

# A Comparison of the Validity of Measurement Methods for the General English Proficiency through Dictation and Read-Aloud Performances

Anonymous ACL submission

## Abstract

This paper compares three measurement methods for the general proficiency of learners of English as a second language (GEP). If students' GEP can be measured on course materials frequently, for instance, at the beginning and end of a semester, English teachers can confirm students' levels of learning achievement. So far, English teachers have two options for GEP measurement: calculating scores for read-aloud or those of dictation performance. This study expands an option to measure GEP using both dictation and read-aloud performances. When comparing the three types of measurement methods, the experimental results suggest that GEP should be measured by calculating dictation and read-aloud performances.

## 1 Introduction

Evaluating general English proficiency (GEP), which includes listening, speaking, reading, and writing, is an essential task for teachers of English as a second language. One of the goals of GEP evaluation is to investigate a learner's learning outcome by comparing their GEP at the beginning and end of a semester. GEP evaluation can be conducted using English tests such as the Test of English as a Foreign Language and Test of English for International Communication (TOEIC), because the use of these tests reduces teachers' time and effort in test preparation.

However, when these tests are used as a classroom-based assessment, that is, an evaluation of learners' performance by teachers, there are three limitations. First, the test content is irrelevant to learners' classes. Second, the test duration requires a couple hours more than the class period. Finally, test fees are expensive, as learners must take a test at least twice a semester.

A solution to these limitations is to introduce computer-assisted language testing (Noijons, 1994;

Suvorov, 2013). Here, GEP is measured by calculating scores for a learner's read-aloud and dictation performance.

Previous research has been classified into two categories. One examined the correlation of GEP with dictation or read-aloud performances (Irvine, Atai, and Oller, Jr. 1974; Iino, Yabuta, and Thomas 2011; Kanzaki 2015; Leeming and Wong 2016). The other developed a measurement method for GEP based on dictation or read-aloud performances (Kotani and Yoshimi 2021a; Kotani and Yoshimi 2021b). Kotani and Yoshimi (2021a) and Kotani and Yoshimi (2021b) measured GEP using dictation performance and read-aloud performance, respectively.

This study expands a teacher's option by providing the third GEP measurement method of using both dictation and read-aloud performances. Previous research (Kotani and Yoshimi 2021a/b) did not examine the extent to which the measurement performance can be improved by measuring GEP based on the third method.

The goal of this study is to determine an effective GEP-measurement method by comparing different patterns of sub-proficiencies. Hence, the research question is as follows:

- Which is the highest GEP-measurement performance among a dictation-based method, a read-aloud method, and a dictation and read-aloud-based method?

These three methods are compared not only regarding the measurement accuracy but also the ease of measurement, specifically, the cost of developing and administering a method.

The contributions of the present study are to (1) investigate an effective GEP-measurement method as a classroom-based assessment alternative to GEP tests, (2) examine the validity of GEP-measurement methods, namely, a dictation-based method, a read-aloud method, and both dictation and read-aloud methods, and (3) verify the

84 robustness of a GEP-measurement method against  
85 the English-language-class size of training data.

## 86 **2 Collection of Dictation and Read- 87 Aloud Data**

88 Data instances comprised sentences transcribed by  
89 a learner, two types of dictation performance scores,  
90 speech sound pronounced by a learner, three types  
91 of read-aloud performance scores, five types of  
92 linguistic features extracted from reference  
93 sentences in a text material, and the learners'  
94 English test scores. The dictation and read-aloud  
95 data included 750 instances (50 learners' tasks for  
96 15 sentences).

### 97 **2.1 Participants**

98 The study participants were 50 English learners.  
99 This number was determined to mimic a large  
100 English class that included learners at different  
101 proficiency levels. The use of class-size training  
102 data reveals the possibility of teachers developing  
103 a GEP-measurement method using training data  
104 compiled in the class.

105 Participants were not randomly chosen. Those  
106 who satisfied the following conditions participated  
107 in the experiment: their first language was Japanese,  
108 and they were students at universities in the area  
109 where this study was conducted (28 men and 22  
110 women; mean age, 20.8 years; standard deviation  
111 [*SD*], 1.3). The participants were paid a fee for the  
112 experiment.

