

1 A practical primer on processing semantic property norm data

2 Erin M. Buchanan¹, Simon De Deyne², & Maria Montefinese^{3,4}

3 ¹ Harrisburg University of Science and Technology

4 ² The University of Melbourne

5 ³ University of Padova

6 ⁴ University College London

7 Author Note

8 Erin M. Buchanan <https://orcid.org/0000-0002-9689-4189>; Simon De Deyne
9 <https://orcid.org/0000-0002-7899-6210>; Maria Montefinese
10 <https://orcid.org/0000-0002-7685-1034>. We would like to thank the editor and two
11 anonymous reviewers for their helpful comments in shaping this manuscript.

12 Correspondence concerning this article should be addressed to Erin M. Buchanan, 326
13 Market St., Harrisburg, PA 17101. E-mail: ebuchanan@harrisburgu.edu

Abstract

14

15 Semantic property listing tasks require participants to generate short propositions (e.g.,
16 *<barks>*, *<has fur>*) for a specific concept (e.g., dog). This task is the cornerstone of the
17 creation of semantic property norms which are essential for modelling, stimuli creation, and
18 understanding similarity between concepts. However, despite the wide applicability of
19 semantic property norms for a large variety of concepts across different groups of people, the
20 methodological aspects of the property listing task have received less attention, even though
21 the procedure and processing of the data can substantially affect the nature and quality of
22 the measures derived from them. The goal of this paper is to provide a practical primer on
23 how to collect and process semantic property norms. We will discuss the key methods to
24 elicit semantic properties and compare different methods to derive meaningful
25 representations from them. This will cover the role of instructions and test context, property
26 pre-processing (e.g., lemmatization), property weighting, and relationship encoding using
27 ontologies. With these choices in mind, we propose and demonstrate a processing pipeline
28 that transparently documents these steps resulting in improved comparability across
29 different studies. The impact of these choices will be demonstrated using intrinsic (e.g.,
30 reliability, number of properties) and extrinsic measures (e.g., categorization, semantic
31 similarity, lexical processing). This practical primer will offer potential solutions to several
32 longstanding problems and allow researchers to develop new property listing norms
33 overcoming the constraints of previous studies.

34

Keywords: semantic, property norm task, tutorial

35 A practical primer on processing semantic property norm data

36 Semantic properties are assumed to be, entirely or in part, the building blocks of
37 semantic representation - the knowledge we have of the world - by a variety of theories (e.g.,
38 Collins & Quillian, 1969; Jackendoff, 1992, 2002; Minsky, 1975; Norman & Rumelhart, 1975;
39 Saffran & Sholl, 1999; Smith & Medin, 1981) and computational models (Caramazza,
40 Laudanna, & Romani, 1988; Farah & McClelland, 1991; Humphreys & Forde, 2001). Within
41 this perspective, the meaning of a concept is conceived as a distributed pattern of semantic
42 properties, which convey multiple types of information (Cree & McRae, 2003; Plaut, 2002;
43 Rogers et al., 2004). For example, the concept HORSE can be described by encyclopedic
44 (<*is a mammal*>), visual (<*is furry*>, <*has legs*>, <*has a tail*>, <*has a mane*>),
45 functional (<*used for racing*>), and motor (<*gallops*>) information. Given the relevance of
46 semantic properties in shaping theories of semantic representation, researchers have
47 recognized the value of collecting semantic property production norms. Typically, in the
48 property generation task, participants are presented with a set of concepts and are asked to
49 list the properties they think are characteristic for each concept meaning. Generally, in this
50 task, the concepts are called *cues*, and the responses to the cue are called *features*¹. While
51 the method is most frequently used to study the semantic representations of concrete
52 concepts and categories (McRae, Cree, Seidenberg, & McNorgan, 2005; Rosch & Mervis,
53 1975; Smith, Shoben, & Rips, 1974), it has also been used for other types of concepts,
54 corresponding to verbs (Vinson & Vigliocco, 2008), events, and abstract concepts (Lebani,
55 Lenci, & Bondielli, 2016; Recchia & Jones, 2012; Wiemer-Hastings & Xu, 2005).

56 On the one hand, many studies adopted the property generation task itself to make
57 inferences about word meaning and its computation (Recchia & Jones, 2012; Santos,
58 Chaigneau, Simmons, & Barsalou, 2011; Wiemer-Hastings & Xu, 2005; Wu & Barsalou,
59 2009). On the other hand, researchers employed the property listing task in order to provide

¹Throughout this article, features will be distinguished from cues using angular brackets and italic font.

60 other researchers with a tool of standardized word stimuli and relative semantic measures.
61 Indeed, based on data obtained from the property production task, it is then possible to
62 calculate numerous measures and distributional statistics both at the feature and the
63 concept level. For example, these feature data can be used to determine the semantic
64 similarity/distance between concepts, often by calculating the feature overlap or number of
65 shared features between concepts (Buchanan, Valentine, & Maxwell, 2019; McRae et al.,
66 2005; Montefinese, Vinson, & Ambrosini, 2018; Montefinese, Zannino, & Ambrosini, 2015;
67 Vigliocco, Vinson, Lewis, & Garrett, 2004), or how different types (Kremer & Baroni, 2011;
68 Zannino et al., 2006a) and dimensions of feature informativeness, such as, distinctiveness
69 (Duarte, Marquié, Marquié, Terrier, & Ousset, 2009; Garrard, Lambon Ralph, Hodges, &
70 Patterson, 2001), cue validity (Rosch & Mervis, 1975), relevance (Sartori & Lombardi, 2004),
71 semantic richness (Pexman, Hargreaves, Siakaluk, Bodner, & Pope, 2008), and significance
72 (Montefinese, Ambrosini, Fairfield, & Mammarella, 2014) are distributed across concepts.

73 Efficient ways to collect data online have boosted the availability of large feature listing
74 data sets. These semantic feature norms are now available across different languages: Dutch
75 (De Deyne et al., 2008; Ruts et al., 2004), English (Buchanan, Holmes, Teasley, & Hutchison,
76 2013; Buchanan et al., 2019; Devereux, Tyler, Geertzen, & Randall, 2014; Garrard et al.,
77 2001; McRae et al., 2005; Vinson & Vigliocco, 2008), German (Kremer & Baroni, 2011),
78 Italian (Catricalà et al., 2015; Kremer & Baroni, 2011; Montefinese, Ambrosini, Fairfield, &
79 Mammarella, 2013; Zannino et al., 2006b), Portuguese (Marques, Fonseca, Morais, & Pinto,
80 2007), and Spanish (Vivas, Vivas, Comesaña, Coni, & Vorano, 2017) as well as for blind
81 participants (Lenci, Baroni, Cazzolli, & Marotta, 2013). However, these norms vary
82 substantially in the procedure of data collection and their pre-processing, and this does not
83 facilitate performing cross-language comparisons and, thus, making inferences about how
84 semantic representations are generalizable across languages.

85 First, there is a lack of agreement in the instructions provided to the participants.

86 Indeed, while some studies use an open-ended verbal feature production (Buchanan et al.,
87 2013, 2019; De Deyne et al., 2008; Montefinese et al., 2013) where participants can list the
88 features related to the concept with any kind of semantic relation, other studies use a
89 constrained verbal feature production (Devereux et al., 2014; Garrard et al., 2001) where
90 participants were instructed to use specific semantic relations between cue concept and
91 features, such as, for example, *<is ...>*, *<has ...>*, *<does ...>*, *<made of ...>*, and so
92 forth. Moreover, authors could instruct the participants to produce a single word as a
93 feature instead of a multiple-word description. This latter case could also determine a
94 problem on subsequent coding steps that affect the identification of pieces of information.
95 For example, if the participant listed the feature *<has four wheels>* for the concept CAR,
96 there is no consensus if this feature should be divided into *<has wheels>* and *<has four*
97 *wheels>*, under the assumption that the participant provided two pieces of information, or
98 rather if it should be considered as a unique feature. Second, some authors gave a time limit
99 to provide the features descriptions (Kremer & Baroni, 2011; Lenci et al., 2013; Marques et
100 al., 2007) or a limited number of features to be listed (De Deyne et al., 2008), with a possible
101 influence on a number of feature-based measures (e.g., semantic richness or distinctiveness).

