

# Brazilian Portuguese Lexicon Manual - LexPorBR: Psycholinguistic Corpus

*Version: Alpha2*

Lyon, 6th October 2015.

## Introduction

The main objective of the Brazilian Portuguese Lexicon - LexPorBR<sup>1</sup> is to deliver a psycholinguistic corpus of the Brazilian Portuguese (BP) language that provides maximum metalinguistic and psycholinguistic information about BP words. The Brazilian Portuguese Lexicon is a free and open corpus consulted in a simple and dynamic interface on the internet. The results are presented in an organized and hierarchical way, containing metalinguistic and psycholinguistic information about the words or groups of word searched.

## Psycholinguistic Corpora

Psycholinguistic corpora are used 1) to control, selection, and manipulation of words and specific criteria in psycholinguistic experiments and 2) in computational linguistic analysis on distribution and lexical behavior (Baayen, 2001). Examples of psycholinguistic corpora are: French – *Lexique*<sup>2</sup> (B. New, Pallier, Brysbaert, & Ferrand, 2004; B. New, Pallier, Ferrand, & Matos, 2001), Spanish – *BuscaPalabras* (Davis & Perea, 2005), English – MRC<sup>3</sup> (Coltheart, 1981), Dutch, English, French, German, and Spanish - ClearPOND<sup>4</sup> (Marian, Bartolotti, Chabal, & Shook, 2012), Cyrillic, Dutch, English, and German - CELEX<sup>5</sup> (Baayen, Piepenbrock, & van Rijn, 1995). For example, these corpora have been used in megastudies investigating the psycholinguistic behavior on word and pseudoword recognition in the

---

<sup>1</sup> <http://www.lexicodoportugues.com/>

<sup>2</sup> <http://www.lexique.org/>

<sup>3</sup> <http://www.psych.rl.ac.uk/>

<sup>4</sup> <http://clearpond.northwestern.edu/>

<sup>5</sup> <http://celex.mpi.nl/>

*English Lexicon Project* (Balota et al., 2007; Boris New, Ferrand, Pallier, & Brysbaert, 2006), *French Lexicon Project* (Ferrand et al., 2010), *Dutch Lexicon Project* (Keuleers, Diependaele, & Brysbaert, 2010), and *British Lexicon Project* (Keuleers, Lacey, Rastle, & Brysbaert, 2012). These corpora are used in the control, selection, and manipulation of words and its characteristics to create psycholinguistic experiments in a number of studies and specific researches (Gimenes & New, 2015), as well as in the development and simulation of linguistic models and computational linguistic analysis (Schreuder & Baayen, 1995).

## Brazilian Portuguese Lexicon

The Brazilian Portuguese Lexicon was born from an idea noted on a post-it in late 2012 when I was beginning my Ph.D in psycholinguistics and neurosciences, in Lyon, France. My Ph.D project investigates the morphological representation and processing of verbal inflection in BP, French, and bilinguals with BP as first language and French as second language. For the French experiments, the stimuli were selected in the *Lexique* corpus (B. New et al., 2004), which offers a range of crucial information for the experiments and result analysis (surface frequency, grammatical category, number of letters, number of neighbors, inverted form, CVCV structure, etc.). In early 2013, when we started to prepare the BP experiments, we faced the complete lack of a BP word-based psycholinguistic corpus. Looking for a corpus that feed our needs, we found the *Linguateca*<sup>6</sup> website (Santos & Bick, 2000) which comprises several European and Brazilian Portuguese corpora. However, we could not find any BP word-based corpus with metalinguistic and psycholinguistic information appropriated for the strict creation of psycholinguistic experiments in BP. It was when I wrote on a post-it “make the Brazilian Portuguese Lexicon”. Currently, the Brazilian Portuguese Lexicon presents the homepage as shown in **Figure 1**.

---

<sup>6</sup> <http://www.linguateca.pt/>

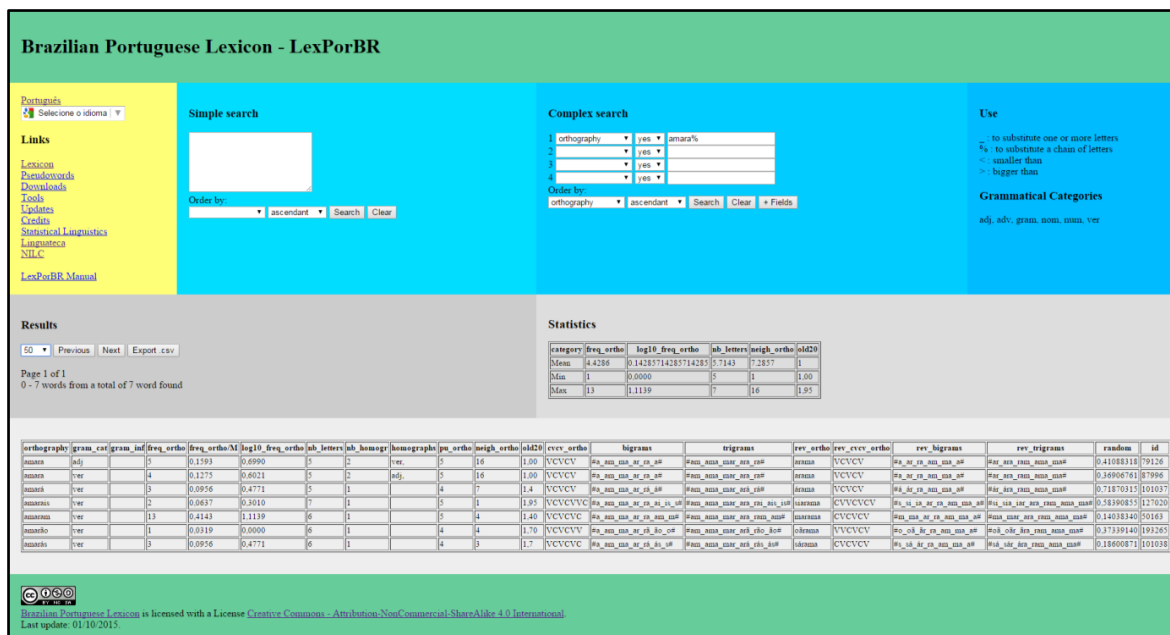


Figure 1: Homepage of the Brazilian Portuguese Lexicon.