### 113 **2.2 Data Collection Procedures**

114 The dictation task proceeded as follows: First, the  
115 50 learners listened to sentences read aloud by a  
116 voice actor (woman, 35 years old) who was a native  
117 speaker of American English and transcribed them  
118 sentence-by-sentence. Subsequently, the learners  
119 subjectively judged their ease of dictation (see  
120 Section 3.1).

121 The read-aloud task was performed as follows:  
122 First, the learners listened to a reference speech  
123 sound by the native speaker. Subsequently, they  
124 read a sentence aloud and subjectively judged the  
125 ease of reading aloud (see Section 3.2). Their read-  
126 aloud durations were recorded to calculate the  
127 speech rates.

128 Learners received three points of instruction: 1)  
129 Each sentence could be listened to or read twice, if  
130 necessary; 2) Each task should be completed at a  
131 speed natural for the learner; 3) It was forbidden to

132 read fast or slowly or to return and revise a sentence  
133 after moving on to the next sentence.

### 134 **2.3 Text Material**

135 Two types of texts were selected from those  
136 distributed by the International Phonetic  
137 Association (1999) and Deterding (2006). As these  
138 texts include basic English sounds, an analysis of  
139 the learners' dictation and read-aloud performance  
140 of these texts would reveal what types of English  
141 sounds influenced their listening and pronunciation.

142 These texts featured two of Aesop's Fables: The  
143 North Wind and the Sun (Text I) and The Boy Who  
144 Cried Wolf (Text II). Texts I and II contained five  
145 and ten sentences, respectively. Text I failed to  
146 encompass certain sounds, such as the initial and  
147 medial /z/ and syllable initial /θ/. However, Text II  
148 included these missing sounds.

### 149 **2.4 General English Proficiency**

150 GEP was determined using participants' TOEIC  
151 listening and reading test scores obtained in the  
152 current or previous year. The reasons for this choice  
153 were as follows: The test scores were strongly  
154 correlated with the GEP test results, specifically,  
155 the Language Proficiency Interview developed at  
156 the Foreign Service Institute of the U.S.  
157 Department of State (Educational Testing Service  
158 1998), and that this test has no dictation or read-  
159 aloud sections.

## 160 **3 Features for Regression**

### 161 **3.1 Dictation Performance**

162 The criteria for evaluating dictation performance  
163 comprised two indexes: learners' subjective  
164 judgment of their ease with dictation (EASE-D)  
165 and dictation accuracy (ACC-D).

166 EASE-D was scored using a five-point Likert  
167 scale for the learners' subjective judgment (1 = easy,  
168 2 = somewhat easy, 3 = average, 4 = somewhat  
169 difficult, and 5 = difficult).

170 ACC-D was calculated by dividing the  
171 Levenshtein edit distance between a given  
172 reference and a transcribed sentence by the number  
173 of characters in a longer sentence than the other.  
174 The Levenshtein edit distance reflects the  
175 differences between the two sentences because of  
176 the substitution, deletion, or character insertion.

### 177 3.2 Read-Aloud Performance

178 The criteria for evaluating read-aloud performance  
179 comprised three indices: learners' subjective  
180 judgment of the ease of reading aloud (EASE-R),  
181 read-aloud accuracy (ACC-R), and speech rate in  
182 words per minute (RATE-R).

183 The EASE-R was determined by the learner's  
184 subjective judgment on a five-point Likert scale.

185 ACC-R was calculated by dividing the number  
186 of words correctly read aloud by the number of  
187 words in the corresponding sentence. A native  
188 English speaker evaluated learners' reading aloud  
189 word-by-word, but not phoneme-by-phoneme,  
190 using a binary decision (correct or incorrect  
191 pronunciation). The evaluator was trained to  
192 replicate interviews and meetings but was not  
193 familiar with the English spoken by learners.  
194 Before the assessment, the evaluator read the  
195 reference texts.

196 RATE-R was calculated by dividing the number  
197 of words by the duration of reading aloud.

### 198 3.3 Linguistic Features

199 In this study, linguistic features included sentence  
200 length, mean word length, number of multiple-  
201 syllable words, and word difficulty.

202 Sentence length (Chall and Dial 1948) was  
203 defined as the number of words in a sentence.

204 The mean word length (Chall and Dial 1948)  
205 was derived by dividing the number of syllables by  
206 the number of words in the sentence. The number  
207 of syllables in a word (Stenton 2013) was counted  
208 using the following steps: count the vowels in the  
209 word, subtract any silent vowels, and subtract one  
210 vowel from every diphthong.