102 Because the feature listing task is a verbal task and language is very productive (i.e.,
103 the same feature can be expressed in many different ways), few features will be listed in
104 exactly the same way across participants. To be able to derive reliable quantitative measures,
105 nearly all studies specify a series of pre-processing steps to group verbal utterances about the
106 same underlying conceptual property together. The main problem is that there is no
107 agreement about how to code/pre-process data derived from the feature listing task.
108 Recoding features is sometimes done in manually (McRae et al., 2005) whereas others use
109 semi-automatic procedures, especially for larger datasets (Buchanan et al., 2019). Further
110 points of debate are related to the inclusion/exclusion of certain types of responses. For
111 example, unlike previous semantic norms (McRae et al., 2005; Montefinese et al., 2013; Vivas
112 et al., 2017), Buchanan et al. (2019) included idiosyncratic features (features produced only

113 by one or a few number of participants) if they were in the top listed features, ambiguous
114 words (words with multiple meanings), and created a special coding for affixes of the root
115 words. Moreover, they discarded stop words, such as, the, an, of, and synonyms were treated
116 as different entries.

117 While hand-coding features leads to features that concise, easily interpretable, and
118 highly predictive of semantic behavior, the increasing scale of recent studies and more
119 powerful natural language processing techniques make automatic procedures an attractive
120 alternative for assistance in processing language data. Moreover, building standard
121 automatic procedures to process feature-listing data would not only add transparency to the
122 process but would also reduce human errors and allow a generalization of the data across
123 languages. For the first time, in this study, we propose an automatic procedure to code the
124 raw feature data derived from a semantic feature listing task. The next sections provide a
125 tutorial on how raw feature data might be processed to a more compact feature output. The
126 tutorial is written for *R* and is fully documented, such that users can adapt it to their
127 language of choice (<https://github.com/doomlab/FLT-Primer>). Figure 1 portrays the
128 proposed set of steps including spell checking, lemmatization, exclusion of stop words, and
129 final processing in a multi-word sequence approach or a bag of words approach. After
130 detailing these steps, the final data form will evaluated and compared to previous norms to
131 determine the usefulness of this approach.

132 **Materials and Data Format**

133 You can load the entire set of libraries for this tutorial as shown below using
134 `dependencies.R` found online².

²A `packrat` project compilation is available on GitHub for reproducibility (Ushey, McPherson, Cheng, Atkins, & Allaire, 2018), and this manuscript was written in Rmarkdown with `papaja` (Aust & Barth, 2017).

```
library(here)
library(dplyr)
#Spelling
library(hunspell)
library(tidytext)
library(stringi)
#Lemmatization
library(koRpus)
library(koRpus.lang.en)
library(tokenizers)
#Stopwords
library(stopwords)
```

135 The data can then be imported with `importData.R`. Additionally, the answers from
136 participants may need to be normalized into lowercase for consistency.

```
# Importing the raw feature lists
X <- read.csv("../raw_data/tidy_words.csv", stringsAsFactors = F)
## Lower case to normalize
X$feature_response <- tolower(X$feature_response)
```

137 The data for this tutorial includes 16,544 unique concept-feature responses for 226
138 concepts from Buchanan et al. (2019). The concepts were taken from McRae et al. (2005),
139 Vinson and Vigliocco (2008), and Bruni, Tran, and Baroni (2014). The concepts include 185
140 nouns, 25 verbs, and 16 adjectives. The concepts were both abstract and concrete, and to
141 describe the concepts, the concreteness ratings collected by Brysbaert, Warriner, and
142 Kuperman (2014) can be used. In their study, they asked participants to rate words on a
143 scale ranging from 1 - abstract (language-based) - to 5 - concrete (experience-based) -
144 concepts. Nouns were rated as most concrete: $M = 4.59$ ($SD = 0.52$), followed by adjectives:
145 $M = 3.78$ ($SD = 0.81$), and verbs: $M = 3.57$ ($SD = 0.79$). The feature listing data consist of
146 a text file where concept-feature observation is a row and each column is a variable. An
147 example of these raw data are shown in Table 1, where the `cue` column is the cue, and the
148 `feature_response` column denotes a single participant's response. The original data can be
149 found at <https://osf.io/cjyzw/>.

150 The data was collected using the instructions provided by McRae et al. (2005),
151 however, in contrast to the suggestions for consistency detailed above (Devereux et al., 2014),
152 each participant was simply given a large text box to include their answer. Each answer
153 includes multiple embedded features, and the tutorial proceeds to demonstrate potential
154 processing addressing the additional challenges in unstructured data of this nature. Figure 1
155 portrays the suggested data processing steps. With structured data entry for participants
156 (e.g., asking participants to type one feature on each line), the multi-word sequence step
157 would be implemented within the data collection design, rather than post-processing. This
158 tutorial presents the more difficult scenario to be applicable to more data collection methods.

159 Spelling

160 The first step (see Figure 1) in processing the features consists of identifying and
161 replacing spelling mistakes. Spell checking can be automated with the `hunspell` package in
162 *R* (Ooms, 2018) using `spellCheck.R`. Each `feature_response` can be checked for
163 misspellings across an entire column of answers, which is in the `X` dataset. Because
164 participants were recruited in the United States, we used the American English dictionary.
165 The `hunspell` vignettes provide details on how to import your own dictionary for
166 non-English languages. The choice of dictionary should also normalize between multiple
167 variants of the same language, for example, the `"en_GB"` would convert to British English
168 spellings.

```
# Extract a list of words
tokens <- unnest_tokens(tbl = X, output = token, input = feature_response)
wordlist <- unique(tokens$token)
# Spell check the words
spelling.errors <- hunspell(wordlist)
spelling.errors <- unique(unlist(spelling.errors))
spelling.sugg <- hunspell_suggest(spelling.errors, dict = dictionary("en_US"))
```

169 The result from the `hunspell()` function is a list object of spelling errors for each row

170 of data. For example, when responding to APPLE, a participant wrote *<fruit, grocery store,*
171 *orchard, red, green, yelloe, good with peanut butter, good with caramell>*, and the spelling
172 errors were denoted as *<yelloe>* and *<caramell>*. After checking for errors, the
173 `hunspell_suggest()` function was used to determine the most likely replacement for each
174 error. For *<yelloe>*, both *<yellow>* and *<yell>* were suggested, and *<caramel>* and
175 *<camel>* were suggested for *<caramell>*.

176 Answers are provided in the most probable order, therefore, the first suggestion is
177 selected as the correct answer. These answers are compiled into a spelling dictionary, which
178 is saved for reproducibility and can be used to manually check the validity of the suggestions
179 in a final (optional) step. In addition to the hunspell dictionary, an auxiliary dictionary with
180 pre-coded error responses and corrections could also be added at this stage to catch any false
181 positives by adding entries to the `spelling.dict`. For example, by examining
182 `spelling.dict`, we found entries that would need to be corrected: *tast* became *tacit*, *frends*
183 became *fends*, and *musles* became *mules*. Since the spelling dictionary is saved this will
184 facilitate the additional step of manually examining the output for incorrect suggestions and
185 to add their own corrections. This file could then be reloaded and used in the step below to
186 provide adjusted spelling corrections. Other paid alternatives, such as Microsoft's Bing Spell
187 Check, can be a useful avenue for datasets that may contain brand names (i.e, *apple* versus
188 *Apple*) or slang terms and provides context sensitive corrections (e.g., keeping *Apple* as a
189 response to computer, but not as a response to green).

```
# Pick the first suggestion
spelling.sugg <- unlist(lapply(spelling.sugg, function(x) x[1]))
spelling.dict <- as.data.frame(cbind(spelling.errors,spelling.sugg))
spelling.dict$spelling.pattern <- paste0("\\b", spelling.dict$spelling.errors, "\\b")
# Write out spelling dictionary
write.csv(x = spelling.dict, file = "../output_data/spelling.dict.csv",
         fileEncoding = "utf8", row.names = F)
```