## Development

In early 2014, we started the development of the Brazilian Portuguese Lexicon in four stages: 1) construction of the corpus with words, metalinguistic, and psycholinguistic information, 2) construction of the webpage in HTML, 3) importation of the corpus to a MySQL database on the internet, and 4) programming in PHP of the operation of the Brazilian Portuguese Lexicon. Additional pages were created: updates, credits, downloads, tools, and links. Then, we developed the BP pseudoword generator engine and the statistical linguistics tools.

## Updates

15/01/2013 - looking for a BP psycholinguistic corpus, we found the *Linguateca* website which concentrates a number of Portuguese corpora, but no BP psycholinguistic corpus. We decided to make the Brazilian Portuguese Lexicon as a BP word-based metalinguistic and

psycholinguistic corpus with free and open access on the internet. Skills to be developed: R<sup>7</sup>, HTML<sup>8</sup>, MySQL<sup>9</sup>, PHP<sup>10</sup>, Java<sup>11</sup>, and CSS<sup>12</sup>.

03/21/2013 - pre-selection in *Linguateca* website of two BP corpora: 1) Corpus Brasileiro<sup>13</sup> (1 billion words, 3.2 GB) and 2) *Núcleo Interdisciplinar de Linguística Computacional de São Carlos* (henceforth NILC)<sup>14,15</sup> (32 million words, 49 MB). After discussion with the researchers responsible for these corpora, we concluded the NILC would be the best corpus for the development of the Brazilian Portuguese Lexicon, following the criteria: 1) number of words (32 million) consistent with other psycholinguistic corpora (*Lexique*, CELEX, ClearPOND) (Brysbaert & New, 2009), 2) number and size of files (13 files, total size 49 MB), 3) organization of the corpus in .txt files separated by grammatical categories, 4) organization of files in two columns (orthography and frequency) and delimited by tabulation, and 5) publications and resources already developed by the NILC.

08/14/2013 – development of a pilot corpus with only the verbs from NILC, accounting about 80,000 forms. We used the R program for the development of 10 columns of information: 1) orthography, 2) orthographic frequency, 3) frequency per million of words, 4) log10 orthographic frequency, 5) number of letters, 6) grammatical category, 7) grammatical information, 8) reversed orthographic form, 9) CVCV structure, and 10) CVCV reversed structure. Construction of the pilot webpage in a localhost with the XAMPP<sup>16</sup> program containing the Apache, MySQL, PHP, and Perls preinstalled modules. Configuration of the phpMyAdmin<sup>17</sup> for the importation of the pilot corpus saved in .csv format to a MySQL database. We used the Notepad++<sup>18</sup> software to program the HTML/PHP interface page between the user and MySQL database.

10/28/2013 – pilot version of the Brazilian Portuguese Lexicon with two search engines: 1) simple search and 2) complex search. The simple search consists of a text area where users

---

<sup>7</sup> <http://www.r-project.org/>

<sup>8</sup> <http://en.wikipedia.org/wiki/HTML>

<sup>9</sup> <http://www.mysql.com/>

<sup>10</sup> <http://www.php.net/>

<sup>11</sup> <http://www.java.com/>

<sup>12</sup> [http://en.wikipedia.org/wiki/Cascading\\_Style\\_Sheets](http://en.wikipedia.org/wiki/Cascading_Style_Sheets)

<sup>13</sup> <http://corpusbrasileiro.pucsp.br/cb/Inicial.html>

<sup>14</sup> <http://www.nilc.icmc.usp.br/nilc/index.php>

<sup>15</sup> <http://www.linguateca.pt/aceso/corpus.php?corpus=SAOCARLOS>

<sup>16</sup> <http://www.apachefriends.org/>

<sup>17</sup> [http://www.phpmyadmin.net/home\\_page/index.php](http://www.phpmyadmin.net/home_page/index.php)

<sup>18</sup> <http://notepad-plus-plus.org/>

can enter multiple words as a list. The complex search consists of four fields for criteria insertion about the words to be searched. Each search engine has a “Search” button to start the search and present the results, and a “Clear” button to delete the current data in the fields. Definition of the Brazilian Portuguese Lexicon pages: Lexicon, Pseudowords, Downloads, Tools, Updates, Credits, Statistical Linguistics, Linguateca, and NILC.

11/30/2013 - programming of Java and PHP algorithms to keep the data in the HTML fields after search. We inserted two fields for results organization, one with the arrangement criterion and another for ascendant or descendant order. We inserted the “+ Fields” button in the complex search engine for the use of eight fields. Choice of the <http://www.biz.nf/> free web server host following the criteria: 1) space of 250 MB, 2) MySQL 5 database, 3) PHP 4/5 support, 4) 5,000 MB transfer, 5) free hosting, 6) free domain <http://portugueselexicon.co.nf>, 7) POP3/SMTP webmail, and 8) FTP control transfer. Importation of the pilot corpus in .csv format to a MySQL database and transfer of all other pages created in HTML/PHP to <http://portugueselexicon.co.nf/> by FTP with the FileZilla<sup>19</sup> program.

12/12/2013 - given the MySQL recognizes the symbols underline “\_” to replace a letter and percentage “%” to replace a chain of letters, this information was added to the tips on the homepage. Programming in PHP for the recognition of the symbols: greater than “>” and less than “<” for numerical search. The Brazilian Portuguese Lexicon development was divided in three versions: 1) Alpha, 2) Beta and 3) Delta. The Alpha version (2014) is the first version of the Brazilian Portuguese Lexicon, providing a pure orthographic corpus. The Beta version (scheduled for 2015) will provide phonological, syllabic, and lemma data. Finally, the Delta version (scheduled for 2016) will provide specific features such as: morphological information, syntactic information, age of acquisition, among others.

01/07/2014 - download on the *Linguateca* website of the 13 files of the NILC<sup>20</sup> in .txt format separated by grammatical categories (6 form files: adjectives, adverbs, grammatical, nouns, numerals, and verbs; 7 lemma files: adjectives, adverbs, grammatical, nouns, proper names, numerals, and verbs). We compared of the total number of words and forms to the data provided in *Linguateca*. Column with the grammatical categories (gram\_cat) for each file (adjectives, adverbs, grammatical, nouns, proper names, numerals, and verbs). Column with the word type (form or lemma). Transformation of all words in lowercases and sum of

---

<sup>19</sup> <https://filezilla-project.org/>

<sup>20</sup> <http://www.linguateca.pt/acesso/contabilizacao.php>

repeated forms. Column with an identification number (id) of the word, according to the organization of the corpus by frequency in descending order and an alphabetic a-z order, so the identification number (id) also becomes the word position in the lexicon and corpus.

