

## Dear Author

Here are the proofs of your article.

- You can submit your corrections **online**, via **e-mail** or by **fax**.
- For **online** submission please insert your corrections in the online correction form. Always indicate the line number to which the correction refers.
- You can also insert your corrections in the proof PDF and **email** the annotated PDF.
- For **fax** submission, please ensure that your corrections are clearly legible. Use a fine black pen and write the correction in the margin, not too close to the edge of the page.
- Remember to note the **journal title**, **article number**, and **your name** when sending your response via e-mail or fax.
- **Check** the metadata sheet to make sure that the header information, especially author names and the corresponding affiliations are correctly shown.
- **Check** the questions that may have arisen during copy editing and insert your answers/corrections.
- **Check** that the text is complete and that all figures, tables and their legends are included. Also check the accuracy of special characters, equations, and electronic supplementary material if applicable. If necessary refer to the *Edited manuscript*.
- The publication of inaccurate data such as dosages and units can have serious consequences. Please take particular care that all such details are correct.
- Please **do not** make changes that involve only matters of style. We have generally introduced forms that follow the journal's style.
- Substantial changes in content, e.g., new results, corrected values, title and authorship are not allowed without the approval of the responsible editor. In such a case, please contact the Editorial Office and return his/her consent together with the proof.
- If we do not receive your corrections **within 48 hours**, we will send you a reminder.
- Your article will be published **Online First** approximately one week after receipt of your corrected proofs. This is the **official first publication** citable with the DOI. **Further changes are, therefore, not possible.**
- The **printed version** will follow in a forthcoming issue.

### Please note

After online publication, subscribers (personal/institutional) to this journal will have access to the complete article via the DOI using the URL:

<http://dx.doi.org/10.3758/s13428-013-0326-1>

If you would like to know when your article has been published online, take advantage of our free alert service. For registration and further information, go to:

<http://www.springerlink.com>.

Due to the electronic nature of the procedure, the manuscript and the original figures will only be returned to you on special request. When you return your corrections, please inform us, if you would like to have these documents returned.

**Metadata of the article that will be visualized in OnlineFirst**

1	Article Title	<b>EsPal: One-stop shopping for Spanish word properties</b>
2	Article Sub- Title	
3	Article Copyright - Year	<b>Psychonomic Society, Inc. 2013 (This will be the copyright line in the final PDF)</b>
4	Journal Name	Behavior Research Methods
5	Corresponding Author	Family Name <b>Duchon</b>
6		Particle
7		Given Name <b>Andrew</b>
8		Suffix
9		Organization Basque Center on Cognition, Brain, and Language
10		Division
11		Address Donostia, Spain
12		e-mail a.duchon@bcbl.eu
13		Author
14	Particle	
15	Given Name <b>Manuel</b>	
16	Suffix	
17	Organization Universitat of València	
18	Division	
19	Address Valencia, Spain	
20	e-mail	
21	Author	
22		Particle
23		Given Name <b>Nuria</b>
24		Suffix
25		Organization Universitat Pompeu Fabra
26		Division
27		Address Barcelona, Spain
28		e-mail
29		Author
30	Particle	

31		Given Name	<b>Antonia</b>
32		Suffix	
33		Organization	Universitat de Barcelona
34		Division	
35		Address	Barcelona, Spain
36		e-mail	
<hr/>			
37		Family Name	<b>Carreiras</b>
38		Particle	
39		Given Name	<b>Manuel</b>
40		Suffix	
41		Organization	Basque Center on Cognition, Brain, and Language
42	Author	Division	
43		Address	Donostia, Spain
44		Organization	IKERBASQUE. Basque Foundation for Science
45		Division	
46		Address	Bilbao, Spain
47		e-mail	
<hr/>			
48		Received	
49	Schedule	Revised	
50		Accepted	
<hr/>			
51	Abstract	<p>This article introduces EsPal: a Web-accessible repository containing a comprehensive set of properties of Spanish words. EsPal is based on an extensible set of data sources, beginning with a 300 million token written database and a 460 million token subtitle database. Properties available include word frequency, orthographic structure and neighborhoods, phonological structure and neighborhoods, and subjective ratings such as imageability. Subword structure properties are also available in terms of bigrams and trigrams, biphones, and bisyllables. Lemma and part-of-speech information and their corresponding frequencies are also indexed. The Web site enables users either to upload a set of words to receive their properties or to receive a set of words matching constraints on the properties. The properties themselves are easily extensible and will be added over time as they become available. It is freely available from the following Web site: <a href="http://www.bcbl.eu/databases/espal/">http://www.bcbl.eu/databases/espal/</a>.</p>	
<hr/>			
52	Keywords separated by ' - '	<input type="checkbox"/>	
<hr/>			
53	Foot note information		

# EsPal: One-stop shopping for Spanish word properties

Andrew Duchon · Manuel Perea · Nuria Sebastián-Gallés ·

Antonia Martí · Manuel Carreiras

© Psychonomic Society, Inc. 2013

**Abstract** This article introduces EsPal: a Web-accessible repository containing a comprehensive set of properties of Spanish words. EsPal is based on an extensible set of data sources, beginning with a 300 million token written database and a 460 million token subtitle database. Properties available include word frequency, orthographic structure and neighborhoods, phonological structure and neighborhoods, and subjective ratings such as imageability. Subword structure properties are also available in terms of bigrams and trigrams, biphones, and bisyllables. Lemma and part-of-speech information and their corresponding frequencies are also indexed. The Web site enables users either to upload a set of words to receive their properties or to receive a set of words matching constraints on the properties. The properties themselves are easily extensible and will be added over time as they become available. It is freely available from the following Web site: <http://www.bcbl.eu/databases/espal/>.

## Keywords

A. Duchon (✉) · M. Carreiras  
Basque Center on Cognition, Brain, and Language, Donostia, Spain  
e-mail: a.duchon@bcbl.eu

M. Perea  
Universitat of València, Valencia, Spain

N. Sebastián-Gallés  
Universitat Pompeu Fabra, Barcelona, Spain

A. Martí  
Universitat de Barcelona, Barcelona, Spain

M. Carreiras  
IKERBASQUE. Basque Foundation for Science, Bilbao, Spain

Researchers from a wide range of disciplines (e.g., neuroscience, artificial intelligence, psychology, linguistics, and education, among others) who work in the interdisciplinary area of language research (e.g., language acquisition, language processing, language learning, bilingualism, and computational linguistics) need quick and efficient access to information about specific properties of words. For example, word frequency is a dominant factor in accounting for visual word recognition speed as measured by lexical decision times (Forster & Chambers, 1973; Monsell, 1991) and eye fixation durations during reading (Rayner, 2009). Unsurprisingly, reading behavior as measured by, for example, lexical decision, naming, fixation times, and so on is affected by a wide range of other properties of words, including orthographic neighborhood (Carreiras, Perea, & Grainger, 1997; Grainger, 1990), syllable frequency (Carreiras, Alvarez, & de Vega, 1993; Carreiras & Perea, 2004; Perea & Carreiras, 1998), and imageability (James, 1975), to cite just a few examples. Similarly, with regard to other fields that employ linguistic stimuli, such as memory research, it has been shown that word frequency plays a role in short-term memory (Hulme et al., 1997) and syllable length in working memory (Gathercole & Baddeley, 1990).

Given the wide range of word properties that can affect language and cognitive processing, it is desirable to have a single, integrated, and updateable source of data. For Spanish, there are now a variety of databases available, but some are based on a relatively small number of tokens (Davis & Perea, 2005; Sebastián-Gallés, Martí, Carreiras, & Cuetos, 2000; Taulé, Martí, & Recasens, 2008), while others provide information about a limited number of variables (Alonso, Fernandez, & Díez, 2011; Cuetos-Vega, González-Nosti, Barbón-Gutiérrez, & Brysbaert, 2011; Davies, 2005; Marian, Bartolotti, Chabal, & Shook, *in press*). EsPal (Español Palabras, meaning simply “Spanish words”) is a Web-based repository available at <http://www.bcbl.eu/databases/espal/> that has been designed to fill this gap, providing information on a comprehensive set of

67 word properties from corpora with hundreds of millions of  
68 words.

69 The most similar effort is the Syllabarium (Duñabeitia,  
70 Cholin, Corral, Perea, & Carreiras, 2010), which is a Web-  
71 based tool accessing a database containing information on  
72 word frequencies and syllable frequencies by token and syl-  
73 lable position. Standalone software packages are also avail-  
74 able for Spanish and other languages that provide subsets of  
75 the properties in EsPal (Davis, 2005; Davis & Perea, 2005;  
76 New, Pallier, Brysbaert, & Ferrand, 2004; Perea et al., 2006).  
77 However, given the size of the corpora (discussed below),  
78 some of the calculations for some of the properties take up  
79 to a week on a standard PC, so a precomputed set of properties  
80 is preferred. With EsPal, the back-end processing for the word  
81 and subword properties is conducted using a multistep pro-  
82 gram written in Java, which precomputes not only basic  
83 properties of word frequency and form, but also orthographic  
84 structure and neighborhoods, phonological structure and  
85 neighborhoods, lemma and part-of-speech properties, and  
86 subword structure properties related to letter bigrams and tri-  
87 grams, bisyllables, and biphones. In addition, other data such  
88 as a word’s subjective ratings (e.g., familiarity, imageability,  
89 etc.) can be easily attached to the data and made searchable.