190 As noted, data was collected with a large text box, allowing participants to list

191 multiple features to the target cue. Participants often used extra spacing, tabs or other
 192 punctuation to denote separate answers to the cue. The `unnest_tokens()` function from
 193 `tidytext` can be used to split their answers into separate response lines and `trimws()` to
 194 remove all extra white spaces (De Queiroz et al., 2019).

```
# Parse features
tokens <- unnest_tokens(tbl = X, output = token,
  input = feature_response, token = stringr::str_split,
  pattern = " |\\, |\\.|\\.|\\;")
tokens$token <- trimws(tokens$token,
  which = c("both", "left", "right"),
  whitespace = "[ \\t\\r\\n]")
```

195 To finalize our data cleaning, we can remove blank lines, and use
 196 `stri_replace_all_regex()` is used to replace the spelling errors with their corrections
 197 from the `stringi` package (Gagolewski & Tartanus, 2019). If the `spelling.dict` output file
 198 was manually edited, it can be (re)loaded here with `read.csv` to update with the adjusted
 199 spelling corrections³. The spell checked dataframe is then output to a comma delimited file
 200 to preserve each workflow step.

```
# Remove empty features
tokens <- tokens[!tokens$token == "", ]
tokens$corrected <- stri_replace_all_regex(str = tokens$token,
  pattern = spelling.dict$spelling.pattern,
  replacement = spelling.dict$spelling.sugg,
  vectorize_all = FALSE)

# Rename columns
tokens <- tokens %>%
  rename(feature = corrected) %>%
  select(cue, feature)

# Write processed file
write.csv(x = tokens, file = "../output_data/spellchecked.features.csv",
  fileEncoding = "utf8", row.names = F)
```

³For transparency, the updated csv file should be renamed, which also practically keeps one from overwriting their adjustments if they rerun their code. The csv should be loaded as `spelling.dict` to continue with the code below.

201 Lemmatization

202 The next step groups different word forms that share the same lemma. The process of
 203 lemmatizing words uses a trained dictionary to convert all tokens part of a lexeme set (i.e.,
 204 all words forms that have the same meaning, *am*, *are*, *is*) to a common lemma (i.e., *be*)⁴.
 205 Lemmatization is performed using the `TreeTagger` program (Schmid, 1994) and
 206 implemented through the `koRpus` package in *R* (Michalke, 2018). `TreeTagger` is a trained
 207 tagger designed to annotate part of speech and lemma information in text, and parameter
 208 files are available for multiple languages. We will create a unique set of tokenized words to
 209 lemmatize to speed computation, as shown in `lemmatization.R`.

```
# Open the spell checked data
X <- read.csv("../output_data/spellchecked.features.csv", stringsAsFactors = F)
# Extract the list of updated tokens
tokens <- unnest_tokens(tbl = X, output = word, input = feature)
cuelist <- unique(tokens$cue)
```

210 The `treetag()` function calls the installation of `TreeTagger` to provide part of speech
 211 tags and lemmas for each token. Importantly, the `path` option should be the directory of the
 212 `TreeTagger` installation.

```
# Create a dataframe for lemmas
tokens.tagged <- data.frame(doc_id=character(),
                           token=character(),
                           wclass=character(),
                           lemma=character(),
                           stringsAsFactors=FALSE)
# Loop over cues and create lemmas + POS tags
```

⁴We mainly focus on lemmatization and do not proceed stemming the word because it introduces additional ambiguity. More specifically, stemming involves processing words using heuristics to remove affixes or inflections, such as *ing* or *s*. The stem or root word may not reflect an actual word in the language, as simply removing an affix does not necessarily produce the lemma. For example, in response to AIRPLANE, `<flying>` can be easily converted to `<fly>` by removing the *ing* inflection. However, this same heuristic converts the feature `<wings>` into `<w>` after removing both the *s* for a plural marker and the *ing* participle marker.

```

for (i in 1:length(cuelist)){
  temp.tag <- suppressWarnings(
    suppressMessages(
      treetag(c(X$feature[X$cue == cuelist[i]], "NULL"),
        treetagger="manual", format="obj",
        TT.tknz=FALSE, lang="en", doc_id = cuelist[i],
        # These parameters are based on your computer
        TT.options=list(path="-/TreeTagger", preset="en"))))
  temp.tag <- temp.tag@TT.res %>%
    mutate_if(is.factor, as.character)
  tokens.tagged <- tokens.tagged %>%
    bind_rows(temp.tag %>%
      select(doc_id, token, wclass, lemma))
}

```

213 This function returns a tagged corpus object, which can be converted into a dataframe
 214 of the token-lemma information. TreeTagger will return <unknown> for unknown values and
 215 @card@ for numbers, and these values were replaced with the original token. Table 2
 216 portrays example results from TreeTagger.

```

tokens.tagged <- tokens.tagged %>%
  rename(cue = doc_id, feature = token, pos = wclass)
# Clean up unknown lookups
tokens.tagged$lemma[tokens.tagged$lemma == "<unknown>"] <- tokens.tagged$feature[tokens.tagged$lemma == "<unknown>"]
tokens.tagged$lemma[tokens.tagged$lemma == "@card@"] <- tokens.tagged$feature[tokens.tagged$lemma == "@card@"]
tokens.tagged$lemma <- tolower(tokens.tagged$lemma)
# Write processed file
write.csv(x = tokens.tagged, file = "../output_data/lemmatized.features.csv",
  fileEncoding = "utf8", row.names = F)

```

217 Stopwords

218 As shown in Figure 1, the next stage of processing would be to exclude stopwords, such
 219 as *the*, *of*, *but*. The `stopwords` package (Benoit, Muhr, & Watanabe, 2017) includes a list of
 220 stopwords for more than 50 languages. At this stage, the `feature` (original tokens, not
 221 lemmatized) or `lemma` (lemmatized tokens) column can be used depending on researcher

222 selection. This code is included in `stopWordRemoval.R`. Within the `filter` command, we
 223 have excluded all lemmas in the stopword list provided by the `stopwords` library. Using
 224 `stopwords(language = "en", source = "snowball")`, one can view the stopword list
 225 and edit it for their own needs.

```
# Open the lemmatized data
X <- read.csv("../output_data/lemmatized.features.csv", stringsAsFactors = F)
# Remove punctuation and stopwords from lemmas
X$lemma <- gsub("\\\\-", " ", X$lemma)
X$lemma <- gsub("^$|\\002", NA, trimws(X$lemma))
X.nostop <- X %>%
  filter(!grepl("[[:punct:]]", lemma)) %>%
  filter(!lemma %in% stopwords(language = "en", source = "snowball")) %>%
  filter(!is.na(lemma))
# Write processed file
write.csv(x = X.nostop, file = "../output_data/nostop.lemmas.csv",
         fileEncoding = "utf8", row.names = F)
```

226 Multi-word Sequences

227 Multi-word sequences are often coded to mimic a Collins and Quillian (1969) semantic
 228 network, where words are nodes and edges are labelled with relations such as “is-a” or
 229 “has-a”. Some instructions specify the use of specific relation types (Devereux et al., 2014;
 230 Garrard et al., 2001), in which case pre-encoded the following step can be omitted. A
 231 potential solution for processing unstructured data involves identifying patterns that mimic
 232 “is-a” and “has-a” strings. Examples of such an approach is the Strudel model (Baroni,
 233 Murphy, Barbu, & Poesio, 2010) in which meaningful relations are grouped together using a
 234 small set of highly specific regular expressions. An examination of the coding in McRae et al.
 235 (2005) and Devereux et al. (2014) indicates that the feature tags are often adverb-adjective
 236 (*<usually-sweet>*), verb-noun (*<made-wood>*), or verb-adjective-noun
 237 (*<requires-lighting-source>*) sequences. Using `TreeTagger` on each concept’s answer set, we
 238 can obtain the parts of speech in context for each lemma. With `dplyr` (Wickham, Francios,