211 The number of multiple-syllable words in a  
212 sentence (Fang 1966) was derived using the  
213 formula  $\sum_{i=1}^N (S_i - 1)$ , where  $N$  denotes the  
214 number of words in the sentence and  $S_i$  denotes the  
215 number of syllables in the  $i$ -th word. This  
216 subtraction derivation ignores the single-syllable  
217 words.

218 Word difficulty (Kiyokawa 1990) was defined as  
219 the rate of words not listed in Kiyokawa's basic  
220 vocabulary list relating to the total number of  
221 words in the sentence.

222 The speech rate was defined as the number of  
223 words read aloud by the native speaker in one  
224 minute.

### 225 4 Measurement of GEP with Dictation 226 and/or Read-Aloud Performances

227 The measurement methods were developed using  
228 support vector regression, with GEP as the  
229 dependent variable. The independent variables  
230 were dictation performance scores, read-aloud  
231 performance scores, and linguistic features.

232 Support vector regression was conducted using  
233 the function "svm()" defined in the "e1071"  
234 package of the software environment R (Meyer  
235 2021). The radial basis function was set as a type  
236 of kernel function, and the other parameters of  
237 "svm()" were set as default.

238 The measurement methods were evaluated using  
239 a leave-one-out cross-validation test. The  
240 training/test data consisted of 750 instances.

241 A correlation analysis was performed between  
242 the measured and observed GEPs. The significance  
243 threshold was adjusted for multiple testing based  
244 on the false discovery rate (FDR) (Benjamini and  
245 Hochberg 1995). A statistically significant  
246 correlation was examined to answer the research  
247 question.

248 To address the research question, three types of  
249 measurement methods were developed: dictation  
250 performance scores, read-aloud performance  
251 scores, and dictation and read-aloud performance  
252 scores. In addition to each type of test score, these  
253 methods use the linguistic features of  
254 dictation/read-aloud materials. The research  
255 question was answered by testing the equality  
256 between the statistically significant correlation  
257 coefficients in the chi-square tests.

### 258 5 Experimental Results and Discussion

259 The mean, minimum, and maximum GEP of the 50  
260 learners were 607.7, 295, and 900, respectively,  
261 and the  $SD$  was 184.45.

262 Table 1 shows the means and  $SD$ s of the  
263 dictation and read-aloud performance scores. Table  
264 2 shows the means and  $SD$ s of the linguistic  
265 difficulty of sentences in the text material.

266 Table 3 shows the correlation coefficients  
267 between the measured and observed GEPs in the  
268 cross-validation tests. Here, D&R refers to a  
269 measurement method using dictation and read-  
270 aloud, D represents a method using dictation, and  
271 R denotes the method using read-aloud. When the  
272 correlation coefficient was significantly different  
273 from zero, it was marked with an asterisk in all  
274 three types of measurement methods.

275 Table 4 shows the results of the chi-square tests  
 276 for the equality of correlations among the three  
 277 measurement methods. Bold chi-square values  
 278 indicate significant differences between the  
 279 correlation coefficients.

280 Table 3 shows the values of correlation  
 281 coefficients in a descending order: D&R > D > R.  
 282 Table 4 indicated the statistical significance of pairs  
 283 of correlation coefficients in the descending order:  
 284 D&R > D, D&R > R, and D > R. The measurement  
 285 method using D&R demonstrated the strongest  
 286 correlation. That is, the results suggest that D and  
 287 R are complementary for measuring GEP.

288 The significant difference in D > R suggests that  
 289 spelling is more associated with TOEIC than  
 290 pronunciation. In D, learners output grapheme  
 291 strings, while they output phoneme strings in R.  
 292 The former needs more sophisticated language  
 293 ability because the errors in spelling can be more  
 294 clearly identified, and learners use the visual,  
 295 auditory, and haptic (kinesthetic and tactile) senses  
 296 (Dobie 1986). Hence, the correlation result, or D >  
 297 R, can be considered evidence that D is more  
 298 associated with TOEIC than R.

299 Therefore, this study suggested that GEP should  
 300 be measured with a method using D&R because of  
 301 the strength of correlation, that is, D&R > D > R.  
 302 However, if teachers must decrease the time for test  
 303 administration and/or to reduce preparation tasks  
 304 for test materials, a measurement can also be  
 305 developed with only D instead of using D and R.