90 The second important factor of EsPal is the capacity to  
91 apply the exact same processing to different corpora. A num-  
92 ber of studies have shown that, across many languages, word  
93 frequencies derived from movie subtitle corpora provide a  
94 better account for various psycholinguistic effects  
95 (Brysbaert, New, & Keuleers, 2012; Cai & Brysbaert, 2010;  
96 Cuetos-Vega et al., 2011; Dimitropoulou, Duñabeitia, Avilés,  
97 Corral, & Carreiras, 2010; Keuleers, Brysbaert, & New, 2010;  
98 New, Brysbaert, Veronis, & Pallier, 2007). However, proper-  
99 ties from written corpora have in the past been more common  
100 and may better predict some phenomena, so it is useful to have  
101 different sources of data available for researchers, depending  
102 on their goals. EsPal currently fulfills this goal by applying the  
103 same processing to both a corpus based on movie subtitles and  
104 one based on written text (fiction, nonfiction, and Web pages).

105 Finally, the Spanish-speaking community is diverse, and  
106 EsPal is constructed to be able to accommodate this diver-  
107 sity, at least in terms of phonological representation. Stan-  
108 dard Castilian Spanish spoken on mainland Spain dif-  
109 fers in a number of dimensions from the Spanish spoken in  
110 the Canary Islands and in Latin America (which itself is  
111 quite diverse). EsPal therefore also allows the user to choose  
112 which phonological representation is used, for example, to  
113 derive properties related to phonological neighborhoods.

114 In the remainder of this article, we describe the collection  
115 and preprocessing of the written and subtitle databases  
116 currently available in EsPal; how we calculate orthographic  
117 and phonological properties, subword properties, lemma  
118 and part-of-speech properties; and the source of the  
119 subjective ratings data.

**Written corpus collection and preprocessing** 120

Written corpus collection 121

122 The EsPal Written Corpus is derived from a wide selection  
123 of texts collected from the Web or available in digital  
124 format. Table 1 provides a listing of percentages in terms  
125 of word tokens across the different sources and genres. We  
126 grouped them into nine subsets according to their content:  
127 academic, culture, law, philosophy, literature, news, politics,  
128 society, and the Spanish Wikipedia. All these texts had to  
129 meet the requirements of being freely available and not  
130 subject to copyright. Most documents were gathered from  
131 Web sites featuring a variety of linguistic styles, including  
132 formal, colloquial, and specialized language.

133 The academic texts are mainly Ph.D. theses selected from  
134 a wide range of scientific fields: anthropology, architecture,  
135 art, biology, law, economics, electronics, philology, philo-  
136 sopy, physics, history, humanities, engineering, mathemat-  
137 ics, medicine, psychology, chemistry, telecommunications,  
138 and veterinary science. The set of culture texts is composed  
139 of news about cultural events from several newspapers and  
140 blogs of opinion about films. Legal texts include mainly  
141 rulings by the High Court of Justice of several autonomous  
142 regions in Spain, as well as news from the judiciary field as  
143 it appeared in popular newspapers (*El Mundo, El País, and*  
144 *El Periódico*). The literary texts come from several Web  
145 pages containing works with expired copyrights (*bdigital,*  
146 *biblioteca\_ignoria, libroteca, logos, and scribd*). These  
147 works are both texts written in Spanish and translations into  
148 Spanish. The news is from the EFE Agency from January,  
149 February, and March 2000. The politics set contains news  
150 texts referring to Spain’s 2007 autonomic elections,  
151 speeches by the Spanish President during 2008, and docu-  
152 ments taken from political party Web sites. The society set is  
153 composed of Web texts about religion, abortion, and psy-  
154 chology. Finally, the Web data are from the whole Spanish  
155 Wikipedia, circa February 2009.

**Table 1** Percentage of terms by source type in the EsPal written corpus t1.1

Source type	Percent of terms	t1.2
Academics	1.8 %	t1.3
Culture	0.2 %	t1.4
Law	1.0 %	t1.5
Philosophy	1.1 %	t1.6
Literature	22.5 %	t1.7
News	8.7 %	t1.8
Politics	16.0 %	t1.9
Society	4.7 %	t1.10
Web/Wikipedia	43.9 %	t1.11

156 The whole corpus underwent a process of cleaning to  
 157 eliminate the metadata usually present in this type of texts.  
 158 This process was both automatic and manual and was ex-  
 159 tremely time consuming.

160 **Written corpus preprocessing**

161 Before the data were incorporated into EsPal, all the text  
 162 was first parsed using the FreeLing part-of-speech tagger  
 163 (Padró, Collado, Reese, Lloberes, & Castellón, 2010) to  
 164 output into a file one term per line with its lemma and its  
 165 part of speech. The parsing resulted in a total of  
 166 309,530,600 terms (no punctuation was included). A “term”  
 167 could be one or more words and included dates (*17 de julio*  
 168 *de 1990* [“July 17, 1990”]), proper nouns (*Congreso de los*  
 169 *Estados Unidos* [“United States Congress”]), or phrases  
 170 (*por ejemplo* [“for example”]). These terms were then  
 171 imported into a raw\_sequence table in EsPal, with one word  
 172 per row (i.e., multiword terms were separated) and columns  
 173 for the lemma and the part-of-speech tag (e.g., [http://](http://nlp.lsi.upc.edu/freeling/doc/tagsets/tagset-es.html)  
 174 [nlp.lsi.upc.edu/freeling/doc/tagsets/tagset-es.html](http://nlp.lsi.upc.edu/freeling/doc/tagsets/tagset-es.html)). If the  
 175 word came from a multiword term, then the word itself  
 176 was used as the lemma. In this manner, the part-of-speech  
 177 tag is maintained for the word within its larger context—for  
 178 example, *de* [“of”] will have lemma statistics as a date and a  
 179 proper noun (among others) in addition to being a preposi-  
 180 tion. In the lemma processing section below, we describe  
 181 further lemma information available for words. The word  
 182 and lemma were changed to all lowercase using the Java  
 183 string function toLowerCase with the “es” locale. This table  
 184 had a total of 325,773,444 rows. Subsequent processing of  
 185 the contents of the raw\_sequence table is described later.

186 **Subtitle corpus collection and pre-processing**

187 **Subtitle corpus collection**

188 A total of 100,659 Spanish subtitle files were originally pro-  
 189 vided by the [www.opensubtitles.org](http://www.opensubtitles.org) Web site including meta-  
 190 data about the file (such as author and total downloads). The  
 191 Internet Movie Database (IMDb) ID was also supplied, by  
 192 which genre, director, and cast information can be obtained.  
 193 Subtitle file formats contain an index number, the start and stop  
 194 time for which the subtitle is to be shown on screen in milli-  
 195 seconds, and the text of the subtitle, all of which were stored in  
 196 the subtitles table of the database. Movies account for 65.6 %  
 197 of the files, with the remainder from television episodes. A  
 198 given show can be labeled with more than one genre, so the  
 199 words in a subtitle file can be double counted, but across all  
 200 such counts, by genre, 22.0 % of the words are from dramas,  
 201 10.9 % from comedies, 10.3 % from thrillers, 7.7 % from crime  
 202 shows, 7.4 % from action shows, 7.3 % from romances, 5.8 %

from mysteries, and 5.5 % from adventure shows, and the 203  
 remaining are in 13 other genres, accounting for less than 204  
 5 % each. Similarly, the source show can contain more than 205  
 one language, and across such counts, 52.6 % of the words are 206  
 from English language shows, followed by French (8.5 %) and 207  
 Spanish (5.5 %). No limits were put on the date of the source, 208  
 since the subtitles themselves, uploaded by users of the Web 209  
 site, are of recent origin. However, given the metadata main- 210  
 tained about the source of the words, a variety of subcorpora 211  
 are possible whose properties might be more appropriate 212  
 depending on the psycholinguistic question being asked. 213

Subtitle corpus preprocessing 214

For a proper parsing of text, complete sentences are needed. 215  
 However, a single subtitle instance could have two speakers 216  
 (usually denoted by a dash [“-”] at the beginning of each of 217  
 their statements), or a single speaker’s statement could con- 218  
 tinue in the next subtitle instance (usually denoted by ellip- 219  
 ses [“...”] at then end). Therefore, a second stage of 220  
 processing was run to fill a statements table with strings 221  
 that were, at a first approximation, single statements (which 222  
 could contain multiple sentences). At this stage, subtitles 223  
 were removed that contained metadata (such as the author of 224  
 the subtitles or translations of the credits); all HTML mark- 225  
 ings were removed; and contents within brackets (often 226  
 indicating sounds) were also removed. 227

Each statement was submitted individually to FreeLing 228  
 (Padró et al., 2010) for part-of-speech tagging and lemma- 229  
 tization. In this case, the lowercased word, lowercased lem- 230  
 ma, and part-of-speech tag were stored directly in the 231  
 raw\_sequence table, along with the file ID, IMDb movie 232  
 ID, statement index, and within-statement index. Thus, the 233  
 provenance, or origin, of every word can be traced back, 234  
 enabling further analyses, which we will be reporting in the 235  
 future. In the end, words from 98,339 distinct files and 236  
 40,444 unique movies are present in this table. 237

**Word selection and frequency processing** 238

The raw\_sequence table holds every individual word token 239  
 from the source. The count of each unique word type is 240  
 accumulated in a second table, raw\_words. Every word type 241  
 in this table is checked against the criteria below. Those that 242  
 do not pass the criteria are marked as rejected. The word had 243  
 to appear in at least one of these publicly available sets of 244  
 Spanish words: OpenOffice,<sup>1</sup> AGME,<sup>2</sup> or SemEval.<sup>3</sup> For 245  
 future comparison, we also allowed words present in other 246

<sup>1</sup> <http://wiki.services.openoffice.org/wiki/Dictionaries>.

<sup>2</sup> <http://www.cic.ipn.mx/~sidorov/agme/>.

