

1                   **Detecting dementia using linguistic analysis:**  
2                   **Terry Pratchett's Discworld tells a more personal story**

3  
4  
5  
6                   Melody Pattison<sup>a</sup>, Ahmet Begde<sup>b,c</sup>, Thomas DW Wilcockson<sup>b\*</sup>

7  
8  
9                   <sup>a</sup>. School of English, Communication and Philosophy, Cardiff University, Cardiff, UK

10                  <sup>b</sup>. School of Sport, Exercise and Health Sciences, Loughborough University,  
11                  Loughborough, UK

12                  <sup>c</sup>. Department of Psychiatry, Oxford University, Oxford, UK

13 \* Correspondence

14 Thomas Wilcockson,

15 School of Sport, Exercise and Health Sciences,

16 Loughborough University,

17 Loughborough, LE11 3TU, UK

18 Email: t.wilcockson@lboro.ac.uk

19  
20 **Running title:** Detecting dementia using linguistic analysis

21  
22 **Keywords:** Dementia; Alzheimer's disease; Linguistics; Linguistic analysis; Terry Pratchett;

23 **Detecting dementia using linguistic analysis:**  
24 **Terry Pratchett's Discworld tells a more personal story**

25  
26 **Abstract**

27 Dementia, characterised by cognitive decline, significantly impacts language abilities. While  
28 the risk of dementia increases with age, it often manifests years before clinical diagnosis.  
29 Identifying early warning signs is crucial for timely intervention. Previous research has  
30 demonstrated that changes in language, such as reduced vocabulary diversity and simpler  
31 sentence structures, may be observed in individuals with dementia. This study investigates the  
32 potential of linguistic analysis to detect early signs of cognitive decline by examining the  
33 writing of Sir Terry Pratchett, a renowned author diagnosed with Posterior Cortical Atrophy  
34 (PCA), a form of dementia caused by Alzheimer's disease. This study analysed 33 Discworld  
35 novels by Terry Pratchett, comparing linguistic features (e.g., vocabulary size, lexical diversity,  
36 word class distribution) before and after a potential turning point identified through analysis of  
37 adjective type-token ratios (TTR). A significant decrease in lexical diversity (TTR) was  
38 observed for nouns and adjectives in later works. Total word count increased, while lexical  
39 diversity decreased, suggesting a shift towards simpler language. This shift coincided with a  
40 decrease in adjective TTR below a defined threshold, occurring approximately ten years before  
41 Pratchett's formal diagnosis. These findings suggest that subtle changes in linguistic patterns,  
42 such as decreased lexical diversity, may precede clinical diagnosis of dementia by a  
43 considerable margin. This research highlights the potential of linguistic analysis as a valuable  
44 tool for early detection of cognitive decline. Further research is needed to validate these  
45 findings in larger cohorts and explore the specific linguistic markers associated with different  
46 types of dementia.

## Detecting dementia using linguistic analysis:

### Terry Pratchett's Discworld tells a more personal story

47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
66  
67  
68  
69  
70  
71  
72  
73  
74  
75  
76  
77  
78  
79  
80

#### 1 Introduction

Dementia is an umbrella term for neurodegeneration that causes cognitive decline, the most common of which is Alzheimer's disease<sup>1</sup>. While the risk of developing dementia becomes more common with age, it is not an inevitable part of ageing<sup>2</sup>. Around a third of people over 85 may experience dementia, but it can also occur in younger individuals, though this is fortunately rare. The cognitive decline observed in people with Alzheimer's disease is caused by amyloid beta and tau proteins<sup>3</sup>, which are both naturally occurring in the brain, but in Alzheimer's disease become toxic and start causing damage to the brain<sup>4</sup>. The underlying reason why amyloid beta and tau become toxic in some people and not others is currently under debate, however, identifying those people with signs of Alzheimer's-related cognitive decline as early as possible would enable interventions to be utilised to delay or even prevent some of the damage cause by the toxic proteins.

People with dementia may first notice they have an issue when they experience increased episodes of confusion or issues with memory or language<sup>5</sup>. However, Alzheimer's pathology likely begins many years and perhaps decades before the onset of symptoms<sup>6</sup>. Indeed, research has shown that there are earlier warning signs of dementia which may be too subtle for a patient to be aware of, for example, problems with attention<sup>7</sup>. Further, research suggests that it is currently possible to predict who will experience dementia 12 years prior to formal diagnosis<sup>8</sup>. Begde, et al.<sup>8</sup> observed that reduced complex visual processing speed is significantly associated with a higher likelihood of a future dementia diagnosis and risk/protective factors. Therefore, it may be possible to use tasks which measure these visual and attentional difficulties to help identify at-risk people before their cognitive decline worsens.

