

# MorphoLex: A derivational morphological database for 70,000 English words

Claudia H. Sánchez-Gutiérrez<sup>1</sup> · Hugo Mailhot<sup>2,3</sup> · S. Hélène Deacon<sup>4</sup> · Maximiliano A. Wilson<sup>3</sup>

Published online: 9 November 2017  
© Psychonomic Society, Inc. 2017

**Abstract** Most of the new words a reader will find are morphologically complex. Also, theoretical models of language processing propose that morphology plays an important role in visual word processing. Nevertheless, studies on the subject show contradicting results that are difficult to reconcile. One factor that may explain this is the lack of a sizeable and reliable morphological database. As a consequence, there are enormous methodological differences in the way the values for morphological variables are calculated across studies. We present a sizeable and freely available database with six new variables for affixes and three for roots for 68,624 words from the English Lexicon Project. We further studied by means of regression models the influence of these new variables on the lexical decision latencies of 4,724 morphologically complex nouns that included one root and one suffix. Results showed that root frequency and suffix length had a facilitatory effect, whereas the percentage of more frequent words in the morphological family of the suffix had an inhibitory effect on

latencies. After controlling for collinearity, root family size, suffix family size, suffix P\*, and suffix frequency also had facilitatory effects. These results shed new light on the importance of suffix length and the frequency of the lexical competitors of the family of a suffix. This database represents a valuable resource for studies on the effect of morphology in visual word processing in English and can be found at <https://github.com/hugomailhot/MorphoLex-en>.

**Keywords** Morphology · Psycholinguistic variables · Lexical decision · Database · Visual word recognition

About 60 % of the new words a reader will encounter are morphologically complex (Angelelli, Marinelli, & Burani, 2014), in that they will have at least two morphemes (i.e., the smallest meaning-bearing linguistic units). These morphemes might be either roots (e.g., *happy*) or affixes (e.g., *un-*, *-ness*). Affixes can take the form of a prefix, if it comes before the root (e.g., *UNhappy*), or a suffix, when it comes after the root (e.g., *happiNESS*). Theoretical models of language processing propose that morphemes play an important role in word recognition and there is widespread evidence that this is the case (for a recent review, see Amenta & Crepaldi, 2012). However, speculations of just how morphological processing occurs present considerable divergences.

Many factors may explain the inconsistencies present in the relevant literature, which we will thoroughly describe in the next paragraphs. Arguably, one element is of paramount importance: the methodological differences in the way the values of the morphological variables are calculated among studies (Amenta & Crepaldi, 2012). This makes thorough interstudy comparisons difficult and highlights the need for a sizeable and reliable morphological database in which all these variables are computed in a unified manner. In order to overcome

---

**Electronic supplementary material** The online version of this article (<https://doi.org/10.3758/s13428-017-0981-8>) contains supplementary material, which is available to authorized users.

---

✉ Claudia H. Sánchez-Gutiérrez  
chsanchez@ucdavis.edu

✉ Maximiliano A. Wilson  
maximiliano.wilson@fmed.ulaval.ca

<sup>1</sup> Department of Spanish and Portuguese, University of California, Davis, CA, USA

<sup>2</sup> Department of Computer Science, University of California, Davis, CA, USA

<sup>3</sup> Centre de recherche CERVO - CIUSSS de la Capitale-Nationale, 2601, de la Canardière, bureau F-2445, Québec - G1J 2G3, Canada

<sup>4</sup> Dalhousie University, Halifax, Nova Scotia, Canada

this limitation, the first aim of this study is to present MorphoLex, a database of derivational morphological variables for each complex word in the complete English Lexicon Project (ELP, <http://elexicon.wustl.edu/>; Balota et al., 2007). The morphological variables included in this database are: (1) affix and root frequency, (2) affix and root family size, (3) percentage of words more frequent in the family size (PFMF) for affixes and roots, (4) P value of affix productivity, (5) P\*-value of affix productivity, and (6) affix length.

