. Indeed,
320 previous studies show that task difficulty and the associated increased cognitive load have been
321 shown to affect speech production, including altering lexical diversity and syntactic complexity
322 [60] as well as disfluency and pause usage [61]. The question in the FR task asked about the
323 participant’s typical routine, which is a topic that is typically static, repetitive, and familiar,
324 requiring few cognitive resources to recall effectively. Pauses produced in this context thus likely
325 resemble typical conversational speech, and would match the documented usage of “uh” fillers as
326 indicators of word-searching behavior [24,62,63]. As a result, cognitive impairment may play a
327 larger role in disrupting this search, leading to greater differences between the HC and MCI groups.
328 This is contrasted with the three other tasks, where required cognitive demands are greater to
329 varying degrees. In these tasks, filled “uh” pauses may still be produced as indicators of searching
330 behavior, but instead of lexical search being the driving force behind pause production as in the
331 spontaneous speech task this searching behavior may be instead more directed towards searching
332 for the correct answer (in NR tasks) or searching for a new object to describe (in the PD task).
333 Thus, while post-filled pause lexical frequency may still be affected by cognitive impairment, the
334 HC group may also use more common words as a result of the increased cognitive load [60],
335 resulting in a less noticeable and thus non-significant difference between groups.

336 In a similar vein, the delayed NR task was arguably the most difficult task administered in this
337 study, requiring participants to recall information presented after being distracted in the interim
338 with other speech tasks while also being under a time limit [64]. Delayed recall activates broader
339 neurological pathways, including activation of the left parahippocampal gyrus, the entorhinal and
340 perirhinal cortices, and both the anterior and posterior hippocampal regions, contrasting with
341 immediate recall which is handled primarily with short-term memory processes in the posterior
342 hippocampus and dorsolateral prefrontal cortex [65,66]. This may explain why post-“um” latencies
343 were significantly different between groups only in this task. “Um” fillers have been shown to
344 signal longer delays than other pauses [24,63], and the delays decrease with skill [23] and increase
345 with cognitive load [67]. Delayed narrative recall performance is known to be significantly
346 affected due to cognitive impairment [68], and so MCI participants likely needed to utilize more
347 cognitive resources to perform the task, leading to increased delays compared to the HC group.

348 While previous studies have shown significant differences between MCI and HC in numbers of
349 pauses, both filled and unfilled [27,28] and especially using the “uh” utterance [69], our results
350 did not show significant differences between the two groups in total count of pauses when
351 examining the post-hoc pairwise tests for task comparison that were not also significantly different
352 due to age effects. While our sample size is relatively limited, the observation of significant age
353 effects may suggest that previous identification of increased numbers of pauses in these studies
354 may be due to the MCI group being older than the HC group. Future examination into total pause
355 counts should take age into account as a possible covariate.

356 *Other lexical differences between tasks*

357 More significant differences between groups were found in the immediate NR task than any others,
358 and when multiple tasks exhibited the same differences in a single measure the *Boruta* feature

359 selection algorithm identified the immediate NR task as the most predictive. The MTLD measure
360 of lexical diversity as well as the Coleman-Liau index both identify a tendency of MCI participants
361 to use simpler language in this task, as does the finding of reduced mean word length and syllables
362 per word. The high usefulness of the immediate NR task has been previously shown in prior work
363 developing parsimonious assessment composites for detecting MCI [68], where it was determined
364 to be one of the most useful metrics for identifying MCI among the neuropsychological tasks
365 administered by ADNI. This may be due to this task requiring coordination of working memory
366 and attention, two areas known to be degraded in MCI but not entirely disrupted.

367 Proper noun usage decreased in both the immediate and delayed NR tasks, indicating that while
368 healthy participants are able to recall names and places, those with MCI are less likely to recall
369 this information. As both were impaired, it is difficult to determine whether the deficit was due to
370 encoding or retrieval. Previous study by Mueller et al. [70] showing that proper noun recall in
371 delayed but not immediate NR is associated with preclinical AD suggests that retrieval is the main
372 factor affected in proper noun recall, however it is important to note that our sample did have
373 biomarker-confirmed AD pathology and may have included other etiologies including vascular
374 dementia.

375 Mean Yngve depth, a measure of syntactic complexity based on hierarchical relationships between
376 clauses in produced speech, was found to be significantly shallower in picture descriptions of the
377 Cookie Theft image produced by MCI participants. Picture description tasks elicit hierarchical
378 language as participants describe the relationships between objects in a picture or photograph. The
379 Cookie Theft picture in particular contains many relationships between objects, with for example
380 the boy engaging with both the cookie jar and the stool and the man both washing the dishes and
381 not attending to them resulting in them overflowing. The reduction of mean Yngve depth in MCI

382 may be a result of a reduced focus on how objects in a scene relate to each other, instead resulting
383 in simpler descriptions of these objects. Previous studies have shown that cognitive impairment
384 affects mean Yngve depth in NR tasks [19,71], but PD tasks have not been previously shown to
385 elicit decreased syntactic complexity. In our sample, mean Yngve depth did not appear to differ
386 between groups for either NR task but did differ in the PD task. It is possible that the updated
387 Cookie Theft image utilized in our study [34] encouraged either increased interrelated descriptions
388 in healthy controls or a decrease in MCI participants. Additionally, the significant differences were
389 only observed in ANOVAs and depth did not significantly correlate with any individual cognitive
390 assessment (**Figure 2**). While this could be explained due to other mediating effects including age,
391 it is also possible that Yngve depth captures differences not identified in any of the assessments
392 used in our correlations. The PD task also prompted decreased word lengths, Coleman-Liau
393 indices, and noun counts in MCI participants, indicating less precise and simpler language in this
394 task, with this finding supported by previous studies [14,72].

395 *Use of speech behavior to predict impairment status*

396 After performing feature selection on the entire group of available features, a random forest
397 utilizing the 16 selected features produced a model with an AUC of 0.828 and a DOR of 13.52.
398 While specificity was excellent (94.59%), sensitivity for detecting impairment was only 43.59%.
399 Nonetheless, if the task is to screen for impairment, the use of speech-based markers alone perform
400 well at reducing numbers of non-impaired individuals. When only the top two post-pause features
401 are used, only marginal decreases in sensitivity and specificity are observed. With high specificity
402 and AUC, the results of this model suggest that post-pause metrics of latency and lexical frequency
403 would be an excellent addition to any machine learning model that utilizes speech behavior and

404 seeks to categorize healthy individuals in some capacity, but the low sensitivity precludes their use
405 as a diagnostic tool.