239 Henry, Muller, & Rstudio, 2019), new columns are added to tagged data to show all bigram
 240 and trigram sequences. All adverb-adjective, verb-noun, and verb-adjective-noun
 241 combinations are selected, and any words not part of these multi-word sequences are treated
 242 as unigrams. Finally, the `count()` function is used to tabulate the final count of n-grams
 243 and their frequency (`multiwordSequences.R`).

```
# Open the no stop words data
X <- read.csv("../output_data/nostop.lemmas.csv", stringsAsFactors = F)
# Combine lemmas and POS
X <- X %>%
  mutate(two.words = paste(lemma, lead(lemma), sep = " "),
         three.words = paste(lemma, lead(lemma),
                             lead(lemma, n = 2L), sep = " "),
         two.words.pos = paste(pos, lead(pos), sep = "."),
         three.words.pos = paste(pos, lead(pos),
                                  lead(pos, n = 2L), sep = "."))
# Patterns
adverb.adj <- grep("\\badverb.adj", X$two.words.pos)
verb.nouns <- grep("\\bverb.noun", X$two.words.pos)
verb.adj.nouns <- grep("\\bverb.adjective.noun", X$three.words.pos)
# Use combined and left over lemmas
X$combined.lemmas <- NA
X$combined.lemmas[c(adverb.adj, verb.nouns)] <- X$two.words[c(adverb.adj, verb.nouns)]
X$combined.lemmas[verb.adj.nouns] <- X$three.words[verb.adj.nouns]
X$combined.lemmas[-c(verb.nouns, verb.nouns+1, verb.adj.nouns,
                    verb.adj.nouns+1, verb.adj.nouns+2)] <- X$lemma[-c(verb.nouns, verb.nouns+1,
                                                                    verb.adj.nouns, verb.adj.nouns+1,
                                                                    verb.adj.nouns+2)]
# Create cue-lemma frequency
multi.words <- X %>%
  filter(!is.na(combined.lemmas)) %>%
  group_by(cue) %>%
  count(combined.lemmas)
# Write processed file
write.csv(x = multi.words, file = "../output_data/multi.nostop.lemmas.csv",
         fileEncoding = "utf8", row.names = F)
```

244 This procedure produces appropriate output, such as FINGERS *<have fingernails>*
 245 and COUCHES *<have cushions>*. One obvious limitation is the potential necessity to

246 match this coding system to previous codes, which were predominately hand processed.
247 Further, many similar phrases, such as the ones for ZEBRA shown below may require
248 flexible regular expressions to ensure that the different codings for *<is a horse>* are all
249 combined together, as shown in Table 3.

250 Bag of Words

251 To be able to evaluate the role of identifying multi-word sequences, we now describe an
252 approach where this information is not retained. This bag of words approach simply treats
253 each token as a separate feature to be tabulated for analysis. After stemming and
254 lemmatization, the data can be processed as single word tokens into a table of frequencies for
255 each cue word. The resulting dataframe is each cue-feature combination with a total for each
256 feature from `bagOfWords.R`. Table 4 shows the top ten most frequent responses to ZEBRA
257 given the bag of words approach.

```
# Open the no stop words data
X <- read.csv("../output_data/nostop.lemmas.csv", stringsAsFactors = F)
# Create cue-lemma frequency
bag.words <- X %>%
  group_by(cue) %>%
  count(lemma)
# Write processed file
write.csv(x = bag.words, file = "../output_data/bag.nostop.lemmas.csv",
         fileEncoding = "utf8", row.names = F)
```

258 Descriptive Statistics

259 The finalized data now represents a processed set of cue-feature combinations with
260 their frequencies for analysis. The data from Buchanan et al. (2019) was collected over
261 multiple years with multiple sample sizes. The sample size for each cue was then merged
262 with the finalized cue-feature information to control for differences in potential maximum

263 frequency. Table 5 includes descriptive statistics for the processed cue-feature set.

264 **Number of response types.** First, the number of cue-feature combinations was
265 calculated by taking the average number of cue-feature listings for each cue. Therefore, the
266 total number of features listed for ZEBRA might be 100, while APPLE might be 45, and
267 these values were averaged. More cue-feature combinations are listed for the multi-word
268 approach, due to differences in combinations for some overlapping features as shown in Table
269 3. The large standard deviation for both approaches indicates that cues have a wide range of
270 possible features listed. For example for the cue ZEBRA, we find a total of 196 features,
271 whereas for APPLE we find 134 features. We expect that the number of different response
272 tokens is a function of the number of times a cue was presented in the study. To investigate
273 this relation, we calculated the correlation provided represents the relation between sample
274 size for a cue and the number of features listed for that cue. These values are high and
275 positive, indicating that the number of unique features increases with each participant.

276 **Idiosyncratic responses.** Potentially, many of the cue-feature combinations could
277 be considered idiosyncratic. The next row of the table denotes the average number of
278 cue-feature responses listed by less than 10% of the participants. This percent of responses is
279 somewhat arbitrary, as each researcher has determined where the optimal criterion should be.
280 For example, McRae et al. (2005) used 16% or 5/30 participants as a minimum standard,
281 and Buchanan et al. (2019) recently used a similar criteria. Many cue-features are generated
282 by a small number of participants, indicating that these are potentially idiosyncratic or part
283 of long tailed distribution of feature responses with many low frequency features. The
284 advantage to the suggested data processing pipeline and code provided here is the ability of
285 each researcher to determine how many low-frequency features should be included.

286 **Response strength.** The next two lines of Table 5 indicate cue-feature combination
287 frequencies, such as the number of times ZEBRA <stripes> or APPLE <red> were listed by
288 participants. The percent of responses is the frequency divided by sample size for each cue,

289 to normalize over different sample sizes present in the data. These average frequency/percent
290 can be seen as a measure of response strength and were calculated for each cue, and then
291 averaged over all cues. The correlation represents the average response strength for each cue
292 related to the sample size for that cue. These frequencies are low, matching the results for a
293 large number of idiosyncratic responses. The correlation between frequency of response and
294 sample size is positive, indicating that larger sample sizes produce items with larger
295 frequencies.

296 Additionally, the correlation between response strength and sample size is negative,
297 suggesting that larger sample sizes are often paired with more items with smaller response
298 strengths. Figure 2 displays the correlations for the average cue-frequency responses and the
299 response strength by sample size. It appears that the relationship between sample size and
300 percent is likely curvilinear, rather than linear. The size of the points indicates the
301 variability (standard deviation of each cue word's average frequency or percent). Variability
302 appears to increase linearly with sample size for average frequency, however, it is somewhat
303 mixed for average percent. These results may imply a necessity to discuss common sample
304 sizes for data collection ($ns \sim 30$) to determine the optimal sample size for an appropriate
305 body of data for each cue word.

306 **Internal Comparison of Approach**

307 In this section, we show that the bag of words approach approximates the data from
308 McRae et al. (2005), Vinson and Vigliocco (2008), and Buchanan et al. (2019), thus
309 comparing data processed completely through code to datasets that were primarily hand
310 coded. These datasets were recoded in a bag of words approach, and the comparison between
311 all three is provided below. The multi-word sequence approach would be comparable if one
312 or more datasets used the same structured data collection approach or with considerable
313 hand coded rules for feature combinations. The data from open ended responses, such as the

314 Buchanan et al. (2019), could potentially be compared in the demonstrated multi-word
 315 sequence approach, if the raw data from other such projects were available.