We based our computations on the ELP because it is the largest freely available English database that provides values of psycholinguistic variables and standardized behavioral data in the frame of an online search engine. Some morphological variables, such as base frequency or morphological family size, can be calculated from the CELEX database (Baayen, Piepenbrock, & Gullikers, 1995). However, the Hyperspace Analogue to Language (HAL) frequency counts from the ELP are based on a corpus of over 130 million tokens (Burgess & Livesay, 1998), as compared to the 17.9 million tokens corpus on which the CELEX is based. This difference in the magnitude of the corpus represents a direct impact on the calculation of the morphological indices. For instance, if some words from a morphological family are missing from a corpus, morphological family size and root and affix frequency counts will be lacking. Additionally, words that are hapaxes (i.e., that appear only once) in smaller corpora may appear more than once in corpora that include more varied sources. This will also influence the computation of affix productivity measures (i.e., P and P\*). Indeed, P and P\* respectively represent the probability that a given affix will be encountered in a hapax and the probability that a hapax contains a certain affix. Thus, for our purposes, a large corpus such as the ELP is better suited to calculate the precise morphological indices that we present here.

The second aim of this study is to illustrate how these morphological variables can be used in morphological processing experiments, by entering them as predictors of lexical decision (LD) reaction times (RTs) for the suffixed nouns from the ELP in a series of regression models. This second part of the study does not pretend to provide evidence in favour or against any specific model of morphological processing but rather to present an example of the type of analyses that can be performed with the newly calculated variables. Thus the results of the regression study will be discussed only in terms of the effects of the variables, without entering into theoretical considerations or implications of those results for processing models. Indeed, our aim is mainly to offer a research tool that will enhance the replicability of future studies and provide a single source of data to morphological processing researchers. In the next paragraphs, we present the selected morphological variables in the context of studies that have assessed their effects.

It has been repeatedly demonstrated that the higher the lexical (i.e., whole-word) frequency of monomorphemic words, the faster they are processed in LD and word-naming tasks (Oldfield & Wingfield, 1965; Scarborough, Cortese, & Scarborough, 1977; Whaley, 1978). Complex words, however, present a more intricate picture when it comes to frequency effects (McCormick, Brysbaert, & Rastle, 2009; Niswander-Klement & Pollatsek, 2006). As they are composed of multiple morphemes, not only their whole-word frequency, but also the frequencies of each one of their morphemes, may play a distinct role in their processing and potentially interact with each other (Andrews, Miller, & Rayner, 2004; Baayen, Dijkstra, & Schreuder, 1997; New, Brysbaert, Segui, Ferrand, & Rastle, 2004).

A common interpretation of frequency effects in complex word processing in LD tasks is the consideration of morpheme frequency effects as evidence for decomposition, and whole-word frequency effects as proof for direct whole-word processing (Colé, Beauvillain, & Segui, 1989; Ford, Davis, & Marslen-Wilson, 2010; Taft, 1979). Concretely, in the frame of parallel dual route models of lexical access (Burani & Laudanna, 1992; Marcolini, Traficante, Zoccolotti, & Burani, 2011; Vannest, Polk, & Lewis, 2005), evidence shows that low frequency complex words are processed faster when decomposed into their morphemic constituents previous to lexical access while higher frequency words are processed faster through whole word access (Lehtonen, Niska, Wande, Niemi, & Laine, 2006; Stemberger & MacWhinney, 1986). Alegre and Gordon (1999), for example, using a LD paradigm, established a frequency threshold at six tokens per million, beyond which words would be accessed faster directly in their whole-word representation. Other authors have even claimed that only newly encountered words need to be accessed via their morphemic constituents (Caramazza, Laudanna, & Romani, 1988), whereas all complex words that have been heard or read at least once are easily accessed as single lexical units.