Dementia also has a negative effect on both speech and writing<sup>9,10</sup>. Therefore measuring the first signs of decline in these functions may also provide an early biomarker for dementia. In early-stage Alzheimer's, researchers have observed impairments in both producing and understanding words and sentences<sup>11</sup>. These studies have also revealed a breakdown in the complex network of knowledge that gives meaning to objects and words, a system known as semantic memory<sup>12</sup>. Such research studies have used standardised tests like word fluency and picture naming to assess language abilities<sup>13</sup>. These tests have shown that factors such as how

81 often a word is used, how familiar a person is with a word, and the age at which a word was  
82 learned can all affect language performance<sup>14</sup>. Therefore, by looking for changes in how  
83 someone uses language, then this may provide an early warning sign for dementia. The  
84 complexity of sentence structure, as measured by factors like the number of clauses per  
85 utterance, decreases with age in both spoken and written language<sup>15</sup>. Older adults struggle more  
86 with complex sentence structures, such as those with left-branching clauses, compared to  
87 younger adults<sup>16</sup>. A longitudinal study by Kemper, et al.<sup>17</sup> revealed a significant decline in  
88 grammatical complexity in older adults, particularly in those with dementia. Additionally,  
89 Bates, et al.<sup>18</sup> demonstrated that grammatical production is impaired in individuals with  
90 Alzheimer's. While patients do not typically produce overt grammatical errors, they encounter  
91 difficulties in selecting the most appropriate grammatical forms to convey their intended  
92 meaning. This suggests that the underlying issue may lie in accessing and retrieving the optimal  
93 fit between meaning and grammatical expression.

94  
95 Overall, linguistic changes are to be expected as people age. However, these changes become  
96 more profound within people with cognitive decline. If a patient's writing history is available,  
97 then linguistic analysis techniques could be used to supplement clinical assessments or as a  
98 standalone early detection tool. Recent studies have done exactly this by measuring individual  
99 writers' publications over their careers to analyse how their language use has evolved. Garrard,  
100 et al.<sup>19</sup> studied the works of Iris Murdoch, a renowned English author who was diagnosed with  
101 Alzheimer's posthumously. Her final novel, published shortly before her diagnosis, is widely  
102 considered to exhibit signs of cognitive decline. While Garrard found minimal differences in  
103 overall structure and syntax, they observed significant and consistent variations in lexical  
104 diversity and word choice between the final book and control books from earlier in Murdoch's  
105 career. These results provide evidence that Alzheimer's may indeed be measured using  
106 linguistic analysis. Le, et al.<sup>20</sup> explored this further by including additional authors, numbers  
107 of books, and improved analysis techniques. Le et al analysed two authors believed to have  
108 Alzheimer's disease during their careers, Iris Murdoch and Agatha Christie, as well as P.D.  
109 James to act as a control participant, who published until the age of 88 without experiencing  
110 evidence of cognitive decline. They included twenty of Murdoch's twenty-six novels,  
111 published between ages 35 and 76, sixteen of Christie's novels written between ages 28 and  
112 82, and fifteen of the novels of P.D. James. They then produced an analysis of the novels at the  
113 lexical level, using a variety of measures, including vocabulary size, lexical repetition, lexical  
114 specificity, word-class deficits, and fillers. Type-token ratio (TTR) calculates the proportion of

115 unique words to the total word count, and the word-type introduction rate (WTIR), which  
116 measures the rate at which new words are introduced in the text, calculated every 10,000 words.  
117 Lexical repetition, while intentional repetition can be a stylistic device, an increasing rate of  
118 repeated words may suggest a limited vocabulary or difficulty accessing words. To examine  
119 this, they conducted two analyses: a global analysis and a local analysis. Lexical specificity is  
120 calculated by the frequency of indefinite nouns and high-frequency, low-imagery verbs in each  
121 text. A higher proportion of these generic words suggests lower overall lexical specificity.  
122 Word class deficit (WCD) is an analysis of the distribution of word classes across each text,  
123 examining both the total number of words and the number of unique words. This allows for  
124 identification of potential deficits or overreliance on specific word classes and to measure the  
125 vocabulary size of open classes. Filler words are a measure of the proportion of interjections  
126 and filler words. While these words often appear in dialogue, fiction authors strive for natural-  
127 sounding conversations. However, this measure may be influenced by stylistic choices rather  
128 than cognitive decline and should be interpreted with care. Le, et al observed that TTR and  
129 WTIR were associated with cognitive decline and a decline in vocabulary led to an increase in  
130 repetitions in content words, and a word-class deficit can be seen in noun-token proportion,  
131 with a compensatory increase in verb-token proportion. They also observed a deficit in noun  
132 tokens that is significantly correlated with a rise in verb and pronoun tokens. Syntactic-  
133 complexity results were also found to fluctuate in a relatively wider range. Interestingly, they  
134 also report that deficits in Murdoch's writing appeared in Murdoch's late 40s and early 50s,  
135 which suggests that language deficits are observed many years before a formal diagnosis and  
136 indicates that Alzheimer's disease has a long preclinical period. Therefore, linguistic analysis  
137 would appear to show promise in identifying whether an author has experienced cognitive  
138 decline and may even indicate when the preclinical phase of dementia has begun.