01/08/2014 – column with the orthographic frequency per million of words (ortho\_freq/M) by calculating  $[1000000 * \text{ortho\_freq} / \text{total\_freq}]$ . Column with the natural log of the orthographic frequency. Column with the number of letters of the form. Exclusion of all forms with more than 30 letters and numerals, but 0-1 and 1st-9th. Column with the number of homographic forms. Column with the grammatical categories of the homographic forms.

01/10/2014 – splitting of words into letters, processing of vowels “V” and consonants “C”, column with the CVCV structure (CVCV\_ortho). The letters were also classified as punctuation “P”, numbers “N”, accents “A”, and symbols “S”. Column with word bigrams. Column with word trigrams.

18/01/2014 - development of an algorithm to calculate the orthographic uniqueness point (pu\_ortho) and creation of a column for it. Column with the number of orthographic neighbors (ortho\_neigh) (Coltheart’s N) (Coltheart, Davelaar, Jonasson, & Besner, 1977) and column with the Orthographic Levenshtein Distance for the 20 closest words (old20) (Yarkoni, Balota, & Yap, 2008). These functions are available in the R package “vwr”<sup>21</sup> developed by Emmanuel Keuleers<sup>22</sup>. Column with a random number between 0 and 1 with eight digits of precision.

01/28/2014 - creation of four columns with the reversed forms from: orthography (rev\_ortho), CVCV\_ortho (rev\_CVCV\_ortho), bigrams (rev\_bigrams), and trigrams (rev\_trigrams). The Brazilian Portuguese Lexicon – Alpha has 21 columns of metalinguistic and psycholinguistic information: 1) orthography, 2) gram\_cat, 3) gram\_inf, 4) ortho\_freq, 5) ortho\_freq/M, 6) log10\_ortho\_freq, 7) nb\_letters, 8) nb\_homogr, 9) homographs, 10) pu\_ortho, 11) ortho\_neigh, 12) old20, 13) CVCV\_ortho, 14) bigrams, 15) trigrams, 16) rev\_ortho, 17) rev\_CVCV\_ortho, 18) rev\_bigrams, 19) rev\_trigrams, 20) random, and 21) id.

02/05/2014 - the word-based psycholinguistic corpus Brazilian Portuguese Lexicon has 21 columns of information and 215,175 rows of BP words. This table has a size of 45 MB in .csv format. This file was split into 36 files of about 1.5 MB in .csv format. Then the .csv files

<sup>21</sup> <http://cran.r-project.org/web/packages/vwr/index.html>

<sup>22</sup> <http://crr.ugent.be/members/emmanuel-keuleers>

were saved in UTF-8 coding (they had ANSI coding) in order to avoid problems with accents, symbols, and special characters. Each file was imported via phpMyAdmin to our internet server host in a way that each new imported file inflated the existing one.

02/12/2014 - development of a module for limitation and navigation of results. The number of words to be presented may be 50, 100, 200, or 500. It was included two buttons (“Previous” and “Next”) to navigate between the result pages. Presentation of four general information about the search: 1) total number of words found, 2) total number of pages, 3) range of words presented, and 4) page displayed. Inclusion of the button "Export .csv" to export all the results of the current search to a downloadable .csv file.

02/18/2014 - development of a basic statistical module from the current search (Davis, 2005), presenting the following information: 1) mean, 2) maximum, and 3) minimum values of the following categories: 1) ortho\_freq, 2) log10\_ortho\_freq, 3) nb\_letters, 4) ortho\_neigh, and 5) old20.

02/23/2014 - development of a BP pseudoword generator engine from bigrams or trigrams. Unlike other pseudoword generators (Keuleers & Brysbaert, 2010; Mota & Resende, 2013), we accounted the overall frequencies bigrams and trigrams, the frequency of bigrams and trigrams according to their position in word and grammatical category. We obtained two tables as databases to generate BP pseudowords, one with the bigrams and another with the trigrams, with 1) general, 2) by position in the word, and 3) by grammatical category.

03/05/2014 – we created a BP pseudoword generator engine where the user must enter four fields: 1) number of letters of the pseudowords to be generated, 2) number of pseudowords to be generated, 3) grammatical category that these pseudowords should belong (all, adj, adv, gram, nom, num, ver), and 4) criterion for the construction of the pseudowords (bigrams or trigrams). The BP pseudoword generator engine builds pseudowords simultaneously in both directions, from left to right and from right to left, beginning with a bigram or trigram like “#xx” or #xx”. Following to the number of letters, the engine concatenates new bigrams or trigrams which share as much orthographic information as possible with the previous bigram or trigram (1 letter for bigrams and 2 letters for trigrams).

03/18/2014 - insertion of four columns in pseudoword results: 1) pseudoword grammatical category, 2) pseudoword frequency calculated as the sum of the bigrams or trigrams

frequencies that build the pseudoword, 3) log<sub>10</sub> of the pseudoword calculated frequency, and 4) number of letters of the pseudoword. We translated the homepage with the simple and complex searches to English.

03/25/2014 - registration of a specific internet domain for the Brazilian Portuguese Lexicon ([www.lexicodoportuguês.com](http://www.lexicodoportuguês.com)) in HostGator<sup>23</sup> and redirection of this domain to the server where the Brazilian Portuguese Lexicon is hosted <http://portugueselexicon.co.nf>. Inauguration of the Brazilian Portuguese Lexicon on March 25th, 2014.

03/21/2015 - development of the Statistical Linguistics page with various open and free tools and resources in HTML/PHP to linguistic and statistical analysis: F1, F2, F', minF', Hartley test, normalization between 0 and 1, reverse word, Hamming distance, Levenshtein distance, orthographic neighbors (Coltheart's N), average Levenshtein distance, relative entropy, word frequency, Zipf's distribution, etc.

09/20/2015 – writing and availability of the *Manual do Léxico do Português Brasileiro - Alfa 2* in Brazilian Portuguese and writing and availability of the Brazilian Portuguese Lexicon - Alpha 2 Manual in English. We translated and implemented to English all pages of the Brazilian Portuguese Lexicon. We implemented the Google Translate gadget in all pages of the Brazilian Portuguese Lexicon for its translation to the various languages available on Google Translate. It is suggested to use the site in English in a way that the results presented in Portuguese will not be translated.