Nonetheless, this idea of low whole-word frequency as a predictor of favored processing through the decompositional route has been challenged by McCormick, Brysbaert, et al. (2009), among others. They found equivalent facilitation for pairs of words where the primes were highly frequent (e.g., national-NATION), less frequent (e.g., notional-NOTION), or even pseudomorphological (corner-CORN). Such results point to an automatic decomposition process for all complex words or morphologically structured nonwords, independent of whole-word frequency (Longtin & Meunier, 2005; McCormick, Rastle, & Davis, 2009; Rastle & Davis, 2008; Rastle, Davis, & New, 2004).

These studies are all based on the opposition between the frequency of a derived word (e.g., *national*) and the frequency of its lexical base (e.g., *nation*). But, additionally, cumulative root frequency, as the sum of the frequencies of all the words

that share the same root (e.g., *human*, *humanity*, *humanist*, *humanism*, etc.), seems to play a role in the processing of complex words (Burani & Thornton, 2003; Caramazza et al., 1988; Luke & Christianson, 2011; Taft & Forster, 1975; Taft & Forster, 1976). When whole-word frequency is controlled (and generally maintained in the low ranges), only the words that have a high frequency root are decomposed, whereas those that include a low frequency root are processed as whole-words (Hay, 2001). The underlying rationale is that low frequency words are processed faster through a decompositional route, unless their constituent morphemes are too infrequent, in which case, the whole-word route would result, again, in faster processing. It is noteworthy that some authors have not found such a robust cumulative root frequency effect in the processing of complex words (Schreuder & Baayen, 1997; Sereno & Jongman, 1997), or have even found it to be inhibitory (Baayen, Tweedie, & Schreuder, 2002; Baayen, Wurm, & Aycocock, 2007). Also, Bradley (1979) found cumulative root frequency effects in both frequent and infrequent words, which indicates that morphological effects do not necessarily depend on the frequency of the whole word.

The incongruent results obtained in these studies have been tentatively interpreted as a reflection of different types of morphological structures. Colé et al. (1989), for instance, noticed that cumulative root frequency effects could only be observed in suffixed but not in prefixed words. This idea was further developed by Beauvillain (1996), as she analyzed eye-movement patterns on complex words in a semantic relatedness task. The results showed that cumulative frequency affected first fixation times in suffixed words, whereas it only had an effect on the second fixation for prefixed words, which indicates that cumulative root frequency effects are influenced by the morphological structure of complex words.

Another interpretation of the contradicting results observed in the study of base and root frequency effects states that only words that include highly productive affixes will be accessed through their morphemic constituents and, thus, present root frequency effects (Bertram, Schreuder, & Baayen, 2000). While Ford et al. (2010) provide tentative evidence for this idea, a methodological issue arises with the concept of productivity. Indeed, as Plag (2006) argues, researchers are lacking a common framework to study morpheme productivity effects, which refrains from the development of much-needed studies in this area. In this context, most empirical work has followed Baayen's definition of affix productivity as the probability that an affix appears in a hapax (i.e., a word that appears only once in a corpus) (Baayen, 2009; Baayen & Lieber, 1991). Baayen and Renouf (1996) established two complementary measures to quantify it: the P and P\* values. The former approximates the likelihood of the morpheme to appear in a hapax, while the latter approximates the morpheme's likelihood in the set of all hapaxes in a given corpus. Baayen et al. (2007) have been the only authors, to date, to

introduce the P value as a continuous measure of affix productivity in a regression design while no study, to the best of our knowledge, has manipulated P\* values. It is interesting to point here that, in Kuperman, Bertram, and Baayen (2010), the measure used for affixal productivity was the family size of the affix, instead of P or P\*. As our database includes P, P\* and family size indices for all affixes, the three variables can be used as predictors and can be compared in terms of their effects on lexical processing.