139  
140 The current research further explores the idea of using linguistic analysis in dementia by  
141 studying the works of Sir Terry Pratchett. Terry Pratchett was an English author, humourist,  
142 and satirist, best known for his Discworld series of 41 comic fantasy novels published between  
143 1983–2015. Terry Pratchett was diagnosed with Posterior Cortical Atrophy (PCA) in  
144 December 2007. This diagnosis came at a time when he was still actively writing and  
145 publishing his beloved Discworld series. Despite the challenges posed by his condition,  
146 Pratchett continued to write and advocate for dementia awareness until his passing in 2015.  
147 PCA is a rare form of Alzheimer's disease that primarily affects visual processing and spatial  
148 awareness<sup>21</sup>. It affects areas in the back of the brain responsible for spatial perception, complex

149 visual processing, spelling, and calculation<sup>22</sup>. Given Terry Pratchett's prolific writing career  
150 and the fact he continued writing after his diagnosis, a linguistic analysis of his novels could  
151 provide valuable insights into the potential early signs of cognitive decline. By comparing his  
152 earlier works to his later ones, particularly those written closer to his dementia diagnosis, it  
153 would be possible to identify subtle changes in linguistic patterns, such as decreased lexical  
154 diversity, increased reliance on simpler sentence structures, and a decline in the use of specific  
155 and descriptive language. Such an analysis could contribute to our understanding of the  
156 linguistic markers of Alzheimer's disease and potentially aid in the development of early  
157 detection tools.

## 158 **2 Method**

159 The methodology of this study enabled the analysis of the lexical diversity of Terry Pratchett's  
160 writing pre- and post-PCA diagnosis, and explored whether this measure can be used as a  
161 predictor of dementia. In order to explore this lexical diversity, SketchEngine<sup>23,24</sup> was used to  
162 establish the TTR of content lexical items, as well as vocabulary repetitions, in 33 out of the  
163 41 of Terry Pratchett's 'Discworld' novels. Several titles were excluded from the analysis due  
164 to them being either shorter than the other full-length novels<sup>†1</sup>, or because they are part of his  
165 titles for younger readers<sup>†2</sup>. They were excluded due to the fact that for SketchEngine to  
166 accurately measure TTR, it is important for texts to be of roughly the same length and/or aimed  
167 at similar levels of readership.

168

169 Word classes of interest included nouns, verbs, adjectives, and adverbs, in addition to the  
170 numbers of unique lemmas to determine the overall vocabulary size. SketchEngine was chosen  
171 to retrieve this information due to its ability to work with a large sample of text and for the  
172 researchers to be able to specify lemmas and word classes of interest for analysis.

173

174 In order to retrieve the information, plain text files of the 33 Terry Pratchett books of interest  
175 were imported into SketchEngine<sup>†3</sup>. Analysis of TTR to measure the WCD and overall  
176 vocabulary size, as well as analysis of repetitions (lines where an individual word is repeated

---

<sup>†1</sup> 'Eric' (1990), 'The Last Hero' (2001).

<sup>†2</sup> 'The Amazing Maurice and His Educated Rodents' (2001), 'The Wee Free Men' (2003), 'A Hat Full of Sky' (2004), 'Wintersmith' (2006), 'I Shall Wear Midnight' (2010), 'The Shepherd's Crown' (2015).

<sup>†3</sup> For ethical considerations, it should be noted that all books have been individually purchased by at least one of the three authors of this paper.

177 within ten surrounding words) was then conducted for each book. To calculate the TTR in  
178 relation to vocabulary size, the number of individual lemmas was divided by the total number  
179 of words for each book. To calculate the TTR to determine the WCD, the same method was  
180 applied to each word class. Most of the tokens could be retrieved automatically; however, some  
181 manual analysis of values was required in some instances, such as where SketchEngine  
182 returned the character 'Death' and the noun 'death' as the same value when considering the  
183 top ten repetitions of each book. Percentages and TTRs were also calculated manually. These  
184 data were then statistically analysed to compare the linguistic measures before and after  
185 Pratchett's PCA diagnosis in 2007. Thus, 29 books<sup>†4</sup> were analysed as pre-diagnosis, and 4  
186 books<sup>†5</sup> as post-diagnosis.

187  
188 The advantages of using WCD, vocabulary size and repetitions can be exemplified through the  
189 work of Le<sup>25</sup> and Le, et al.<sup>20</sup>, as mentioned in the introduction of this paper. Despite using a  
190 different method to extract lexical tokens, their findings indicated that a loss of TTR in  
191 vocabulary size, an increase in repetitions, and a deficit in nouns, alongside other measures,  
192 could reliably predict cognitive linguistic decline in the works of Murdoch and Christie, which  
193 were not evident in the works of James. It is therefore our intention to replicate such linguistic  
194 analysis as part of the present study.

195

## 196 **2.1 Data availability**

197 The data that support the findings of this study are available from the corresponding author,  
198 upon reasonable request.

199

## 200 **2.2 Statistical Analyses**

201 Descriptive statistics were calculated to summarise the key linguistic features used in  
202 Pratchett's books. Independent t-tests were used to compare linguistic measures released before  
203 and after the dementia diagnosis. Linear regression analyses were conducted to investigate the  
204 relationship between various TTR types and age. Additionally, Receiver Operating  
205 Characteristic (ROC) curve analyses were performed to evaluate the accuracy of TTR measures

---

†4 From *'The Colour of Magic'* (1983) to *'Thud!'* (2005).

†5 From *'Making Money'* (2007) to *'Raising Steam'* (2013).