Performance score	<i>n</i>	Mean	<i>SD</i>
EASE-D	750	4.22	0.77
ACC-D	750	0.44	0.19
EASE-R	750	3.03	0.91
ACC-R	750	0.95	0.06
RATE-R	750	100.66	27.39

307 Table 1: Descriptive statistics of the dictation  
 308 and read-aloud performances

Linguistic features	<i>n</i>	Mean	<i>SD</i>
Sentence length	15	21.93	7.57
Mean word length	15	1.26	0.11
Number of multiple-syllable words	15	5.93	2.84
Word difficulty	15	0.26	0.11
Speech rate	15	178.44	17.41

310 Table 2: Descriptive statistics of the linguistic  
 311 difficulty of the sentences

Measurement methods	<i>r</i>	<i>t</i>	<i>df</i>	<i>p</i>
D&R	0.80*	36.13	748	< 0.05
D	0.75*	31.17	748	< 0.05
R	0.59*	19.78	748	< 0.05

313 Table 3: Correlation coefficients of the three  
 314 measurement methods

Measurement methods	<i>chi sq.</i>	<i>df</i>	<i>p</i>	<i>FDR</i>
D&R > D	<b>4.89</b>	1	0.03	0.05
D&R > R	<b>65.79</b>	1	< 0.02	0.02
D > R	<b>34.78</b>	1	< 0.03	0.03

316 Table 4: Chi-square tests for equality among the  
 317 three measurement methods

## 318 6 Conclusion

319 This study determined which GEP-measurement  
 320 method achieved the best performance. The three  
 321 GEP-measurement methods were developed using  
 322 dictation and/or read-aloud performance scores as  
 323 well as the linguistic features of the dictation/read-  
 324 aloud materials. These methods were compared  
 325 respecting the measurement accuracy and ease of  
 326 measurement.

327 The experimental results suggested that GEP  
 328 should be measured with the dictation and read-  
 329 aloud-based method, as the measured GEP had the  
 330 strongest correlation with the observed GEP.  
 331 However, if teachers must decrease testing time  
 332 and/or preparation tasks for test materials, the  
 333 dictation-based method can also be utilized.

334 Future research should examine what  
 335 combinations of dictation performance (EASE-D  
 336 and ACC-D) and read-aloud performances (EASE-  
 337 R, ACC-R, and RATE-R) can achieve the best  
 338 measurement performance. How the measurement  
 339 depends on learners' GEP should also be  
 340 investigated.

## 341 Acknowledgments

## 342 References

343 Yoav Benjamini and Yosef Hochberg.1995.  
 344 Controlling the False Discovery Rate: A Practical  
 345 and Powerful Approach to Multiple Testing. *Journal*  
 346 *of the Royal Statistical Society Series B*  
 347 *(Methodological)*, 57(1):289-300.  
 348