While base and root frequency effects have been widely studied, affix frequency has been generally neglected. This lack of attention to suffixes and prefixes is compelling, as the assumption that complex words can be accessed through their morphological constituents necessarily implies that their affixes should also play a central role (Burani, Dovetto, Thornton, & Laudanna, 1997). Interestingly, Burani and Thornton (2003) evaluated the effect of root versus suffix frequency in three LD experiments. High frequency suffixes slowed down decision times for pseudowords, indicating that pseudowords that include a frequent suffix are more difficult to reject than pseudowords that comprise low-frequency suffixes. The inverse effect with low-frequency Italian words was not observed. This indicates that highly frequent suffixes do not ensure faster lexical access through morphological decomposition, especially when they are part of words that include low-frequency roots. Root frequency, on the other hand, did exert a significant facilitative effect on RTs, independent of suffix frequency. Thus, the high frequency of a suffix only seems to favor morpheme-based lexical access when the root to which it is attached is highly frequent too.

Burani and Thornton (2003) come to interpret these data in the light of a difficulty to disentangle the unique contribution of token (i.e., the sum of frequency counts for each word that shares a common morpheme) and type frequency (i.e., the total number of different words that include that morpheme). Indeed, a high frequency suffix is also, generally, a suffix that appears in many different words. For instance, a highly frequent suffix in Italian could be attached to as many as 650 roots (Burani & Thornton, 2003). Thus, the fast recognition of a suffix does not, by itself, facilitate access to a specific word but rather to an unmanageable amount of competing lexical forms. In such case, the decomposition route does not result in faster processing than whole-word access, which may explain why suffix frequency effects are not observed in words that include a low-frequency root and a high-frequency suffix.

A similar discussion also arose concerning root frequency effects. Indeed, it has been hypothesized that roots that pertain to large families will also tend to have a higher cumulative root frequency. Thus, what has traditionally been interpreted as a root frequency effect could actually be driven by the size of the morphological family of the root (Bertram, Schreuder, et al., 2000; Schreuder & Baayen, 1997). However, Ford et al. (2010) found independent effects of root frequency and

morphological family size, stating that those are different variables that independently contribute to lexical access in complex words.

Morphological family size effects – as the effect of the number of words that share the same root – have been consistently found in LD tasks (Balling & Baayen, 2008; Bertram, Baayen, & Schreuder, 2000; De Jong, Schreuder, & Baayen, 2000; Ford et al., 2010; Moscoso del Prado Martín, Bertram, Häikiö, Schreuder, & Baayen, 2004). This indicates that words from a larger family are easier to access than words from smaller families, as the root is more easily accessed when it is activated through a broader network of words that share it.

Another important aspect of morphological family size effects is that they are modulated by the relative frequency of a given word compared to the rest of the members of the family (Colé et al., 1989; Meunier & Segui, 1999). Specifically, Meunier and Segui (1999) compared the reaction times in a LD task with complex words that had many higher frequency members in their morphological family and complex words that had few of them. The results showed that those with less high frequency competitors were responded to significantly faster than the others. This indicates that, once all morphological family members are activated, access to the target is easier when it has less higher frequency competitors. This effect has been found in auditory (Meunier & Segui, 1999) as well as written (Colé et al., 1989) modalities of stimuli presentation.

Finally, while variables based on type or token frequency counts have been studied repeatedly, affix length (i.e., the number of letters in an affix) has been less explored. However, the available results on affixal length effects are more consistent across studies when compared with those on the effects of frequency counts, which present numerous inconsistencies. Indeed, Laudanna and Burani (1995) proposed that words that have affixes that are more salient will more probably be decomposed into morphemes than words that have affixes that are less salient. One of the variables that significantly increases affixal salience is its length, as a longer affix will be more visually noticeable than a shorter one. Interestingly, Kuperman et al. (2010) confirmed this hypothesis by showing that whole-word frequency effects are less noticeable in words that include longer suffixes than in words with shorter ones. This observation indicates that affixal length is a relevant variable to include in studies that aim at assessing the interaction between whole-word and morphemic factors in word processing.