206 in distinguishing between pre- and post-diagnosis phases. All statistical analyses were  
 207 performed using SPSS (version 29), with the significance level set at  $p < 0.05$ .

### 208 3 Results

209 The analysis included 33 books, with 29 books written before the dementia diagnosis and 4  
 210 book written after the diagnosis. Independent t-tests were conducted to compare linguistic  
 211 features between these periods (see Table 1).

212 Table 1: Comparisons of Linguistic Features Pre- and Post-Dementia Diagnosis

<b>Linguistic Features</b>	<b>Pre-Diagnosis Mean (SD) (n=29)</b>	<b>Post-Diagnosis Mean (SD) (n=4)</b>	<b>p-value</b>
<b>Total Nouns</b>	21,819 (4,042)	29,195 (1,696)	0.001*
<b>Type Token Ratio (nouns)</b>	0.197 (0.025)	0.170 (0.008)	0.044*
<b>Total Verbs</b>	20,015 (3,758)	26,405 (2,398)	0.003*
<b>Type Token Ratio (verbs)</b>	0.078 (0.014)	0.067 (0.005)	0.138
<b>Total Adjectives</b>	5,422 (759)	7,394 (426)	<0.001*
<b>Type Token Ratio (adjective)</b>	0.241 (0.024)	0.213 (0.010)	0.028*
<b>Total Adverbs</b>	7,206 (1,356)	10,082 (1,227)	<0.001*
<b>Type Token Ratio (adverb)</b>	0.079 (0.017)	0.063 (0.004)	0.076
<b>Total Words</b>	94,379 (16,646)	129,226 (10,325)	<0.001*
<b>TTR</b>	0.071 (0.010)	0.063 (0.004)	0.114
<b>Lemmas</b>	6,604 (623)	8,166 (250)	<0.001*

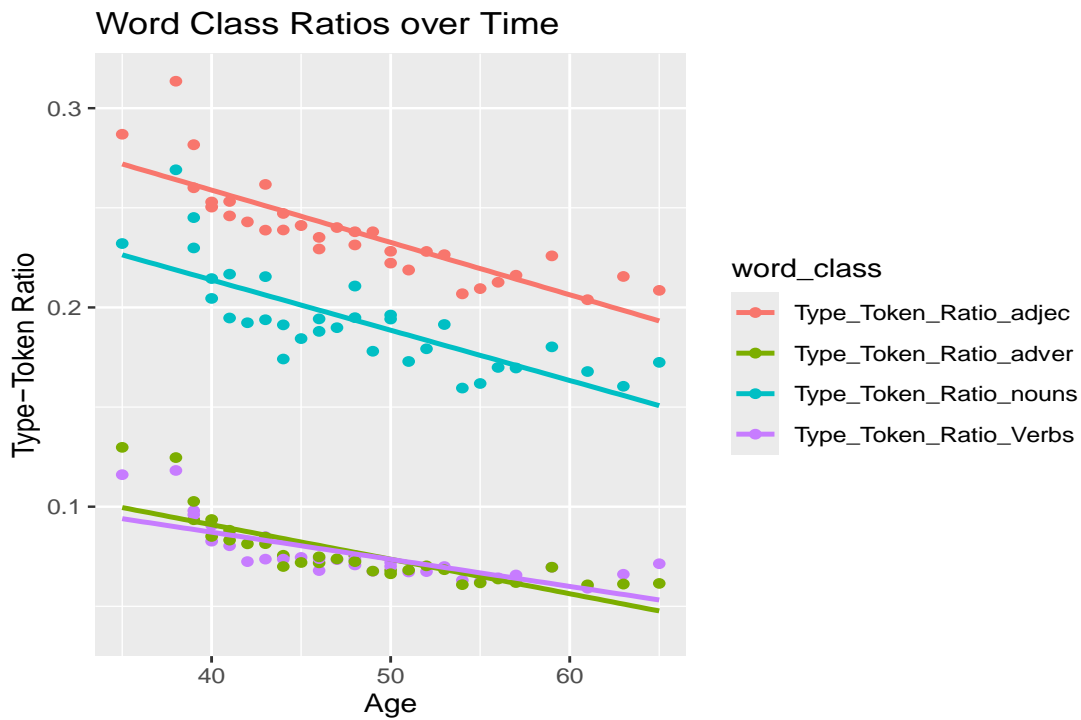
213 Note: \* indicates statistical significance ( $p < 0.05$ ), TTR = Total-Type Token Ratio

214 The analysis showed significant changes in several linguistic features following Pratchett's  
 215 dementia diagnosis. Books written after diagnosis showed significantly higher total word  
 216 counts across all word classes (nouns:  $p=0.001$ ; verbs:  $p=0.003$ ; adjectives:  $p<0.001$ ; adverbs:



217  $p < 0.001$ , see Table 1) and overall text length ( $p < 0.001$ ). However, this increase was  
218 accompanied by a decrease in lexical diversity, as shown by significant lower type-token ratios  
219 for nouns ( $p = 0.044$ ) and adjectives ( $p = 0.028$ ). No significant changes were observed in the  
220 type-token ratios for verbs and adverbs, or the overall TTR, suggesting that some aspects of  
221 linguistic complexity remained stable after the diagnosis. While total word usage increased,  
222 the variety of words used (measured by type-token ratios) either decreased or remained  
223 unchanged across different word classes. The total number of unique word forms (lemmas)  
224 increased significantly ( $p < 0.001$ ).

