

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25

Using Latent Semantic Analysis to grade brief summaries:

A study exploring texts at different academic levels.

Ricardo Olmos & José A. León

Universidad Autónoma de Madrid, Spain

&

Guillermo Jorge-Botana & Inmaculada Escudero

Universidad Nacional de Educación a Distancia

Running Head: LSA to grade summaries

Please send correspondence to:

Ricardo Olmos Albacete

Departamento de psicología social y metodología

Facultad de Psicología

Universidad Autónoma de Madrid

Campus de Cantoblanco

28049-Madrid, Spain

e-mail: ricardo.olmos@uam.es

Telephone: +34 91 497 85 86

26 Abstract

27

28 In this study we propose an integrated method to automatically evaluate very brief
29 summaries (around 50 words) using the computational tool Latent Semantic Analysis
30 (LSA). The method proposed is based on a regression equation calculated with a corpus
31 of a hundred summaries (the training sample), and is validated on a different sample of
32 summaries (validation sample). The equation incorporates two parameters extracted
33 from LSA: (1) the semantic similarity of the summary, measured using the Summary–
34 expert summaries method (Landauer et al., 1998; León et al., 2006; Olmos et al., 2009)
35 and (2) the vector length (Redher et al., 1998). The study is based on a sample of 786
36 summaries by students at four academic levels. All of these students summarized either
37 an expository or a narrative text; their summaries were then evaluated by four graders
38 on a scale of 0-10. The results support three ideas. First, that incorporating both
39 parameters into the method is more successful than the traditional cosine measure. The
40 reliability of LSA for evaluating summaries rises above the 0.80 level for the expository
41 text. Second, that LSA shows practically the same level of sensitivity as the human
42 graders to the quality of the summaries at different academic levels. Third, that the
43 method overcomes a serious limitation of LSA: its difficulties evaluating very brief
44 texts (Redher et al., 1998; Wiemer-Hastings et al., 1999).

45

46 **Keywords:** Latent Semantic Analysis, assessment summaries, academic levels,
47 university students, Secondary students, Primary students, vector length, expository
48 text, narrative text.

49

50 Introduction

51

52 In recent decades educators have often used multiple choice tests to evaluate
53 comprehension of material. There are undoubtedly advantages to multiple choice
54 testing, including speed of assessment and the possibility of evaluating many different
55 aspects in a short time-frame, low cost, objective reliability measures (test-retest,
56 Cronbach's alpha, etc.) and relatively simple analysis of the psychometric properties of
57 items. This form of assessment has its limitations, however: comprehension would be
58 more superficial than for a student challenged with an open-ended question (Millis,
59 Magliano, Wiemer-Hastings, Todaro & McNamara, 2007; Shapiro & McNamara,
60 2000). The cognitive demands of a multiple choice test are those of recognition rather
61 than recall, so the student's learning strategy does not demand deep understanding of the
62 text (Far, Pritchard & Smitten, 1990). Causal relationships, drawing of inferences and
63 the form of expressing ideas would not be assessed by this type of tasks.

64

65 Exposing students to open-ended tasks, on the other hand, offers a means of
66 evaluating deeper understanding of the material. From the constructivist perspective,
67 building explanations or producing written material such as summaries gives rise to and
68 improves comprehension of texts more than multiple choice testing (Graesser, Lu,
69 Jackson, Mitchell, Ventura, Olney & Lowerse, 2004). Several studies have
70 demonstrated the importance of knowing how to summarize succinctly in understanding
71 and learning, and have shown it is a good measure of comprehension processes (Brown,
72 Bransford, Ferrara & Campione, 1983; Wade-Stein & E. Kintsch, 2004).

73 Summaries play a major role in research on comprehension of texts and its
74 assessment. For some authors a text has not been understood if the reader cannot
75 summarize it (Palinscar & Brown, 1984). However, it is often taken for granted that
76 students learn to summarize as they move to higher-level studies. In countries like
77 Spain, however, courses teaching summarizing skills are uncommon. Whilst it is true
78 that the actual concept of a summary is somewhat imprecise, there might be general
79 agreement that the process of producing a good summary implies understanding a text,
80 identifying the main ideas and transmitting them succinctly (E. Kintsch, Caccamise,
81 Franzke, Johnson & Dooley, 2007; León et al., 2006). For authors such as van Dijk and
82 Kintsch (1983) summarizing involves the capacity to generalize, synthesize and write
83 coherently. It thus goes far beyond reading, since it implies profound comprehension
84 of what is read. In their model of comprehension, summarizing is essential to
85 understanding, since it involves extracting the main content of what is read, and at the
86 same time eliminating superficial details. Kintsch himself (2002) used the LSA model
87 in an attempt to find the phrases that best summarize the content of a text. To achieve
88 this aim he sought structures (titles, subheadings or paragraphs) that best represent the
89 information contained in the text (the abstract information from the macrostructure).

90

91 Summarizing, then, involves establishing relationships between important
92 concepts, and presenting them in a coherent, organized manner. The information must
93 be restructured, further abstracting it from the content of the text. The summary allows
94 us easier access to factual and conceptual knowledge in memory. Summarized texts
95 allow us to build on information in the classroom much more and better than by simply
96 rereading a text. It allows students to formulate more pertinent questions, and the
97 teacher to evaluate the extent to which the material was understood. In this line we

98 subscribe to a popular eighties school of thought (see review by Bransford, Brown &
99 Cocking, 2000) that learning to summarize is a central aspect of the comprehension
100 process, so that reliably evaluating a summary is key to knowing whether a student has
101 a deep understanding of a text.

102

103 The demands placed on teachers make it very difficult to find time to evaluate
104 essays and summaries and give feedback to each of the students individually; so many
105 educators favor multiple choice tests (Wade-Stein & E. Kintsch, 2004). There are tools
106 currently available which can evaluate texts reliably (see review by Dikli, 2006). These
107 tools, although they are not apt for awarding a final grade, might be valuable for
108 monitoring a student's progress, or to provide the student with longitudinal information
109 regarding their level of ability. The task of introducing these tools in the classroom is
110 extremely complex, and depends on many factors (e.g. schools infrastructure, budget,
111 subject matter and access to technology) and also on cooperation between many
112 individuals (e.g. educational psychologists and teachers). However, one of these
113 contributions - Automated Essay Scoring - is already far advanced enough to take a step
114 forward.

115

116 This paper was driven by two factors: on the one hand the need to introduce
117 summaries into the classroom to synthesize the key information from syllabus units
118 covered, and on the other the possibility of extracting reliable automated evaluations of
119 these summaries using a computational tool known as *Latent Semantic Analysis* (LSA).
120 These two aspects have normally been impossible to reconcile, since LSA only provides
121 reliable evaluations using more extensive texts - mostly over 250 words (Rehder,
122 Schreiner, Wolfe, Laham, Landauer & Kintsch, 1998). In fact, some authors indicate

123 that LSA has special difficulty in analyzing texts between two and sixty words
124 (Wiemer-Hastings, Wiemer-Hastings & Graesser, 1999) - the range of the summaries
125 used in our research. In this study, then, we propose an LSA-based method that provides
126 reliable evaluations of very brief summaries. The method we propose combines
127 information on semantic similarity commonly provided by LSA, together with
128 information on the extent of knowledge LSA has of the terms represented in a semantic
129 space. This means that the method combines the cosine measure with vector length
130 information. We will begin by looking at the importance of summaries in the classroom
131 together with the need for evaluation of open-ended responses as a complement to
132 multiple choice exams. In fact, this paper attempts to analyze a combination of ratings
133 from LSA cosine and vector length measures on the one hand, and human graders'
134 evaluations of content and coherence on the other. These types of measures will be
135 applied in the assessment of brief summaries by students at different academic levels
136 and using two types of texts (a narrative text and an expository text). In terms of
137 academic levels, we took 238 students from 6th grade, 192 students from 8th grade, 198
138 students from 10th grade, and lastly 158 university students. In addition, these open-
139 ended responses will be compared with standard multiple choice scores obtained by the
140 same students on the same texts.