316 Cosine similarity is often used as a measure of semantic similarity, indicating the
 317 feature overlap between two sets of cue-feature lists. For each concept or cue it provides an
 318 estimate of similarity based using a vector consisting of features with magnitudes
 319 corresponding to their frequency. The formula is identical to a Pearson product correlation
 320 when the vectors are centered to have mean zeros. First, matching feature (i) frequencies of
 321 cues A and B are multiplied and then summed, and this value is divided by products of the
 322 vector length of A and B for all features:

$$\frac{\sum_{i=1}^n A_i \times B_i}{\sqrt{\sum_{i=1}^n A_i^2} \times \sqrt{\sum_{i=1}^n B_i^2}}$$

323 As all of the frequencies are positive, these values can range from 0 (no overlap) to 1
 324 (perfect overlap). Two cosine values can be derived from the Buchanan et al. (2019) data:
 325 the raw cosine, which included all features as listed and the cosine for lemmatized responses.
 326 Each cue in the sample data for this project was compared to the corresponding cue in the
 327 Buchanan et al. (2019). The example participant responses provided in this tutorial are a
 328 subset of the Buchanan et al. (2019) data, and therefore, if the participant responses were
 329 processed in an identical fashion, the cosine values would be nearly 1. Additionally, if the
 330 processing detailed here matches the hand coding in the Buchanan et al. (2019), the overlap
 331 with the McRae et al. (2005) and Vinson and Vigliocco (2008) should be similar. These
 332 values were: original feature cosine = .54-.55, and lemmatized⁵ features = .66-.67. However,
 333 all previous datasets have been reduced by eliminating idiosyncratic features at various
 334 points, and therefore, we might expect that noise in the data would reduce the average

⁵These results were lemmatized by creating a lookup dictionary from the features listed in the Buchanan et al. (2019) norms.

335 cosine values.

336 Table 6 shows the role of using a cut-off for low-frequent or idiosyncratic responses by
337 calculating the cosine values when using varying cut-offs or stopword filtering. On the left,
338 the cosine values with stopwords are provided for both the original feature listed (i.e., no
339 lemmatization) and the lemmatized features. The right side of the table includes the cosine
340 values once stopwords have been removed. The removal of stopwords increases the match
341 between sets indicating how removing these terms can improve prediction. When stop words
342 were excluded, cosine values indicated somewhat comparable set of data, with lower values
343 for McRae et al. (2005) than previous results in the original feature sets. These values
344 portray that the data processed entirely in R produces a comparable set of results, albeit
345 with added noise of small frequency features.

346 External Comparison of Approach

347 The MEN dataset (Bruni et al., 2014) contains cue-cue pairs of English words rating
348 for similarity by Amazon Mechanical Turk participants for stimuli taken from the McRae et
349 al. (2005) feature norms. In their rating task, participants were shown two cue-cue pairs and
350 asked to select the more related pair of the two presented. Each pair was rated by 50
351 participants, and thus, a score of 50 indicates high relatedness, while a score of 0 indicates
352 no relatedness. The ratings for the selected set of cues provided in this analysis was 2 - 49
353 with an average rating of 25.79 ($SD = 12.00$). The ratings were compared to the cosine
354 calculated between cues using the bag of words method with and without stopwords. The
355 correlation between bag of words cosines with stopwords and the MEN ratings was $r = .54$,
356 95% CI [.42, .63], $N = 179$, indicating fair agreement between raters and cosine values. The
357 agreement between ratings and bag of word cosine values was higher when stopwords were
358 excluded, $r = .70$, 95% CI [.61, .76].

Discussion

Semantic feature listing tasks are used across various disciplines and are likely to remain an important source of information about the subjective meaning of concepts. In this article we have outlined a workflow to process large datasets where features consist of unstructured short propositions derived from written language. The advantage to this workflow is two-fold. First, science practices are shifting to open procedures and practices (Nosek et al., 2015), and reproducible research is key (Peng, 2011). Second, automated processing provides faster data analysis than hand-coded systems, and the ability to examine how processing steps affect results. We have shown that the automated procedure provides a comparable set of results to the hand-coded systems from Buchanan et al. (2019), McRae et al. (2005), and Vinson and Vigliocco (2008). The addition of specialized lemmas and other word exclusions (i.e., *<sometimes>*, *<usually>*, *<lot>* or idiosyncratic features) would provide more reduction, and thus, more overlap between hand and automated processing. Further, the automated data processing showed positive correlations with external subjective ratings of cue-cue relatedness in the MEN dataset. We suggest the workflow shown in Figure 1 and the suggested *R* code can provide a framework for researchers to use on their own data. In closing, the use of automated procedures will depend on specific use cases and cannot entirely replace careful human annotation (e.g. in the case of spell-checking). It can, however, greatly facilitate such checking.

Extending the approach. An attractive property of the subjective feature listing task is that it results in transparent representations. As a result, many researchers have taken additional steps to group specific types of knowledge together, depending on semantic relations (e.g., taxonomy relations) or their mapping onto distinct brain regions (Fairhall & Caramazza, 2013). Typically this involves applying a hand-crafted coding scheme, which requires a substantial effort. One of the common ontologies is the one developed by Wu and Barsalou (2009). The ontology is structured as a hierarchical taxonomy for coding categories

385 as part of the feature listing task. It has been used in several projects, notably the McRae et
386 al. (2005). Examples of the categories include taxonomic (synonyms, subordinates), entity
387 (internal components, behavior, spatial relations), situation (location, time), and
388 introspective properties (emotion, evaluation). Coding ontology may be best performed
389 systematically with look-up rules of previously decided upon factors, however, clustering
390 analyses may provide a potential avenue to explore categorizing features within the current
391 dataset. One limitation to this method the sheer size of the idiosyncratic features as
392 mentioned above, and thus, features smaller in number may be more difficult to group.

393 Potentially, a simple ontology can be mapped using an approach similar to Strudel
394 (structured dimension extraction and labeling, Baroni et al., 2010). Strudel is a corpus-based
395 semantic model wherein cue words are found in a large text corpus and matched to nouns,
396 verbs, and adjectives that appear near a concept. Using specific patterns of expected feature
397 listing, Baroni et al. (2010) were able build a model of English concepts and their properties
398 that aligned with semantic feature production norms. From this model, they were able to
399 cluster properties based on their lexical patterns. For example, if a sentence included the
400 phrase *fruit, such as an apple*, this lexical pattern would be classified as *such_as+right*,
401 indicating that the concept (apple) was found to the right of the property (fruit) with the
402 phrase such as connecting them. Using clustering, Baroni et al. (2010) were able to assign
403 four ontology labels to properties: part, category, location, and function. Using these results,
404 we can match 2279 of the bag of words features (5%). These features were predominately
405 parts (39.7), followed by function (30.7), location (24.2), and category (5.4). Table 7
406 indicates ten of the most frequent cue-feature pairs for each ontology label, excluding
407 duplicate features across cues. An examination of the top results indicates coherent labels
408 (parts: ZEBRA <*stripe*>, location: SHOE <*foot*>, and category: FURNITURE <*table*>);
409 however, there are also a few mismatches (location: SCISSORS <*cut*>, function: LEAF
410 <*green*>). This model represents an area in which one might begin to automate the labeling
411 process, likely combined with other pre-defined rule sets. Taxonomic labeling often

412 represents a large time demand on a researcher or team and by shifting the burden of the
413 taxonomic labeling to a semi-automated process, this time may be reduced. With the
414 addition of ontology labels to property norm data, theoretical questions about semantic
415 representation can be explored (Jones & Golonka, 2012; Santos et al., 2011).

416 **Some limitations.** So far we have not investigated to what extent the automatic
417 procedure leads to equally good representations for different types of concepts. More
418 specifically, abstract concepts tend to have a larger number of features. This result can be
419 explained by the larger context-variability of these concepts, but could also reflect to the
420 level of detail in the specific ontologies used to code these features (Recchia & Jones, 2012).
421 Pooling together features might improve the quality of the final representation, especially for
422 these types of concepts. Potentially, this might require additional steps in which features are
423 not only grouped based on surface properties but might also benefit from grouping
424 synonymous words. Within this framework, the properties could be added within a lookup
425 dictionary to further promote an open and transparent coding for data processing.