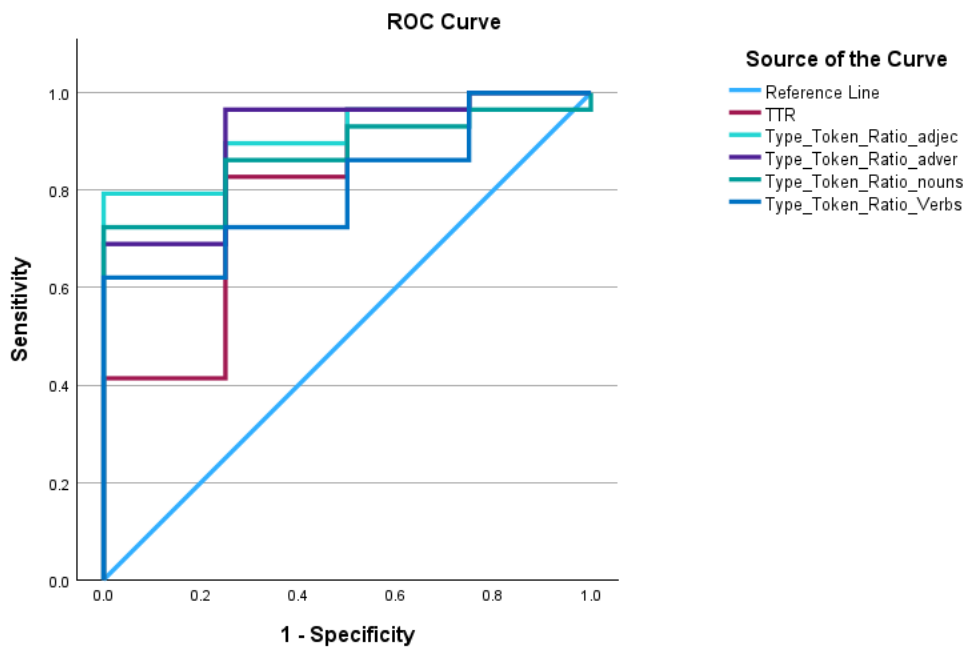
225 Figure 1 shows the changes in TTR across different word classes (adjectives, adverbs, nouns,  
226 and verbs) in Pratchett's writing over time. All word classes demonstrated a declining trend in  
227 TTR with age, indicating a general decrease in lexical diversity. Linear regression analyses  
228 were conducted to examine the relationship between various TTR and age. All models showed  
229 significant predictive relationships ( $p < 0.001$ ). The TTR for adjectives emerged as the  
230 strongest predictor ( $F(1,31) = 73.101$ ,  $p < 0.001$ ), followed by adverb TTR ( $F(1,31) = 53.694$ ,  
231  $p < 0.001$ ). Noun TTR ( $F(1,31) = 45.728$ ,  $p < 0.001$ ) and verb TTR ( $F(1,31) = 39.413$ ,  $p <$   
232  $0.001$ ) also showed significant relationships with age. The overall TTR demonstrated a  
233 comparatively lower but significant, predictive relationship ( $F(1,31) = 31.407$ ,  $p < 0.001$ ).  
234 These findings suggest that lexical diversity might be a strong indicator of age-related linguistic  
235 changes.



236

237 Figure 1: Type/token ratio over time in Pratchett’s Discworld books

238 Receiver Operating Characteristic (ROC) curve analysis was performed to assess the accuracy  
 239 of various TTR measures in detecting linguistic changes associated with dementia. All TTR  
 240 measures demonstrated significant predictive ability ( $p < 0.05$ ), with AUC values ranging from  
 241 0.80 to 0.91.



242

243 Figure 2. Receiver Operating Characteristic (ROC) curve analysis for the TTR variables

244 The TTR for adjectives showed the highest classification accuracy (AUC = 0.91, 95% CI: 0.80-  
 245 1.02,  $p < 0.001$ ), with a cut-off value of 0.227 showing 76% sensitivity and 100% specificity.  
 246 This was followed by adverb TTR (AUC = 0.90, 95% CI: 0.76-1.04,  $p < 0.001$ ), which  
 247 demonstrated 96% sensitivity and 75% specificity at a cut-off of 0.616. Noun TTR also showed  
 248 strong predictive performance (AUC = 0.87, 95% CI: 0.73-1.00,  $p < 0.001$ ), with 86%  
 249 sensitivity and 75% specificity at a cut-off of 0.172.

250 Verb TTR and overall TTR show relatively lower but significant accuracy (both AUC = 0.80,  
 251  $p \leq 0.015$ ). Verb TTR achieved highest specificity (100%) but lower sensitivity (62%) at a cut-  
 252 off of 0.717, while overall TTR showed balanced performance with 82% sensitivity and 75%  
 253 specificity at a cut-off of 0.065.

