



# The Portuguese Vocabulary Profile: Pilot Study

Shintaro Torigoe

Tokyo University of Foreign Studies, Tokyo, Japan.  
Gunma Prefectural Women's University, Gunma, Japan.  
torigoe.shintaro@tufs.ac.jp

## Abstract

This paper reports on a pilot study from the *Portuguese Vocabulary Profile* project. In this pilot study, a vocabulary list for learners of Portuguese was developed by analysing learner corpora, an approach inspired by CEFR-based wordlists, such as the *English Vocabulary Profile*. A draft wordlist was constructed from two learner corpora of L2 Portuguese, the *Corpora do PLE* and the *Corpus de PEAPL2*. The draft wordlist was then compared to the *LMCPC*, a wordlist derived from a million-word native speaker corpus, in order to investigate differences between learners and native speakers and to identify aspects of the wordlist needing improvement. The pilot study indicated that the use of Portuguese by the Intermediate and Advanced learner is quite different from that of native speakers and that learner's language use was affected by data collection tasks and learning environments.

## 1 Introduction

The aim of the pilot study reported in this paper was to construct the *Portuguese Vocabulary Profile*, a vocabulary list for elementary, intermediate, and advanced learners of Portuguese. The list was derived by analysing learner corpora, an approach inspired by the *English Vocabulary Profile* (Capel 2010, 2012). In Japan, an increasing number of Portuguese language textbooks have been published since 2010 (cf. Torigoe & Yamada 2015), yet there are no guidelines for appropriate tasks, grammar, or vocabulary. In response to this situation, the author proposed developing a prototype wordlist by adopting the *Common European Framework of Reference for Languages (CEFR)* levels (Council of Europe 2001) as proficiency criteria and by using learner corpora as the data source. Such a list could be expected to improve the quality of Portuguese language textbooks in Japan and also contribute to the development of textbooks for Intermediate learners (resource which are currently lacking in the Japanese market).

## 2 Background

As mentioned above, there is an increasing number of Portuguese language textbooks in Japan, especially since 2010. Most of these are intended for beginners; however, the goals of the various books differ. Moreover, within these textbooks there is no specification of the amount of vocabulary included, or the criteria adopted for vocabulary selection.

In Japan, there are robust national guidelines, *the Courses of Studies* (Ministry of Education, Culture, Sports, Science and Technology-Japan 2009), which define detailed goals for the learning and teaching of the majority of subjects including English. The guidelines, however, simply describe the goals of foreign languages other than English as they ‘should follow the objectives and contents of English instruction’ (ibid, p.90). This effectively means that the decisions about the teaching of these subjects are delegated to language teachers and textbook writers.

For English language learning, there are examples. One approach to establishing vocabulary goals is based on the language use of native speakers (NS), with corpora used as the data sources. For example, Tono (2013) reports that the top 100 words of the spoken subcorpora of the *British National Corpus* account for 67 per cent of the total word frequency, the top 1000 words account for 87 per cent, and the top 2000 words account for 92 per cent. A vocabulary target of 1000 words is equivalent to the goal specified for English language learning in the 8th grade in Japan, and 2000 words is equivalent to the expectation for 10th grade. In Japan, there are some learning books for English vocabulary based on corpus information, such as Ishikawa (2006) or Tono (2008).

For Portuguese language learning, some corpus-based teaching materials are also available. They are mainly dictionaries; for example, *Dicionário Verbo Língua Portuguesa* (2008), based on the *Corpus de Referência de Português Contemporâneo (CRPC)*<sup>\*</sup>; *Dicionário Gramatical de Verbos Portugueses* (2007), based on the *Português Fundamental*<sup>†</sup>; and *A Frequency Dictionary of Portuguese* (Davis & Preto-Bay 2008), based on the *Corpus do Português*<sup>‡</sup>. In Japan, there are some examples of corpus-based teaching materials for Portuguese, such as Aires & Iyanaga (2012), a vocabulary book based totally on the *LMCPC* (a subcorpus of the *CRPC*); and Ichinose et al. (2015), a Portuguese-Japanese dictionary that uses the *Corpus do Português* as one reference for important words.

### 2.1 Does native speakers’ language use correspond to learners’ use?

All the resources identified in the previous section are examples of teaching materials based on native speaker corpora. In second language acquisition research, however, there is a suggestion that NSs’ use of language may not always correspond to learners’ use. The author agrees with this idea and, therefore, adopts data sources that depict learner use of Portuguese as the basis for developing the new vocabulary list. In other words, learner corpora are used to describe learner use. The *Common European Framework of Reference for Languages (CEFR)* (Council of Europe 2001) is also used to provide criteria for learner proficiency.