426 **Compliance with Ethical Standards**

427 *Funding:* This work was supported by the European Union's Horizon 2020 research
428 and innovation programme under the Marie Skłodowska-Curie Grant Agreement No. 702655
429 and by the University of Padua (SID 2018) to MM.

430 *Ethical Approval:* All procedures performed in studies involving human participants
431 were in accordance with the ethical standards of the institutional and/or national research
432 committee (include name of committee + reference number) and with the 1964 Helsinki
433 declaration and its later amendments or comparable ethical standards.

434 *Conflict of Interest:* The authors declare that they have no conflict of interest.

References

- 435
- 436 Aust, F., & Barth, M. (2017). papaja: Create APA manuscripts with R Markdown.
437 Retrieved from <https://github.com/crsh/papaja>
- 438 Baroni, M., Murphy, B., Barbu, E., & Poesio, M. (2010). Strudel: A corpus-based semantic
439 model based on properties and types. *Cognitive Science*, *34*(2), 222–254.
440 doi:10.1111/j.1551-6709.2009.01068.x
- 441 Benoit, K., Muhr, D., & Watanabe, K. (2017). stopwords: Multilingual stopword lists.
442 Retrieved from <https://cran.r-project.org/web/packages/stopwords/index.html>
- 443 Bruni, E., Tran, N. K., & Baroni, M. (2014). Multimodal distributional semantics. *Journal*
444 *of Artificial Intelligence Research*, *49*, 1–47. doi:10.1613/jair.4135
- 445 Brysbaert, M., Warriner, A. B., & Kuperman, V. (2014). Concreteness ratings for 40
446 thousand generally known English word lemmas. *Behavior Research Methods*, *46*(3),
447 904–911. doi:10.3758/s13428-013-0403-5
- 448 Buchanan, E. M., Holmes, J. L., Teasley, M. L., & Hutchison, K. A. (2013). English
449 semantic word-pair norms and a searchable Web portal for experimental stimulus
450 creation. *Behavior Research Methods*, *45*(3), 746–757. doi:10.3758/s13428-012-0284-z
- 451 Buchanan, E. M., Valentine, K. D., & Maxwell, N. P. (2019). English semantic feature
452 production norms: An extended database of 4436 concepts. *Behavior Research*
453 *Methods*, *51*(4), 1849–1863. doi:10.3758/s13428-019-01243-z
- 454 Caramazza, A., Laudanna, A., & Romani, C. (1988). Lexical access and inflectional
455 morphology. *Cognition*, *28*(3), 297–332. doi:10.1016/0010-0277(88)90017-0
- 456 Catricalà, E., Della Rosa, P. A., Plebani, V., Perani, D., Garrard, P., & Cappa, S. F. (2015).

- 457 Semantic feature degradation and naming performance. Evidence from
458 neurodegenerative disorders. *Brain and Language*, *147*, 58–65.
459 doi:10.1016/J.BANDL.2015.05.007
- 460 Collins, A. M., & Quillian, M. R. (1969). Retrieval time from semantic memory. *Journal of*
461 *Verbal Learning and Verbal Behavior*, *8*(2), 240–247.
462 doi:10.1016/S0022-5371(69)80069-1
- 463 Cree, G. S., & McRae, K. (2003). Analyzing the factors underlying the structure and
464 computation of the meaning of chipmunk, cherry, chisel, cheese, and cello (and many
465 other such concrete nouns). *Journal of Experimental Psychology: General*, *132*(2),
466 163–201. doi:10.1037/0096-3445.132.2.163
- 467 De Deyne, S., Verheyen, S., Ameel, E., Vanpaemel, W., Dry, M. J., Voorspoels, W., &
468 Storms, G. (2008). Exemplar by feature applicability matrices and other Dutch
469 normative data for semantic concepts. *Behavior Research Methods*, *40*(4), 1030–1048.
470 doi:10.3758/BRM.40.4.1030
- 471 De Queiroz, G., Hvitfeldt E, Keyes O, Misra K, Mastny T, Erickson J, ... Silge J. (2019).
472 tidytext: Text mining using 'dplyr', 'ggplot2', and other tidy tools. Retrieved from
473 <https://cran.r-project.org/web/packages/tidytext/index.html>
- 474 Devereux, B. J., Tyler, L. K., Geertzen, J., & Randall, B. (2014). The Centre for Speech,
475 Language and the Brain (CSLB) concept property norms. *Behavior Research*
476 *Methods*, *46*(4), 1119–1127. doi:10.3758/s13428-013-0420-4
- 477 Duarte, L. R., Marquié, L., Marquié, J. C., Terrier, P., & Ousset, P. J. (2009). Analyzing
478 feature distinctiveness in the processing of living and non-living concepts in
479 Alzheimer's disease. *Brain and Cognition*, *71*(2), 108–117.
480 doi:10.1016/j.bandc.2009.04.007

- 481 Fairhall, S. L., & Caramazza, A. (2013). Category-selective neural substrates for person- and
482 place-related concepts. *Cortex*, *49*(10), 2748–2757. doi:10.1016/j.cortex.2013.05.010
- 483 Farah, M. J., & McClelland, J. L. (1991). A computational model of semantic memory
484 impairment: Modality specificity and emergent category specificity. *Journal of*
485 *Experimental Psychology: General*, *120*(4), 339–357. doi:10.1037/0096-3445.120.4.339
- 486 Gagolewski, M., & Tartanus, B. (2019). stringi: Character string processing facilities.
487 Retrieved from <https://cran.r-project.org/web/packages/stringi/index.html>
- 488 Garrard, P., Lambon Ralph, M. A., Hodges, J. R., & Patterson, K. (2001). Prototypicality,
489 distinctiveness, and intercorrelation: Analyses of the semantic attributes of living and
490 nonliving concepts. *Cognitive Neuropsychology*, *18*(2), 125–174.
491 doi:10.1080/02643290125857
- 492 Humphreys, G. W., & Forde, E. M. (2001). Hierarchies, similarity, and interactivity in
493 object recognition: "category-specific" neuropsychological deficits. *The Behavioral and*
494 *Brain Sciences*, *24*(3), 453–476.
- 495 Jackendoff, R. (1992). *Semantic structures*. Boston, MA: MIT Press.
- 496 Jackendoff, R. (2002). *Foundations of language (brain, meaning, grammar, evolution)*.
497 Oxford, UK.: Oxford University Press.
- 498 Jones, L. L., & Golonka, S. (2012). Different influences on lexical priming for integrative,
499 thematic, and taxonomic relations. *Frontiers in Human Neuroscience*, *6*, 205.
500 doi:10.3389/fnhum.2012.00205
- 501 Kremer, G., & Baroni, M. (2011). A set of semantic norms for German and Italian.
502 *Behavior Research Methods*, *43*(1), 97–109. doi:10.3758/s13428-010-0028-x
- 503 Lebani, G. E., Lenci, A., & Bondielli, A. (2016). You are what you do: An empirical

- 504 characterization of the semantic content of the thematic roles for a group of Italian
505 verbs. *Journal of Cognitive Science*, *16*(4), 401–430. doi:10.17791/jcs.2015.16.4.401
- 506 Lenci, A., Baroni, M., Cazzolli, G., & Marotta, G. (2013). BLIND: A set of semantic feature
507 norms from the congenitally blind. *Behavior Research Methods*, *45*(4), 1218–1233.
508 doi:10.3758/s13428-013-0323-4
- 509 Marques, J. F., Fonseca, F. L., Morais, S., & Pinto, I. A. (2007). Estimated age of acquisition
510 norms for 834 Portuguese nouns and their relation with other psycholinguistic
511 variables. *Behavior Research Methods*, *39*(3), 439–444. doi:10.3758/BF03193013
- 512 McRae, K., Cree, G. S., Seidenberg, M. S., & McNorgan, C. (2005). Semantic feature
513 production norms for a large set of living and nonliving things. *Behavior Research
514 Methods*, *37*(4), 547–559. doi:10.3758/BF03192726
- 515 Michalke, M. (2018). koRpus: An R package for text analysis. Retrieved from
516 <https://cran.r-project.org/web/packages/koRpus/index.html>
- 517 Minsky, M. (1975). A framework for representing knowledge. In P. H. Winston (Ed.), *The
518 psychology of computer vision* (pp. 211–277). Winston, NY: McGraw Hill.
- 519 Montefinese, M., Ambrosini, E., Fairfield, B., & Mammarella, N. (2013). Semantic memory:
520 A feature-based analysis and new norms for Italian. *Behavior Research Methods*,
521 *45*(2), 440–461. doi:10.3758/s13428-012-0263-4
- 522 Montefinese, M., Ambrosini, E., Fairfield, B., & Mammarella, N. (2014). Semantic
523 significance: a new measure of feature salience. *Memory & Cognition*, *42*(3), 355–369.
524 doi:10.3758/s13421-013-0365-y
- 525 Montefinese, M., Vinson, D., & Ambrosini, E. (2018). Recognition memory and featural
526 similarity between concepts: The pupil's point of view. *Biological Psychology*, *135*,