141

142 **LSA in educational tasks**

143

144 Latent Semantic Analysis (LSA) is a computational technique that contains a
145 mathematical representation of language. During the last twenty years its capacity to
146 simulate aspects of human semantics has been widely demonstrated (Landauer &
147 Dumais, 1997). LSA is based on three fundamental ideas (Steyvers & Griffiths, 2007):

148 (1) To begin to simulate human semantics of language we first obtain an occurrence
149 matrix of terms by document, (2) the dimensionality of this matrix is reduced using
150 singular value decomposition (SVD), a mathematical technique that effectively makes
151 the tool a latent semantic space, and (3) any word or text is represented by a vector in
152 this new latent semantic space. Since each word is a vector, a text is the sum of the
153 words that comprise it, i.e. another vector.

154

155 To evaluate the semantic similarity between two texts, we need only extract the
156 cosine formed between the two vectors. When the cosine is close to zero the semantic
157 similarity is null. When the cosine is close to one the semantic similarity is very high.
158 Formally, the cosine is defined

159 as:

$$160 \text{Cos} = \frac{\sqrt{\sum_{i=1}^k x_i y_i}}{\sqrt{\sum_{i=1}^k x_i^2} \sqrt{\sum_{i=1}^k y_i^2}}$$

161

162

163 where the numerator contains the scalar product between the k coordinates of the first
164 and second word (x and y respectively), and the denominator contains the product of the
165 word x and word y vector lengths.

166

167 There are other possible LSA-derived measures that could be used, such as the
168 dot product or Euclidean distance between the two vectors, or the length of an
169 individual vector. A vector can be thought of as a position within an n-dimensional
170 space. The value of a vector is represented as a series of coefficients, each coefficient
171 representing a value (or distance) along a particular dimension in the n-dimensional

172 space (Reder, Schreiner, Wolfe, Laham, Landauer & Kintsch, 1998). The vector length
173 formula can be defined as:

174
$$VectorLength = \sqrt{\sum_{i=1}^k x_i^2}$$

175 that is, the square root of the sum of squares of the k coordinates that represent the word
176 x .

177

178 The vector length also informs us of the knowledge LSA has of this text. A
179 graphical representation of how LSA works can be seen in Figure 1. There are several
180 points to note regarding the diagram. First, that if we semantically compare the text
181 “lush forest” with “tropical jungle” the cosine is close to one, since both vectors have
182 very similar directions due to their similar meanings. However, these two vectors have
183 a cosine close to zero with the text “modern building”, since the vectors produce a
184 practically orthogonal angle. Second, since the vector length for “modern building” is
185 greater than that of the other two vectors, we can assume that LSA has more knowledge
186 about this text than about the other two. In short, after LSA creates a semantic space in
187 which it can represent the texts vectorially, we can calculate the semantic relationship
188 between two texts using the cosine formed by the two vectors, and measure the
189 knowledge that LSA has of a text by the length of the vector that represents it. Vector
190 length contains the quantity of summary elaboration, and the cosine contains the
191 quantity of semantic similarity.

192

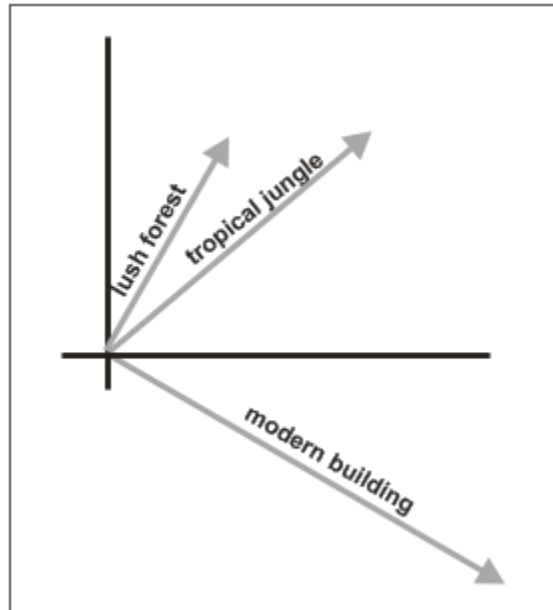


Figure 1. Graphical example of LSA representing three texts (vectors)

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

LSA assessment methods normally only use the former to attempt evaluations, ignoring vector length (Foltz et al., 1999; Landauer & Dumais, 1997; Landauer et al., 1998; Wade-Stein & Kintsch, 2004). Nonetheless, some studies have analyzed vector length in the assessment. Redher et al. (1998) used a multiple regression to demonstrate that the combination of the semantic similarity (cosine) and vector length explain the correlation of LSA with human criteria, and that in fact vector length is the factor responsible for the greater part of the variance. In this study the authors examine these and other measures derived from the cosine or the vector length, such as the Euclidean distance or the scalar product, to test their efficiency in automatic evaluation of essays. They found that in the regression the strongest predictors were vector length and cosine, although alone (i.e. not in combination with the others) the most powerful measure was the scalar product. Other research has also covered alternative measures such as Euclidean distance in the automatic grading of summaries (Jorge-Botana, León, Olmos & Escudero, 2010; Olmos, León, Jorge-Botana & Escudero, 2009).

211 LSA has been used abundantly in the field of education. For example, it has
212 been used to evaluate online comprehension using verbal protocols while asking
213 students to read (usually self-explanations) (McNamara, Levinstein & Boonthum, 2004;
214 Millis et al., 2007). LSA is a satisfactory tool for capturing the meaning of these verbal
215 protocols. In these studies LSA reveals different kinds of learning strategies when the
216 students read, or rather when they discuss what they are reading. It detects, for
217 example, that some tend to paraphrase what they just read, whilst others usually relate it
218 to other phrases read previously, or to their prior knowledge. Since these different
219 strategies generally imply different levels of comprehension, LSA can be used quite
220 successfully either to predict the level of comprehension, or to evaluate the
221 predominance of reading strategies or give appropriate feedback to students, coaching
222 them towards better use of strategies (McNamara, Boonthum, Levinstein & Millis,
223 2007; Millis et al., 2007).

224

225 Another educational application of LSA is with computer tutors (Graesser,
226 Chipman, Haynes & Olney, 2005; Wade-Stein, E. Kintsch, 2004). Graesser et al.
227 (2005) created a tool called AutoTutor, using a computer to hold conversations with
228 students in natural language. Students are presented with common problems from the
229 curriculum script, and using an animated agent AutoTutor gives them feedback until
230 they manage to give a satisfactory response to each problem. Another computer tutor is
231 Summary Street by Wade-Stein and E. Kintsch (2004). This tool coaches during the
232 process of writing a summary. The basic idea underlying this tool is how to teach
233 students to summarize. These computer tutors teach students to resolve problems (e.g.
234 summarizing) without individual attention from a teacher - something totally impossible
235 if we consider the lack of educational resources available today. In simple terms, what

236 LSA does is compare what students write with texts incorporated into the tools (ideal
237 summaries, main topics, key words, etc.) using the cosine measure. Thresholds are set
238 such that if the cosines rise above them we assume the student response covers the
239 pertinent aspects. If the cosine does not reach the threshold, the computer tutor gives
240 clues to help the student incorporate the missing information. The central feature is the
241 dynamic interaction between student and machine. If a stimulating environment is
242 combined with good task design, the results show a notable improvement in student
243 responses (Graesser, Penumatsa, Ventura, Cai & Hu, 2007; E. Kintsch et al., 2007).

244

245 A third example of an LSA application in the field of education is based on
246 automatic assessors (Foltz, Laham & Landauer, 1999; Landauer, Foltz & Laham, 1998).
247 For example, Foltz et al. (1999) created the Intelligent Essay Assessor (IEA). These
248 methods are based on using the cosine to compare student essays with a source text.
249 One very common method is when the source text consists of an expert summary
250 (normally by a grader or teacher), thus creating what they call a “golden summary”
251 (Landauer et al., 1998; León et al., 2006). These tools automatically provide an essay
252 score, sometimes offering impressive results (e.g. above .80), proving as reliable as the
253 expert judges (trained graders or teachers) themselves. Another similar application by
254 the French authors Dessus & Lemaire is the APEX system (2002). This system, based
255 on LSA, provides texts for students to summarize, and then evaluates the summary.
256 Other more innovative procedures include the so-called EssayAid (Kakkonen &
257 Sutinen, 2011) for semi-automatic essay evaluation, or evaluation procedures that
258 combine LSA with n-grams (Monjurul Islam & Latiful Hoque, 2012).