<sup>3</sup> [http://www.lsi.upc.edu/~nlp/semeval/msacs\\_download.html](http://www.lsi.upc.edu/~nlp/semeval/msacs_download.html).



247 recent Spanish corpora projects (Alonso et al., 2011;  
 248 Cuetos-Vega et al., 2011). In addition, we included a large  
 249 number of Spanish first names, surnames, and place names  
 250 from publicly available Web sites. Rejection criteria were  
 251 that words could not be longer than 30 characters,<sup>4</sup> contain a  
 252 nonletter (which excluded hyphenated words), have more  
 253 than 3 characters in a row of the same character, nor contain  
 254 non-Spanish characters that is, outside of a–z, áéíóúñü.  
 255 Words that passed these filters were placed into the word\_  
 256 data table with their counts.<sup>5</sup> Table 2 contains the final  
 257 counts of word types and word tokens for the two corpora.

258 The word\_data table contains all the information about  
 259 each word and, thus, what can be searched for simultaneou-  
 260 sly via the Web interface. We will be presenting the various  
 261 properties available for each word with its column name in  
 262 bold italics. For each *word*, we store the count (*cnt*), the  
 263 frequency per million (*frq*),  $\log_{10}(cnt+1)$  (*log\_cnt*), and  
 264  $\log_{10}(frq + 1/N)$ , where  $N$  = millions of words in the  
 265 database (*log\_frqN*), which has been shown to be a fruitful  
 266 way to compare frequencies across corpora (Brysbaert et al.,  
 267 2011).

268 Subtitle corpus contextual diversity processing

269 Recent work has found that the number of different contexts  
 270 in which a word occurs can be more informative than the  
 271 token frequency (Adelman, Brown, & Quesada, 2006;  
 272 Brysbaert & New, 2009; Dimitropoulou et al., 2010;  
 273 Keuleers et al., 2010; Perea, Soares, & Comesaña, in  
 274 press). The original EsPal subtitles database described  
 275 above uses all the files available, so some shows are multi-  
 276 ply represented. Therefore, EsPal provides a third database  
 277 of properties (subtitles\_cdm) that are based on the number  
 278 of different movies (IMDb IDs) that the word appears in. In  
 279 this database, *cnt* refers to the count of different movies, and  
 280 *frq* is equal to the percent of movies (i.e.,  $100 * cnt/40,444$ ).  
 281 We also explored using the count of different subtitle files,  
 282 with the expectation that this would have some relationship  
 283 to popularity (e.g., there are almost 300 versions of *Lord of*  
 284 *the Rings: Return of the King*) and, therefore, provide word  
 285 frequencies that were better predictors of certain psycholin-  
 286 guistic variables. However, in all the cases we have explored  
 287 to date, the contextual diversity based on the number of  
 288 movies has given slightly better results.

<sup>4</sup> A cutoff was made for processing and memory considerations. Out of over 460 million tokens in the raw subtitle data set, only 735 tokens have a length greater than 30.

<sup>5</sup> The system is designed such that at this stage, it would also have been possible to further reduce the words by removing accents or tildes and collapsing the counts across the subsequent word forms. Some psycholinguistic research questions, such as studies focused on stress assignment (e.g., Shelton, Gerfen, & Gutiérrez-Palma, 2011), might benefit from this type of frequency data. However, the first version of these sources has the actual form of the word.

**Table 2** Counts of word types and word tokens in each corpus t2.1

Corpus	Word types	Word tokens	t2.2
Written	277,771	307,772,547	t2.3
Subtitles	244,983	462,611,693	t2.4

**Orthographic properties processing** 289

Orthographic structure 290

The basics of the orthographic structure are, of course, 291  
 present in the *word* column itself. In addition, the number 292  
 of letters (*num\_letters*) and whether or not there are repeat- 293  
 ed letters (*rep\_letters*) within the word (0 = false, and 1 = 294  
 true) are stored. A straightforward consonant–vowel struc- 295  
 ture (*orth\_cv\_structure*) was also created by replacing each 296  
 vowel character (*a,e,i,o,u*, with or without accents, but not 297  
*y*) with “V” and all other characters with “C.” Note, how- 298  
 ever, that there are certain limitations to this simple heuris- 299  
 tic, especially with regard to the letters *y* and *h*. 300

Orthographic neighborhoods 301

Orthographic neighborhood size affects a large number of 302  
 psycholinguistic phenomena (Carreiras et al., 1997; Davis, 303  
 Perea, & Acha, 2009; Grainger, 1990; Yarkoni, Balota, & 304  
 Yap, 2008). For EsPal, each word was compared with all 305  
 other words in the same source in order to provide an array 306  
 of neighborhood properties. For single-change substitution, 307  
 addition, deletion, and transpose letter neighbors, data are 308  
 provided such as the list of neighbors and the frequency of 309  
 the highest frequency neighbor. The average edit distance 310  
 (Levenshtein distance) of the 20 closest words (no matter 311  
 how far) is also provided (*Lev\_N*). Another way to compare 312  
 a word with all the others is the character in the word at 313  
 which it is no longer like any other word (*orth\_uniq\_point*), 314  
 which is a factor in reading studies (e.g., Miller, Juhasz, & 315  
 Rayner, 2006). If the word is completely unique, a second- 316  
 ary orthographic uniqueness point (*orth\_sec\_uniq\_point*) is 317  
 determined as well in case the uniqueness is simply due to, 318  
 for example, the plural form of the word. Table 3 contains 319  
 all of the measures available concerning orthographic 320  
 neighborhoods. 321

**Phonological properties processing** 322

Phonological structure 323

Spanish is a relatively transparent language, so syllable and 324  
 phonological structure can be derived from the orthography 325  
 in a rule-based fashion. To derive the syllable structure, we 326

t3.1 **Table 3** Orthographic neighborhood variable names and meanings

t3.2	Variable name	Variable meaning
t3.3	<i>N</i>	Number of substitution neighbors
t3.4	<i>NHF</i>	Number of higher frequency substitution neighbors
t3.5	<i>frq_hf_s</i>	Frequency of the highest frequency substitution neighbor
t3.6	<i>hf_s</i>	Highest frequency substitution neighbor
t3.7	<i>hf_s_list</i>	List of substitution neighbors in descending frequency, with the place of the word itself marked by "OOOOOO"
t3.8	<i>P</i>	Number of positions with substitution neighbors
t3.9	<i>PHF</i>	Number of positions with higher frequency substitution neighbors
t3.10	<i>avg_frq_Ns</i>	Average frequency of substitution neighbors
t3.11	<i>N_TL</i>	Number of transposed-letter neighbors
t3.12	<i>frq_hf_tl</i>	Frequency of the highest frequency transposed-letter neighbor
t3.13	<i>hf_tl</i>	Highest frequency transposed-letter neighbor
t3.14	<i>hf_tl_list</i>	List of transposed-letter neighbors in descending frequency, with the place of the word itself marked by "OOOOOO"
t3.15	<i>N_A</i>	Number of addition-letter neighbor
t3.16	<i>frq_hf_A</i>	Frequency of the highest frequency addition-letter neighbor
t3.17	<i>hf_A</i>	Highest frequency addition-letter neighbor
t3.18	<i>hf_A_list</i>	List of addition-letter neighbors in descending frequency, with the place of the word itself marked by "OOOOOO"
t3.19	<i>N_D</i>	Number of deletion-letter neighbors
t3.20	<i>frq_hf_D</i>	Frequency of the highest frequency deletion-letter neighbor
t3.21	<i>hf_D</i>	Highest frequency deletion-letter neighbor
t3.22	<i>hf_D_list</i>	List of substitution neighbors in descending frequency, with the place of the word itself marked by "OOOOOO"
t3.23	<i>orth_uniq_point</i>	The character in the word at which it is no longer like any other word
t3.24	<i>orth_sec_uniq_point</i>	If the word is unique, then the uniqueness point with the last letter removed
t3.25	<i>Lev_N</i>	Average Levenshtein distance of the 20 closest words (OLD20)

327 implemented, with some minor changes, the rules in  
 328 Silabeador TIP (Hernández-Figueroa, Rodríguez-Rodríguez,  
 329 & Carreras-Riudavets, 2009) to obtain orthographic syllable  
 330 boundaries (*orth\_syll\_structure*). The most notable change  
 331 was the addition of the onset, nucleus, and coda information  
 332 being stored for each character. From this information, the  
 333 number of syllables (*num\_syll*) and the position of the syllable  
 334 with the accent was also derived (*syll\_accent*).

335 The phonetic transcription of the word (*phon\_structure*)  
 336 was derived using a Java implementation of the rules in the  
 337 SAGA project (Nogueiras & Mariño, 2009) taking advan-  
 338 tage, when necessary, of the syllabification described above.  
 339 For example, the letter *t* is phonetically transcribed as *t* (*toro*  
 340 {"bull"} → *toro*), except when it is syllable final (*etnia*  
 341 ["ethnicity"] → *eDnja*). The codes were modified to be a  
 342 single character and are shown in Table 4. From this  
 343 information, the number of phonemes (*num\_phon*), the  
 344 initial phoneme (*init\_phon*), and the phonetically based  
 345 CV structure (*phon\_cv\_structure*) were derived.<sup>6</sup>

<sup>6</sup> Note that exceptions to the rules have not been implemented, and we are investigating other methods by which to derive phonetic transcriptions.