## Alpha Version

Considering the large amount of metalinguistic and psycholinguistic information that can and will be computed, implemented, and made available in the Brazilian Portuguese Lexicon, its development was divided in three versions: 1) Alpha (2014), 2) Beta (2015), and 3) Delta (2016). Currently, the Brazilian Portuguese Lexicon is in the Alpha version, inaugurated on March 25<sup>th</sup>, 2014. The main characteristic of the Brazilian Portuguese Lexicon Alpha version is that it provides a pure orthographic corpus with all information provided being strictly computed from orthographic data of the BP words. Further, the Beta version will provide

---

<sup>23</sup> <http://hostgator.com/>



information about: 1) phonology, 2) syllables, and 3) lemma association. Delta version will provide: 1) morphological information, 2) syntactic information, 3) allomorphy information, and as far as possible 4) Reaction Time measures from the recognition of a large number of BP words and pseudowords, following the Lexicon Projects (Balota et al., 2007; Ferrand et al., 2010; Keuleers et al., 2010, 2012).

## Creative Commons License



The Brazilian Portuguese Lexicon from Gustavo Lopez Estivalet<sup>24</sup> is licensed with a License Creative Commons - Attribution-NonCommercial-ShareAlike 4.0 International<sup>25</sup>. Based on the work available in <http://www.linguateca.pt/aceso/corpus.php?corpus=SAOCARLOS>. Permissions beyond the scope of this license may be available at [http://www.lexicodoportugues.com/credits\\_en.php](http://www.lexicodoportugues.com/credits_en.php).

## NILC/São Carlos and Linguateca

The Brazilian Portuguese Lexicon was developed from the *Núcleo Interinstitucional de Linguística Computacional de São Carlos* corpus (NILC) (Inter-Institutional Computational Linguistics Center) (Pinheiro & Aluísio, 2003) located at the *Instituto de Ciências Matemáticas e de Computação de São Carlos* (Institut of Mathematics and Computation) (ICMC/São Carlos)<sup>26</sup>, at the University of São Paulo in São Carlos (USP/São Carlos)<sup>27</sup>. The lists of forms and lemmas divided in grammatical categories were downloaded from the *Linguateca* website, which also provides a series of information about the NILC, as quantitative and statistical data<sup>28</sup>, corpus origin<sup>29</sup>, and especially, the forms<sup>30</sup> and the lemmas<sup>31</sup> files in .txt format, separated by grammatical categories.

<sup>24</sup> [http://www.researchgate.net/profile/Gustavo\\_Estivalet](http://www.researchgate.net/profile/Gustavo_Estivalet)

<sup>25</sup> <http://creativecommons.org/licenses/by-nc-sa/4.0/>

<sup>26</sup> <http://www.icmc.usp.br/Portal/>

<sup>27</sup> <http://www.saocarlos.usp.br/>

<sup>28</sup> [http://www.linguateca.pt/aceso/desc\\_corpus.php?corpus=SAOCARLOS](http://www.linguateca.pt/aceso/desc_corpus.php?corpus=SAOCARLOS)

<sup>29</sup> <http://www.linguateca.pt/aceso/NILCsaocarlos.html>

<sup>30</sup> <http://www.linguateca.pt/aceso/contabilizacao.php#listaPosSAOCARLOS>

<sup>31</sup> <http://www.linguateca.pt/aceso/contabilizacao.php#listaLemasSAOCARLOS>

“All the material that we provide is not restricted to any group and was authorized (under the provided terms) by the respective authors or copyright holders. From resource to resource the conditions are different and are specified in its specific documentation. The tools created by the *Linguateca* are available under the GNU<sup>32</sup>.”<sup>33</sup>

## Lexique

The Brazilian Portuguese Lexicon was inspired by the French psycholinguistic corpus *Lexique* (B. New et al., 2004, 2001). The *Lexique* has already offered data on French words to a number of studies and researches, being a great example of psycholinguistic corpus. This corpus exemplifies the features and utilities that a psycholinguistic corpus should and can offer as resources for research in psycholinguistics and computational linguistics. A detailed description of the *Lexique* can be found in the user’s manual<sup>34</sup>.

## R Program and Packages

The Brazilian Portuguese Lexicon was developed with the R program with the original data imported from .txt files and each new column being created and computed through certain functions and algorithms. The number of orthographic neighbors (Coltheart’s N) (Coltheart et al., 1977) and the Orthographic Levenshtein Distance from de 20 closest words (OLD20) (Yarkoni et al., 2008) were calculated with the functions “colheart.N” and “old20” available in the “vwr” R package developed by Emmanuel Keuleers. Functions from the “languageR”<sup>35</sup> R package developed by Harald Baayen<sup>36</sup> were also used in the development of the Brazilian Portuguese Lexicon.

---

<sup>32</sup> <http://www.gnu.org/copyleft/gpl.html>

<sup>33</sup> <http://www.linguateca.pt/FAQ/#faq1.8>

<sup>34</sup> <http://www.lexique.org/docLexique.php>

<sup>35</sup> <https://cran.r-project.org/web/packages/languageR/index.html>

<sup>36</sup> <http://www.sfs.uni-tuebingen.de/~hbaayen/>

## LexPorBR - Alpha

### Conventions

For the use of the Brazilian Portuguese Lexicon, some conventions were determined to perform the searches and understand the results.

- Grammatical categories: **adj** - adjective, **adv** – adverb, **gram** - grammatical, **nom** - noun, **num** - numeral, **prop** – proper name, **ver** - verb.
- The CVCV structures have the letters: **V** - vowel, **C** - consonant, **P** - punctuation, **N** - number, **A** - accentuation, **S** - symbol.
- The wildcard used are: “<” greater than, “>” less than, “\_” replaces a letter, “%” replaces a chain of letters.
- Order of the results presentation: **ascendant** - presents the results in ascending order, **descendant** - presents the results in descending order.
- Buttons: **Search** - performs the search and displays the results, **Clear** - clears all data from the fields, + **Fields** - redirects the user to a page with more fields for complex research.
- Select yes/no: **yes** - considers the criterion, **not** – disregards the criterion.

### Columns

The Brazilian Portuguese Lexicon - Alpha has 215,175 rows with different lexical entries and 21 columns with different metalinguistic and psycholinguistic information. Thus, each row of the Brazilian Portuguese Lexicon contains a word and each column contains an information about this word.