254 Table 2. Area Under the Curve results for each TTR variable

Measure	AUC (95% CI)	Cut-off	Sensitivity	Specificity	p-value
TTR-Adjectives	0.91 (0.80-1.02)	0.227	0.76	1.00	<0.001*
TTR-Adverbs	0.90 (0.76-1.04)	0.616	0.96	0.75	<0.001*
TTR-Nouns	0.87 (0.73-1.00)	0.172	0.86	0.75	<0.001*
TTR-Verbs	0.80 (0.62-0.98)	0.717	0.62	1.00	0.001*
Overall TTR	0.80 (0.55-1.04)	0.065	0.82	0.75	0.015*

255 Note: AUC = Area Under the Curve; CI = Confidence Interval; \* indicates statistical  
 256 significance ( $p < 0.05$ )

257

258 As TTR for adjectives showed the highest diagnostic accuracy, we can use the cut-off value of  
 259 0.227 to identify when in Pratchett's writing his TTR for adjectives started to fall below this  
 260 cut-off score. It was found that eleven of Pratchett's works were found to have a TTR for  
 261 adjectives lower than 0.227, with the earliest published being *The Last Continent* (Discworld  
 262 22), which was published in May 1998 – 9 years and 7 months before his formal diagnosis. All  
 263 books published after this date were found to have a TTR for adjectives of less than 0.227  
 264 whilst all books published before this date were found to have a TTR for adjectives more than  
 265 0.227. The only outliers are Discworlds 23 (*Carpe Jugulum*) and 25 (*The Truth*) which both  
 266 had scores of 0.228, which do not meet the cut-off despite being published after Discworld 22,  
 267 but are only 0.001 outside the cut-off.

268 Because Discworld 22 was the first book to fall below the TTR for adjectives cut-offs,  
 269 independent t-tests were performed on pre and post Discworld 22 to compare linguistic features  
 270 between these periods (see Table 3).

271 Table 3: Comparisons of Linguistic Features Pre- and Post-Discworld 22

<b>Linguistic Features</b>	<b>Pre-Discworld 22 Mean (SD) (n=20)</b>	<b>Post-Discworld 22 Mean (SD) (n=13)</b>	<b>p-value</b>
<b>Total Nouns</b>	20,179 (3,548)	26,613 (2,769)	<0.001*
<b>Type Token Ratio (nouns)</b>	0.206 (0.024)	0.175 (0.013)	<0.001*
<b>Total Verbs</b>	18,388 (3,240)	24,484 (2,312)	<0.001*
<b>Type Token Ratio (verbs)</b>	0.083 (0.015)	0.067 (0.004)	<0.001*
<b>Total Adjectives</b>	5,098 (677)	6,528 (677)	<0.001*
<b>Type Token Ratio (adjective)</b>	0.251 (0.021)	0.217 (0.009)	<0.001*
<b>Total Adverbs</b>	6,595 (1,153)	9,032 (1,042)	<0.001*
<b>Type Token Ratio (adverb)</b>	0.085 (0.017)	0.065 (0.004)	<0.001*
<b>Total Words</b>	87,197 (14,372)	116,152 (12,303)	<0.001*
<b>TTR</b>	0.074 (0.010)	0.065 (0.004)	<0.001*
<b>Lemmas</b>	6,337 (511)	7,496 (583)	<0.001*

272 Note: \* indicates statistical significance ( $p < 0.05$ ), TTR = Total-Type Token Ratio

273 The analysis showed significant changes in several linguistic features following Pratchett's  
 274 publication of Discworld 22. Discworld books written after Discworld 22 showed significantly  
 275 higher total word counts across all word classes (nouns:  $p < 0.001$ ; verbs:  $p < 0.001$ ; adjectives:  
 276  $p < 0.001$ ; adverbs:  $p < 0.001$ , see Table 3) and overall text length ( $p < 0.001$ ). However, this  
 277 increase was accompanied by a decrease in lexical diversity, as shown by significant lower  
 278 type-token ratios for nouns:  $p < 0.001$ ; verbs:  $p < 0.001$ ; adjectives:  $p < 0.001$ ; and adverbs:

279 p<0.001. Overall TTR (p<0.001) and lemmas also significantly differed (p<0.001), indicating  
280 overall that linguistic features differed on all aspects pre and post Discworld 22.

#### 281 **4 Discussion**

282 The current study aimed to explore the potential of linguistic analysis as a tool for early  
283 detection of cognitive decline, specifically focusing on the case of Terry Pratchett and his  
284 diagnosis of PCA due to Alzheimer's disease. Our analysis of Pratchett's Discworld series  
285 revealed significant changes in linguistic patterns over time. The most notable finding was a  
286 significant decrease in lexical diversity, as measured by TTR, for adjectives and nouns in  
287 Pratchett's later works. This suggests a decline in vocabulary richness and a reliance on simpler  
288 language structures. While the overall TTR remained relatively stable, the decrease in lexical  
289 diversity within specific word classes indicates a subtle but significant change in linguistic  
290 style. These findings align with previous research on the linguistic markers of Alzheimer's  
291 disease and other forms of dementia<sup>20</sup>. Studies have shown that individuals with dementia often  
292 exhibit reduced vocabulary diversity, simpler sentence structures, and increased reliance on  
293 clichés and formulaic language<sup>26</sup>. Our analysis of Pratchett's works suggests that similar  
294 linguistic changes could potentially also be observed in individuals with PCA due to  
295 Alzheimer's disease.