259

260 More recent studies have introduced the possibility of evaluating bridging
261 inferences based on text cohesion, measuring this cohesion as the similarity between
262 adjacent phrases using the cosine with LSA. The tool is known as Interactive Strategy
263 Trainer for Active Reading and Thinking (iSTART) created by Bellissens, Jeuniaux,
264 Duran McNamara (2010). This tool was conceived so that with the help of LSA
265 students can improve their reading strategies. In fact they promote what they call Self-
266 Explanation Reading Training (SERT).

267

268 Lastly, we do not want to give the impression that LSA is only applied to the
269 evaluation of essays. Its versatility makes it a fertile resource in more complex tasks -
270 for example LSA can be used to adapt texts to the participant's level of ability, so that
271 they make the best use of the texts they learn with (Wolfe, Schreiner, Rehder, Laham,
272 Foltz, Kintsch & Landauer, 1998). LSA is not only applicable to the field of written
273 tasks, as we see for example in its successful application to reasoning task analysis for
274 complex problems (Quesada, Kintsch & Gomez, 2001).

275

276 This study focuses on the latter type of applications, more specifically on
277 presenting a method that allows us to evaluate summaries reliably. LSA-based
278 evaluations have normally been applied to relatively long essays (over 200 words), but
279 few have tackled the use of LSA to evaluate brief summaries of only fifty words. When
280 texts contain fewer than 1000 words, it is only natural that the summaries are shorter
281 than those used in most LSA applications (Wade-Stein & Kintsch, 2004; Landauer et
282 al., 1998, Rehder et al., 1998). For a detailed up-to-date review of research on automatic
283 essay evaluation with LSA see Haley (2009). In our case, the mean number of words
284 across the 396 summaries of a narrative text (which contained 402 words) was 39.15 – a

285 ratio of 1 word in the summary for every 10 in the text. The mean number of words in
286 the 388 summaries of an expository text (500 words) which was also applied was 28.22,
287 giving a ratio of approximately 1 word in the summary for every 18 in the text.

288

289 Objectives

290

291 The aim of this study was to use a LSA-based computational method to reliably
292 evaluate especially brief summaries (maximum fifty words). The method incorporates
293 the essential information from the latent semantic space: (1) A measure of semantic
294 similarity using the cosine and (2) a measure of the vector length or extent of
295 knowledge about the text. With this general aim we sought four goals. First, to try to
296 obtain reliable evaluations (at least > 0.70) combining both kinds of essential
297 information derived from the latent semantic space (cosine and vector length, analyzing
298 their effects separately as well as in combination) compared to human graders and for
299 each type of text. Second, to take Pearson correlations between LSA measures (cosine
300 and vector length) vs. human graders' measures (content and coherence), related to
301 summary length (number of words). Third, to analyze whether LSA is sensitive to the
302 quality of brief summaries by students at different academic levels, as trained human
303 graders are. And fourth, to analyze whether positive correlations exist between
304 summaries and multiple choice scores.

305

306 Method

307

308 **Participants.** 786 students from four grade levels took part in this study. Student ages
309 ranged from 10 to 23 years. The youngest students were from 6th grade (a total of 238

310 students), followed by 192 students from 8th grade, 198 students from 10th grade, and
311 lastly 158 second year psychology undergraduates.

312

313 **Material.** Two texts were applied in this study. A narrative text (*The Carob Tree*
314 *Legend*) was 402 words long and its comprehension requires only general knowledge.
315 The expository text (*The Strangler Trees*) was 500 words long and understanding it also
316 requires general knowledge.

317

318 **The Spanish LSA corpus.** The generalist corpus used in this study belongs to the
319 University of Colorado. The corpus contains material from on-line encyclopedias,
320 newspapers, textbooks and several Internet sources. In total the corpus has 2,059,234
321 documents (i.e. paragraphs) and 1,661,954 different terms. A semantic space with 337
322 dimensions was used. The method used was document to document.

323

324 **Procedure.** Each student read the text at their own pace in a classroom. Before reading
325 the text they were told that it was important to understand the text in order to answer a
326 set of questions. After reading, students were allowed 15 minutes to write a summary
327 (maximum fifty words). Finally, they were asked to answer a set of multiple choice
328 comprehension questions. The 786 summaries used in this research, 396 of them were
329 summaries of a narrative text, and the remaining 390 were summaries of an expository
330 text.

331

332 **The four judges' evaluations.** To obtain the expert judges' evaluations with which to
333 later compare those of LSA, four doctorate students received four sessions of instruction
334 in evaluating summaries on a scale from 0 to 10. The evaluation criteria followed were

335 taken from research by León & the Reading Literacy Research Group. Whilst grading
 336 the summaries, the judges took two aspects into consideration. First the content of the
 337 summaries was evaluated on a scale of 0 (no content) to 4 (all key content). Both the
 338 narrative and the expository texts contained four main ideas that had to be considered in
 339 the summaries (see León et al., 2006). Each main idea counted as one point. Secondly
 340 coherence was evaluated on a scale of 0 (incoherent) to 6 (highly coherent). To assess
 341 coherence the organization, causal relationships, use of connectives, extent of
 342 conceptualization of the summary and lack of redundancy were analyzed. The judges
 343 carried out the evaluations independently and without knowing the student's academic
 344 level.

345

346 **The proposed method**

347

348 To implement our method we used a database comprising 107 summaries of
 349 narrative text and 93 expository summaries, distributed across four grade levels. The
 350 sample used to adjust the method is called the training sample, and allows us to
 351 calculate the way we obtain the scores with LSA, although it is not used to evaluate the
 352 reliability of the method. Table 1 shows scores from the 200 summaries in the training
 353 sample used to implement the method.

354

		Educational level				Total
		6th grade	8th grade	10th grade	University	
Type of text	Narrative	27	29	32	19	107
	Expository	29	26	24	14	93
	Total	56	55	56	33	200

355

356 Table 1. Training sample of summaries in the narrative and expository text

357

358 Each of these summaries was graded independently by each of the four judges
359 on a scale of 0 to 10, awarding up to four points for content and six points for coherence
360 of the summary. Blind scoring was used, in other words the graders were unaware of
361 the student's academic level. An average score was obtained using the four graders'
362 scores. After this a regression line was calculated, where the dependent variable was
363 the graders' average score and the independent variables were vector length and
364 semantic similarity.

365

366 The regression equation for the narrative text was:

367

$$368 \quad \text{NarrativeScore} = \beta_0 + \beta_1 * \text{Vectorlength} + \beta_2 * \text{Similarity}$$

369

370 And the regression equation for the expository text was:

371

$$372 \quad \text{ExpositoryScore} = \beta_0 + \beta_1 * \text{Vectorlength} + \beta_2 * \text{Similarity}$$

373

374 where β_0 is the constant, β_1 is the coefficient for vector length and β_2 is the
375 coefficient for semantic similarity.

376

377 Once the regression calculations are done, we took a new sample of summaries
378 (the validation sample). This time there were 289 summaries of the narrative text and
379 297 summaries of the expositive text. Table 2 shows the distribution of summaries for
380 the validation sample used to check the method.

381

382

383

384

		Educational level				Total
		6th grade	8th grade	10th grade	University	
Type of text	Narrative	92	69	68	60	289
	Expository	90	68	74	65	297
	Total	182	137	142	125	586

385

386

Table 2. Validation sample of summaries of the narrative and expository text

387

388

389

390

391

392

393

Summaries were again graded independently by the four graders, using the same scale from 0 to 10. An average was taken of the four judges' scores to obtain a single grade. These grades were used to assess the reliability of the scores awarded by LSA using the regression equations. The independent validation sample was used to avoid overfitting, to make it easier to generalize to new summaries.