A search example from the complex research with **cat\_gram=yes-ver** and **nb\_letters=yes->3<8** can be seen in **Figure 2**. This search presents all the words that have the grammatical category defined as “verb”, less than 8 letters, and more than 3 letters.

orthography	gram_cat	gram_inf	freq_ortho	freq_ortho/Al	log10_freq_ortho	nb_letters	nb_homogr	homographs	pu_ortho	height_ortho	old20	cvcc_ortho	bigrams	trigrams	rev_ortho	rev_cvcc_ortho	rev_bigrams	rev_trigrams	random	id	
case	ver	37	1.1792	1.5682	1.00	CVCV	4	35	1	1.00	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.67386608	28665	
cast	ver	4	0.1275	0.6021	1.00	CVCV	4	17	1	1.00	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.94087864	89116	
cabo	ver	84	2.6771	1.9243	1.00	CVCV	4	32	1	1.00	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.66035062	17673	
caha	ver	8	0.2550	0.9031	1.00	CVCV	4	55	2	1.00	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.65841162	63940	
cato	ver	5	0.1593	0.6990	1.00	CVCV	4	41	3	1.00	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.64417913	79982	
cava	ver	8	0.2550	0.9031	1.00	CVCV	4	47	3	1.00	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.78947595	63950	
cave	ver	2	0.0637	0.3010	1.00	CVCV	4	32	1	1.00	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.01989358	130366	
cear	ver	8	0.2550	0.9031	1.00	CVCV	4	16	1	1.00	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.40265500	63956	
ceci	ver	1	0.0319	0.0000	1.00	CVCV	4	5	1	1.70	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.42591075	204072	
ceda	ver	17	0.5418	1.2304	1.00	CVCV	4	26	1	1.00	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.93177629	44012	
cede	ver	167	5.3223	2.2227	1.00	CVCV	4	18	3	1.00	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.63774699	11590	
cedi	ver	10	0.3187	1.0000	1.00	CVCV	4	11	1	1.35	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.44488334	47412	
cedo	ver	23	0.7330	1.6617	1.00	CVCV	4	19	1	1.00	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.77500498	75522	
ceei	ver	4	0.1275	0.6021	1.00	CVCV	4	2	1	1.85	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.89151235	89165	
cega	ver	13	0.4143	1.1139	1.00	CVCV	4	23	1	1.00	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.48261355	50413	
cego	ver	10	0.3187	1.0000	1.00	CVCV	4	3	1	1.00	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.74791504	57414	
ceia	ver	33	1.0517	1.5185	1.00	CVCV	4	27	1	1.00	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.63512325	40389	
ceio	ver	1	0.0319	0.0000	1.00	CVCV	4	29	1	1.00	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.13375229	204161	
ceir	ver	1	0.0319	0.0000	1.00	CVCV	4	19	1	1.00	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.64279096	204335	
cf	ver	20	0.6374	1.3010	1.00	CVCV	4	2	2	1.75	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.04254472	40356	
chai	ver	1	0.0319	0.0000	1.00	CVCV	4	4	1	1.55	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.05094832	204509	
chia	ver	6	0.1912	0.7782	1.00	CVCV	4	12	1	1.00	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.28125504	73356	
cia	ver	68	2.1671	1.8325	1.00	CVCV	4	30	1	1.00	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.53433187	20134	
ciara	ver	11	0.3506	1.0414	1.00	CVCV	4	14	1	1.00	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.81248379	54845	
cias	ver	2	0.0637	0.3010	1.00	CVCV	4	17	1	1.00	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.28997654	130722	
ciãd	ver	1	0.0319	0	1.00	CVCV	4	9	1	1.55	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.41265629	205189	
cie	ver	5	0.1593	0.6990	1.00	CVCV	4	18	3	1.00	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.12407612	80063	
ciui	ver	1	0.0319	0	1.00	CVCV	4	5	1	1.7	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.16318023	205378	
ciua	ver	618	19.6955	2.7910	1.00	CVCV	4	30	3	1.00	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.45828409	4245	
cite	ver	11	0.3506	1.0414	1.00	CVCV	4	14	1	1.05	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.23306359	44857	
ciro	ver	55	1.7528	1.7404	1.00	CVCV	4	34	1	1.00	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.81889330	22891	
ciun	ver	2	0.0637	0.3010	1.00	CVCV	4	4	1	1.80	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.99817464	130851	
ciun	ver	1	0.0319	0	1.00	CVCV	4	10	2	1.1	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.69128345	205808	
côa	ver	1	0.0319	0	1.00	CVCV	4	8	1	1.5	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.27219436	205862	
coar	ver	2	0.0637	0.3010	1.00	CVCV	4	4	1	1.00	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.60956044	130937	
coas	ver	6	0.1912	0.7782	1.00	CVCV	4	15	1	1.00	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.44505048	73414	
coça	ver	4	0.1275	0.6021	1.00	CVCV	4	33	1	1.00	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.10540561	89374	
codi	ver	4	0.1275	0.6021	1.00	CVCV	4	3	1	1.5	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.1638345	206060	
coe	ver	41	1.3067	1.6128	1.00	CVCV	4	23	3	1	1.00	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.63565076	27090
cola	ver	39	1.2429	1.5911	1.00	CVCV	4	50	1	1.00	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.05251178	27863	
cole	ver	31	0.9880	1.4914	1.00	CVCV	4	36	1	1.00	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.85920788	31700	
coli	ver	3	0.0956	0.4771	1.00	CVCV	4	20	1	1.00	CVCV	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	VCVC	Ec_ee_aa_ee_ee	Ec_ee_aa_ee_ee	0.16136673	103093	

Figure 2: Search result example.

It follows below the name, abbreviation, and description of each one of the 21 columns of information presented in the results obtained from a search in the Brazilian Portuguese Lexicon, Alpha.

**Orthography** (ortho): orthographic form of the word in lower case (except for proper names), respecting the specific accentuation of each word<sup>37</sup>.

**Grammatical category** (gram\_cat): grammatical category of the word (adj, adv, gram, nom, num, prop, and ver).

**Grammatical information** (gram\_inf): grammatical information about the word (e.g. singular/plural, masculine/feminine, past/present/future, 1/2/3 person, etc.).

**Orthographic frequency** (ortho\_freq): number of times that the word appears in NILC (around 32 million words).

**Orthographic frequency per million** (ortho\_freq/M): number of times the word appears in NILC in 1 million words. Standard value used to express the frequency of a word.

**Natural logarithm of the orthographic frequency** (log10\_ortho\_freq): natural logarithm of the orthographic frequency. Logarithmic values are used to linearize the behavior of the frequencies of words in a corpus.

**Number of letters** (nb\_letters): number of letters of the word.

**Number of homographs** (nb\_homogr): number of homographic words. Words that have the same orthography or accentuation differences, but belong to different grammatical categories.