Two phonetic representations were derived, one for 346  
 Castilian Spanish and one for Latin American Spanish. 347  
 Although this is a complex topic and pronunciation varies 348  
 dramatically within and between countries (Moreno & 349  
 Mariño, 1998), for this introduction of EsPal, the only 350  
 difference between these two representations are that *z* and 351  
*c* (followed by *e* or *i*) are transcribed as *T* in Castilian and *s* 352  
 in Latin American Spanish. However, the software and Web 353  
 site are capable of accommodating any number of phonetic 354  
 representations, and more accurate representations can be 355  
 added over time. In the database and Web site output, these 356  
 columns and the neighborhood columns described below 357  
 are prepended by either *es* or *sa* for Castilian and 358  
 Latin American Spanish, respectively, depending on 359  
 which representation is chosen. 360

Phonological neighborhoods 361

With a single-character representation of the phonemes of 362  
 each word, we can use exactly the same neighborhood 363  
 processing as was used for the orthographic neighborhoods. 364  
 However, in the spoken word recognition literature, slightly 365  
 different variables are typically investigated, so the proper- 366  
 ties provided are different from those for the orthographic 367



t4.1 **Table 4** Phonetic transcription codes used in EsPal

t4.2	SAGA code	EsPal code	Sound
t4.3	p	p	voiceless bilabial plosive
t4.4	b	b	voiced bilabial plosive
t4.5	t	t	voiceless dental plosive
t4.6	d	d	voiced dental plosive
t4.7	k	k	voiceless velar plosive
t4.8	g	g	voiced velar plosive
t4.9	m	m	voiced bilabial nasal
t4.10	n	n	voiced alveolar nasal
t4.11	N	N	voiced velar nasal (preceding a velar consonant)
t4.12	J	J	voiced palatal nasal
t4.13	tS	C	voiceless palatal affricate
t4.14	f	f	voiceless labiodental fricative
t4.15	T	T	voiceless interdental fricative
t4.16	s	s	voiceless alveolar fricative
t4.17	z	z	voiced alveolar fricative (preceding a voiced consonant)
t4.18	jj	H	voiced palatal fricative
t4.19	x	x	voiceless velar fricative
t4.20	l	l	voiced alveolar lateral
t4.21	L	L	voiced lateral palatal
t4.22	rr	R	voiced alveolar trill
t4.23	j	j	palatal semivowel
t4.24	w	w	labiovelar semivowel
t4.25	B	B	voiced bilabial approximant
t4.26	D	D	voiced dental approximant
t4.27	G	G	voiced velar approximant
t4.28	r	r	simple vibrating voiced alveolar
t4.29	a	a	open central vowel
t4.30	e	e	front half vowel
t4.31	i	i	front closed vowel
t4.32	o	o	half rounded back vowel
t4.33	u	u	closed rounded back vowel

368 neighborhoods. Table 5 contains a listing of those phonological neighborhood variables currently available.

370 Subword processing

371 Infralexical, or subword, features are known to influence  
 372 lexical decision and naming times (Carreiras et al., 1993;  
 373 Carreiras & Perea, 2004). The processing was very similar  
 374 for bigrams, trigrams, biphones, and bisyllables, but for  
 375 exposition we will describe only bigram processing. A  
 376 new table `bigram_raw` is created to hold for each bigram-  
 377 word-position combination the sum of word token frequen-  
 378 cies (*frq*) and word type counts from the `word_data` table.  
 379 For instance, when the word *casa* (“house”) is encountered,  
 380 it is found to contain three bigrams (*ca*, *as*, *sa*) with

positions 1, 2, and 3, respectively.<sup>7</sup> An entry is made in the `bigram_raw` table for each of these bigrams at their positions, and the frequency per million (*frq*) of *caso* is added to the token frequency column and 1 is added to the type count column. When the word *caso* (“case”) is encountered, *ca* at position 1 and *as* at position 2 have their token frequency and type count columns incremented by the frequency per million of *caso* and 1, respectively; and a new entry for *so* is made at position 3.

After information from all the words was added to the `bigram_raw` table, each word was reanalyzed to obtain properties of its bigrams. For example, across the entire word *casa*, we can sum or average, in terms of token frequency or type count, its three bigram frequencies. These sums and averages can also either respect the position of the bigram or not (e.g., *ca* at position 1 vs. at any position). Thus, there are eight bigram values that are available for each word as a whole.

For a given word, EsPal also provides each bigram’s token frequency and type count, either for the bigram in that position only or for the bigram in that position found anywhere in a word. So *caso* has three nonzero bigram data sets, and the first data set has the token frequency and type count of *ca* at position 1 and of *ca* at any position. Bigram and trigram data are calculated for words with up to 20 characters. Similar processing is done for biphones on the basis of the phonetic structure (*phon\_structure*) up to 20 phonemes and for bisyllables on the basis of the individual syllables in the orthographic syllable structure (*orth\_syll\_structure*) up to eight syllables.

To provide this large amount of infralexical information, we created a systematic method for deriving property names. Property name affixes are added for each n-gram length (bigram [*B*] or trigram [*T*]), and for each n-gram modality (orthographic [*O*], phonemic [*P*], syllabic [*S*]). So, bigram = *BO*; trigram = *TO*; biphone = *BP*; and bisyllable = *BS*. The system is designed to be extensible, so any other combination of interest could be added. Currently, the frequency per million (*frq*) is used and denoted by *F* in the variable name, but the count (*cnt*) could also be used, as well as the log of either. We can add such versions of the calculations as they are requested. Eight variables are made for each length–type combination. These have combinations that are position sensitive (*pos\_*) or independent (*abs\_*) sums (*S*) or means (*M*) of the token frequency (*tok\_*) or type count (*type\_*). The previous code is then appended to the property name. For example, the position-independent mean of biphone token frequencies is *abs\_tok\_MBPF*.

<sup>7</sup> Note that it is common to have markers for the beginning and end of words as well; for example, *casa* would also produce the bigrams *\_c* and *a\_* and the trigrams *\_ca* and *sa\_*. This information will be available in a subsequent version of the database.

t5.1 **Table 5** Phonological neighborhood variable names and meanings

t5.2	Variable name	Variable meaning
t5.3	<i>NP</i>	Number of phonological neighbors (all kinds)
t5.4	<i>NPHF</i>	Number of higher frequency phonological neighbors
t5.5	<i>frq_hfp</i>	Frequency of the highest frequency phonological neighbor
t5.6	<i>hfp</i>	Phonological neighbor with the highest frequency
t5.7	<i>hfp_list</i>	List of phonological neighbors in descending frequency, with the place of the word itself marked by "OOOOOO"
t5.8	<i>pf</i>	Number of phonemes/positions with phonological neighbors
t5.9	<i>pf_hf</i>	Number of phonemes/positions with higher frequency phonological neighbors
t5.10	<i>avg_frq_Np</i>	Average frequency of phonological neighbors
t5.11	<i>phon_uniq_point</i>	Phoneme position at which it is no longer like any other word. Set to 0 if greater than the number of phonemes (i.e., it is subsumed by some other word and not unique)
t5.12	<i>homoph</i>	Number of other word entries with the same <i>phon_structure</i>
t5.13	<i>homoph_list</i>	List of homophones in descending frequency

429 **Lemma and part-of-speech processing**

430 While word-form frequencies have tended to dominate anal-  
 431 yses, the lemma and part-of-speech frequencies may also  
 432 influence behavior (Baayen, Dijkstra, & Schreuder, 1997;  
 433 Taft, 1979). To set the values for the lemma and part-of-  
 434 speech properties, we return to the raw\_sequence table.  
 435 Counts were made of every unique combination of word,  
 436 lemma, and part-of-speech tag, rejecting combinations where  
 437 the lemma contains non-Spanish characters or is too long (>  
 438 255 characters). For the written database, there were  
 439 388,270 word-lemma-code types, and for the subtitles  
 440 database, there were 404,394 word-lemma-code types.  
 441 Since there was more than one row per word, these data  
 442 were stored in a separate lemma\_data table for search-  
 443 ing (cf. Brysbaert et al., 2012).

444 For each word, EsPal gives the percentage of occurrences  
 445 with each lemma-code combination. For example, the word  
 446 *caso* most often appears as a common masculine singular noun  
 447 (“case”) but can also appear as a conjunction (*caso de que*  
 448 [“if”]), an adverb (*en todo caso* [“in any case”]), a preposition  
 449 (*en caso de* [“in case of”]), a verb (*yo me caso* [“I marry”]), as  
 450 well as a proper noun and URL. Similarly, for each lemma,  
 451 EsPal gives the percentage of occurrences with each word-  
 452 code combination. For example, the lemma *caso*, besides  
 453 occurring with the previous parts of speech, also occurs with

the masculine plural noun *casos*. The variable *percent\_word* 454  
 gives the percentage of each word (by *\_type* or *\_tok*) that has 455  
 that word-lemma-code, and *percent\_lemma* gives the per- 456  
 centage of each lemma (by *\_type* or *\_tok*) that has that word- 457  
 lemma-code. For example, for the word-lemma-code combi- 458  
 nations with *caso* as the word, *percent\_word\_type* = 16.76 % 459  
 in the written database, since *caso* appears with six different 460  
 lemma-code combinations, and the *percent\_word\_tok* for the 461  
 masculine singular noun lemma-code = 81.5 %, and for the 462  
 simple preposition = 5.6 %. 463

The part-of-speech tags are also expanded to allow search- 464  
 ing and organization of results. The part-of-speech informa- 465  
 tion includes *Category, Type, Degree, Appreciative,* 466  
*Diminutive, Person, Mode, Tense, Form, Gender, Number,* 467  
*Function, Possessor, and Politeness.* A full list for Spanish 468  
 can be found on the FreeLing Web site,<sup>8</sup> which shows, for 469  
 example, how the different attributes of an adjective are 470  
 specified. 471