### 2.2 The *Common European Framework of Reference for Languages*

The *Common European Framework of Reference for Languages* (Council of Europe 2001) provides versatile guidelines for language learning in multilingual and multicultural contexts. The *CEFR* is not based on dogmatic notions of ‘what learners should learn’, but, rather, on actual learner

\* [www.clul.ul.pt/pt/investigacao/183-reference-corpus-of-contemporary-portuguese-crpc](http://www.clul.ul.pt/pt/investigacao/183-reference-corpus-of-contemporary-portuguese-crpc)

† [www.clul.ul.pt/pt/recursos/84-spoken-corpus-qportugues-fundamental-pfq-r](http://www.clul.ul.pt/pt/recursos/84-spoken-corpus-qportugues-fundamental-pfq-r)

‡ [www.corpusdportugues.org/](http://www.corpusdportugues.org/)

behaviour, that is ‘what learners can do’. The *CEFR* is applied very widely, including in learner assessment by teachers, learner self-assessment, textbook writing, syllabus development, etc.

The *CEFR* assesses listening, speaking, reading, writing, and oral interaction, using six proficiency levels: A1 and A2 (‘basic user’), B1 and B2 (‘independent user’), and C1 and C2 (‘proficient user’). According to Tono (2013), C-level language learners can perform as well as, or even better than, native speakers. For example, a C-level learner ‘[c]an hold his/her own in formal discussion of complex issues, putting an articulate and persuasive argument, at no disadvantage to native speakers’ (Council of Europe 2001, 78), although, of course, it is difficult to compare these two groups objectively.

For further details of the *CEFR*, such as its history, the examples of its implementation, and related problems, see Council of Europe (2001).

### 2.3 CEFR-based wordlist

The *English Vocabulary Profile* (Capel 2010, 2012) and the *CEFR-J Wordlist* (Tono 2013) are two instances of vocabulary lists that have been constructed with reference to learner corpora, which have been divided into subcorpora according to the *CEFR* proficiencies.

According to Capel (2010), the *EVP* is a wordlist initially organised using the headwords of the *Cambridge Learner Dictionary*, which is based on a native speaker corpus (the *Cambridge International Corpus*). However, to enhance its validity, the vocabulary in the *EVP* was compared with a *CEFR*-based learner corpus (the *Cambridge Learner Corpus*) and textbook corpora, and was checked by native speakers. A noteworthy feature of the *EVP* is that the classification of vocabulary into the *CEFR* proficiency levels was done not by lemma, but by meaning.

The *CEFR-J* (Tono 2013) is a modified version of the *CEFR* and was developed to reflect the particular context of English teaching in Japan. It has more proficiency classifications within the A and B levels than the original *CEFR*. The *CEFR-J Wordlist* (ibid) is derived from grading corpora of English textbooks in East Asian countries by *CEFR* proficiencies. Tono (ibid) compares the *EVP* and the *CEFR-J Wordlist* quantitatively (Table 1).

<i>CEFR</i> proficiency levels	Words in <i>EVP</i>	Words in <i>CEFR-J Wordlist</i>
A1	601	1000
A2	925	1000
B1	1429	2000
B2	1711	2000
C1	2300	n/a

**Table 1:** Comparison between the English Vocabulary Profile and the *CEFR-J* wordlist

Although Tono does not compare the content of the two wordlists, he shows that the *CEFR-J Wordlist* has more words in each proficiency level from A1-B2, on the other hand, the *CEFR-J Wordlist* does not have the C-levels.

For Portuguese language learning, there is currently no wordlist available that is based on learner corpora. However, two learner corpora are available that are composed of subcorpora classified by *CEFR* proficiencies (see 4.1 below). These may represent useful source material for developing a wordlist for learners for learners of Portuguese.

## 3 Research questions

This paper describes a pilot study that involved using *CEFR*-based learner corpora to construct a wordlist for learners of Portuguese. The research questions are below.

1. How much vocabulary should learners learn?
2. Does native speaker's use correspond to learner's use?

The author first developed the new wordlist from learner corpora and then examines the degree of similarity between this list and a wordlist derived from a NS corpus.