406 *Limitations*

407 The main limitations of this study were the limited sample size, unequal sample size between
408 impaired and non-impaired groups, the racial and ethnic makeup of the sample being primarily
409 non-Hispanic White, and the lack of biomarker data available in this early sample. As MCI is a
410 heterogeneous disorder, it is possible that separation of MCI subtypes (e.g., MCI due to AD, MCI
411 due to Lewy body or vascular etiologies) will lead to more detailed characterization using
412 neurobehavioral markers such as speech. Additional recruitment and biomarker examination of
413 existing and new participants is ongoing, with increased recruitment efforts targeting impaired as
414 well as racially and ethnically diverse populations.

415 We also observed age differences between groups. Instead of reducing our healthy control sample
416 to minimize age effects through age-matching, we chose to instead include age as a covariate
417 within our analyses due to the small size of the impaired group. Our future research will
418 incorporate age-matching, as we aim to recruit a sufficient number of biomarker-confirmed
419 impaired participants to support an evenly sized control group. Additionally, some significant
420 measures exhibited standard deviations that exceeded their means, indicating a high degree of
421 variability or skewed distributions for these measures; however, these typically occurred only
422 when all tasks were considered. Thus, this phenomenon is not entirely unexpected as each task
423 was shown to perform differently with respect to each measure and between groups. The only
424 significant task-specific measure with this property was the increased post-pause adverb usage in
425 the free response task.

426 In this study, only a handful of fillers were examined (“uh”, “um”, and “er”). While many other
427 fillers exist including “you know”, “like”, “okay”, and “right”, particularly among younger
428 individuals, we opted to not include these filler words in our analyses as these were both relatively
429 uncommon in our sample as well as difficult to code in automatic processing; each of these
430 additional filler words can also appear as normal non-filler words, thus requiring manual
431 classification for each of these instances. However, it is possible that the inclusion of these more
432 uncommon filler words would have revealed additional differences between groups. In addition,
433 while the filler “er” was examined in all analyses, no significant differences were found between
434 groups for the usage of this filler nor for words following it. As a result, it is not referenced
435 elsewhere in this paper. Further, this study only examined English-speaking participants: these
436 findings may not generalize to other languages, particularly in the type of frequency of filler words.
437 Nonetheless, this study identified compelling patterns in post-pause speech behavior between MCI
438 and HC in English-speaking populations, and future studies will examine non-native English
439 speakers as well as Spanish speakers to determine whether this behavior translates to other
440 languages and contexts.

441 **CONCLUSIONS**

442 The increase of post-pause lexical frequency in the absence of task demands and post-filler latency
443 in tasks with high difficulty were observed in MCI participants, with both able to accurately predict
444 impairment status with an AUC of 79.1%. Future research will examine likely causes of pause
445 production, comparing lexical or word-finding pauses with those driven by task demands or
446 cognitive load, as well as neural correlates of speech degradation.

447 **FUNDING**

448 Work on this study was supported by grants from the National Institute on Aging (R01 AG071514,
449 R01 AG069765, and R01 NS101483), the Alzheimer’s Association (AARF-22-923592), the
450 Evelyn F. McKnight Brain Research Foundation (FP00006751), and the Harry T. Mangurian
451 Foundation.

452 **ACKNOWLEDGEMENTS**

453 The authors have no acknowledgements to report.

454 **CONFLICT OF INTEREST**

455 The authors have no conflict of interest to report

456 **DATA AVAILABILITY**

457 The data supporting the findings of this study are available on request from the corresponding
458 author.

REFERENCES

- [1] Forbes-McKay KE, Venneri A (2005) Detecting subtle spontaneous language decline in early Alzheimer's disease with a picture description task. *Neurol Sci* **26**, 243–254.
- [2] Ahmed S, Haigh A-MF, de Jager CA, Garrard P (2013) Connected speech as a marker of disease progression in autopsy-proven Alzheimer's disease. *Brain* **136**, 3727–3737.
- [3] Kavé G, Goral M (2016) Word retrieval in picture descriptions produced by individuals with Alzheimer's disease. *J Clin Exp Neuropsychol* **38**, 958–966.
- [4] Fraser KC, Meltzer JA, Rudzicz F (2015) Linguistic features identify Alzheimer's disease in narrative speech. *Journal of Alzheimer's Disease* **49**, 407–422.
- [5] Yamada Y, Shinkawa K, Nemoto M, Ota M, Nemoto K, Arai T (2022) Speech and language characteristics differentiate Alzheimer's disease and dementia with Lewy bodies. *Alz & Dem Diag Ass & Dis Mo* **14**,.
- [6] Yuan J, Cai X, Bian Y, Ye Z, Church K (2021) Pauses for Detection of Alzheimer's Disease. *Front Comput Sci* **0**,.
- [7] Boschi V, Catricalà E, Consonni M, Chesi C, Moro A, Cappa SF (2017) Connected Speech in Neurodegenerative Language Disorders: A Review. *Frontiers in Psychology* **8**, 269.
- [8] Petersen RC (2016) Mild Cognitive Impairment. *Continuum (Minneap Minn)* **22**, 404–418.
- [9] Themistocleous C, Eckerström M, Kokkinakis D (2020) Voice quality and speech fluency distinguish individuals with Mild Cognitive Impairment from Healthy Controls. *PLOS ONE* **15**, e0236009.
- [10] Mueller KD, Kosciak RL, Hermann BP, Johnson SC, Turkstra LS (2018) Declines in connected language are associated with very early mild cognitive impairment: Results from the Wisconsin Registry for Alzheimer's Prevention. *Frontiers in Aging Neuroscience* **9**, 1–14.
- [11] Goodglass H, Kaplan E, Weintraub S (2001) *BDAE: The boston diagnostic aphasia examination*, Lippincott Williams & Wilkins Philadelphia, PA.
- [12] Wechsler D (1987) Wechsler memory scale-revised. The psychological corporation. *New York*.
- [13] Craft S, Newcomer J, Kanne S, Dagogo-Jack S, Cryer P, Sheline Y, Luby J, Dagogo-Jack A, Alderson A (1996) Memory improvement following induced hyperinsulinemia in alzheimer's disease. *Neurobiology of Aging* **17**, 123–130.
- [14] Beltrami D, Gagliardi G, Rossini Favretti R, Ghidoni E, Tamburini F, Calzà L (2018) Speech Analysis by Natural Language Processing Techniques: A Possible Tool for Very Early Detection of Cognitive Decline? *Front Aging Neurosci* **0**,.
- [15] Roark B, Mitchell M, Hollingshead K (2007) Syntactic complexity measures for detecting mild cognitive impairment. In *Proceedings of the Workshop on BioNLP 2007 Biological, Translational, and Clinical Language Processing - BioNLP '07* Association for Computational Linguistics, Prague, Czech Republic, p. 1.
- [16] Clarke N, Barrick TR, Garrard P (2021) A Comparison of Connected Speech Tasks for Detecting Early Alzheimer's Disease and Mild Cognitive Impairment Using Natural Language Processing and Machine Learning. *Frontiers in Computer Science* **3**,.
- [17] Taler V, Kousaie S, Sheppard C (635775264000000000) Lexical access in mild cognitive impairment. *The Mental Lexicon* **10**, 271–285.
- [18] Fraser KC, Lundholm Fors K, Eckerström M, Öhman F, Kokkinakis D (2019) Predicting MCI Status From Multimodal Language Data Using Cascaded Classifiers. *Frontiers in Aging Neuroscience* **11**, 1–18.
- [19] Fraser KC, Meltzer JA, Graham NL, Leonard C, Hirst G, Black SE, Rochon E (2014) Automated classification of primary progressive aphasia subtypes from narrative speech transcripts. *Cortex* **55**, 43–60.