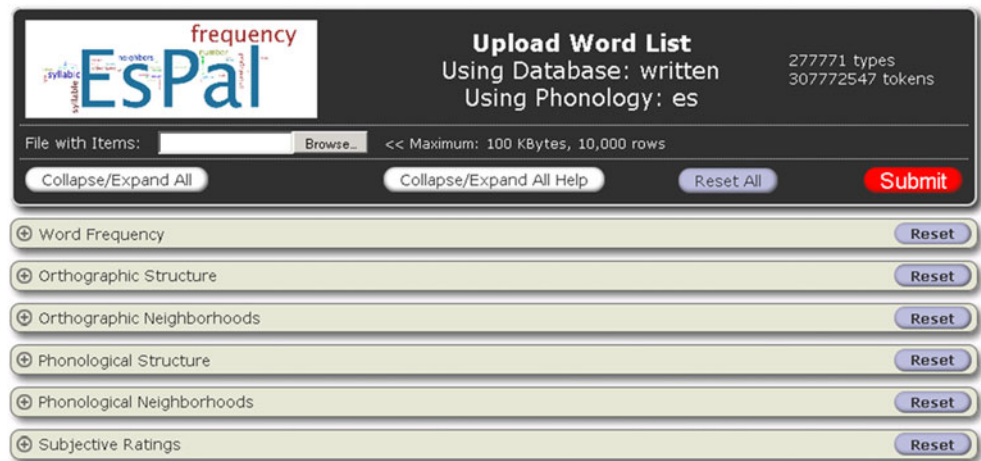
Some of the lemma information is also added to the word\_ 472  
 data table—namely, information about the most common part 473  
 of speech associated with the word (the “maximum lemma”) 474  
 and the “lemma frequency” of the word, which is based on the 475  
 sum of the counts of all the words that have the same lemma as 476  
 any of the lemmas of the word (Keuleers et al., 2010). For the 477  
 maximum lemma of a word, EsPal provides the lemma itself 478  
 (*max\_lem\_lemma*), the detailed part-of-speech code 479  
 (*max\_lem\_code*), and the percentage of all the word’s tokens 480  
 with that code (*max\_lem\_perc*), the category (*max\_lem\_cat*), 481  
 and the percentage as that category (*max\_lem\_cat\_sum\_perc*). 482  
 So for example, in the subtitles database, the word *caso* 483  
 mentioned above appears 90.15 % as a common mas- 484  
 culine singular noun and 90.55 % as a noun overall (the 485  
 additional appearances probably labeled as a proper 486



**Fig. 1** Screenshot of the choice of database and phonology from the EsPal Web site

<sup>8</sup> <http://nlp.lsi.upc.edu/freeling/doc/tagsets/tagset-es.html>.

**Fig. 2** Screenshot of the Word to Properties page where one can upload a list of words to receive the properties of the types shown

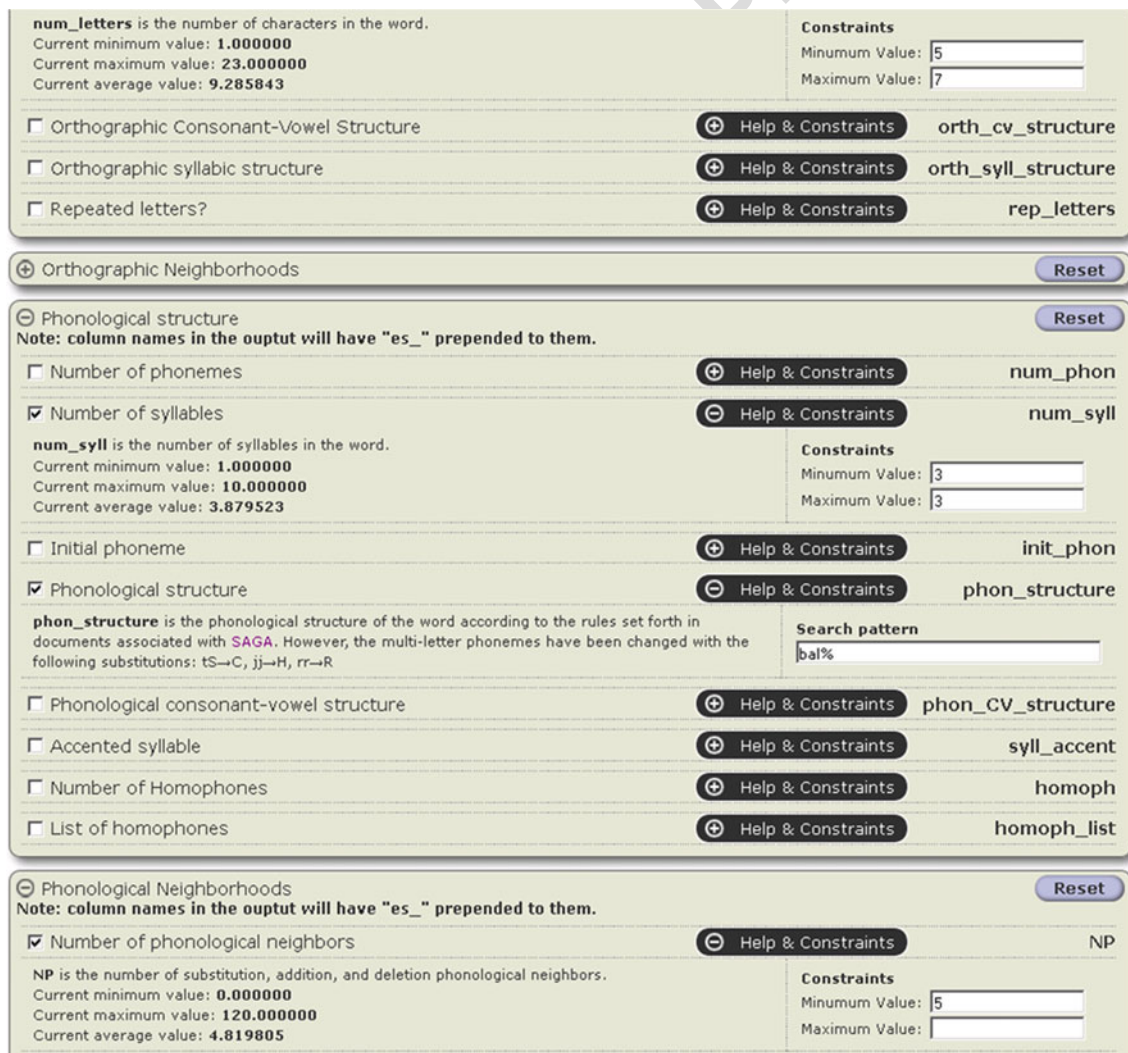


487 noun). For the lemma frequencies, EsPal makes avail-  
 488 able the  $\log(\text{count}+1)$ , as well as the  $\log^2(\text{cnt} + 1)$ ,  
 489 which Keuleers et al. (2010) found helped account for  
 490 more variance in lexical decision times in Dutch.

**Subjective ratings**

491

Subjective ratings, such as the imageability of the thing that a  
 492 word refers to, also modulate the process of lexical access  
 493



**Fig. 3** Screenshot of the Constraints to Words page where a variety of constraints have been applied



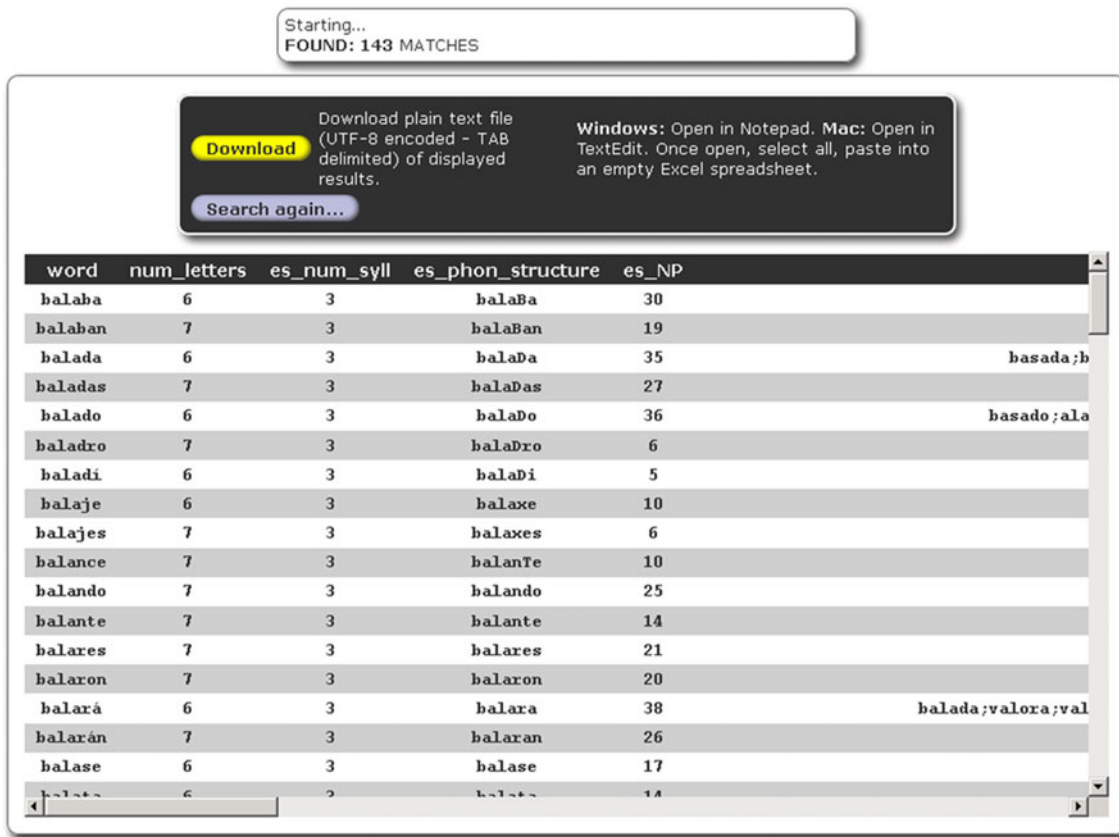


Fig. 4 Screenshot of the results of the query from Fig. 3. All checked properties are returned, which here include *num\_letters*, *es\_num\_syll*, *es\_phon\_structure*, *es\_NP* (number of phonological neighbors), and *es\_hfp\_list* (list of higher frequency phonological neighbors)