## 4 Methodology

Although this study was inspired by the *EVP* (Capel 2010, 2012), it adopts different methodology. The *EVP* was initially derived from the headwords of a dictionary, and then these were compared to a learner corpus. The wordlist developed in the current study, in contrast, draws on learner corpora first in order to obtain more learner-centric data. Furthermore, the words in the present study are not selected based on simple frequency alone, but using additional statistical techniques.

### 4.1 Learner data

As the source of the new wordlist, the author uses two learner written corpora: the *Corpora do Português Língua Estrangeira (PLE)*<sup>§</sup>, collected by the Linguistic Centre of the University of Lisbon; and the *Corpus de Produções Escritas de Aprendentes de PL2 (PEAPL2)*<sup>\*\*</sup>, collected by the General and Applied Linguistics Centre of the University of Coimbra. The writing samples in both corpora were obtained using the same writing task. Summaries of the two corpora are shown below (Table.2).

	<i>PLE</i>	<i>PEAPL2</i>
Total words	Approx. 70,500	Approx. 120,000
Number of subcorpora	471	546
Number of mother tongues of the informants	27	39
Learning environment	Studying in home country	Studying at University of Coimbra

**Table 2:** Summaries of the *PLE* and the *PEAPL2*

All the informants for both corpora responded to a can-do self-assessment questionnaire for placement, and then each informant was classified by *CEFR* proficiency. Whereas in the *PEAPL2* the informants were classified into five levels (A1, A2, B1, B2, and C1), in the *PLE* the informants are classified into just three levels (A1-A2, B1-B2, and C1-C2). To allow both these corpora to be used for the present study, the author integrated the A1 and A2 levels, as well as the B1 and B2 levels, of the *PEAPL2*. Following this, the author integrated the *PLE* and the re-classified *PEAPL2* into a final, three-level learner corpora by integrating the A1-A2, the B1-B2 and the C1-C2 levels from both corpora. Before integrating the corpora, the author examined the proximity between the subcorpora by clustering<sup>††</sup>, in order to validate the integration. The analysis confirmed the proximity between subcorpora of the same proficiency level, with the exception of the B1-B2 level of the *PEAPL2* (Figure 1). This seemed to be because of the larger size of this subcorpus compared to the others. Despite this finding, the author did not discard the B1-B2 subcorpus of the *PEAPL2* from this study.

<sup>§</sup> [www.clul.ul.pt/pt/recursos/314-corpora-of-ple](http://www.clul.ul.pt/pt/recursos/314-corpora-of-ple)

<sup>\*\*</sup> [www.uc.pt/fluc/rcpl2](http://www.uc.pt/fluc/rcpl2)

<sup>††</sup> The author used *Seagull-Stat*, a statistical Add-in for Microsoft Excel<sup>®</sup>.  
[www.7b.biglobe.ne.jp/~hayakari/](http://www.7b.biglobe.ne.jp/~hayakari/)

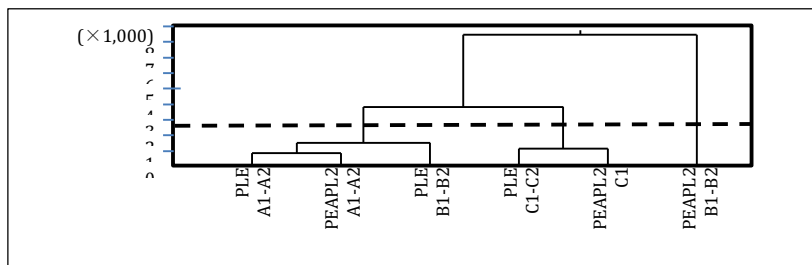


Figure 1: Result of clustering subcorpora

## 4.2 Results of the learner corpus analysis

The results of the learner corpus analysis are shown below (Table 3).

	A1-A2	B1-B2	C1-C2	Total
Words	59932	101758	16793	178483
Lemmas	2198	3961	1286	4865
Including foreign words and orthographical mistakes				

Table 3: Number of words and lemmas in integrated subcorpora

Across the whole learner corpus, the top 100 words accounted for 72 per cent of the total text, the top 1000 words accounted for 95 per cent, and the top 2000 words accounted for 98 per cent. Restricting the analysis to the A-level subcorpora, the top 500 words accounted for 95 per cent of the total text in the remaining subcorpus, the top 600 words account for 96 per cent, and the top 700 account for 97 per cent. Given these results, the author set the cut off points for each proficiency level at the top 500th, 1000th, and 1500th words.