- [20] Meteyard L, Quain E, Patterson K (2014) Ever decreasing circles: Speech production in semantic dementia. *Cortex* **55**, 17–29.
- [21] Goldman-Eisler F (1968) Psycholinguistics: Experiments in spontaneous speech.
- [22] Gilden D, Mezaraups T (2022) Laws for pauses. *Journal of Experimental Psychology: Learning, Memory, and Cognition* **48**, 142–158.
- [23] Womack K, McCoy W, Ovesdotter Alm C, Calvelli C, Pelz JB, Shi P, Haake A (2012) Disfluencies as Extra-Propositional Indicators of Cognitive Processing. In *Proceedings of the Workshop on Extra-Propositional Aspects of Meaning in Computational Linguistics*, Morante R, Sporleder C, eds. Association for Computational Linguistics, Jeju, Republic of Korea, pp. 1–9.
- [24] Clark HH, Fox Tree JE (2002) Using uh and um in spontaneous speaking. *Cognition* **84**, 73–111.
- [25] Shriberg E (2001) To ‘errrr’ is human: ecology and acoustics of speech disfluencies. *Journal of the International Phonetic Association* **31**, 153–169.
- [26] Toth L, Hoffmann I, Gosztolya G, Vincze V, Szatloczki G, Banreti Z, Pakaski M, Kalman J (2017) A Speech Recognition-based Solution for the Automatic Detection of Mild Cognitive Impairment from Spontaneous Speech. *Current Alzheimer Research* **14**, 130–138.
- [27] Vincze V, Szatloczki G, Tóth L, Gosztolya G, Pákáski M, Hoffmann I, Kálmán J (2020) Telltale silence: temporal speech parameters discriminate between prodromal dementia and mild Alzheimer’s disease. *Clinical Linguistics & Phonetics* **00**, 1–16.
- [28] Egas-López JV, Balogh R, Imre N, Hoffmann I, Szabó MK, Tóth L, Pákáski M, Kálmán J, Gosztolya G (2022) Automatic screening of mild cognitive impairment and Alzheimer’s disease by means of posterior-thresholding hesitation representation. *Computer Speech & Language* 101377.
- [29] Wang T, Hong Y, Wang Q, Su R, Ng M, Xu J, Yan N (2021) Identification of Mild Cognitive Impairment Among Chinese Based on Multiple Spoken Tasks. *Journal of Alzheimer’s Disease* **82**, 1–20.
- [30] Besser LM, Chrisphonte S, Kleiman MJ, O’Shea D, Rosenfeld A, Tolea M, Galvin JE (2023) The Healthy Brain Initiative (HBI): A prospective cohort study protocol. *PLOS ONE* **18**, e0293634.
- [31] Morris JC (1993) The Clinical Dementia Rating (CDR): Current version and scoring rules. *Neurology* **43**, 2412–2412.
- [32] Weintraub S, Besser L, Dodge HH, Teylan M, Ferris S, Goldstein FC, Giordani B, Kramer J, Loewenstein D, Marson D, Mungas D, Salmon D, Welsh-Bohmer K, Zhou XH, Shirk SD, Atri A, Kukull WA, Phelps C, Morris JC (2018) Version 3 of the Alzheimer Disease Centers’ Neuropsychological Test Battery in the Uniform Data Set (UDS). *Alzheimer Disease and Associated Disorders* **32**, 10–17.
- [33] Bolognani SAP, Miranda MC, Martins M, Rzezak P, Bueno OFA, de Camargo CHP, Pompeia S (2015) Development of alternative versions of the Logical Memory subtest of the WMS-R for use in Brazil. *Dement Neuropsychol* **9**, 136–148.
- [34] Berube S, Nonnemacher J, Demsky C, Glenn S, Saxena S, Wright A, Tippett DC, Hillis AE (2019) Stealing Cookies in the Twenty-First Century: Measures of Spoken Narrative in Healthy Versus Speakers With Aphasia. *Am J Speech Lang Pathol* **28**, 321–329.
- [35] Heuer S (2016) The influence of image characteristics on image recognition: a comparison of photographs and line drawings. *Aphasiology* **30**, 943–961.
- [36] Hux K, Frodsham K (2023) Speech and Language Characteristics of Neurologically Healthy Adults When Describing the Modern Cookie Theft Picture: Mixing the New With the Old. *Am J Speech Lang Pathol* **32**, 1110–1130.
- [37] Radford A, Kim JW, Xu T, Brockman G, McLeavey C, Sutskever I (2022) Robust Speech Recognition via Large-Scale Weak Supervision.
- [38] Robert J, Webbie M, Jean-philippe S, Ramdasan A, Lee C, Campbell J, McFarland R, McMellen J, Orby P (2023) PyDub.