494 (Balota, Cortese, Sergent-Marshall, Spieler, & Yap, 2004). For  
 495 EsPal, 6,500 words were selected (mostly nouns and verbs,  
 496 although some nouns could be considered also adjectives). The  
 497 words corresponded to those with the highest frequencies in the  
 498 Alameda and Cuetos (1995) and the Juilland and Chang-  
 499 Rodríguez (1964) word frequency lists. Nouns with gender  
 500 (e.g., *niña*, *niño*) and number (e.g., *corte* and *cortes*) inflections  
 501 were generally both included for evaluation. We decided to  
 502 include both since, in many cases, the different usages with  
 503 the two gender forms or the two number forms hold different  
 504 semantic features. For instance, the word *corte* suggests more  
 505 clearly the action of cutting than does the plural form *cortes*. In  
 506 addition, each form involves different semantic meanings:  
 507 *Cortes* is a term that can be used to refer to the parliament of  
 508 Spain (*Cortes Generales*), while *corte* is linked more to the  
 509 royalty. On the other hand, some nouns could also be consid-  
 510 ered adjectives; for example, the word *rojo* (“red”) can refer to  
 511 the color itself (as well as a Communist) or be used as an  
 512 adjective. Finally, we have included nonreflexive and reflexive  
 513 verbal forms when the two are common, such as *aplicar* [“to  
 514 apply/attach”] and *aplicarse* [“to apply oneself/work hard”],  
 515 because there are important semantic differences between them.

516 From the 6,500 words, we created 130 questionnaires of  
 517 100 words each. This way, each word appeared in a different

position in two questionnaires and was embedded in a different 518  
 context of other words. Then we created three forms for each 519  
 of the 130 questionnaires, so that each word was evaluated on a 520  
 scale of 1–7 for three different values: concreteness, familiarity, 521  
 and imageability. Subjective ratings were obtained in two 522  
 different time windows. The first wave was obtained in 523  
 1998–1999 and corresponds to the data appearing in 524  
 LEXESP (Sebastian-Gallés et al., 2000). The questionnaires 525  
 were answered by undergraduates from 12 different Spanish 526  
 universities, including Universitat Autònoma de Barcelona, 527  
 Universidad Autónoma de Madrid, Universitat de Barcelona, 528  
 Universidad Complutense de Madrid, Universidad de 529  
 Granada, Universidad de Oviedo, Universidad de La Laguna, 530  
 Universitat Rovira i Virgili, Universitat de València, 531  
 Universidad de Santiago de Compostela, Universidad de 532  
 Málaga, and Universidad de Salamanca. Due to the random 533  
 sampling, not all words were equally evaluated, and around 534  
 2,000 words in each dimension did not reach the min- 535  
 imum of 30 responses. In a second wave (taking place 536  
 between 2007 and 2009), an additional set of under- 537  
 graduate students from the Universitat de Barcelona and 538  
 Universidad de La Laguna answered new questionnaires 539  
 so that a minimum of 30 responses for each word were 540  
 finally reached. The data present in EsPal are the 541

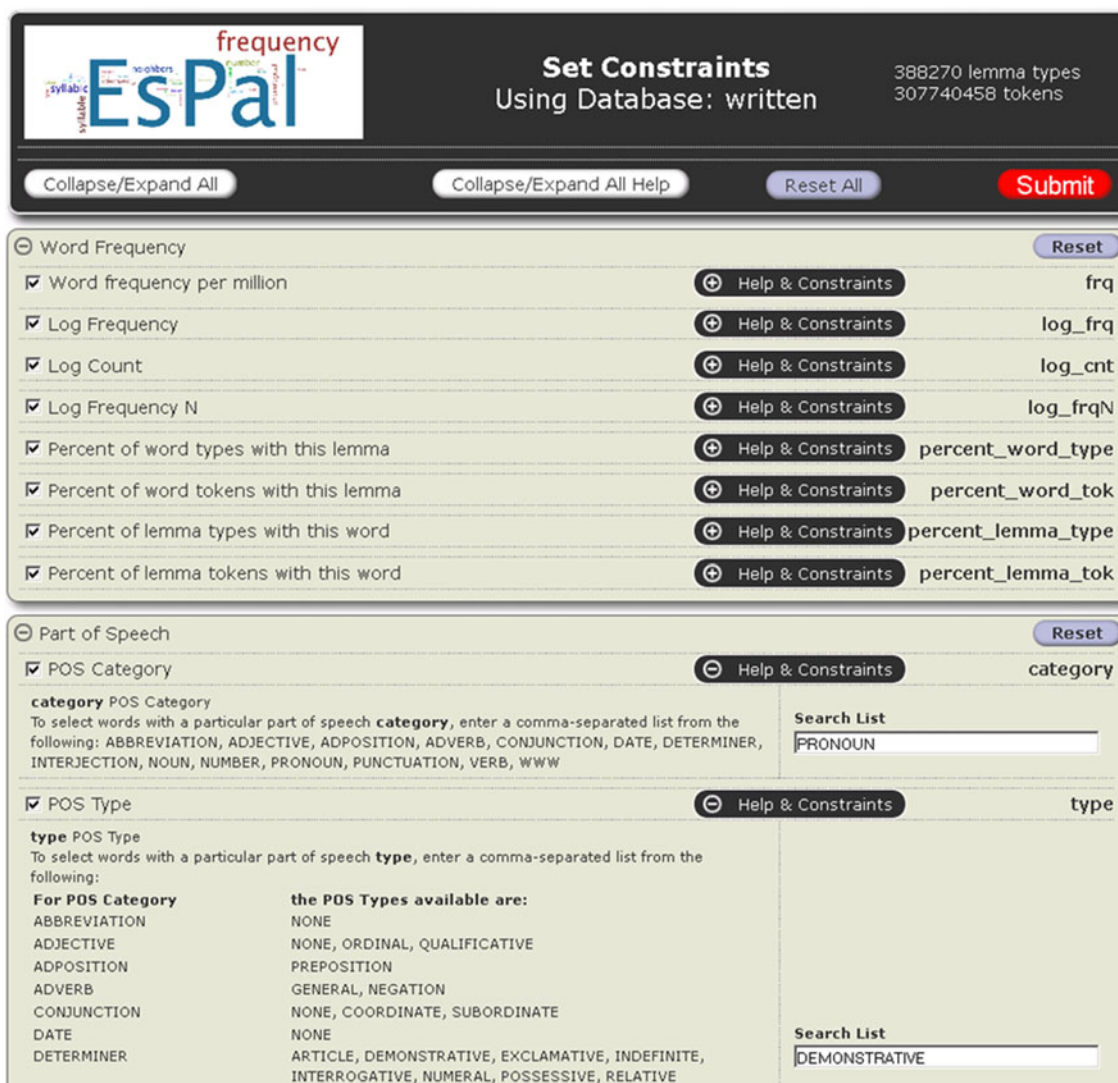


Fig. 5 Screenshot of the POS Constraints to Words page from which frequency properties will be returned for all demonstrative pronouns

542 average ratings for over 6,400 words from at least 30  
 543 participants and from at least 2 universities.

544 **EsPal Web site**

545 The EsPal Web site can be accessed at [http://www.bcbl.eu/  
 546 databases/espal/](http://www.bcbl.eu/databases/espal/). When the user first goes to the Web site

(Fig. 1), the user must first choose a source database and 547  
 phonology via the radio buttons. There are then four ways to 548  
 obtain information from EsPal. The user can upload a file of 549  
 words, one word per line, to receive chosen properties on 550  
 those words, or the user can set constraints on the properties 551  
 to receive a list of words having those constraints. These 552  
 two actions can be performed on data in either the word\_ 553  
 data table or the lemma\_data table. 554

t6.1 **Table 6** Frequency correlations: Correlations of frequency—that is,  $\log(\text{count} + 1)$ —between different corpora (number of common words)

t6.2	EsPal-written	EsPal-subtitles	EsPal-subtitles CDM	LEXESP (B-Pal)	Oral frequency	
t6.3	EsPal-subtitles	.693 (193,757)				
t6.4	EsPal-subtitles CDM	.713 (193,757)	.977 (244,947)			
t6.5	LEXESP (B-Pal)	.855 (30,277)	.794 (28,727)	.799 (28,727)		
t6.6	Oral frequency	.700 (65,388)	.655 (62,271)	.649 (62,271)	.827 (18,723)	
t6.7	SUBTLEX-ESP	.663 (88,303)	.938 (93,949)	.936 (93,949)	.777 (20,316)	.725 (44,374)



555 For example, clicking “Words to Properties” brings the  
 556 user to a Web page (Fig. 2) where he or she can upload a file  
 557 of words, one word per line, and within each of the sub-  
 558 panels, choose which properties to receive. Clicking  
 559 “Submit” brings the user to the results page, which contains  
 560 a table of the results, as well as a button to download a file  
 561 containing the results (returned in the order of the original  
 562 file). Instructions on the page specify how best to convert  
 563 the downloaded file into a spreadsheet program.