## 4.3 Determining final for word placement

Through the analysis described in the previous section, the raw frequency of the words was obtained. However, dividing the wordlist into proficiency levels by simply using the cut off points may not be adequate. For this reason, additional statistical techniques were used to determine the final word placement..

First, the top 500 words of the A-level subcorpus were automatically classified as the Elementary level words, because statistical analysis confirmed that there was a high correlation ( $r=.90$ ) between the top 500 words of the A-level subcorpus and the whole corpus.

Second, to classify the Intermediate and Advanced level words, the remaining 1000 words from the top 1500 words of the whole corpora were analysed using the chi-square test. The words found to be significant for the A-level were added to the existing list of 500 Elementary words; similarly, those significant for the B- and C-levels were added to the existing Intermediate and Advanced lists respectively.

Third, the remaining word (which had not been significant at any proficiency level according to the chi-square test) were classified. Those which had been produced within the whole corpus more than 10 times were classified as Intermediate words, and those produced less than 10 times were classified as Advanced words.

Finally, Arabic numerals ( $1, 2, 3, \dots$ ), alphabet letters ( $a, b, c, \dots$ ), foreign words not naturalized in Portuguese (*new, espanhã*), and some orthographical mistakes (*tambem, portugues*) were all excluded.

## 4.4 Native speaker baseline data: The LMCPC

For the native speaker baseline data, the author chose the *Léxico Multifuncional Computorizado do Português Contemporâneo (LMCPC)*, a wordlist derived from the *CORLEX* (a subcorpus of the *CRPC*). This list was selected because of its balance and the number of words included. The *LMCPC* is frequency list of more than 25000 words. The rank and frequency of the words in the *LMCPC* were annotated manually on the learner wordlist obtained as a result of the steps described in in 4.1-4.3.

## 5 Findings

The final wordlist derived from the two learner corpora is presented in the Appendix. In this section, findings related to this list are detailed.

The correlation (mutual information score) between the word frequencies in the learner corpora and in the *LMCPC* was examined. For the elementary level vocabulary, high mutual information score was obtained ( $r=.98$ ); lower scores were obtained for the intermediate ( $r=.30$ ) and the advanced vocabulary ( $r=.28$ ).

Reviewing the final wordlist qualitatively, the author noted that there were some words that appeared frequently in the learner corpora but much less frequently in the NS corpus. One obvious example concerns the bias in the learner corpora towards European Portuguese vocabulary. For example, *autocarro* (bus), *comboio* (train), and *desporto* (sport) all appeared in the top 500 words of the final word list; meanwhile, their counterparts in Brazilian Portuguese, *ônibus*, *trem*, and *esporte*, did not even appear even in the top 1500 words. Similarly, the frequency of *você* (you)<sup>\*\*</sup>, one of the most frequent words in Brazilian Portuguese, was not as highly ranked in the learner use wordlist ( $n=40$ , 267<sup>th</sup>) as in the native speaker corpus. A second type of bias observed in the learner corpus was towards words related to academic life, for example, *universidade* (university), *faculdade* (faculty, college), *república* (a dormitory system in Coimbra), *doutoramento* (doctoral programme), etc. One further bias noted within the final wordlist was caused by the tasks used for data collection. For example, the words *coreano* (Korean), *Edimburgo* (Edinburgh), *Bucareste* (Bucharest), which seemed to have been elicited in the self-presentation task<sup>§§</sup> responses of A1-A2 learners, all appeared in the Elementary level vocabulary list.

Finally, some intuitively basic words such as numerals (*sete* [seven], *nove* [nine]) and months (*janeiro* [January], *outubro* [October]) appeared in the Intermediate or Advanced level. These words would normally be included in textbooks for beginners.

## 6 Discussion

Let us now consider the second research question: Does native speakers' use of language correspond to learners' use? The high correlation between the frequencies of the Elementary words in the learner corpora and the NS baseline corpus (the *LMCPC*) suggests that certain basic vocabulary is common to both L2 learners and NS. On the other hand, one interpretation of the low correlation between the frequencies of the Intermediate and Advanced words in the learner corpora and the NS baseline corpus suggests that the most important words for these groups of learners may differ; if so, this may directly answer the second research question.

---

<sup>\*\*</sup> In Brazilian Portuguese usage, the third-person singular *você* is almost exclusively used to indicate interlocutor, whereas in European Portuguese, *você* and the second-person singular *tu* are optional depending on the relationship with interlocutor.

<sup>§§</sup> One of the optional tasks used for corpus data collection.