- [39] Bittner R, Humphrey EJ, Bello J (2016) *PySOX: Leveraging the Audio Signal Processing Power of SOX in Python*, Proceedings of the 17th International Society for Music Information Retrieval Conference Late Breaking and Demo Papers, New York, NY, USA.
- [40] Honnibal M, Montani I, Van Landeghem S, Boyd A (2020) spaCy: Industrial-strength Natural Language Processing in Python.
- [41] Speer R (2022) rspeer/wordfreq: v3.0.
- [42] Hansen L, Olsen LR, Enevoldsen K (2023) TextDescriptives: A Python package for calculating a large variety of metrics from text. *JOSS* **8**, 5153.
- [43] Coleman M, Liau TL (1975) A computer readability formula designed for machine scoring. *Journal of Applied Psychology* **60**, 283–284.
- [44] Brysbaert M, New B (2009) Moving beyond Kucera and Francis: a critical evaluation of current word frequency norms and the introduction of a new and improved word frequency measure for American English. *Behav Res Methods* **41**, 977–990.
- [45] Yngve VH (1960) A Model and an Hypothesis for Language Structure. *Proceedings of the American Philosophical Society* **104**, 444–466.
- [46] Covington MA, McFall JD (2010) Cutting the Gordian Knot: The Moving-Average Type–Token Ratio (MATTR). *Journal of Quantitative Linguistics* **17**, 94–100.
- [47] McCarthy PM, Jarvis S (2010) MTL, vocd-D, and HD-D: A validation study of sophisticated approaches to lexical diversity assessment. *Behavior Research Methods* **42**, 381–392.
- [48] Bortfeld H, Leon SD, Bloom JE, Schober MF, Brennan SE (2001) Disfluency rates in conversation: effects of age, relationship, topic, role, and gender. *Lang Speech* **44**, 123–147.
- [49] Vallat R (2018) Pingouin: statistics in Python. *JOSS* **3**, 1026.
- [50] Nasreddine ZS, Phillips NA, Bédirian V, Charbonneau S, Whitehead V, Collin I, Cummings JL, Chertkow H (2005) The Montreal Cognitive Assessment, MoCA: A brief screening tool for mild cognitive impairment. *Journal of the American Geriatrics Society* **53**, 695–699.
- [51] Cahn-Hidalgo D, Benabou R, Kewin S (2019) Validity, Reliability, and Psychometric properties of Cognivue®, a quantitative assessment of cognitive impairment. *The American Journal of Geriatric Psychiatry* **27**, S212.
- [52] Brandt J (1991) The hopkins verbal learning test: Development of a new memory test with six equivalent forms. *Clinical Neuropsychologist* **5**, 125–142.
- [53] Galvin JE, Tolea MI, Moore C, Chrisphonte S (2020) The Number Symbol Coding Task: A brief measure of executive function to detect dementia and cognitive impairment. *PLOS ONE* **15**, e0242233.
- [54] Kleiman MJ, Galvin JE (2021) The Vulnerability Index: A weighted measure of dementia and cognitive impairment risk. *Alzheimer's & Dementia: Diagnosis, Assessment & Disease Monitoring* **13**, e12249.
- [55] Galvin JE, Kleiman MJ, Chrisphonte S, Cohen I, Disla S, Galvin CB, Greenfield KK, Moore C, Rawn S, Riccio ML, Rosenfeld A, Simon J, Walker M, Tolea MI (2021) The Resilience Index: A Quantifiable Measure of Brain Health and Risk of Cognitive Impairment and Dementia. *Journal of Alzheimer's Disease* **84**, 1729–1746.
- [56] Ke G, Meng Q, Finley T, Wang T, Chen W, Ma W, Ye Q, Liu T-Y (2017) LightGBM: A Highly Efficient Gradient Boosting Decision Tree. In *Advances in Neural Information Processing Systems* Curran Associates, Inc.
- [57] Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R, Dubourg V, Vanderplas J, Passos A, Cournapeau D, Brucher M, Perrot M, Duchesnay E (2011) Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research* **12**, 2825–2830.

- [58] Keany E (2020) BorutaShap: A wrapper feature selection method which combines the Boruta feature selection algorithm with Shapley values.
- [59] Akiba T, Sano S, Yanase T, Ohta T, Koyama M (2019) Optuna: A Next-generation Hyperparameter Optimization Framework. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* Association for Computing Machinery, New York, NY, USA, pp. 2623–2631.
- [60] Lee J (2019) Task Complexity, Cognitive Load, and L1 Speech. *Applied Linguistics* **40**, 506–539.
- [61] Müller C, Großmann-Hutter B, Jameson A, Rummer R, Wittig* F (2001) Recognizing Time Pressure and Cognitive Load on the Basis of Speech: An Experimental Study. In *User Modeling 2001*, Bauer M, Gmytrasiewicz PJ, Vassileva J, eds. Springer, Berlin, Heidelberg, pp. 24–33.
- [62] Swerts M (1998) Filled pauses as markers of discourse structure. *Journal of Pragmatics* **30**, 485–496.
- [63] Degand L, Gilquin G, Meurant L, Simon AC (2019) *Fluency and Disfluency Across Languages and Language Varieties*, Presses universitaires de Louvain.
- [64] Earles JL, Kersten AW, Berlin Mas B, Miccio DM (2004) Aging and Memory for Self-Performed Tasks: Effects of Task Difficulty and Time Pressure. *The Journals of Gerontology: Series B* **59**, P285–P293.
- [65] Libby LA, Hannula DE, Ranganath C (2014) Medial temporal lobe coding of item and spatial information during relational binding in working memory. *J Neurosci* **34**, 14233–14242.
- [66] Kühn S, Gallinat J (2014) Segregating cognitive functions within hippocampal formation: a quantitative meta-analysis on spatial navigation and episodic memory. *Hum Brain Mapp* **35**, 1129–1142.
- [67] Russo M, Bendazzoli C, Defrancq B (2017) *Making Way in Corpus-based Interpreting Studies*, Springer.
- [68] Kleiman MJ, Barenholtz E, Galvin JE (2021) Screening for Early-Stage Alzheimer’s Disease Using Optimized Feature Sets and Machine Learning. *Journal of Alzheimer’s Disease* **81**, 355–366.
- [69] Yuan J, Bian Y, Cai X, Huang J, Ye Z, Church K (2020) Disfluencies and Fine-Tuning Pre-Trained Language Models for Detection of Alzheimer’s Disease. In *Interspeech 2020 ISCA*, pp. 2162–2166.
- [70] Mueller KD, Kosciak RL, Du L, Bruno D, Jonaitis EM, Kosciak AZ, Christian BT, Betthausen TJ, Chin NA, Hermann BP, Johnson SC (2020) Proper names from story recall are associated with beta-amyloid in cognitively unimpaired adults at risk for Alzheimer’s disease. *Cortex* **131**, 137–150.
- [71] Roark B, Mitchell M, Hosom J-P, Hollingshead K, Kaye J (2011) Spoken Language Derived Measures for Detecting Mild Cognitive Impairment. *IEEE Transactions on Audio, Speech, and Language Processing* **19**, 2081–2090.
- [72] Mueller KD, Hermann B, Mecollari J, Turkstra LS (2018) Connected speech and language in mild cognitive impairment and Alzheimer’s disease: A review of picture description tasks. *Journal of Clinical and Experimental Neuropsychology* **40**, 917–939.