394

395

396

397

398

399

400

401

402

403

404

405

406

The two LSA measures, vector length and semantic similarity (cosine), were obtained as follows. Given that in LSA each document is represented by a vector, the vector length is simply calculated as the length of each summary vector. In the equation the vector length component represents how detailed the summary is - the greater the vector length the more detail, and the more familiar or relevant words appear in the semantic space. The measure of similarity is somewhat more difficult to obtain. For this reason we use a well-known method, habitually used in Automated Essay Scoring with LSA: the Summary-expert summaries method (Dikli, 2006; Foltz et al., 1999; Kintsch et al., 2007; León et al., 2006; Landauer & Dumais, 1997; Landauer et al., 1998; Olmos, et al., 2009). The Summary-expert summaries method consists of assessing student summaries by comparing them with an expert summary (Landauer, Laham & Foltz, 1998). It is conceived as a method that can capture how semantically similar a student summary is to other summaries written by experts, usually known as

407 'golden summaries'. For the present study, six summaries written by experts were
408 chosen as the standard. With this method, LSA scores a student summary as follows:
409 LSA computes cosines between the student summary and each of the six expert
410 summaries. The final score for the student summary is the average of these six cosines.

411

412 To obtain a measure of similarity we had to compare the summary with a source
413 text. This method uses ideal summaries, in other words summaries written by experts
414 containing the essential information from the text (all relevant details and very strong
415 coherence). To this end, six teachers with expertise in comprehension were asked to
416 write a summary of no more than fifty words, both of the expository text and of the
417 narrative text. To obtain a measure of semantic similarity the cosine between the
418 student summary and each expert summary was calculated. Since there are six cosines,
419 one for each expert summary, the average cosine is taken as the final measure of
420 semantic similarity. Both vector length and semantic similarity values were obtained
421 automatically for all 786 summaries, as described. What we refer to in the regression
422 equation as *similarity* will be called the *Expert Method* henceforth, since the similarity
423 is calculated using the expert's method.

424

425 **Data analysis.** The data was analyzed in four stages. First the regression equation with
426 the ordinary least square estimation method that best fits the judges' scores was obtained
427 using the training sample, as described in the previous section. Secondly, the reliability
428 (Pearson correlations) between LSA and judges was calculated using the validation
429 sample. Third, an ANOVA was used to study the sensitivity of judges and LSA to
430 differences in the quality of summaries from different grade levels. Lastly, the
431 reliability (Pearson correlations) between LSA measures (cosine and vector length) vs.

432 human graders' measures (content and coherence), related to summary length (number
433 of words), and the relationship between summary and multiple choice scores.

434

435 **Results**

436 **Regression equations**

437

438 To predict the average judges' scores a training sample was used, comprising
439 107 summaries of the narrative text and 93 summaries of the expository text chosen
440 completely at random. These samples were used to find two regression equations, later
441 applied to predict the judges' grades for the validation sample. The regression line to
442 predict the grades for the summaries of the narrative text was:

443

$$444 \quad \text{NarrativeScore} = -1.62 + 5.76 * \text{Vectorlength} + 11.26 * \text{ExpertMethod}$$

445

446 Where the coefficient of Vector Length was statistically significant ($t(1)=4.98, p < .05$),

447 as too was the Expert Method coefficient of semantic similarity ($t(1)=6.34, p < .05$).

448 The tolerance between the two variables was .80.

449

450 The regression line obtained for the expository text was:

451

$$452 \quad \text{ExpositoryScore} = -4.19 + 10.18 * \text{Vectorlength} + 15.61 * \text{ExpertMethod}$$

453

454 Once again, the coefficient of Vector Length was statistically significant

455 ($t(1)=8.22, p < .05$) as was the coefficient associated with Expert Method ($t(1)=8.63, p$

456 $< .05$). The tolerance between the two variables was .92.

457

458 To obtain a grade for a new summary (narrative or expository), then, we had
459 only to calculate the vector length of the summary and a measure of semantic similarity
460 using the expert method.

461

462 We then substitute into the corresponding equation to obtain the grade for the
463 summary. Both equations obtained positive coefficients, informing us that the greater
464 the summary detail (as measured by vector length) the greater the summary grade, and
465 the greater the semantic similarity between the student summary and the experts'
466 summaries the greater the summary grade.

467

468 Lastly, we examined the impact of vector length on reliability, compared to the
469 traditional method using only the measure of similarity (Expert method) (Foltz et al.,
470 1999; Landauer et al., 1998; Wade-Stein & E. Kintsch, 2004). It was found that on
471 including the vector length information, the change in the proportion of variance
472 explained was statistically significant in the case of the narrative summaries ($F(1,104) =$
473 $24.81, p < .05$) as well as in the case of the expository summaries ($F(1,90) = 67.51, p <$
474 $.05$). In other words, with these brief summaries there was a substantial drop in
475 reliability if we do not incorporate vector length (see following section).

476

477 **A) Reliability between the four human graders themselves**

478

479 Before analyzing the reliability of LSA, we evaluated reliability between the
480 four judges themselves. This reliability varied between .78 and .86 for the narrative

481 text, and between .83 and .88 for the expository text - all highly reliable and statistically
482 significant.

483

484 **B) Reliability between LSA and human graders' scores.**

485

486 The LSA-grader reliabilities were calculated using the validation sample. This
487 sample was set aside to avoid overfitting of reliabilities, and thus allow generalization
488 of results to other summaries. The validation sample consisted of 289 summaries of the
489 narrative text and 297 of the expository text. The LSA grades for this new sample of
490 summaries were obtained by simply using the equations presented in the previous
491 section.

492

493 Table 3 shows the reliability of LSA with individual judges and with the average
494 judges' scores in both texts (calculated with Pearson's correlation). For the narrative
495 text the reliability of LSA ranged from .60 to .67 for the individual judges, and reached
496 .68 with the average judges' scores. As for the expository text, the LSA-grader
497 reliability ranged from .76 to .78, reaching .82 with the average judges' scores. For the
498 narrative text the reliability scores were fairly high; in the expository text they were
499 high.

500

Text	Human Grader				
	Grader 1	Grader 2	Grader 3	Grader 4	Average grader
Narrative text	.61**	.67**	.60**	.63**	.68**
Expository text	.76**	.78**	.77**	.78**	.82**

501 ** p<.01.

502

503

504

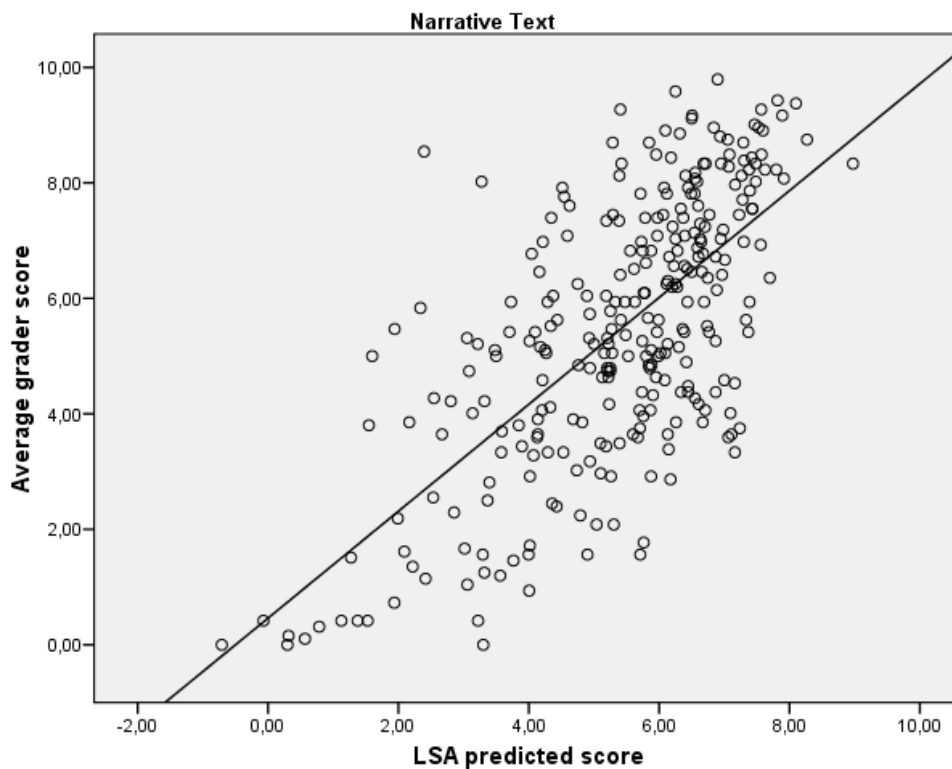
Table 3. LSA-grader reliability for each text and human grader

505 As for the LSA-grader reliability comparison we found significant differences
506 between LSA-human graders' reliability for the expository text ($r = .82, p < .01$) and
507 LSA-human graders' reliability for the narrative text ($r = .68, p < .01$). Reliability for the
508 expository text was thus significantly higher ($p < .01$).