Another, perhaps more probable, explanation of the low correlations observed for Intermediate and Advanced words, is that these results may have been due to the cut off points selected in developing the learner word list as well as particular characteristics of the corpora used. In terms of the cut off points, compared to the *English Vocabulary Profile* and the *CEFR-J Wordlist* (see **Table 1**), just 500 words at each proficiency level are too few; this could mean that the level of difficulty of some words, which should really be classified as Elementary, is overestimated and the words are placed into the Intermediate or Advanced levels. More words within each proficiency level are needed due to the size of both the *PLE* and the *PEAPL2*. In terms of the characteristics of the corpora used, the discrepancy between the C-level word frequencies for learners and NS, as well as the biases observed in the previous section, seem to be influenced by the data collection task and the informants' learning environment.

For these reasons, the first research question still requires further study: How much vocabulary should learners learn?

## 7 Concluding remarks

This paper reports on a pilot study designed to develop a prototype of the *Portuguese Vocabulary Profile*, a wordlist for learners of Portuguese and derived from *CEFR*-based learner corpora. As the wordlist reported in this paper is the pilot version, opportunities for further improvement can be identified as well as some limitations resulting from the corpora used. First, the new wordlist was based solely on written corpora, meaning that there was no consideration of learners' use of spoken language. Secondly, constrained by the proficiency classifications used in the *PLE*, this study adopted a three-level classification; however, it would, of course, be more desirable to follow the original, six-level *CEFR* structure. The use of a three-level classification is also in contrast to the recent trends favouring the use of more proficiency levels; for example, the *CEFR-J* (Tono 2013) adopts 12 proficiency levels in order to reflect the distribution of learners of English in Japan. Finally, the points mentioned in the previous sections (that is, the absence of L2 Brazilian Portuguese words within the learner corpora, the biases caused by characteristics of the corpora used, and the adequateness of the cut off points selected for the word list) are also areas to be improved.

In conclusion, some areas for future study are indicated. First, the wordlist developed in this pilot study should be compared with additional corpus-based wordlist other than the *LMCPC*. Second, as Capel (2010) did for the *English Vocabulary Profile*, the learner wordlist should be compared with textbook corpora and should also be assessed by native speakers. Third, a list of *n*-grams should also be obtained and examined through the methods outlined above. The more sophisticated wordlist is expected to help improve the quality of both textbooks for Portuguese language learners and learner dictionaries in Japan.

## References

- Aires, P., & Iyanaga, S. (2010/2012). *Verbos Fundamentis do Português*. Kyoto: Kyoto University of Foreign Studies.
- Baba, J. (2010). *Comparison of French Textbooks Published in Japan and France*. Tokyo: Tokyo University of Foreign Studies.
- Capel, A. (2010). A1-B2 Vocabulary: Insights and Issues Arising from the English Profile Wordlists Project. *English Vocabulary Journal*, 1-11.
- Capel, A. (2012). Completing the English Vocabulary Profile: C1 and C2 Vocabulary. *English Vocabulary Journal*, 5(1), 1-14.

- Council of Europe. (2001). *Common European Framework of Reference for Languages: Learning, Teaching, Assessment*. Cambridge: Cambridge University Press.
- Davis, M., & Preto-Bay, A. R. (2008). *A Frequency Dictionary of Portuguese: Core Vocabulary for Learners*. New York: Routledge.
- Ichinose, A., Toida, H., Hayashida, M., & Yoshino, T. (2015). *Shogakukan Dicionário da Língua Portuguesa*. Tokyo: Shogakukan.
- Ishikawa, S. (2006). *Vital-series CORPUS 4500 words*. Tokyo: Bun-Ei-Do.
- Tono, Y. (2008). *Favourite Corpus 1800 words*. Tokyo: Tokyo Syoseki.
- Tono, Y. (2013). *CEFR-J Guidebook*. Tokyo: Taisyukan.
- Torigoe, S., & Yamada, M. (2015, November 29). Analysis of the Proficiency Goals of Portuguese Language Textbooks Based on the CEFR-J. Paper presented at the 19th annual meeting of The Japan Association of Foreign Language Education. Tokyo, Japan.

## Acknowledgement

This study is supported by *KAKENHI (Grant-in-Aid for Young Scientists (B))* of the Japan Society for the Promotion of Science), number 15K16792.











Table with 14 columns: rank, lemma, English, A1-A2, B1-B2, C1-C2, total, LMPGP rank, LMCPG freq, advanced level, LMPGP rank, LMCPG freq. It lists vocabulary items with their levels and frequencies.