Table 1. Subject characteristics

	Healthy Controls	Mild Cognitive Impairment	p-val
Age	68.28 (7.21)	76.64 (8.90)	0.000
% Female	0.72 (0.45)	0.64 (0.50)	0.478
Years of Education	16.43 (2.65)	17.00 (3.22)	0.920
% Non-White	0.11 (0.31)	0.21 (0.43)	0.159
% Hispanic	0.12 (0.33)	0.21 (0.43)	0.130
MoCA Total	27.05 (2.26)	23.29 (2.73)	0.000
FAQ Total	0.14 (0.45)	0.18 (0.60)	0.750
Numbers Backwards	5.37 (1.33)	3.86 (0.86)	0.000
Trailmaking B	68.11 (24.34)	188.08 (241.85)	0.000
Categorical Fluency	21.21 (5.19)	16.50 (4.03)	0.027
HVLT Immediate	25.32 (3.80)	20.07 (5.05)	0.004
HVLT Delayed	9.52 (1.84)	6.21 (3.89)	0.001
Craft Story Immediate	21.20 (5.68)	18.93 (6.20)	0.198
Craft Story Delayed	19.17 (5.56)	14.50 (8.30)	0.020
Verbal IQ	116.33 (9.44)	111.50 (12.06)	0.008
Vulnerability Index	6.55 (2.02)	8.64 (2.73)	0.078
Resilience Index	183.33 (33.06)	172.54 (32.30)	0.326
Number Symbol Coding Test	45.89 (9.16)	31.36 (6.49)	0.000

All comparisons performed using ANCOVA with age as a covariate, except age which used ANOVA.

MoCA = Montreal Cognitive Assessment; FAQ = Functional Activities Questionnaire; HVLT = Hopkins Verbal Learning Test
 Bold = significant

Table 2.a. Speech measures for the Immediate Narrative Recall task

Speech measure	HC	MCI	Stat (df = 60)	p-val	Test
Total Pause Count	7.325 (3.422)	9.133 (4.809)	2.138	0.144	KW
Unfilled Pause Count	5.975 (2.824)	6.533 (2.588)	1.179	0.278	KW
“Uh” Pause Count	0.707 (1.470)	1.333 (1.633)	3.167	0.075	KW
“Um” Pause Count	0.610 (1.321)	2.200 (3.570)	2.509	0.113	KW
Post-“uh” WF Lexical frequency	0.008 (0.016)	0.004 (0.008)	0.029	0.866	KW
Post-“um” Mean Latency	1.395 (0.471)	1.762 (0.721)	1.970	0.160	KW
Total Word Count	84.378 (26.173)	70.059 (24.517)	3.818	0.055	ANOVA
MTLD	50.422 (15.642)	38.791 (16.528)	6.616	0.013	ANOVA
Mean WF Lexical frequency	0.006 (0.001)	0.007 (0.002)	0.557	0.458	ANOVA
Mean Word Length	4.158 (0.239)	3.944 (0.258)	9.449	0.003	ANOVA
Median Word Length	3.967 (0.165)	3.735 (0.437)	7.772	0.005 ^X	KW
Mean Syllables per Word	1.156 (0.046)	1.117 (0.055)	7.980	0.006	ANOVA
Coleman-Liau Index	6.737 (1.670)	5.271 (2.306)	7.660	0.007	ANOVA
Mean Yngve Depth	4.785 (1.450)	4.739 (1.733)	0.011	0.917	ANOVA
Noun Count	12.733 (3.774)	9.824 (3.226)	7.902	0.007	ANOVA
Proper Noun Count	5.244 (2.542)	2.706 (2.392)	12.694	0.001	ANOVA

KW = Kruskal-Wallis; ANOVA = Analysis of Variance; HC = Healthy Control; MCI = Mild cognitive impairment; MTLD = measure of textual lexical diversity; X = Significant age effects; WF = *wordfreq* lexical frequency

Bold = significant

Table 2.b. Speech measures for the Delayed Narrative Recall task

Speech measure	HC	MCI	F (df = 60)	p-val	Test
Total Pause Count	9.375 (5.138)	9.000 (7.774)	2.168	0.141	KW
Unfilled Pause Count	7.500 (5.208)	6.333 (5.080)	2.34	0.126	KW
“Uh” Pause Count	1.140 (1.656)	1.062 (1.611)	0.123	0.726	KW
“Um” Pause Count	0.698 (1.264)	2.062 (2.720)	4.682	0.030	KW
Post-“uh” WF Lexical frequency	0.003 (0.003)	0.005 (0.004)	1.231	0.267	KW
Post-“um” Mean Latency	1.174 (0.500)	1.886 (0.766)	5.740	0.015	KW
Total Word Count	96.065 (35.668)	79.750 (35.422)	2.492	0.120	ANOVA
MTLD	50.375 (23.418)	45.202 (14.668)	0.683	0.412	ANOVA
Mean WF Lexical frequency	0.007 (0.002)	0.007 (0.001)	0.146	0.704	ANOVA
Mean Word Length	4.109 (0.334)	3.889 (0.226)	5.954	0.018	ANOVA
Median Word Length	3.902 (0.467)	3.500 (0.483)	9.827	0.002	KW
Mean Syllables per Word	1.168 (0.054)	1.137 (0.060)	3.747	0.058	ANOVA
Coleman-Liau Index	5.905 (2.242)	5.128 (2.519)	1.339	0.252	ANOVA
Mean Yngve Depth	4.639 (1.404)	5.156 (1.741)	1.420	0.238	ANOVA
Noun Count	15.848 (5.692)	12.688 (6.488)	3.405	0.070	ANOVA
Proper Noun Count	4.891 (2.915)	1.938 (1.692)	13.026	0.000	ANOVA