509

510 Figures 2 and 3 show *scatter plots*, as a graphical representation of LSA
511 evaluations compared to the average judges' grades. Each point represents a pair of
512 grades: the judges' average and the LSA grade.

513

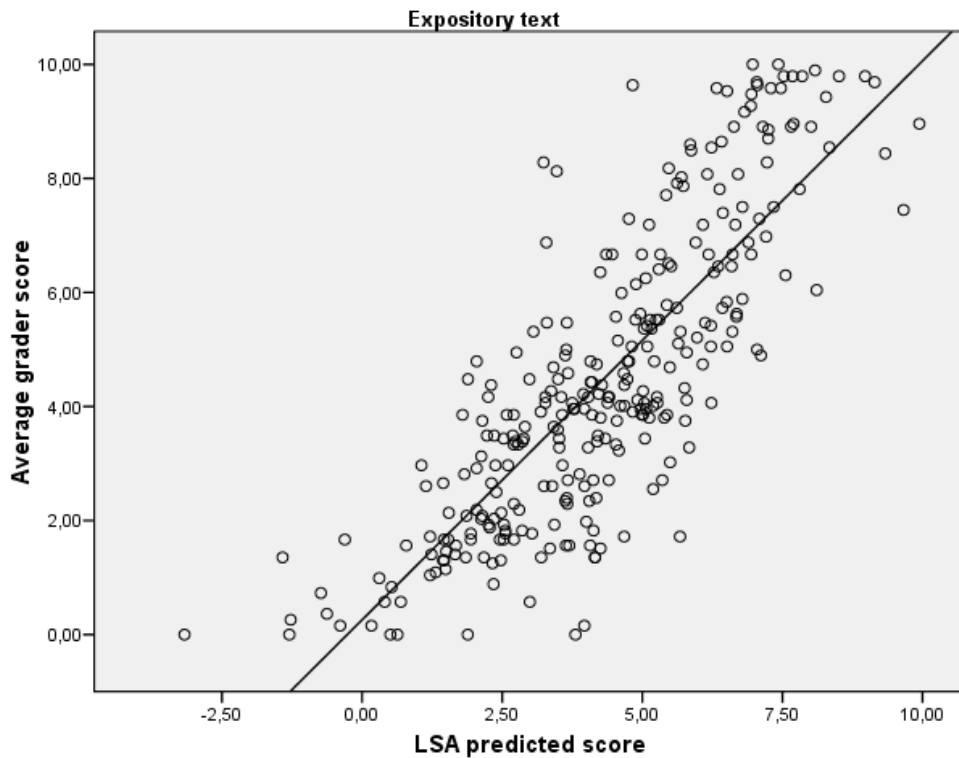


514

515 Figure 2. Average grader scores plotted against LSA-predicted narrative scores.

516

517



518

519

Figure 3. Average grader scores plotted against LSA-predicted expository scores.

520

521

522

523

524

525

526

527

528

529

530

531

532

533

The graphs show the average judges' scores against the scores predicted by LSA using the regression equation above. As we saw from the table of reliability scores, the cluster of points fits the line better for summaries of the expository text than for those of the narrative text. It is clear that when the judges award the summary a low score LSA does the same. LSA also follows suit when the judges give a high grade. Both graphs show linear relationships; there were no summaries whose LSA grade differed markedly from the average judges' grade.

534 **C) Pearson correlations between LSA measures (cosine and vector length) vs.**
 535 **human graders measures (content and coherence)**

536

537 In Table 4 we can see the data analyzed by type of measure, in order to
 538 evaluate the results of the previous section in more detail.

539

540

	LSA	Judge's assessments		
		Content	Coherence	Global (content + coherence)
Narrative summaries	Cosine (expert method)	.61**	.51**	.60**
	Vector length	.56**	.46**	.55**
	Cosine + Vector length	.68**	.56**	.68**
Expository summaries	Cosine (expert method)	.68**	.63**	.67**
	Vector length	.64**	.65**	.66**
	Cosine + Vector length	.82**	.81**	.82**

541

**p<.01.

542

543 Table 4. Table Correlations between LSA and Judge's assessments for each text and type of measure
 544 (cosine, vector length, content, coherence).

545

546

547 Whilst the results obtained are all significant, they highlight differences that
 548 should be borne in mind when considering the type of text. For the expository text there
 549 is a high degree of homogeneity in the correlations obtained between LSA and experts,
 550 which range from .63 to .68. There is a positive effect with the combination of cosine
 551 and vector length, offering a very high correlation of .82. In contrast, the correlations
 552 obtained for the narrative text show greater divergence between LSA cosines and
 553 human evaluations of content and coherence. Whilst the cosine correlates highly with
 554 content (.61), the score for coherence is lower (.51). Vector length correlates less with
 555 both coherence (.46) and content (.56), and the combined effect of cosine and vector
 556 length does not strengthen the result of correlations with human graders. It is clear,
 557 then, that expository text summaries are better evaluated using LSA in all cases.

558

559 There were significant differences in the reliability (Pearson) between human
560 graders (content and coherence) and LSA, but only for the narrative summaries. The
561 combined cosine and vector length LSA method offers significantly higher correlation
562 with human assessment of content ($r=.680$, $p<.01$) than of coherence ($r=.563$, $p<.01$).
563 However, in the expository summaries there are no significant differences in the
564 correlations between LSA and grading of content ($r = .823$, $p<.01$) and LSA and grading
565 of coherence ($r = .806$, $p<.01$).

566

567

568 **D) Correlations between LSA measures (cosine and vector length) vs Human**
569 **grader measures (Content and coherence) related to the length of summary.**

570

571 In this section we have included a new analysis of the measures made by human
572 graders (content, coherence and the two combined) and LSA (cosine and vector length,
573 and the two combined) regarding the length of the summary made by each student, to
574 assess whether scores from LSA and Human graders correlated with the length the
575 summary.

576

577 The results show a significant correlation in every case, which may indicate that
578 summary length affects the score awarded by both expert and LSA. This result is not
579 unexpected as both content and coherence of a summary are favored by longer
580 summaries with more information to evaluate. The same occurs here with LSA,
581 although the correlations are even higher when analyzed using vector length, which is
582 more sensitive to summary length as it measures length as well as familiarity (vector

583 length has been interpreted as the familiarity LSA has with technical words, Rehder et
 584 al., 1998).
 585

Type of text		Number of words in the summary
Narrative	Cosine	.52**
	Vector length	.88**
	Cosine + vector length	.89**
	Content	.64**
	Coherence	.65**
	Content + coherence	.69**
Expository	Cosine	.49**
	Vector length	.83**
	Cosine + vector length	.87**
	Content	.76**
	Coherence	.78**
	Content + coherence	.79**

586 **p<.01.

587

588

589 Table 5. LSA-Human grader reliability for each text and type of measure (cosine, vector length, content,
 590 coherence) related to the length of summary

591

592 In Table 5 we can see that summary length (number of words) affects the
 593 assessment of human graders whether evaluating content or coherence. The mean
 594 number of words in the narrative text summary was 39.15 (S.D.=16.28 (range=88). The
 595 mean number of words in the expository text summary was 22.22 (S.D=14.66,
 596 range=89). Thus we could say that the longer the summary the higher the grade awarded
 597 for both content and coherence. This phenomenon is more pronounced for the
 598 expository than the narrative text. It can also be seen in LSA assessment, but only when
 599 applying vector length, which is more sensitive to summary length. This is not an
 600 unexpected result given the nature of this measure, but is not the case when cosine is
 601 used. The fact that the cosine is not susceptible to this effect may indicate that the two

602 measures are complementary, and that together they can allow us to better evaluate a
603 summary.