296

297 The high predictive accuracy of TTR measures, particularly for adjectives, suggests that  
298 linguistic analysis could be a valuable tool for early detection of cognitive decline. By  
299 identifying subtle changes in language use, it may be possible to detect the onset of dementia  
300 years before a formal diagnosis. This early detection could enable timely interventions and  
301 potentially slow the progression of the disease. The analysis revealed a significant shift in  
302 Pratchett's writing style following the publication of "The Last Continent" (Discworld 22). This  
303 book marks a potential turning point, as it was the first to exhibit a TTR for adjectives below  
304 the established cut-off value of 0.227, a threshold with a high diagnostic accuracy for  
305 potentially detecting linguistic changes related to dementia in the Discworld series. To further  
306 investigate this shift, we compared linguistic features in books published before and after "The  
307 Last Continent." The results revealed significant differences across various measures. Post-  
308 "Last Continent" books exhibited higher word counts across all word classes (nouns, verbs,  
309 adjectives, adverbs), increased overall text length, and a significant decrease in lexical diversity  
310 across all word classes as measured by TTR. Furthermore, the total number of unique lemmas

311 (vocabulary size) also increased significantly. These findings suggest a substantial shift in  
312 Pratchett's writing style after this point, characterised by an increase in word count but a  
313 decrease in lexical diversity and a potential shift towards simpler sentence structures. It should  
314 be noted that this book was published 9 years and 7 months before his formal diagnosis,  
315 indicating a long preclinical period prior to diagnosis. This observation is striking and  
316 highlights the potential for linguistic analysis to identify subtle changes in writing style that  
317 may precede clinical diagnosis of dementia by a considerable margin.

318

319 However, it is crucial to acknowledge the limitations of this study. Firstly, the analysis was  
320 based on a single case study, limiting the generalisability of the findings. Further research is  
321 needed to validate these findings in a larger sample of individuals with dementia, including  
322 those with PCA. Secondly, while PCA and Alzheimer's disease share some common features  
323 and PCA is thought to be caused by Alzheimer's disease, they are distinct conditions with  
324 different clinical presentations<sup>21</sup>. It is possible that the linguistic markers of PCA may differ  
325 from those of Alzheimer's disease. Furthermore, several methodological limitations should be  
326 considered. The precise chronology of Pratchett's writing is uncertain. We do not know the  
327 exact dates of writing for each book, whether he worked on multiple books concurrently, or if  
328 the publication order accurately reflects the writing order. These uncertainties could introduce  
329 potential biases into the analysis. Additionally, the observed changes in Pratchett's writing style  
330 after his diagnosis may not solely be attributed to PCA. Factors such as reduced writing time  
331 due to the disease, potential collaboration with other writers, or significant editorial alterations  
332 could also have contributed to these changes. Finally, it is crucial to emphasise that the  
333 observed changes in Pratchett's writing may not exclusively reflect cognitive decline due to  
334 PCA. Age-related changes in writing style are expected, and it is possible that some of the  
335 observed changes represent natural stylistic evolution rather than disease-related decline.

336

337 Despite these limitations, this study provides valuable insights into the potential of linguistic  
338 analysis as a tool for detecting early signs of cognitive decline. By analysing the works of Terry  
339 Pratchett, we have demonstrated how subtle changes in language use can be indicative of  
340 underlying cognitive impairment. Further research is needed to explore the full potential of  
341 linguistic analysis as a diagnostic tool for dementia, including the development of more  
342 sophisticated analytical methods and the investigation of larger and more diverse datasets.  
343 Future research could explore the use of more advanced linguistic analysis techniques, such as  
344 computational linguistics and machine learning, to identify additional linguistic markers of

345 cognitive decline. By developing more sensitive and accurate diagnostic tools, we may be able  
346 to improve early detection and intervention for dementia.

347

348 In conclusion, our study provides evidence that linguistic analysis can be a valuable tool for  
349 detecting early signs of cognitive decline. By analysing the works of Terry Pratchett, we have  
350 demonstrated how subtle changes in language use can be indicative of underlying cognitive  
351 impairment. The results also emphasises that language deficits may be observed many years  
352 before a formal diagnosis and indicates that Alzheimer's disease has a long preclinical period,  
353 in the case of Terry Pratchett, potentially almost ten years. Further research is now needed to  
354 explore the full potential of linguistic analysis as a diagnostic tool for dementia.

355

### 356 **Funding**

357 This research did not receive any specific grant from funding agencies in the public,  
358 commercial, or not-for-profit sectors.