564 Clicking “Constraints to Words” from the EsPal homepage  
 565 allows the user to set constraints for returned words. The  
 566 example in Fig. 3 shows how the user might search for words  
 567 with five to seven letters and three syllables that start with the  
 568 phonemes “bal” and have at least five phonological neighbors.  
 569 In the written database, this returns 143 words (Fig. 4).

570 Clicking “Words to Lemma and POS Properties” from  
 571 the EsPal homepage allows the user to receive lemma and  
 572 part-of-speech information for a list of words. Starting with  
 573 “POS Constraints to Words” the user can request, for exam-  
 574 ple, the frequencies of all demonstrative pronouns (Fig. 5).

575 **Index comparisons and validity**

576 While the main purpose of this article is to describe the source  
 577 of the word frequency data and how it has been processed and  
 578 made available, readers may wish to note how it compares to  
 579 other corpora with regard to the psycholinguistic data men-  
 580 tioned in the introduction. We compare the three EsPal  
 581 corpora with three other Spanish data sources: LEXESP  
 582 (Sebastián-Gallés et al., 2000), although in the form of B-Pal  
 583 (Davis & Perea, 2005), which did extensive cleaning of the  
 584 data; SUBTLEX-ESP (Cuetos-Vega et al., 2011), which also  
 585 used subtitles (although from different online sources); and  
 586 oral frequency data from Alonso et al. (2011). Table 6 shows  
 587 the overall frequency correlations between each of these sour-  
 588 ces based on the number of words they have in common,  
 589 although better means may be available for such comparisons  
 590 (Brysbaert & Diependaele, *in press*). As one would expect, the  
 591 written databases (EsPal-Written and LEXESP [B-Pal]) are  
 592 most similar to each other, and the subtitle databases are most  
 593 similar to each other, with the oral frequency data somewhere  
 594 in between.

595 Given that the lexical decision times used in the  
 596 SUBTLEX-ESP paper are not yet available and our own  
 597 are still forthcoming, we provide some basic comparisons  
 598 with two other data sets currently available: word-naming  
 599 times (Cuetos & Barbón, 2006) and picture-naming times  
 600 (Cuetos, Ellis, & Alvarez, 1999), which are shown in  
 601 Tables 7 and 8, respectively, along with word length as an  
 602 added factor in the multiple regression, as previous authors  
 603 have done. Among the EsPal corpora, the subtitles CDM  
 604 database performs best and reinforces previous findings in

**Table 7** Word naming: Regression analysis results using word length and the frequency,  $\log(\text{count} + 1)$ , from different corpora on word naming times (Cuetos & Barbón, 2006) t7.1

Factors	Weights	Adjusted $R^2$	t7.2
LEXESP (BPAL)	-8.464 *	.302	t7.3
Length	10.786 ***	( $N = 240$ )	t7.4
Oral frequency	-8.263 **	.301	t7.5
Length	10.774 ***	( $N = 235$ )	t7.6
SUBTLEX-ESP	-8.353 **	.312	t7.7
Length	10.390 ***	( $N = 239$ )	t7.8
EsPal-written	-5.870 †	.298	t7.9
Length	10.700 ***	( $N = 240$ )	t7.10
EsPal-subtitle tokens	-7.430 *	.304	t7.11
Length	10.604 ***	( $N = 240$ )	t7.12
EsPal-subtitle CDM	-9.074 *	.305	t7.13
Length	10.611 ***	( $N = 240$ )	t7.14

Note. All adjusted  $R^2$  have  $ps < .001$ .

†  $p < .1$   
 \*  $p < .05$   
 \*\*  $p < .01$   
 \*\*\*  $p < .001$

other languages. The EsPal subtitle data sets (both by token and CDM) are very similar to the SUBTLEX-ESP and oral frequency data sets with respect to word-naming times and account for slightly more variance than SUBTLEX-ESP with respect to the picture-naming times. 605 606 607 608 609

**Table 8** Picture naming: Regression analysis results using word length and the frequency,  $\log(\text{count} + 1)$ , from different corpora on picture naming times (Cuetos, Ellis, & Alvarez, 1999) t8.1

Factors	Weights	Adjusted $R^2$	t8.2
LEXESP (BPAL)	-61.187 ***	.161	t8.3
Length	9.634 †	( $N = 139$ )	t8.4
Oral frequency	-65.848 ***	.188	t8.5
Length	7.446	( $N = 137$ )	t8.6
SUBTLEX-ESP	-44.897 ***	.118	t8.7
Length	12.244 *	( $N = 138$ )	t8.8
EsPal-written	-39.378 **	.100	t8.9
Length	11.748 *	( $N = 139$ )	t8.10
EsPal-subtitle tokens	-46.61 ***	.123	t8.11
Length	11.84 *	( $N = 139$ )	t8.12
EsPal-subtitle CDM	-59.008 ***	.133	t8.13
Length	11.050 †	( $N = 139$ )	t8.14

Note. All adjusted  $R^2$ s have  $ps < .001$ .

†  $p < .1$   
 \*  $p < .05$   
 \*\*  $p < .01$   
 \*\*\*  $p < .001$

610 EsPal currently provides the properties of two data sources,  
 611 one written and one based on subtitles, with additional infor-  
 612 mation based on the contextual diversity (by movie) of the  
 613 subtitles data. We provide initial evidence that these data  
 614 sources, the latter especially, are comparable to other corpora  
 615 in Spanish in terms of their frequency data helping to predict  
 616 some psycholinguistic phenomena. We should note, however,  
 617 that there are some limitations that researchers should keep in  
 618 mind when using the data contained in EsPal, especially the  
 619 subtitle data. These data are based on a large number of  
 620 amateur translations of media that are most often English,  
 621 not Spanish, in source, and since proper nouns are typically  
 622 not translated (e.g., “John” is not renamed “Juan”), such terms  
 623 will appear with some frequency. We have used publicly  
 624 available lists of “Spanish words” in order to restrict what is  
 625 inserted into our databases, as well as allow comparison with  
 626 other experimental data. Even so, when using EsPal to gener-  
 627 ate Spanish words for an experiment, one should have a native  
 628 speaker, from the same culture as the subjects, cull out these  
 629 perhaps undesirable elements. Nevertheless, our initial vali-  
 630 dation results suggest that despite what pollution may occur  
 631 because of these foreign words, the frequencies given for the  
 632 “true” Spanish words are useful.

### 633 Conclusion

634 EsPal is a free online application that makes available a wide  
 635 range of frequency, orthographic, phonological, and subjec-  
 636 tive information about Spanish words. EsPal provides an  
 637 extensible, ever-improving, and accurate set of data sources  
 638 and analyses. Initial testing of the current data indicates that  
 639 they are at least comparable to extant sources. This system  
 640 may, therefore, assist the research communities of many dis-  
 641 ciplines to accelerate selection of stimuli for their experiments  
 642 and thereby increase the rate of scientific progress.

643 **Acknowledgments** We would like to thank Daniel Diaz for his  
 644 technical help during the initial phases of the project. Our reviewers  
 645 have been extremely helpful as well. This work was partially funded by  
 646 a grant, HUM2007–30271–E/FILO, from the Spanish Ministry of  
 647 Science and Innovation.  
 648  
 649

### 650 References

- 652 Adelman, J. S., Brown, G. D. A., & Quesada, J. F. (2006). Contextual  
 653 diversity, not word frequency, determines word-naming and lex-  
 654 ical decision times. *Psychological Science*, *17*(9), 814–823.  
 655 Alameda, J., & Cuetos, F. (1995). *Diccionario de frecuencias de las*  
 656 *unidades lingüísticas del español*. Oviedo: Servicio de Publicaciones  
 657 de la Universidad de Oviedo.  
 658 Alonso, M. A., Fernandez, A., & Díez, E. (2011). Oral frequency  
 659 norms for 67,979 Spanish words. *Behavior Research Methods*,  
 660 *43*, 449–458.