604

605 **E) Sensitivity to differences between different academic levels**

606

607 Lastly, an ANOVA was carried out on each text to study the capacity of both
608 judges and LSA to detect differences in the quality of summaries from different grade
609 levels. In this case the analysis was carried out using only the validation sample. The
610 results obtained showed that in the narrative text judges detected more differences in the
611 quality of summaries than LSA. For the expository text, however, the results for LSA
612 and human graders were similar.

613

614 Figure 4 shows the average LSA and judges' grades for summaries of the
615 narrative text at each of the four grade levels. Both LSA and the judges detect
616 differences in the quality of summaries ($F(3,285) = 10.96, p < .05$ and $F(3,285) = 36.60,$
617 $p < .05$, for LSA and average judges' scores respectively). However, a *post hoc* test
618 revealed three groups of averages for the judges and only two groups for LSA. The
619 judges detected that the group that summarized with the highest quality was the
620 undergraduate group, followed by 10th grade and then 6th grade and 8th grade students
621 (no difference was found between these last two). LSA detected two groups of
622 averages: the group that summarized with the highest quality was the university
623 students and the other grade levels formed another group of lower-quality summaries
624 (no statistically significant differences were detected between them).

625

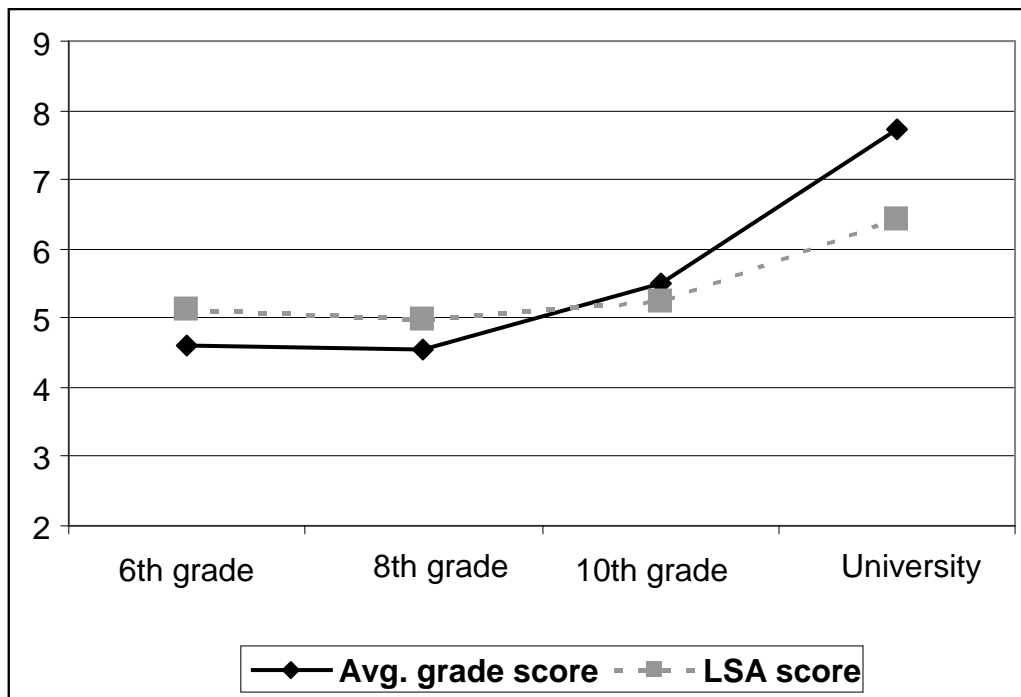


Figure 4. LSA and graders' mean score at each grade level for the narrative text

626

627

628

629

630

631

632

633

634

635

636

637

638

At the same time, Figure 5 shows the average LSA and judges' grades for summaries of the expository text at each of the four grade levels. As for the narrative text, both LSA and the judges detected differences in the quality of summaries ($F(3,293) = 47.88, p < .05$ and $F(3,293) = 89.97, p < .05$, for LSA and average judges' scores respectively). This time the *post hoc* test showed that both LSA and the judges detected two groups of averages: the undergraduate group summarized with higher quality than the other grade levels (no statistically significant difference was found between the averages). In Figure 5 we see the pattern of averages for LSA and judges' average scores are practically identical, corroborating the high reliability shown.

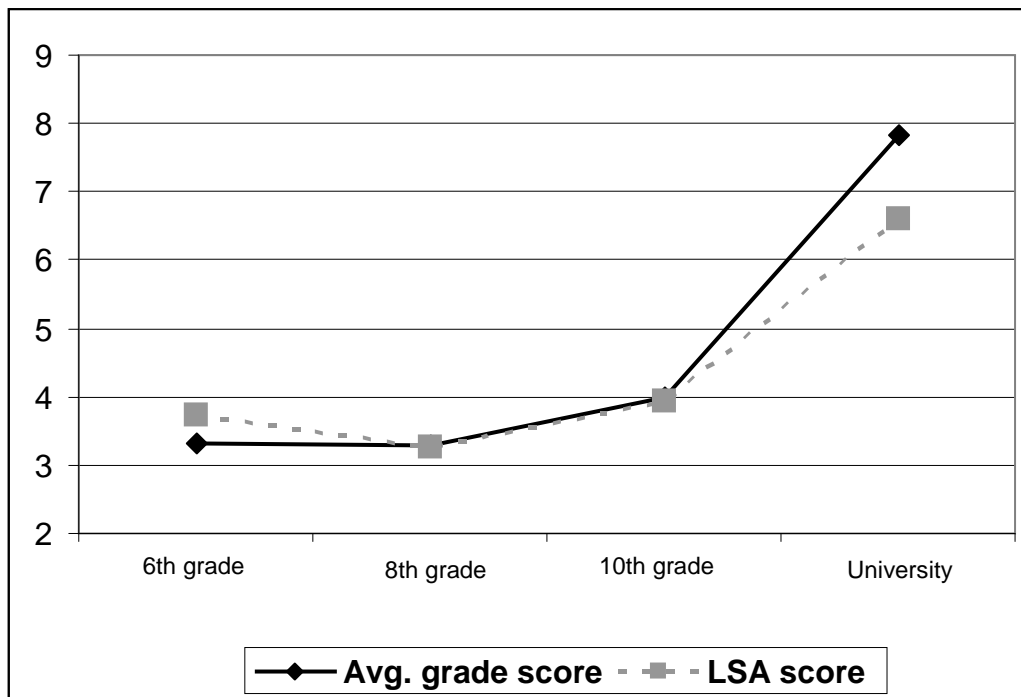


Figure 5. LSA and graders' mean score at each grade level for the expository text

An interesting phenomenon shown in figures 4 and 5 is that LSA scores

overestimate the lowest grade level summaries and underestimate the university level

summaries, while at the intermediate level it matches the human graders' opinions.

This effect is clearest for the narrative text. LSA measures seem less appropriate for

evaluating the lowest and highest scores than human graders

F) Correlations between multiple choice and summary assessments

Lastly, in this section our aim is to investigate the relationship between

summary gradings and multiple choice scores. As shown by the data in Table 6, there

are correlations between multiple choice scores and expert grader scores for summaries,

quite similar for both texts. However, the correlations obtained for LSA are generally

lower than those for experts. LSA scores using vector length vary greatly depending on

655 the type of text, with no significant correlation for the narrative text and a significant
 656 positive correlation similar to those obtained with expert graders in the expository text.

657

658

	Human graders			LSA		
	Content	Coherence	Global (content + coherence)	Cosine	Vector length	Cosine + vector length
Narrative multiple choice test (N=249)	.37**	.45**	.43**	.18**	.11	.18**
Expository multiple choice test (N=243)	.36**	.41**	.39**	.18**	.37**	.34**

659 ** p<.01.

660

661 Table 6. Correlations between multiple choice tests and summary assessments (judges and LSA)
 662 for each text and type of measure (cosine, vector length, content, coherence)

663

664 The results shown in Table 6 are homogeneous for human grader scores on
 665 both types of text and for both content and coherence, when compared with multiple
 666 choice scores. One conclusion that can be drawn from this data is that the graders
 667 scored in a very stable manner, independently of the text type. However, the correlation
 668 between multiple choice scores and LSA cosine and vector length measures were lower
 669 and more discrepant. The vector length measure seems more sensitive to the type of
 670 text studied here.