KW = Kruskal-Wallis; ANOVA = Analysis of Variance; HC = Healthy Control; MCI = Mild cognitive impairment; MTLD = measure of textual lexical diversity; X = Significant age effects; WF = *wordfreq* lexical frequency

Bold = significant

Table 2.c. Speech measures for the Picture Description task

Speech measure	HC	MCI	F (df = 60)	p-val	Test
Total Pause Count	11.561 (5.473)	14.125 (6.195)	0.397	0.531	ANCOVA
Unfilled Pause Count	8.854 (4.709)	10.500 (5.304)	0.018	0.893	ANCOVA
“Uh” Pause Count	1.178 (1.482)	1.471 (1.068)	0.553	0.460	ANOVA
“Um” Pause Count	1.400 (1.959)	2.588 (2.476)	3.916	0.052	ANOVA
Post-“uh” WF Lexical frequency	0.010 (0.016)	0.014 (0.019)	1.720	0.190	KW
Post-“um” Mean Latency	1.498 (0.720)	1.912 (0.752)	2.658	0.103	KW
Total Word Count	134.356 (55.009)	114.176 (43.805)	1.840	0.180	ANOVA
MTLD	38.641 (11.760)	33.141 (12.431)	2.617	0.111	ANOVA
Mean WF Lexical frequency	0.010 (0.002)	0.010 (0.002)	0.934	0.338	ANOVA
Mean Word Length	3.981 (0.143)	3.843 (0.179)	4.638	0.035 ^X	ANCOVA
Median Word Length	3.356 (0.472)	3.000 (0.000)	8.613	0.003	KW
Mean Syllables per Word	1.217 (0.040)	1.191 (0.053)	0.337	0.564 ^X	ANCOVA
Coleman-Liau Index	4.740 (0.884)	3.651 (1.448)	7.069	0.008	KW
Mean Yngve Depth	3.983 (0.863)	3.356 (0.668)	4.913	0.012	ANOVA
Noun Count	32.933 (13.048)	25.941 (9.959)	3.987	0.05	ANOVA
Proper Noun Count	0.002 (0.006)	0.001 (0.003)	0.268	0.606	ANOVA

KW = Kruskal-Wallis; ANOVA = Analysis of Variance; HC = Healthy Control; MCI = Mild cognitive impairment; MTLD = measure of textual lexical diversity; X = Significant age effects; WF = *wordfreq* lexical frequency

Bold = significant

Table 2.d. Speech measures for the Free Response task

Speech measure	HC	MCI	F (df = 60)	p-val	Test
Total Pause Count	12.049 (7.218)	11.867 (7.170)	0.007	0.934	ANOVA
Unfilled Pause Count	8.488 (5.437)	8.000 (4.342)	0.098	0.756	ANOVA
“Uh” Pause Count	1.932 (2.366)	1.933 (1.981)	0.105	0.746	KW
“Um” Pause Count	1.682 (2.197)	4.200 (4.945)	3.281	0.070 ^X	KW
Post-“uh” WF Lexical frequency	0.003 (0.003)	0.009 (0.008)	8.027	0.005	KW
Post-“um” Mean Latency	1.252 (0.391)	1.731 (0.716)	5.250	0.022	KW
Total Word Count	126.000 (71.155)	119.938 (84.409)	0.078	0.781	ANOVA
MTLD	42.294 (24.467)	41.056 (21.867)	0.032	0.859	ANOVA
Mean WF Lexical frequency	0.007 (0.001)	0.007 (0.002)	1.177	0.282	ANOVA
Mean Word Length	3.597 (0.260)	3.572 (0.212)	0.121	0.729	ANOVA
Median Word Length	3.130 (0.341)	3.031 (0.125)	0.000	0.992 ^X	ANCOVA
Mean Syllables per Word	1.152 (0.057)	1.145 (0.043)	0.207	0.650	ANOVA
Coleman-Liau Index	3.446 (1.797)	3.530 (1.239)	0.030	0.864	ANOVA
Mean Yngve Depth	4.867 (1.909)	4.731 (1.369)	0.068	0.795	ANOVA
Noun Count	24.130 (13.016)	19.375 (11.949)	1.649	0.204	ANOVA
Proper Noun Count	0.009 (0.016)	0.012 (0.032)	0.122	0.728	ANOVA

KW = Kruskal-Wallis; ANOVA = Analysis of Variance; HC = Healthy Control; MCI = Mild cognitive impairment; MTLD = measure of textual lexical diversity; X = Significant age effects; WF = *wordfreq* lexical frequency

Bold = significant

Table 3. Measures of predictive ability for models utilizing speech behavior

	Sensitivity	Specificity	PPV	NPV	DOR	AUC
All Significant Features	38.46%	94.59%	71.43%	81.40%	10.94	0.825
Post-Pause Features	41.03%	93.70%	69.57%	81.89%	10.33	0.791

All Significant Features: Total filled pauses (INR), Total “uh” pauses (INR & DNR), Mean *wordfreq* frequency post- “uh” (FR), mean latency post-“um” (DNR), word count (INR), MTLN (INR), mean *wordfreq* frequency (INR), median word length (DNR) Coleman-Liau Index (PD & INR), mean Yngve depth (PD), proper noun ratio (INR), noun count (DNR)

Post-pause Features: Mean *wordfreq* frequency post- “uh” (FR), mean latency post-“um” (DNR)

FR=“Free-response task”; INR=“Immediate Narrative Recall task”; DNR=“Delayed Narrative Recall task”; PD=“Picture Description task”; PPV=“Positive Predictive Value”; NPV=“Negative Predictive Value”; DOR=“Diagnostic Odds Ratio”; AUC=“Area under the ROC curve”