359

### 360 **Competing Interests**

361 The authors report no competing interests.

362

### 363 **References**

- 364 [1] Fymat AL. On dementia and other cognitive disorders. *Clin Res Neurol*. 2019;2(1):1-4.
- 365 [2] Fotuhi M, Hachinski V, Whitehouse PJ. Changing perspectives regarding late-life  
366 dementia. *Nat Rev Neurol*. 2009;5(12):649-658.
- 367 [3] Allsop D, Mayes J. Amyloid  $\beta$ -peptide and Alzheimer's disease. *Essays Biochem*.  
368 2014;56:99-110.
- 369 [4] Breijyeh Z, Karaman R. Comprehensive review on Alzheimer's disease: causes and  
370 treatment. *Molecules*. 2020;25(24):5789.
- 371 [5] Knopman DS. The initial recognition and diagnosis of dementia. *Am J Med*.  
372 1998;104(4):2S-12S.
- 373 [6] Lloret A, Esteve D, Lloret MA, et al. When does Alzheimer's disease really start? The role  
374 of biomarkers. *Int J Mol Sci*. 2019;20(22):5536.
- 375 [7] Wilcockson TD, Mardanbegi D, Xia B, et al. Abnormalities of saccadic eye movements in  
376 dementia due to Alzheimer's disease and mild cognitive impairment. *Aging (Albany NY)*.  
377 2019;11(15):5389.

- 378 [8] Begde A, Wilcockson T, Brayne C, Hogervorst E. Visual processing speed and its  
379 association with future dementia development in a population-based prospective cohort: EPIC-  
380 Norfolk. *Sci Rep.* 2024;14(1):5016.
- 381 [9] Calzà L, Gagliardi G, Favretti RR, Tamburini F. Linguistic features and automatic  
382 classifiers for identifying mild cognitive impairment and dementia. *Comput Speech Lang.*  
383 2021;65:101113.
- 384 [10] Luzzatti C, Laiacona M, Agazzi D. Multiple patterns of writing disorders in dementia of  
385 the Alzheimer type and their evolution. *Neuropsychologia.* 2003;41(7):759-772.
- 386 [11] Szatloczki G, Hoffmann I, Vincze V, Kalman J, Pakaski M. Speaking in Alzheimer's  
387 disease, is that an early sign? Importance of changes in language abilities in Alzheimer's  
388 disease. *Front Aging Neurosci.* 2015;7:195.
- 389 [12] Nebes RD. Semantic memory in Alzheimer's disease. *Psychol Bull.* 1989;106(3):377.
- 390 [13] Maseda A, Lodeiro-Fernández L, Lorenzo-López L, et al. Verbal fluency, naming and  
391 verbal comprehension: three aspects of language as predictors of cognitive impairment. *Aging*  
392 *Ment Health.* 2014;18(8):1037-1045.
- 393 [14] Ferrante FJ, Migeot J, Birba A, et al. Multivariate word properties in fluency tasks reveal  
394 markers of Alzheimer's dementia. *Alzheimers Dement.* 2024;20(2):925-940.
- 395 [15] Burke DM, Shafto MA. Language and aging. In: Craik FIM, Salthouse TA, eds. *The*  
396 *Handbook of Aging and Cognition.* 3rd ed. New York, NY: Psychology Press; 2008:373-443.
- 397 [16] Maxim J, Bryan K. *Language of the elderly: a clinical perspective.* London, UK: Whurr  
398 Pub Ltd; 1994.
- 399 [17] Kemper S, Thompson M, Marquis J. Longitudinal change in language production: effects  
400 of aging and dementia on grammatical complexity and propositional content. *Psychol Aging.*  
401 2001;16(4):600.
- 402 [18] Bates E, Harris C, Marchman V, Wulfeck B, Kritchewsky M. Production of complex  
403 syntax in normal ageing and Alzheimer's disease. *Lang Cogn Process.* 1995;10(5):487-539.
- 404 [19] Garrard P, Maloney LM, Hodges JR, Patterson K. The effects of very early Alzheimer's  
405 disease on the characteristics of writing by a renowned author. *Brain.* 2005;128(2):250-260.
- 406 [20] Le X, Lancashire I, Hirst G, Jokel R. Longitudinal detection of dementia through lexical  
407 and syntactic changes in writing: a case study of three British novelists. *Lit Linguist Comput.*  
408 2011;26(4):435-461.
- 409 [21] Mendez MF, Ghajarian M, Perryman KM. Posterior cortical atrophy: clinical  
410 characteristics and differences compared to Alzheimer's disease. *Dement Geriatr Cogn Disord.*  
411 2002;14(1):33-40.



- 412 [22] Yong KX, Crutch SJ, Schott JM. Posterior cortical atrophy. In: Oxford Textbook of  
413 Neurologic and Neuropsychiatric Epidemiology. Oxford, UK: Oxford University Press;  
414 2020:141.
- 415 [23] Sketch Engine. Sketch Engine website. <http://www.sketchengine.eu>. Published 2003.  
416 Accessed November 15, 2019.
- 417 [24] Kilgarriff A, Baisa V, Bušta J, et al. The Sketch Engine: ten years on. *Lexicography*.  
418 2014;1(1):7-36.
- 419 [25] Le X. Longitudinal Detection of Dementia Through Lexical and Syntactic Changes in  
420 Writing [master's thesis]. Toronto, ON: University of Toronto, Department of Computer  
421 Science; 2010.
- 422 [26] Banovic S, Zunic LJ, Sinanovic O. Communication difficulties as a result of dementia.  
423 *Mater Sociomed*. 2018;30(3):221.