671

672 Conclusions

673

674 Summarizing is an extremely important task in facilitating students'
 675 comprehension of texts (Brown et al., 1983; E. Kintsch et al., 2007; van Dijk & Kintsch,
 676 1983). The study presents a potential system for automated assessment of readers'
 677 summaries. This is an emergent approach to text analysis and automated grading that is
 678 gaining acceptance, as evidenced by intelligent tutors, automated grading systems,
 679 questionnaires, search systems, dialog management, etc.

680

681 Nonetheless, given the complexities of evaluating open-ended responses, it is
682 still too early to declare the availability of a computer-based tool able to carry out these
683 tasks automatically in an effective manner.

684

685 The aims of this study were threefold: (1) To obtain a reliable LSA-based
686 method that combines vector length and the measure of semantic similarity (2) To
687 compare the sensitivity of LSA and judges to differences in quality of summaries by
688 different grade levels, type of text, and correlations with multiple choice test scores (3)
689 To overcome the limitations of LSA with very short texts. The results showed progress
690 toward fulfilling these aims, although in general they are relatively more satisfactory in
691 evaluations of the expository than the narrative summaries studied in this paper.

692

693 On the one hand, we can conclude that the use of information on vector length
694 together with semantic similarity provides an important improvement when evaluating
695 highly conceptualized summaries with LSA. In previous studies where the texts graded
696 were longer (Foltz et al., 1999; Landauer et al., 1998; Wade-Stein & E. Kintsch, 2004),
697 high levels of reliability were achieved without the need for the measure of the vector
698 length. But the limitations of LSA when applied to shorter texts (Rehder et al., 1998;
699 Wiemer-Hastings et al., 1999) call for more information from the latent semantic space
700 to be used. In this case, using semantic similarity with vector length we are better able
701 to simulate evaluation by human graders. The additional information contributes to
702 making LSA a more reliable tool. On the other hand, the results obtained show that
703 with summaries of the expository text reliability is always above .70, and even above
704 .80 when the average judges' scores are used. The results with the narrative text

705 summaries are slightly less impressive, even though the reliability level was always
706 above .60. It's common for the results to be weaker when working with narrative texts
707 (Wolf, 2005), probably because they are less descriptive, less structured, less factual and
708 more metaphorical than the expository texts. However, in an academic context it is
709 more important to obtain reliable results with expository texts, which are prototypical
710 academic texts. The better results for evaluations of the expository text translate into
711 better sensitivity of LSA to differences in quality between summaries from different
712 grade levels, more so than with the narrative text summary. Lastly, we should point out
713 that a generalist corpus was used to obtain these results, the same as those on the
714 University of Colorado official LSA webpage (<http://lsa.colorado.edu>). No *ad hoc*
715 corpus was built to train LSA, and this should be encouraging news for researchers
716 looking to apply and perhaps improve on this and other similar methods, without the
717 need to incorporate their own subject area-specific corpus. We do not aim to generalize
718 the results obtained to all types of narrative or expository texts, but rather to generalize
719 the regression method proposed for any kind of text to be evaluated with the LSA tool.
720 We should note, however, that as indicated the success of the method is greater for
721 expository texts than for narrative texts. The method's improved performance on the
722 expository text could be partly due to it being slightly shorter than the narrative text,
723 offering a certain advantage for the expository text when we demand that the summaries
724 are so concise.

725

726 One general finding from the results obtained is that LSA overestimates the
727 lowest scores from lower grade levels, and underestimates the highest scores (university
728 level), while for intermediate grade levels there is a closer match with the graders'
729 criteria. This phenomenon is much clearer in the narrative text. One possible

730 explanation is that previous knowledge is often greater in the narrative text, and there is
731 greater usage of connectives and transitions that would exaggerate the differences
732 between good and bad readers. Another factor may be the increased importance of
733 implicit knowledge introduced by a greater number of inferences in narrative than
734 expository texts (Graesser, 1981). On the other hand, these differences are perhaps
735 lessened by the fact that in an expository text the summary more closely matches the
736 content of the text, being less reliant on previous knowledge and with less usage of
737 connectives, synonyms, etc. This should be borne in mind for future studies, to confirm
738 whether the finding can be generalized.

739

740 This study also presents some limitations. The computer-based tool described in
741 this paper is not yet ready for classroom implementation, so we must realistically
742 consider what the current findings indicate for classroom prospects, and what still needs
743 to be done to make this tool more useable in classroom applications. In addition, further
744 research is required with more texts in order to form a sufficient basis for the
745 generalization of our findings.

746

747 Nevertheless, the availability of automatic tools that evaluate reliably, help to
748 detect weak or strong points in the summaries, may take pressure off of teachers, at the
749 same time providing the student with assistance in everyday tasks. LSA has become
750 one of the most widely-used computational tools of recent years, and one of the fastest-
751 growing areas of application has been the field of education. Today, LSA is already a
752 reality in some U.S. classrooms. It is gradually finding its way into schools to help
753 improve students' writing and comprehension strategies (Dikli, 2006; E. Kintsch et al.,
754 2007). This tool needs to be complemented with new algorithms to overcome some of

755 its limitations (Kintsch, 2001; Jorge-Botana et al., 2009; Olmos et al., 2009),
756 mathematically optimize the usage of the latent semantic space (Hu, Cai, Wiemer-
757 Hastings, Graesser & McNamara, 2007). Also it is possible to joint several of them
758 such as we implemented a combination between cosine and vector length. But also
759 need it to link it with psychological models such as semantic memory (Denhière,
760 Lemaire, Bellissens & Jhean-Larose, 2007; León, Jorge-Botana, Olmos & Escudero,
761 2010) or with other computational models of language (Steyvers & Griffiths, 2007).
762 Once all of these contributions are added, the potential and the capability of LSA in the
763 educational sector will be far greater.

764

765

766 ACKNOWLEDGEMENTS

767 This work was supported by Grant SEJ2006-09916 from the Spanish Ministry of
768 Science and Technology and PSI 2009-31932 from the Spanish Ministry of Education.

769

770

771 References

772

773 Bellissens, C., Jeuniaux, P., Duran, N. D. & McNamara, D. S. (2010). A text relatedness
774 and dependency computational model: Using Latent Semantic Analysis and Coh-Metrix
775 to predict self-explanation quality. *Studia Informatica Universalis*, 8, 85-125.

776

777 Bransford, J. D., Brown, A. L. & Cocking, R. R. (2000). *How people learn: Brain,*
778 *mind, experience, and school*. National Research Council Commission on Behavioral
779 and Social Sciences and Education. Washington, DC: National Academy Press.

780

781 Brown, A. L., Bransford, J. D., Ferrara, R. A. & Campione, J. C. (1983). Learning,
782 remembering, and understanding. In J. Flavell & E. M. Markman (Eds.), *Handbook of*
783 *child psychology* (4th ed.). *Cognitive development* (Vol. 3, pp. 515- 629). New York:
784 Wiley.

785

786 Dessus, P. & Lemaire, B. (2002). Using production to assess learning: an ILE that
787 Fosters Self-Regulated Learning. In S. A. Cerri, G. Gouardères & F. Paraguaçu (Eds.),
788 *Intelligent Tutoring Systems (ITS 2002)* (pp. 772-781). Berlin : Springer, LNCS 2363.

789

790 Denhière, G., Lemaire, B., Bellissens, C. & Jhean-Larose, S. (2007). A Semantic Space
791 Modeling Children's Semantic Memory. In T. K. Landauer, D. McNamara, S. Dennis &
792 W. Kintsch (Eds.), *The handbook of Latent Semantic Analysis* (pp. 143- 167). Mahwah,
793 NJ: Erlbaum.

794

795 Dikli, S. (2006). An Overview of Automated Scoring of Essays. *Journal of Technology,*
796 *Learning, and Assessment*, 5, 1. Retrieved [date] from <http://www.jtla.org>.

797

798 Far, R., Pritchard, R. & Smitten, B. (1990). A description of what happens when an
799 examinee takes a multiple-choice reading comprehension test. *Journal of Educational*
800 *Measurement*, 27, 209-226.