In sum, conflicting results are found among studies on the effects of different morphological variables on the processing of complex words. The key to answering these debates lies in the comprehensive study of the features of complex words in empirical studies, and this can only be done using a common and sizeable database in which all morphological variables are computed in the same way. That is why we created MorphoLex, a database that offers a comprehensive list of

prefixes, roots, and suffixes in English, as well as morphological values for the 79,672 words included in the Complete ELP. It must be noted here that the variables included in this study focus exclusively on token and type morphological frequency indices as well as affix length, but no measures that require subjective ratings or that are related to semantic or orthographic transparency have been included here. Studies that provide norms for this type of variables are time-consuming endeavours and thus are usually based on smaller scale databases. For example, in a recent study, Davies, Izura, Socas, and Dominguez (2016) collected age of acquisition, imageability and semantic distance norms for 2,204 English words. Collecting these types of norms for the ~70,000 words of the complete ELP is beyond the scope of the present database. Meanwhile, researchers interested in assessing the effects and interactions of objective and subjective morphological variables could combine data from Davies et al.'s study with data from MorphoLex.

### Word segmentation method

Of the 79,672 lexical entries of the complete database of the ELP available online (<http://elexicon.wustl.edu/>; Balota et al., 2007), 68,624 words are segmented into morphemes, with the following codes: << for prefixes, >> for suffixes, and {} for lexical bases. Bases sometimes contain, between curly brackets {}, prefixes and suffixes that are not annotated as such, but are still segmented from the root by a double dash (e.g. {reciproc-ity}, {reciproc-al}, {reciproc-ate}). For example, the word *miscalculations* is segmented as follows: <mis<{calcul-ate}>ion>>s>, where *mis-* is the prefix, *-ion* and *-s* are the suffixes, and *calculate* is the base, that includes the suffix *-ate* and the root *calcul*.

We performed a series of changes on these initial segmentations manually in order to homogenize them for the computation of the new morphological variables and to differentiate roots and affixes in the bases. We first removed inflectional suffixes such as *-s*, *-ing*, or *-ed*, as well as contractions such as *'ll*, *'s*. This first step was done in order to base our calculations of family size on lemmas and not on wordforms. We then normalized the treatment of neoclassical compounds, as several classical morphemes (e.g., *thermo*) appeared as roots in some words (e.g., {thermo}{form}) and as affixes in others (e.g., <thermo<{plastic}>), with no apparent theoretical or practical justification for those different classifications. Without normalization, the frequency counts of such morphemes in the database would be randomly split between their affix and root versions. We annotated as prefixes the classical morphemes that indicate either a position (e.g., *pre-*, *sub-*, *trans-*, *supra-*), a negation (e.g., *non-*), or a quantity (e.g., *ultra-*, *mega-*, *maxi-*, *bi-*, *tri-*, *deca-*), when such morphemes are not the only potential root of the word they are in.

Morphemes that are Latin or Greek nouns (e.g., *-cephal-*, *-thermo-*) were coded as roots, independently of their position in the word (Bauer & Nation, 1993). Thus, while *unicycle* is segmented as <uni<{(cycle)}>, owing to the first morpheme having a meaning of quantity, *unity*, having no other valid candidate roots, is segmented as {un}>ity>. This last example shows that there are still instances of morpheme counts that are split between different morphemic categories. However, the rules we applied greatly reduced the incidence of such cases, and eliminated the often arbitrary categorizations of classical morphemes.