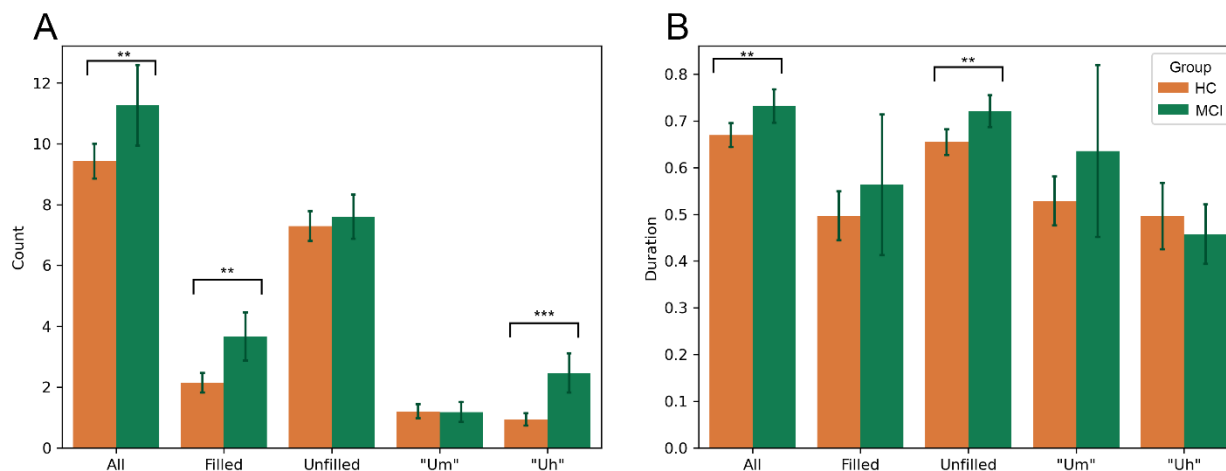


Figure 1. Count and duration of pause types separated by impairment status. (A) Total number of pauses separated by type of pause (all pauses, filled, unfilled, filled with “um”, filled with “uh”) and by impairment status (HC in orange vs MCI in green). All tasks were included in these comparisons. Significant differences between impairment status can be observed when examining All pauses, Filled pauses, and “Uh” pauses, all of which increased in the MCI group, with no differences found for unfilled and “um” pauses when all tasks were included. (B) Total duration in seconds separated by type of pause and impairment status. Unfilled pauses and All pauses were found to significantly increase in the MCI group. Significance is indicated using asterisks, with *** = $p < .001$, ** = $p < .01$, * = $p < .05$.

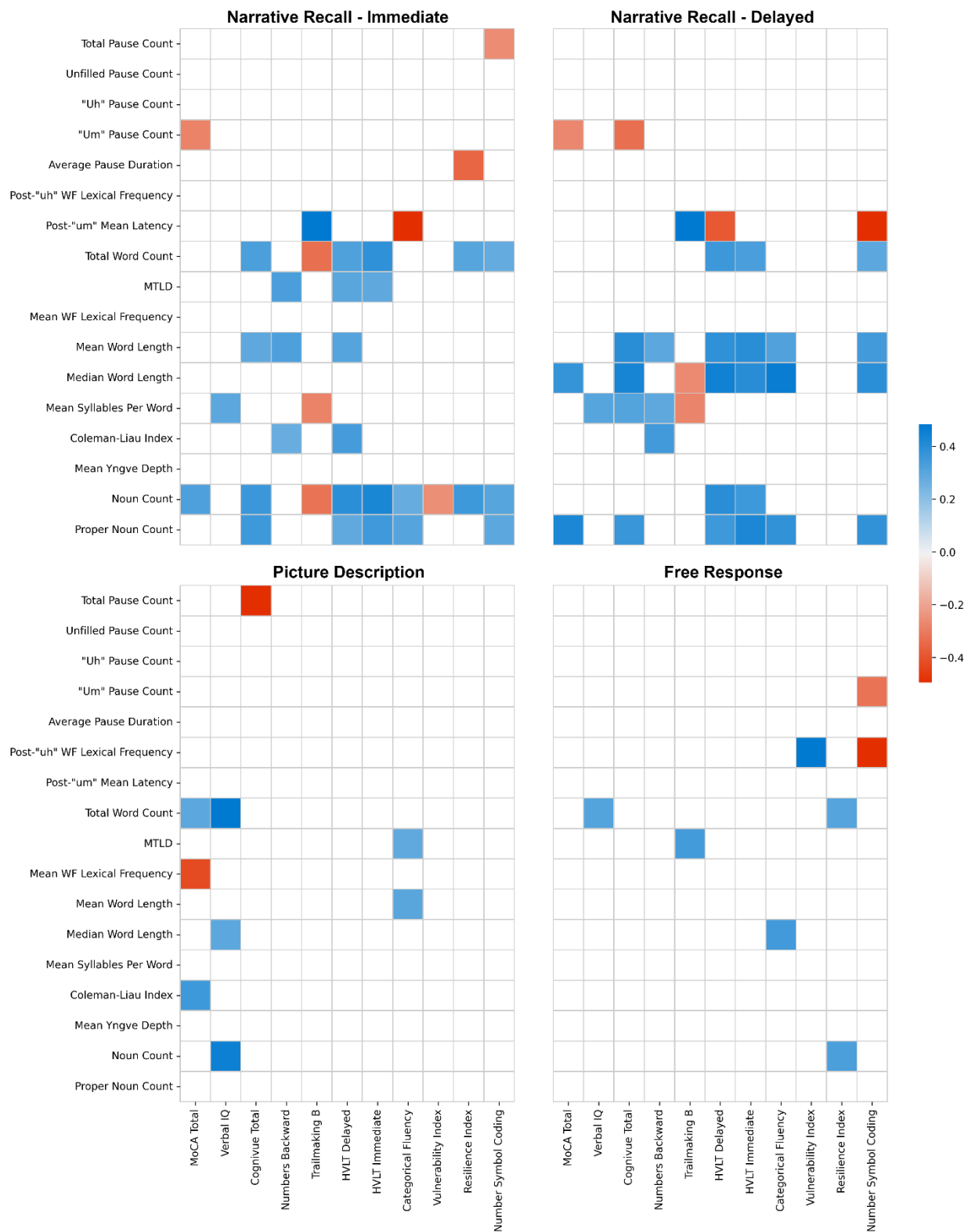


Figure 2. Correlation matrix using Pearson's correlations to compare significant speech metrics with cognitive and neuropsychological assessments for each of the four tasks. Spaces without color are non-significant determined by a p -value above 0.015.