801

802 Foltz, P. W., Laham, D. & Landauer, T. K. (1999). The Intelligent Essay Assessor:
803 Applications to educational technology. *Interactive Multimedia Electronic Journal of*

804 *Computer-Enhanced Learning, 1*. Retrieve June 29, 2004, from
805 <http://knowledgetechnologies.com>.
806
807 Graesser, A. C., Chipman, P., Haynes, B. C. & Olney, A. (2005). AutoTutor: An
808 intelligent tutoring system with mixed-initiative dialogue. *IEEE Transactions in*
809 *Education, 48*, 612-618.
810
811 Graesser, A. C., Lu, S., Jackson, G. T., Mitchell, H, Ventura, M., Olney, A. & Lowerse,
812 M. M. (2004). Autotutor: A tutor with dialogue in natural language. *Behaviour*
813 *Research Methods, Instruments, and Computers, 36*, 180-193.
814
815 Graesser, A. C., Penumatsa, P., Ventura, M., Cai, Z. & Hu, X. (2007) .Using LSA in
816 AutoTutor: Learning through Mixed-initiative Dialogue in Natural Language. In T. K.
817 Landauer, D. McNamara, S. Dennis & W. Kintsch (Eds.). *The Handbook of Latent*
818 *Semantic Analysis* (pp. 243-262). Mahwah, NJ: Erlbaum.
819
820 Haley, D. (2009). Applying latent semantic analysis to computer assisted assessment in
821 the Computer Science domain: a framework, a tool, and an evaluation. PhD thesis, The
822 Open University.
823
824 Hu, X., Cai, Z., Wiemer-Hastings, Graesser, A. C. & McNamara, D.S. (2007).
825 Strengths, limitations, and extensions of LSA. In T.K. Landauer, D. McNamara, S.
826 Dennis & W. Kintsch (Eds.), *The handbook of Latent Semantic Analysis* (pp. 401- 426).
827 Mahwah, NJ: Erlbaum.
828

829 Jorge-Botana, G., Olmos, R. & León J.A. (2009). Using LSA and the predication
830 algorithm to improve extraction of meanings from a diagnostic corpus. *The Spanish*
831 *Journal of Psychology*, 12, 2, 424-440.

832

833 Jorge-Botana, G., León J.A., Olmos, R. & Escudero, I. (2010). Latent Semantic
834 Analysis Parameters for Essay Evaluation using Small-Scale Corpora. *Journal of*
835 *Quantitative Linguistics*, 17, 1, 1-29.

836

837 Kakkonen, T. & Sutinen, E. (2011). EssayAid: Towards a Semi-automatic System for
838 Assessing Student Texts. *International Journal of Continuing Engineering Education*
839 *and Life-Long Learning*, 21, 2-3, 119-139.

840

841 Kintsch, E., Caccamise, D., Franzke, M., Johnson, N. & Dooley, S. (2007). Summary
842 Street: Computer-Guided Summary Writing. In T. K. Landauer, D. McNamara, S.
843 Dennis & W. Kintsch (Eds.). *The handbook of Latent Semantic Analysis* (pp. 263- 277).
844 Mahwah, NJ: Erlbaum.

845

846 Kintsch, W. (2002). On the notions of theme and topic in psychological process models
847 of text comprehension. In M. Louwerse & W. van Peer (Eds.), *Thematics:*
848 *Interdisciplinary Studies* (pp. 157-170). Amsterdam: Benjamins.

849

850 Landauer, T. K. & Dumais, S. T. (1997). A solution to Plato's problem: the Latent
851 Semantic Analysis theory of the acquisition, induction, and representation of
852 knowledge. *Psychological Review*, 104, 211-240.

853

854 Landauer, T. K., Foltz, P. W. & Laham, D. (1998). Introduction to Latent Semantic
855 Analysis. *Discourse Processes*, 25, 259-284.
856

857 León, J. A., Olmos, R., Escudero, I., Cañas, J. J. & Salmerón, L. (2006). Assessing short
858 summaries with human judgments procedure and Latent Semantic Analysis in narrative
859 and expository texts. *Behavior Research Methods, Instruments and Computers* 38, 4,
860 616–627.
861

862 McNamara, D. S., Boonthum, C., Levinstein, I. B. & Millis, K. (2007). Evaluating Self-
863 Explanations in iSTART: Comparing Word-Based and LSA Algorithms. In T. K.
864 Landauer, D. McNamara, S. Dennis, & W. Kintsch (Eds.). *The handbook of Latent*
865 *Semantic Analysis* (pp. 227-242). Mahwah, NJ: Erlbaum.
866

867 McNamara, D. S., Levinstein, I. B. & Boonthum, C. (2004). iSTART: interactive
868 strategy training for active reading and thinking. *Behaviour Research Methods,*
869 *Instruments, and Computers*, 36, 222-233.
870

871 Millis, K., Magliano, J., Wiemer-Hastings, K., Todaro, S. & McNamara, D. S. (2007).
872 Assessing and improving comprehension with Latent Semantic Analysis. In T. K.
873 Landauer, D. McNamara, S. Dennis, & W. Kintsch (Eds.), *The handbook of Latent*
874 *Semantic Analysis* (pp. 207-225). Mahwah, NJ: Erlbaum.
875

876 Monjurul Islam, Md. & Latiful Hoque, A. S. M. (2012). Automated Essay Scoring
877 Using Generalized Latent Semantic Analysis. *Journal of Computers*, 7, 3, 616-626.
878

879 Olmos, R., León, J. A., Jorge-Botana, G. & Escudero, I. (2009). New algorithms
880 assessing short summaries in expository texts using Latent Semantic Analysis.
881 *Behaviour Research Methods, Instruments, and Computers*, 41, 944-950.
882

883 Palincsar, A. S. & Brown, A. L. (1984). Reciprocal teaching of comprehension fostering
884 and comprehension-monitoring activities. *Cognition & Instruction*, 1, 117-175.
885

886 Quesada, J.F, Kintsch, W. & Gomez, E. (2001). A Computational Theory of Complex
887 Problem Solving Using the Vector Space Model (part I): Latent Semantic Analysis,
888 Through the Path of Thousands of Ants. In J. J. Cañas (Ed.) *Cognitive research with*
889 *Microworlds*, (pp. 117-131). Granada (Spain).
890

891 Rehder, B., Schreiner, M. E., Wolfe, B. W., Laham, D., Landauer, T. K. & Kintsch, W.
892 (1998). Using Latent Semantic Analysis to assess knowledge: Some technical
893 considerations. *Discourse Processes*, 25, 337-354.
894

895 Shapiro, A. M. & McNamara, D. S. (2000). The use of Latent Semantic Analysis as a
896 tool for the quantitative assessment of understanding and knowledge. *Journal of*
897 *Educational Computing Research*, 22, 1-36.
898

899 Steyvers, M. & Griffiths, T. (2007). Probabilistic Topic Models. In T. K. Landauer, D.
900 McNamara, S. Dennis, & W. Kintsch (Eds.), *The handbook of Latent Semantic Analysis*
901 (pp. 427-448). Mahwah, NJ: Erlbaum.
902

903 Van Dijk, T. A. & Kintsch, W. (1983). *Strategies of discourse comprehension*. New
904 York: Academic Press.
905
906 Wade-Stein, D. & Kintsch, E. (2004). Summary Street: Interactive computer support for
907 writing. *Cognition and Instruction*, 22, 333-362.
908
909 Wiemer-Hastings, P., Wiemer-Hastings, K. & Graesser, A. (1999). How latent is Latent
910 Semantic Analysis? In *Proceedings of the Sixteenth International Joint Congress on*
911 *Artificial Intelligence* (pp. 932–937). San Francisco: Morgan Kaufmann.
912
913 Wolfe, M. B. W. (2005) Memory for narrative and expository text: Independent
914 influences of semantic associations and text organization. *Journal of Experimental*
915 *Psychology: Learning, Memory and Cognition*, 31, 2, 359-364.
916
917 Wolfe, M. B., Schreiner, M. E., Rehder, B., Laham, D., Foltz, P. W., Kintsch, W. &
918 Landauer, T. K. (1998). Learning from text: Matching readers and text by Latent
919 Semantic Analysis. *Discourse Processes*, 25, 309-336.