**Homographs** (homographs): grammatical categories of the homographic words.

**Orthographic uniqueness point** (pu\_ortho): letter position from which the word is dissociated from the other corpus words. Direction: from left to right.

**Orthographic neighbors** (ortho\_neigh): number of orthographic neighbors according to the Coltheart's N, changing only one letter at a time (Coltheart et al., 1977).

---

<sup>37</sup> We remark that the NILC was computed in 1999, before the Portuguese orthographic reform.

**Orthographic Leveinshtein Distance** (old20): average of the Orthographic Leveinshtein Distance from the 20 closest words calculated from linear regressions (Yarkoni et al., 2008).

**CVCV structure** (CVCV\_ortho): word structure, where consonants are C and vowels V. Also, A for the accentuation, P for punctuation, N for numbers, and S for symbols.

**Bigrams** (bigrams): the word bigrams separated by underline “\_” and bounded by hashmarks “#”. The number of bigrams equals the number of letters of the word plus 1.

**Trigrams** (trigrams): the word trigrams separated by underline “\_” and bounded by hashmarks “#”. The number of trigrams equals the number of letters of the word.

**Reversed orthography** (rev\_ortho): orthographic word reversed from (ortho).

**Reversed CVCV structure** (rev\_CVCV\_ortho): reversed CVCV structure of the word from (CVCV\_ortho).

**Reversed bigrams** (rev\_bigrams): reversed word bigrams separated by underline “\_” and bounded by hashmarks “#” from (bigrams).

**Reversed trigrams** (rev\_trigrams): reversed word trigrams separated by underline “\_” and bounded by hshmarks “#” from (trigrams).

**Random number between 0 and 1** (random): random number between 0 and 1 with eight digits of precision.

**Identification number** (id): word identification number assigned from the corpus organized by decreasing frequency and a-z alphabetical order. The identification number is the position of the word in the lexicon and corpus.

## Simple Search

The simple search engine was developed for the search of specific words or word lists, as shown in **Figure 3**. The user must inserts the orthographic form of the words to perform a search, indeed, wildcard symbols as underline “\_” for a letter and percentage “%” for a chain of letters can be used. The user can enter a list of words separated by different rows. For

example, the user can copy and paste a list of words from a spreadsheet or text editor. The user can choose the category used for the organization and presentation of words and the ascending or descending sense of organization. The “Search” button performs the search and displays the results and the “Clear” button clears the information from the fields.



**Simple search**

amar  
braba  
dormir  
ontem

Order by:

orthography · ascendant · Search Clear

**Figure 3:** Simple search.

## Complex Search

The complex search engine was developed to perform complex searches based on specific criteria of words, as the number of letters, frequency, grammatical category, orthographic neighbors, etc., as shown in **Figure 4**. In the first field, the user must choose the column of information to perform the search. In the second field, the user must choose whether to consider “yes” or disregard “no” the criteria. In the third field, the user must enter the specific criteria to be searched. The wildcard symbols underline “\_” for a letter and percentage “%” for a chain of letters can be used. Still, the symbols greater than “>” and less than “<” may be used for numerical searches. The user can choose the category used for the organization and presentation of words and the ascending or descending sense of organization. The “Search” button performs the search and displays the results and the “Clear” button clears the information from the fields. Initially, the complex search presents four fields for criteria insertion, by clicking on the “+ Fields” button, the user is redirected to another page that presents eight fields for criteria insertion in complex searches.

### Complex search

1	gram_cat	·	yes	·	ver
2	nb_letters	·	yes	·	>3 <8
3		·	yes	·	
4		·	yes	·	

Order by:

orthography	·	ascendant	·	Search	Clear	+ Fields
-------------	---	-----------	---	--------	-------	----------

**Figure 4:** Complex search.

## Results

In the results section (**Figures 2 and 5**), the user finds the search results organized in different rows and the metalinguistic and psycholinguistic information in different columns. It can be also found specific search information, as shown in **Figure 5**: 1) total number of words found in the search, 2) range of words presented, 3) total number of pages, and 4) page number displayed. The user can choose in the upper left selection the number of words displayed on the page and navigate through the pages using the buttons “Previous” and “Next”.

The top right of the results section, as shown in **Figure 5**, presents a series of statistical data established and calculated from the search (Davis & Perea, 2005; Davis, 2005): 1) mean, 2) maximum, and 3) minimum values of the following categories: 1) ortho\_freq, 2) log10\_ortho\_freq, 3) nb\_letters, d) ortho\_neigh, and 4) old20. Finally, the button “Export .csv” exports all the search data to a downloadable .csv file.

Results		Statistics			
50	· Previous	Next	Export .csv		
Page 1 of 1					
0 - 4 words from a total of 4 word found					
category	freq_ortho	log10_freq_ortho	nb_letters	neigh_ortho	old20
Mean	11306.0000	2	5.0000	6.5000	1
Min	5	0.6990	4	2	1.00
Max	44199	4.6454	6	12	1.75

**Figure 5:** Results information and basic statistics.



## Pages

In addition to the main pages of the Brazilian Portuguese Lexicon: Lexicon<sup>38</sup> and Pseudowords<sup>39</sup>, the following other pages were also created to complement the corpus and website: Downloads, Tools, Updates, Credits, and Statistical Linguistics. Downloads<sup>40</sup> offer for download several files from the Brazilian Portuguese Lexicon (corpus.txt, manuals, lists, conventions, bigrams, trigrams, R scripts, etc.). Tools<sup>41</sup> present a series of corpora links, programs, and literature on psycholinguistics, computational linguistics, and statistics. Updates<sup>42</sup> describe the development of the Brazilian Portuguese Lexicon and the updates made over time. Credits<sup>43</sup> present the proposal, the source, and the authors of the Brazilian Portuguese Lexicon; it still describes the references and relevance of the NILC/São Carlos corpus, *Linguateca* website, *Lexique* corpus, R program and packages, and Creative Commons License, ending with the acknowledgments. Finally, Statistical Linguistics is a page that provides various open and free tools and resources, as described below.

## Pseudowords

The BP pseudoword generator engine was developed to create pseudowords based on the structure and frequency of the BP words (Keuleers & Brysbaert, 2010). Unlike other pseudoword generator engines based on the syllabic structure of existing words (Keuleers & Brysbaert, 2010; Mota & Resende, 2013), the BP pseudoword generator engine of the Brazilian Portuguese Lexicon uses the bigrams and trigrams (B. New et al., 2001). All the bigrams and trigrams were computed from all the words of the Brazilian Portuguese Lexicon. The BP pseudowords are generated from the frequency and combination of bigrams or trigrams. We computed the 1) general bigram and trigram frequency, 2) bigram and trigram frequency in function of the position in the word, and 3) bigram and trigram frequency by grammatical categories.