- 527 159–169. doi:10.1016/J.BIOPSYCHO.2018.04.004
- 528 Montefinese, M., Zannino, G. D., & Ambrosini, E. (2015). Semantic similarity between old
529 and new items produces false alarms in recognition memory. *Psychological Research*,
530 *79*(5), 785–794. doi:10.1007/s00426-014-0615-z
- 531 Norman, D. A., & Rumelhart, D. E. (1975). *Explorations in cognition*. San Francisco, CA:
532 Freeman.
- 533 Nosek, B. A., Alter, G., Banks, G. C., Borsboom, D., Bowman, S. D., Breckler, S. J., ...
534 Yarkoni, T. (2015). Promoting an open research culture. *Science*, *348*(6242),
535 1422–1425. doi:10.1126/science.aab2374
- 536 Ooms, J. (2018). The hunspell package: High-Performance Stemmer, Tokenizer, and Spell
537 Checker for R. Retrieved from <https://cran.r-project.org/web/packages/hunspell/>
- 538 Peng, R. D. (2011). Reproducible research in computational science. *Science (New York,*
539 *N.Y.)*, *334*(6060), 1226–7. doi:10.1126/science.1213847
- 540 Pexman, P. M., Hargreaves, I. S., Siakaluk, P. D., Bodner, G. E., & Pope, J. (2008). There
541 are many ways to be rich: Effects of three measures of semantic richness on visual
542 word recognition. *Psychonomic Bulletin & Review*, *15*(1), 161–167.
543 doi:10.3758/PBR.15.1.161
- 544 Plaut, D. C. (2002). Graded modality-specific specialisation in semantics: A computational
545 account of optic aphasia. *Cognitive Neuropsychology*, *19*(7), 603–639.
546 doi:10.1080/02643290244000112
- 547 Recchia, G., & Jones, M. N. (2012). The semantic richness of abstract concepts. *Frontiers in*
548 *Human Neuroscience*, *6*, 315. doi:10.3389/fnhum.2012.00315
- 549 Rogers, T. T., Lambon Ralph, M. A., Garrard, P., Bozeat, S., McClelland, J. L., Hodges, J.

- 550 R., & Patterson, K. (2004). Structure and deterioration of semantic memory: A
551 neuropsychological and computational investigation. *Psychological Review*, *111*(1),
552 205–235. doi:10.1037/0033-295X.111.1.205
- 553 Rosch, E., & Mervis, C. B. (1975). Family resemblances: Studies in the internal structure of
554 categories. *Cognitive Psychology*, *7*(4), 573–605. doi:10.1016/0010-0285(75)90024-9
- 555 Ruts, W., De Deyne, S., Ameel, E., Vanpaemel, W., Verbeemen, T., & Storms, G. (2004).
556 Dutch norm data for 13 semantic categories and 338 exemplars. *Behavior Research*
557 *Methods, Instruments, & Computers*, *36*(3), 506–515. doi:10.3758/BF03195597
- 558 Saffran, E., & Sholl, A. (1999). Clues to the function and neural architecture of word
559 meaning. In P. Hagoort & C. Brown (Eds.), *The neurocognition of language*. Oxford
560 University Press.
- 561 Santos, A., Chaigneau, S. E., Simmons, W. K., & Barsalou, L. W. (2011). Property
562 generation reflects word association and situated simulation. *Language and Cognition*,
563 *3*(1), 83–119. doi:10.1515/langcog.2011.004
- 564 Sartori, G., & Lombardi, L. (2004). Semantic relevance and semantic disorders. *Journal of*
565 *Cognitive Neuroscience*, *16*(3), 439–452. doi:10.1162/089892904322926773
- 566 Schmid, H. (1994). Probabilistic part of speech tagging using decision trees.
567 doi:10.1.1.28.1139
- 568 Smith, E. E., Shoben, E. J., & Rips, L. J. (1974). Structure and process in semantic
569 memory: A featural model for semantic decisions. *Psychological Review*, *81*(3),
570 214–241. doi:10.1037/h0036351
- 571 Smith, E., & Medin, D. L. (1981). *Categories and concepts (Vol. 9)*. Cambridge, MA:
572 Harvard University Press.

- 573 Ushey, K., McPherson, J., Cheng, J., Atkins, A., & Allaire, J. (2018). packrat: A
574 dependency management system for projects and their R rackage dependencies.
575 Retrieved from <https://cran.r-project.org/web/packages/packrat/index.html>
- 576 Vigliocco, G., Vinson, D. P., Lewis, W., & Garrett, M. F. (2004). Representing the meanings
577 of object and action words: The featural and unitary semantic space hypothesis.
578 *Cognitive Psychology*, *48*(4), 422–488. doi:10.1016/j.cogpsych.2003.09.001
- 579 Vinson, D. P., & Vigliocco, G. (2008). Semantic feature production norms for a large set of
580 objects and events. *Behavior Research Methods*, *40*(1), 183–190.
581 doi:10.3758/BRM.40.1.183
- 582 Vivas, J., Vivas, L., Comesaña, A., Coni, A. G., & Vorano, A. (2017). Spanish semantic
583 feature production norms for 400 concrete concepts. *Behavior Research Methods*,
584 *49*(3), 1095–1106. doi:10.3758/s13428-016-0777-2
- 585 Wickham, H., Francios, R., Henry, L., Muller, K., & Rstudio. (2019). dplyr: A grammar of
586 data manipulation. Retrieved from
587 <https://cloud.r-project.org/web/packages/dplyr/index.html>
- 588 Wiemer-Hastings, K., & Xu, X. (2005). Content differences for abstract and concrete
589 concepts. *Cognitive Science*, *29*(5), 719–736. doi:10.1207/s15516709cog0000_33
- 590 Wu, L.-l., & Barsalou, L. W. (2009). Perceptual simulation in conceptual combination:
591 Evidence from property generation. *Acta Psychologica*, *132*(2), 173–189.
592 doi:10.1016/j.actpsy.2009.02.002
- 593 Zannino, G. D., Perri, R., Pasqualetti, P., Caltagirone, C., & Carlesimo, G. A. (2006a).
594 Analysis of the semantic representations of living and nonliving concepts: A
595 normative study. *Cognitive Neuropsychology*, *23*(4), 515–540.
596 doi:10.1080/02643290542000067

- 597 Zannino, G. D., Perri, R., Pasqualetti, P., Caltagirone, C., & Carlesimo, G. A. (2006b).
598 (Category-specific) semantic deficit in Alzheimer's patients: The role of semantic
599 distance. *Neuropsychologia*, *44*(1), 52–61. doi:10.1016/j.neuropsychologia.2005.04.008

Table 1

Example of Data Formatted for Tidy Data

Cue	Participant Answer
airplane	you fly in it its big it is fast they are expensive they are at an airport you have to be trained to fly it there are lots of seats they get very high up
airplane	wings engine pilot cockpit tail
airplane	wings it flys modern technology has passengers requires a pilot can be dangerous runs on gas used for travel
airplane	wings flys pilot cockpit uses gas faster travel
airplane	wings engines passengers pilot(s) vary in size and color
airplane	wings body flies travel