- 349 Jeanne S. Chall and Harold E. Dial. 1948. Predicting  
350 Listener Understanding and Interest in Newscasts.  
351 *Educational Research Bulletin*, 27(6):141-153+168.
- 352 David Coniam. 1991. Reading Aloud Speed as a Factor  
353 in Oral Fluency and General Language Proficiency?.  
354 *Hongkong Papers in Linguistics and Language*  
355 *Teaching*, 14:47-69.
- 356 David Deterding. 2006. The North Wind versus a Wolf:  
357 Short Texts for the Description and Measurement of  
358 English Pronunciation. *Journal of the International*  
359 *Phonetic Association*, 36(2):187-196.
- 360 Ann B. Dobie. 1986. Orthographical Theory and  
361 Practice, or how to Teach Spelling. *Journal of Basic*  
362 *Writing*, 5(2): 41-48.
- 363 Educational Testing Service. 1998. *TOEIC Technical*  
364 *Manual*. Educational Testing Service, Princeton: NJ.
- 365 Irving E. Fang. 1966. The Easy Listening Formula.  
366 *Journal of Broadcasting*, 11(1):63-68.
- 367 Atsushi Iino, Yukiko Yabuta, and Joel Thomas. 2011.  
368 Relationship between Criteria for Reading Aloud  
369 Evaluation and English Proficiency. *Journal of the*  
370 *Chubu English Language Education*, 40:159-166.
- 371 International Phonetic Association. 1999. *Handbook of*  
372 *the International Phonetic Association: A Guide to*  
373 *the Use of the International Phonetic Alphabet*.  
374 Cambridge University Press, Cambridge, UK.
- 375 Patricia Irvine, Parvin Atai, and John W. Oller, Jr. 1974.  
376 Cloze, Dictation, and the Test of English as a  
377 Foreign Language. *Language Learning*, 24(2):245-  
378 252.
- 379 Masaya Kanzaki. 2015. Minimal English Test: Item  
380 Analysis and Comparison with TOEIC Scores.  
381 *Shiken: JALT Testing & Evaluation SIG Newsletter*,  
382 19(2):12-23.
- 383 Semin Kazazoglu. 2013. Dictation as a Language  
384 Learning Tool. *Procedia-Social and Behavioral*  
385 *Sciences*, 70:1338-1346.
- 386 Hideo Kiyokawa. 1990. A Formula for Predicting  
387 Listenability: The Listenability of English  
388 Language Materials 2. *Wayo Women's University*  
389 *Language and Literature*, 24:57-74.
- 390 Katsunori Kotani and Takehiko Yoshimi. 2021a.  
391 Predicting English Proficiency with Read-Aloud  
392 Performance and Linguistic Difficulty of Sentences.  
393 *Proceedings of International Technology, Education*  
394 *and Development Conference*, 1180-1185.
- 395 Katsunori Kotani and Takehiko Yoshimi. 2021b.  
396 Prediction of General ESL Proficiency Considering  
397 Learners' Dictation Performance. *The 3rd ETLTC*  
398 *International Conference on Information and*  
399 *Communications Technology, EDP Sciences*,  
400 *France*.  
401 <https://doi.org/10.1051/shsconf/202110201003>.
- 402 Satsuki Kojima and Soichi Ota. 2012. Shadowing,  
403 Dictation and Reading Aloud: Which is Effective?.  
404 *Journal of The Japan Association of College*  
405 *English Teachers*, 4:29-40.
- 406 Yongeun Lee. 2014. Quantifying English Fluency in  
407 Korean Speakers' Read-Aloud and Picture-Cued  
408 Storytelling Speech. *Linguistic Research*,  
409 31(3):465-490.
- 410 Paul Leeming and Aeris Wong. 2016. Using Dictation  
411 to Measure Language Proficiency: A Rasch  
412 Analysis. *Language Testing and Assessment*, 5(2):1-  
413 25.
- 414 David Meyer, et al. 2021. e1071: *Misc Functions of the*  
415 *Department of Statistics (Formerly: E1071)*, TU  
416 Wien. [https://cran.r-](https://cran.r-project.org/web/packages/e1071/e1071.pdf)  
417 [project.org/web/packages/e1071/e1071.pdf](https://cran.r-project.org/web/packages/e1071/e1071.pdf).
- 418 Jose Noijons. 1994. Testing Computer Assisted  
419 Language Testing: Towards a Checklist for CALT.  
420 *CALICO Journal*, 12(1):37-58.
- 421 John W. Oller, Jr. 1983. *Evidence for a General*  
422 *Language Proficiency Factor: An Expectancy*  
423 *Grammar*. In John W. Oller, Jr. (Ed.), *Issues in*  
424 *Language Testing Research*. Newbury House,  
425 Rowley: MA.
- 426 Anthony Stenton. 2013. The Role of the Syllable in  
427 Foreign Language Learning: Improving Oral  
428 Production through Dual-Coded, Sound-  
429 Synchronised, Typographic Annotations. *Language*  
430 *Learning in Higher Education. Journal of the*  
431 *European Confederation of Language Centres in*  
432 *Higher Education*, 2(1):145-161.
- 433 Ruslan Suvorov and Volker Hegelheimer. 2013.  
434 *Computer-Assisted Language Testing*. In Antony J,  
435 Kunnan (Ed.), *The Companion to Language*  
436 *Assessment*, John Wiley & Sons, Hoboken: NJ.
- 437 Elvis Wagner. 2020. Duolingo English Test, Revised  
438 Version July 2019, *Language Assessment Quarterly*,  
439 17(3):300-315.
- 440 Aeris Wong and Paul Leeming. 2014. Using Dictation  
441 to Measure Language Proficiency. *Language*  
442 *Education in Asia*, 5(1):160-169.
- 443 Anoushe Yazdinejad and Mitra Zeraatpishe. 2019.  
444 Investigating the Validity of Partial Dictation as a  
445 Test of Overall Language Proficiency. *International*  
446 *Journal of Language Testing*, 9(2):44-55.