---

<sup>38</sup> [http://www.lexicodoportugues.com/index\\_en.php](http://www.lexicodoportugues.com/index_en.php)

<sup>39</sup> [http://www.lexicodoportugues.com/pseudowords\\_en.php](http://www.lexicodoportugues.com/pseudowords_en.php)

<sup>40</sup> [http://www.lexicodoportugues.com/downloads\\_en.php](http://www.lexicodoportugues.com/downloads_en.php)

<sup>41</sup> [http://www.lexicodoportugues.com/tools\\_en.php](http://www.lexicodoportugues.com/tools_en.php)

<sup>42</sup> [http://www.lexicodoportugues.com/updates\\_en.php](http://www.lexicodoportugues.com/updates_en.php)

<sup>43</sup> [http://www.lexicodoportugues.com/credits\\_en.php](http://www.lexicodoportugues.com/credits_en.php)

The user must enter four fields in the BP pseudoword generator engine: 1) number of words to be generated, 2) number of letters of the words to be generated, 3) grammatical category that these words must belong (all, adj, adv, gram, nom, num, ver), and 4) criterion for the construction of words (bigrams or trigrams). The BP pseudoword generator engine builds words simultaneously in both directions, from left to right and from right to left, beginning with a bounded bigram or trigram “#xx” or “#xx”. According to the number of letters, the engine concatenates new bigrams or trigrams that share as much orthographic information as possible to the previous bigram or trigram (1 letter for bigrams and 2 letters for trigrams). The BP pseudoword generator engine has two buttons: “Submit” to generate and present the pseudoword results and “Clear” to clear all fields, as shown in **Figure 6**.

**Figure 6:** BP pseudoword generator engine.

In the pseudoword results, as shown in **Figure 7**, it is presented four columns with information about the pseudowords: 1) grammatical category defined by the user, 2) pseudoword frequency calculated from the sum of the frequencies of the pseudoword bigrams or trigrams, 3) log10 of the pseudoword calculated frequency, and 4) number of letters of the pseudoword. It is available the button “Export .csv” to export the results of the BP pseudowords generator engine to a downloadable .csv file.

Results									
<a href="#">Export .csv</a>									
pseudo left-right	category left-right	freq left-right	log_freq left-right	nb_letters left-right	pseudo right-left	category right-left	freq right-left	log_freq right-left	nb_letters right-left
deste	toda	37248	10.5254	5	cadaz	toda	42028	10.6461	5
reste	toda	35294	10.4715	5	prado	toda	41806	10.6408	5
inado	toda	38404	10.5559	5	cazes	toda	29797	10.3022	5
prese	toda	31623	10.3616	5	casse	toda	27245	10.2126	5
caose	toda	19178	9.8615	5	atram	toda	24379	10.1015	5
entes	toda	41396	10.6309	5	prada	toda	32872	10.4004	5
supes	toda	21001	9.9523	5	dente	toda	47546	10.7695	5
estas	toda	37925	10.5434	5	caria	toda	24547	10.1083	5
expla	toda	8401	9.0361	5	dicão	toda	24379	10.1015	5
mando	toda	30824	10.336	5	aprem	toda	16447	9.7079	5

**Figure 7:** BP pseudoword generator engine results.

## Statistical Linguistics

The page Statistical Linguistics<sup>44</sup> in the Brazilian Portuguese Lexicon provides free and open tools and resources for psycholinguistic, linguistic, and statistics that can be consulted directly on the page on the internet. These resources and tools were developed in HTML/PHP: a) F' and minF' - MS, b) minF' - F1.F2, c) Hartley test, d) normalization between 0 and 1, d) reverse words, f) Hamming distance, g) Levenshtein distance, h) orthographic neighbors (Coltheart's N), i) average Levenshtein distance, j) relative entropy, k) word frequency, and l) Zipf's distribution.

## Authors

The Brazilian Portuguese Lexicon is being developed by Gustavo Lopez Estivalet during his Ph.D financed by the program Science without Borders (CsF)<sup>45</sup> from the National Council for Scientific and Technological Development (CNPq)<sup>46</sup> scholarship, Brazil, between 2012 and 2016. Gustavo Lopez Estivalet holds his Ph.D in France in the city of Lyon, in the University Claude Bernard Lyon 1 (UCBL)<sup>47</sup> in the Doctoral School of Neurosciences et Cognition (ED NSCo)<sup>48</sup>, in the Laboratory on Language, Brain, and Cognition (L2C2)<sup>49</sup>, located at the Institute of Cognitive Sciences (ISC)<sup>50</sup>, and being supervised by Prof. Dr. Fanny Meunier<sup>51</sup>, which are funded by the National Council for Scientific Research (CNRS)<sup>52</sup>.

---

<sup>44</sup> [http://www.lexicodoportugues.com/stat\\_ling\\_en.php](http://www.lexicodoportugues.com/stat_ling_en.php)

<sup>45</sup> <http://www.cienciasemfronteiras.gov.br/web/csf>

<sup>46</sup> <http://www.cnpq.br/>

<sup>47</sup> <http://www.univ-lyon1.fr/>

<sup>48</sup> <http://nsco.universite-lyon.fr/>

<sup>49</sup> <http://l2c2.isc.cnrs.fr/fr/>

<sup>50</sup> <http://www.isc.cnrs.fr/>

<sup>51</sup> [http://www.researchgate.net/profile/Fanny\\_Meunier/](http://www.researchgate.net/profile/Fanny_Meunier/)

<sup>52</sup> <http://www.cnrs.fr/>

## Acknowledgements

For the realization and success of the Brazilian Portuguese Lexicon, I thank the National Council for Scientific and Technological Development (CNPq) by the Ph.D scholarship from the program Science without Borders (CsF). I thank my Ph.D supervisor Prof. Dr. Fanny Meunier and Prof. Dr. Michel Hoen<sup>53</sup>, who understood the importance of developing a BP psycholinguistic corpus. I thank the NILC Prof. Dr. Sandra M. Aluísio<sup>54</sup> and Prof. Dr. Maria das Graças Volpe Nunes<sup>55</sup> by the help in materials, information, and discussions on the NILC, as well as support, motivation and recognition of this work. I thank my work colleagues Léo Varnet<sup>56</sup> and Emmanuel Trouche<sup>57</sup> by the discussions on scripts and algorithms for the development of Brazilian Portuguese Lexicon. I thank the internet community users who work with website development and database management for the forum discussions and tutorials available. Finally, I thank Prof. Dr. Mailce Borges Mota, and the best BP Prof. Lise Lopez. Finally, I thank Luanda Lins by understanding the importance of this project and my motivation to do it. Thank you all!