Table 2

Lemma and Part of Speech (POS) Information from TreeTagger

Cue	Feature	POS	Lemma
airplane	is	verb	be
airplane	fast	adverb	fast
airplane	they	pronoun	they
airplane	are	verb	be
airplane	expensive	adjective	expensive
airplane	they	pronoun	they

Table 3

Multi-Word Sequence Examples for Zebra

Cue	Combined Lemmas	N
zebra	horse	27
zebra	horse like	1
zebra	look similar horse	1
zebra	relate horse	2
zebra	resemble small horse	1
zebra	stripe similar horse	1

Table 4

Bag of Words Examples for Zebra

Cue	Lemma	N
zebra	stripe	64
zebra	black	63
zebra	white	61
zebra	animal	54
zebra	horse	32
zebra	africa	28
zebra	zoo	22
zebra	leg	20
zebra	life	20
zebra	eat	17

Table 5

Descriptive Statistics for Multi-word Sequences and Bag-of-words Approaches

Statistics	Multi-Word Sequences			Bag of Words		
	<i>Mean</i>	<i>SD</i>	<i>r</i>	<i>Mean</i>	<i>SD</i>	<i>r</i>
Number of Cue-Features	192.27	99.14	.78	173.44	77.21	.67
Frequency of Idiosyncratic Response	183.29	97.38	.80	160.57	74.26	.69
Frequency of Cue-Feature Response	2.09	3.39	.65	2.70	4.76	.83
Percent of Cue-Feature Response	3.41	5.10	-.64	4.30	4.76	-.62

Note. The correlation (*r*) represents the relation between frequency of response and sample size.

Table 6

Cosine Overlap with Previous Data Collection

Statistic	With Stopwords		No Stopwords	
	Original	Translated	Original	Translated
B Mean	.55	.58	.69	.74
B SD	.16	.16	.16	.15
M Mean	.33	.50	.39	.59
M SD	.15	.13	.18	.13
V Mean	.50	.50	.60	.59
V SD	.18	.18	.18	.19

Note. Translated values are hand coded lemmatization from Buchanan et al. (2019). B: Buchanan et al. (2019), M: McRae et al. (2005), V: Vinson & Vigliocco (2008). *N* values are 226, 61, and 68 respectively.

Table 7

Top Ten Ontology Labels

Parts	Function	Location	Category
brush use	brush hair	scissors cut	flute instrument
lawn grass	river water	snow cold	snow white
snail shell	branch tree	farm land	elephant animal
river stream	chair sit	cabin wood	cabbage green
radio music	leaf plant	rocket space	dagger knife
elephant trunk	kitchen food	breakfast day	apple fruit
zebra stripe	hammer nail	stone rock	hammer tool
river flow	garden flower	bacon pig	lion king
door open	oven cook	shoe foot	cabbage vegetable
dragon fire	leaf green	toy play	furniture table

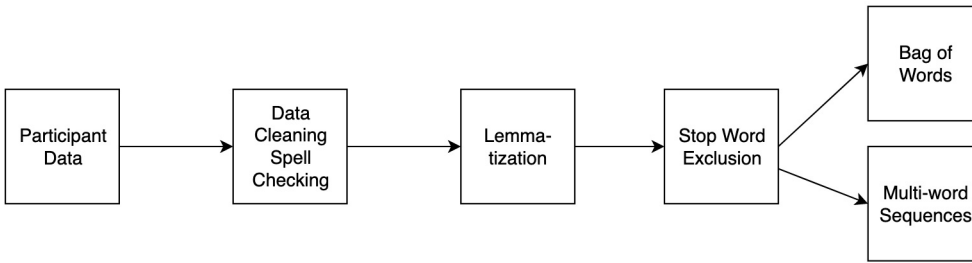


Figure 1. Flow chart illustrating how feature lists are recoded to obtain a standard feature format.

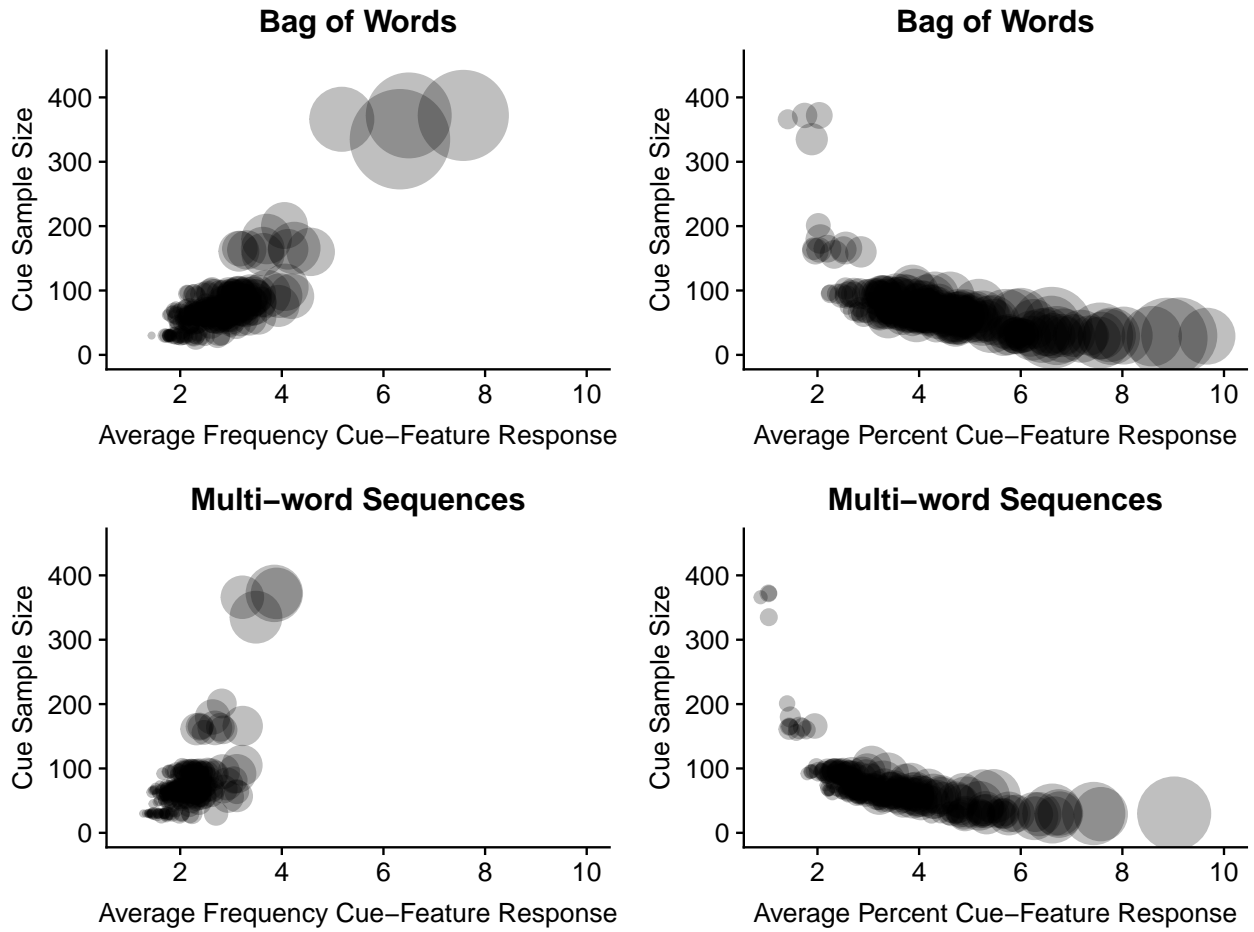


Figure 2. Correlation of sample size with the average cue-feature frequency (left) and percent (right) of response for each cue for both processing approaches. Each point represents a cue word, and the size of the point indicates the variability of the average frequency (left) or percent (right).