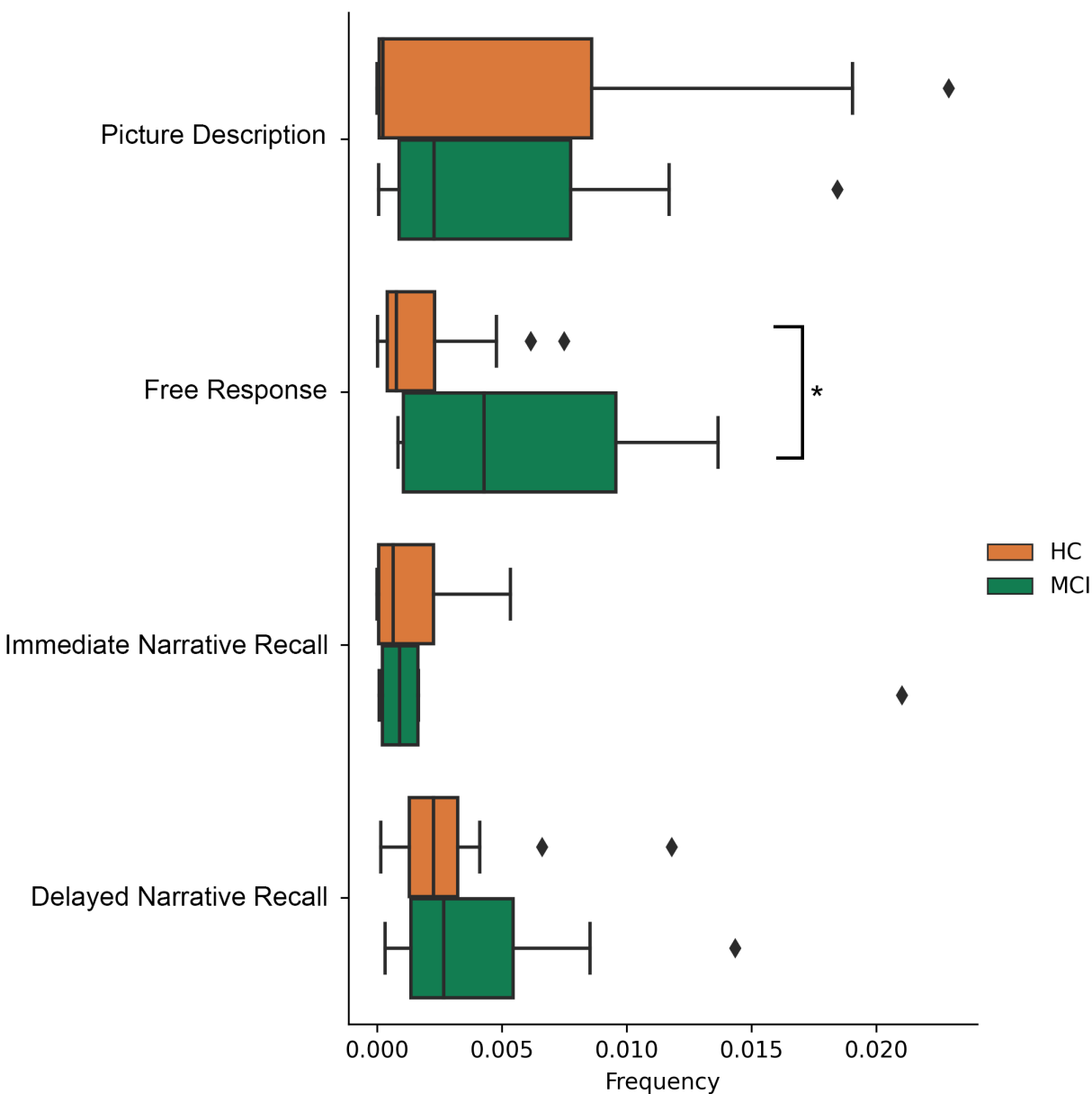


Figure 3. Task comparisons of post-“uh” pause lexical frequency rankings. This measure was used for comparison to highlight differences in pause and post-pause behavior between tasks. In both narrative recall tasks, all groups tended to be searching for comparatively uncommon words following “uh” pauses, while in the picture description task all groups searched for more common words. Only in the free response task were there differences between impairment status, with healthy controls (HC) searching only for more uncommon words and participants with mild cognitive impairment (MCI) producing words of more variable frequency following “uh” pauses.

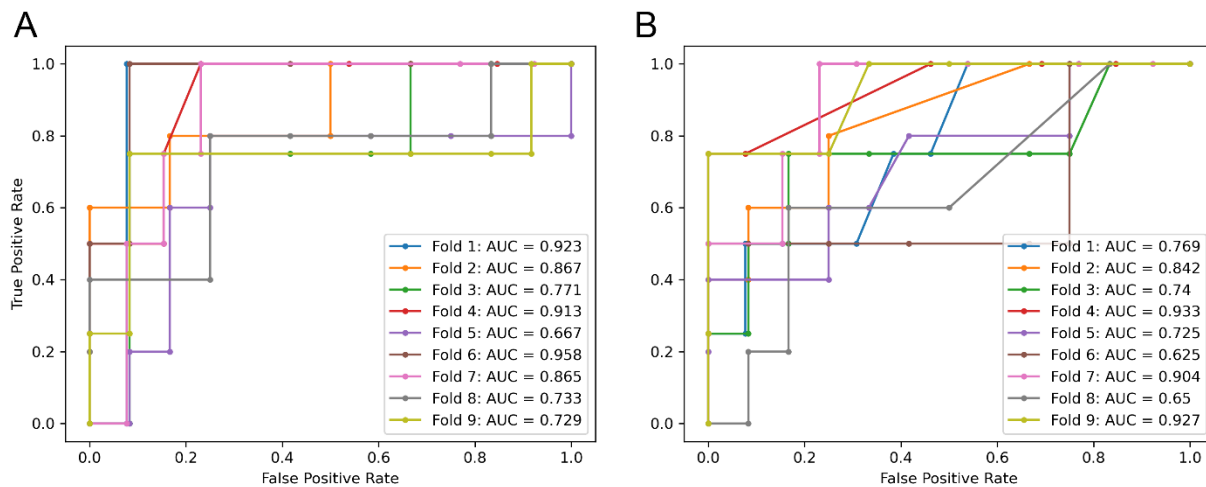


Figure 4. Areas under ROC curves for predictive models. (A) After performing feature selection on all available lexical features, 16 features were identified to best identify impairment status (HC vs MCI). Incorporating these into a repeated stratified 3-fold 3-repeat cross-validation procedure using *LightGBM* gradient boosted machines as the models generated a mean AUC of 0.828. AUCs for individual folds are depicted in multiple colors. (B) The best identified post-pause features (post-“uh” lexical frequency in free response and post-“um” latency in delayed narrative recall) were selected and used as sole features in a *LightGBM* gradient boosted machine, examined using repeated stratified 3-fold 3-repeat cross-validation. This parsimonious model performed similarly to the larger model with an AUC of 0.791. AUCs for individual folds are depicted in multiple colors.