## References

- Baayen, H. R. (2001). *Word Frequency Distributions* (Vol. 18). Dordrecht; Boston; London: Kluwer Academic Publishers.
- Baayen, H. R., Piepenbrock, R., & van Rijn, H. (1995). *The CELEX lexical database. Release 2 [CD-ROM]*. Philadelphia, PA: Linguistic Data Consortium, University of Pennsylvania.
- Balota, D. A., Yap, M. J., Hutchison, K. A., Cortese, M. J., Kessler, B., Loftis, B., ... Treiman, R. (2007). The English Lexicon Project. *Behavior Research Methods*, 39(3), 445–459. doi:10.3758/BF03193014
- 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. doi:10.3758/BRM.41.4.977

<sup>53</sup> [http://www.researchgate.net/profile/Michel\\_Hoen/](http://www.researchgate.net/profile/Michel_Hoen/)

<sup>54</sup> [http://www.researchgate.net/profile/Sandra\\_Aluisio/](http://www.researchgate.net/profile/Sandra_Aluisio/)

<sup>55</sup> [http://www.researchgate.net/profile/Maria\\_Nunes10/](http://www.researchgate.net/profile/Maria_Nunes10/)

<sup>56</sup> [http://www.researchgate.net/profile/Leo\\_Varnet/](http://www.researchgate.net/profile/Leo_Varnet/)

<sup>57</sup> <http://cnrs.academia.edu/EmmanuelTrouche>

- Coltheart, M. (1981). The MRC psycholinguistic database. *The Quarterly Journal of Experimental Psychology Section A*, 33(4), 497–505. doi:10.1080/14640748108400805
- Coltheart, M., Davelaar, E., Jonasson, J. T., & Besner, D. (1977). Access to the internal lexicon. In S. Dornic (Ed.), *Attention and Performance VI* (pp. 535–555). Hillsdale, NJ: Lawrence Erlbaum Associates.
- Davis, C. J. (2005). N-Watch: A program for deriving neighborhood size and other psycholinguistic statistics. *Behavior Research Methods*, 37(1), 65–70. doi:10.3758/BF03206399
- 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. doi:10.3758/BF03192738
- Ferrand, L., New, B., Brysbaert, M., Keuleers, E., Bonin, P., Méot, A., ... Pallier, C. (2010). The French Lexicon Project: Lexical decision data for 38,840 French words and 38,840 pseudowords. *Behavior Research Methods*, 42(2), 488–496. doi:10.3758/BRM.42.2.488
- Gimenes, M., & New, B. (2015). Worldlex: Twitter and blog word frequencies for 66 languages. *Behavior Research Methods*. doi:10.3758/s13428-015-0621-0
- Keuleers, E., & Brysbaert, M. (2010). Wuggy: A multilingual pseudoword generator. *Behavior Research Methods*, 42(3), 627–633. doi:10.3758/BRM.42.3.627
- Keuleers, E., Diependaele, K., & Brysbaert, M. (2010). Practice Effects in Large-Scale Visual Word Recognition Studies: A Lexical Decision Study on 14,000 Dutch Mono- and Disyllabic Words and Nonwords. *Frontiers in Psychology*, 1. doi:10.3389/fpsyg.2010.00174
- Keuleers, E., Lacey, P., Rastle, K., & Brysbaert, M. (2012). The British Lexicon Project: Lexical decision data for 28,730 monosyllabic and disyllabic English words. *Behavior Research Methods*, 44(1), 287–304. doi:10.3758/s13428-011-0118-4
- Marian, V., Bartolotti, J., Chabal, S., & Shook, A. (2012). CLEARPOND: Cross-Linguistic Easy-Access Resource for Phonological and Orthographic Neighborhood Densities. *PLoS ONE*, 7(8), e43230. doi:10.1371/journal.pone.0043230
- Mota, M. B., & Resende, N. (2013). Metodologia da pesquisa em psicolinguística: desenvolvimento de uma ferramenta para a geração automática de pseudoverbos. *Letras de Hoje*, 48(1), 100–107.
- New, B., Ferrand, L., Pallier, C., & Brysbaert, M. (2006). Reexamining the word length effect in visual word recognition: New evidence from the English Lexicon Project. *Psychonomic Bulletin & Review*, 13(1), 45–52. doi:10.3758/BF03193811
- New, B., Pallier, C., Brysbaert, M., & Ferrand, L. (2004). Lexique 2 : A new French lexical database. *Behavior Research Methods, Instruments, & Computers*, 36(3), 516–524. doi:10.3758/BF03195598

- New, B., Pallier, C., Ferrand, L., & Matos, R. (2001). Une base de données lexicales du français contemporain sur internet : LEXIQUE<sup>TM</sup>//A lexical database for contemporary french : LEXIQUE<sup>TM</sup>. *L'année Psychologique*, 101(3), 447–462.  
doi:10.3406/psy.2001.1341
- Pinheiro, G. M., & Aluísio, S. M. (2003). *Corpus NILC: descrição e análise crítica com vistas ao projeto Lacio - Web. Série de Relatórios do Núcleo Interinstitucional de Lingüística Computacional NILC - ICMC - USP*. São Carlos, SP: Universidade Federal de São Carlos - UFSCar.
- Santos, D., & Bick, E. (2000). Providing internet access to Portuguese corpora: the AC/DC project. In M. Gavrilidou, G. Carayannis, S. Markantonatou, S. Piperidis, & G. Stainhauer (Eds.), *Proceedings of the Second International Conference on Language Resources and Evaluation (LREC2000)* (pp. 205–210). Athens, Greece.
- Schreuder, R., & Baayen, R. H. (1995). Modeling Morphological Processing. In L. B. Feldman (Ed.), *Morphological Aspects of Language Processing* (pp. 131–154). Hillsdale, New Jersey: Lawrence Erlbaum Associates, Inc., Publishers.
- Yarkoni, T., Balota, D., & Yap, M. (2008). Moving beyond Coltheart's N: A new measure of orthographic similarity. *Psychonomic Bulletin & Review*, 15(5), 971–979.  
doi:10.3758/PBR.15.5.971