- Baayen, R. H., Dijkstra, T., & Schreuder, R. (1997). Singulars and  
 plurals in Dutch: Evidence for a parallel dual-route model. *Journal of Memory and Language*, *37*(1), 94–117.  
 Balota, D. A., Cortese, M. J., Sergent-Marshall, S. D., Spieler, D. H., &  
 Yap, M. J. (2004). Visual word recognition of single-syllable  
 words. *Journal of Experimental Psychology: General; Journal of Experimental Psychology: General*, *133*(2), 283–316.  
 Brysbaert, M., Buchmeier, M., Conrad, M., Jacobs, A. M., Bölte, J., &  
 Böhl, A. (2011). The word frequency effect: A review of recent  
 developments and implications for the choice of frequency esti-  
 mates in German. *Experimental Psychology*, *58*(5), 412.  
 Brysbaert, M., & Diependaele, K. (in press). Dealing with zero word  
 frequencies: a review of the existing rules of thumb and a sug-  
 gestion for an evidence-based choice. *Behavior Research Methods*.  
 Brysbaert, M., & New, B. (2009). Moving beyond Kučera and Francis:  
 A critical evaluation of current word frequency norms and the  
 introduction of a new and improved word frequency measure for  
 American English. *Behavior Research Methods*, *41*(4), 977–990.  
 Brysbaert, M., New, B., & Keuleers, E. (2012). Adding part-of-speech  
 information to the SUBTLEX-US word frequencies. *Behavior Research Methods*, 1–7.  
 Cai, Q., & Brysbaert, M. (2010). SUBTLEX-CH: Chinese word and  
 character frequencies based on film subtitles. *PLoS One*, *5*(6),  
 e10729. doi:10.1371/journal.pone.0010729  
 Carreiras, M., Alvarez, C. J., & de Vega, M. (1993). Syllable frequency  
 and visual word recognition in Spanish. *Journal of Memory and Language*.  
 Carreiras, M., & Perea, M. (2004). Naming pseudowords in Spanish:  
 Effects of syllable frequency. *Brain and Language*, *90*(1–3), 393–  
 400.  
 Carreiras, M., Perea, M., & Grainger, J. (1997). Effects of the ortho-  
 graphic neighborhood in visual word recognition: Cross-task  
 comparisons. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *23*(4), 857–871.  
 Cuetos, F., & Barbón, A. (2006). Word naming in Spanish. *European Journal of Cognitive Psychology*, *18*(03), 415–436.  
 Cuetos, F., Ellis, A. W., & Alvarez, B. (1999). Naming times for the  
 Snodgrass and Vanderwart pictures in Spanish. *Behavior Research Methods*, *31*(4), 650–658.  
 Cuetos-Vega, F., González-Nosti, M., Barbón-Gutiérrez, A., &  
 Brysbaert, M. (2011). SUBTLEX-ESP: Spanish word frequencies  
 based on film subtitles. *Psicología: Revista de Metodología y Psicología Experimental*, *32*(2), 133–143.  
 Davies, M. (2005). The advantage of using relational databases for  
 large corpora: speed, advanced queries, and unlimited annotation. *International Journal of Corpus Linguistics*, *10*(3), 307–334.  
 Davis, C. J. (2005). N-Watch: A program for deriving neighborhood  
 size and other psycholinguistic statistics. *Behavior Research Methods*, *37*(1), 65–70.  
 Davis, C. J., & Perea, M. (2005). BuscaPalabras: A program for  
 deriving orthographic and phonological neighborhood statistics  
 and other psycholinguistic indices in Spanish. *Behavior Research Methods*, *37*(4), 665–671.  
 Davis, C. J., Perea, M., & Acha, J. (2009). Re (de) fining the orthographic  
 neighborhood: The role of addition and deletion neighbors in lexical  
 decision and reading. *Journal of Experimental Psychology: Human Perception and Performance*, *35*(5), 1550–1570.  
 Dimitropoulou, M., Duñabeitia, J. A., Avilés, A., Corral, J., &  
 Carreiras, M. (2010). Subtitle-based word frequencies as the best  
 estimate of reading behavior: The case of Greek. *Frontiers in Psychology*, *1*(218).  
 Duñabeitia, J. A., Cholin, J., Corral, J., Perea, M., & Carreiras, M. (2010). SYLLABARIUM: An online application for deriving  
 complete statistics for Basque and Spanish orthographic syllables. *Behavior Research Methods*, *42*(1), 118–125.

727 Forster, K. I., & Chambers, S. M. (1973). Lexical access and naming time. 767  
 728 *Journal of Verbal Learning and Verbal Behavior*, 12(6), 627–635. 768  
 729 Gathercole, S. E., & Baddeley, A. D. (1990). The role of phonological 769  
 730 memory in vocabulary acquisition: A study of young children 770  
 731 learning new names. *British Journal of Psychology*, 81(4), 439–454. 771  
 732 Grainger, J. (1990). Word frequency and neighborhood frequency 772  
 733 effects in lexical decision and naming. *Journal of Memory and 773*  
 734 *Language*, 29(2), 228–244. 774  
 735 Hernández-Figueroa, Z., Rodríguez-Rodríguez, G., & Carreras- 775  
 736 Riudavets, F. (2009). *Separador de sílabas del español - 776*  
 737 *Silabeador TIP*. Retrieved from <http://tip.dis.ulpgc.es> 777  
 738 Hulme, C., Roodenrys, S., Schweickert, R., Brown, G. D. A., Martin, 778  
 739 S., & Stuart, G. (1997). Word-frequency effects on short-term 779  
 740 memory tasks: evidence for a reintegration process in immediate 780  
 741 serial recall. *Journal of Experimental Psychology: Learning, 781*  
 742 *Memory, and Cognition*, 23(5), 1217. 782  
 743 James, C. T. (1975). The role of semantic information in lexical 783  
 744 decisions. *Journal of Experimental Psychology. Human 784*  
 745 *Perception and Performance*, 1(2), 130–136. 785  
 746 Juilland, A., & Chang-Rodríguez, E. (1964). *Frequency dictionary of 786*  
 747 *Spanish words*. The Hague: Mouton. 787  
 748 Keuleers, E., Brysbaert, M., & New, B. (2010). SUBTLEX-NL: A new 788  
 749 measure for Dutch word frequency based on film subtitles. 789  
 750 *Behavior Research Methods*, 42(3), 643–650. 790  
 751 Marian, V., Bartolotti, J., Chabal, S., & Shook, A. (in press). 791  
 752 CLEARPOND: Cross-Linguistic Easy-Access Resource for 792  
 753 Phonological and Orthographic Neighborhood Densities. *PLoS ONE*. 793  
 754 Miller, B., Juhasz, B. J., & Rayner, K. (2006). The orthographic 794  
 755 uniqueness point and eye movements during reading. *British 795*  
 756 *Journal of Psychology*, 97(2), 191–216. 796  
 757 Monsell, S. (1991). The nature and locus of word frequency effects in 797  
 758 reading. In D. Besner & G. W. Humphreys (Eds.), *Basic processes 798*  
 759 *in reading: Visual word recognition* (pp. 148–197). Hillsdale, NJ: 799  
 760 Lawrence Erlbaum Associates. 800  
 761 Moreno, A., & Mariño, J. B. (1998). Spanish dialects: Phonetic tran- 801  
 762 scription. *Fifth International Conference on Spoken Language 802*  
 763 *Processing (ICSLP '98)* (pp. 189–192). Sydney, Australia. 803  
 764 New, B., Brysbaert, M., Veronis, J., & Pallier, C. (2007). The use of film 804  
 765 subtitles to estimate word frequencies. *Applied PsychoLinguistics*, 805  
 766 28(4), 661. 806  
 767 New, B., Pallier, C., Brysbaert, M., & Ferrand, L. (2004). Lexique 2: A new 767  
 768 French lexical database. *Behavior Research Methods*, 36(3), 516–524. 768  
 769 Nogueiras, A., & Mariño, J. (2009). *SAGA: Transcriptor fonético de 769*  
 770 *las variedades dialectales del español*. Retrieved from [http://](http://www.talp.upc.edu/index.php/technology/tools/signal-processing-tools/81-saga) 770  
 771 [www.talp.upc.edu/index.php/technology/tools/signal-processing-](http://www.talp.upc.edu/index.php/technology/tools/signal-processing-tools/81-saga) 771  
 772 [tools/81-saga](http://www.talp.upc.edu/index.php/technology/tools/signal-processing-tools/81-saga) 772  
 773 Padró, L., Collado, M., Reese, S., Lloberes, M., & Castellón, I. (2010). 773  
 774 Freeling 2.1: Five years of open-source language processing tools. 774  
 775 *Proceedings of 7th Language Resources and Evaluation 775*  
 776 *Conference*. La Valletta, Malta. 776  
 777 Perea, M., & Carreiras, M. (1998). Effects of syllable frequency and 777  
 778 syllable neighborhood frequency in visual word recognition. 778  
 779 *Journal of Experimental Psychology. Human Perception and 779*  
 780 *Performance*, 24(1), 134–144. 780  
 781 Perea, M., Soares, A. P., & Comesaña, M. (in press). Contextual 781  
 782 diversity is a main determinant of word-identification times in 782  
 783 young readers. *Journal of Experimental Child Psychology*. 783  
 784 Perea, M., Urkia, M., Davis, C. J., Agirre, A., Laseka, E., & Carreiras, 784  
 785 M. (2006). E-Hitz: A word frequency list and a program for 785  
 786 deriving psycholinguistic statistics in an agglutinative language 786  
 787 (Basque). *Behavior Research Methods*, 38(4), 610–615. 787  
 788 Rayner, K. (2009). Eye movements and attention in reading, scene 788  
 789 perception, and visual search. *The Quarterly Journal of 789*  
 790 *Experimental Psychology*, 62(8), 1457–1506. 790  
 791 Sebastián-Gallés, N., Martí, M., Carreiras, M., & Cuetos, F. (2000). 791  
 792 *LEXESP: Léxico Informatizado del Español*. Barcelona: Universitat 792  
 793 de Barcelona. 793  
 794 Shelton, M., Gerfen, C., & Gutiérrez-Palma, N. (2011). The interaction 794  
 795 of subsyllabic encoding and stress assignment: A new examina- 795  
 796 tion of an old problem in Spanish. *Language & Cognitive 796*  
 797 *Processes*, 27(10), 1459–1478. 797  
 798 Taft, M. (1979). Recognition of affixed words and the word frequency 798  
 799 effect. *Memory & Cognition*, 7(4), 263–272. 799  
 800 Taulé, M., Martí, M. A., & Recasens, M. (2008). Ancora: Multilevel 800  
 801 annotated corpora for catalan and spanish. *Proceedings of the 6th 801*  
 802 *International Conference on Language Resources and Evaluation 802*  
 803 *(LREC-2008)*. 803  
 804 Yarkoni, T., Balota, D., & Yap, M. (2008). Moving beyond Coltheart's 804  
 805 N: A new measure of orthographic similarity. *Psychonomic 805*  
 806 *Bulletin & Review*, 15(5), 971–979. 806

## AUTHOR QUERIES

### **AUTHOR PLEASE ANSWER ALL QUERIES.**

- Q1. Keywords are desired. Please provide.
- Q2. Please provide complete bibliographic details of the following reference Brysbaert & Diependaele (in press); Marian et al. (in press); Perea et al. (in press).
- Q3. Please provide the VolumeNumber of this reference item.
- Q4. Please provide the VolumeNumber and PageNumber of this reference item.

UNCORRECTED PROOF