Morphemes inside the bases identified in the ELP segmentations are only separated by dashes (e.g., {calcul-ate}) and not marked as roots or affixes. Thus, we selected for each of the bases at least one root among the segmented morphemes. After this manipulation, all roots inside and outside the bases were marked between parentheses. For example, <dis<{quiet-ude}> became <dis<{(quiet)>ude>, {clean-ly}>ness> became {(clean)>ly}>ness>, and <thermo<{plastic}> became {(thermo)}{(plastic)}. Following this stage, a systematic manual revision was performed in order to solve issues related to allomorphy and to eliminate the few pseudo-derivations (e.g., *corn* in *corner*) that remained from the original segmentation. We listed all the English affixes that also appeared as roots in the database and corrected their segmentations by hand, and did the same for roots that appeared as affixes. For instance, all English location adverbs and prepositions (e.g., *on*, *in*, *after*, *under*, *for*) were marked as prefixes when they preceded a root (e.g., <under<{(score)}>) and were only considered as roots when no other root was identifiable in the word (e.g., {(under)}), following the same criterion used for neoclassical morphemes. Finally, we created a list of allomorphs for all the suffixes (e.g., *-tion/-ation/-ition*), prefixes (e.g., *a-/al-/ac-*) and roots (e.g., *-cephal-/cephalo-, spectacle/spectacul-*) in the database. We then identified a canonical form for each set of allomorphs (e.g., *-ate* is the canonical form of *-ate*, *-uate*, *-cate*, *-iate*). This allowed us to compute morphological variables in which all allomorphs were counted as one single morpheme. Thus, the frequency counts of *-uate*, *-cate*, and *-iate* were added up towards calculating the frequency of the canonical morpheme *-ate*. Allomorphs that were orthographically identical but phonologically distinct (e.g., *reciproc-ity* vs. *reciproc-ate*) were not marked in any way for this study.

Homographs are not differentiated in the database if they pertain to the same morphemic category (i.e., prefix, root, suffix). Therefore, a suffix such as *-ar*, which can form nouns with the meaning of “person who” (e.g., *beggar*, *liar*, *burglar*) or adjectives that mean “pertaining to” (e.g., *solar*, *lunar*, *alveolar*), will be considered as one, independent of the word category in which they are embedded. However, our database does differentiate morphological homographs when they are encountered in different positions (i.e., suffixes and prefixes).

For example, the segment *al* will be counted as a prefix in *alchemy* and as a suffix in *chemical*.

### The new morphological variables

Based on these segmentations, we calculated six new morphological variables for affixes and three for the root. These variables are described in the following paragraphs.

*Morphological family size* is the number of word types in which a given morpheme is a constituent (Baayen, Feldman, & Schreuder, 2006). The family size of a morpheme was calculated by counting all its types in the ELP database. For instance, in the example {*attendance*, *pleasance*, *pleasure*, *appearance*}, the suffix *-ance* has a morphological family size of 3 {*attendance*, *pleasance*, *appearance*}, while the root *-pleas-* has a morphological family size of 2 {*pleasure*, *pleasance*}.

*Summed token frequency* is the summed frequency of all members in the morphological family of a morpheme. Thus, following the above example, the frequency of the root *-pleas-* would be the result of adding the frequency of the word *pleasure* and that of the word *pleasance*. The frequency count used for this calculation was the HAL frequency provided in the ELP.

*Affix productivity* was computed for prefixes and suffixes only, using two measures of productivity: *P* and *P\** (Baayen & Renouf, 1996). Both *P* and *P\** take values between 0 and 1. Values closer to 1 indicate high morphemic productivity, whereas values closer to 0 indicate low productivity. The *P* value for a morpheme in a corpus is defined as:

$$P_{C_m} = \frac{H_{C,m}}{STF_m}$$

where  $H_{C,m}$  is the total of all hapaxes in corpus *C* that contain morpheme *m*, and  $STF_m$  is the summed token frequency of morpheme *m*. We identified as hapaxes all words in the ELP with a HAL frequency value of 0 or 1. This first measure approximates the likelihood that a word containing morpheme *m* (i.e., the affix) is a new word. For example, based on its *P* value, *-ness* is among the 20 most productive suffixes due to the fact that out of the 1,243 words of its morphological family (that sum a total frequency of 181,553), 106 are hapaxes. The suffix *-al*, on the other hand, is less productive because, while its morphological family is composed of 1,431 words (with a total frequency of 4,704,731), only 43 words in the family are hapaxes. Thus, it is more likely that a word ending in *-ness* will be a newly invented word than is the case for words ending in *-al*.

*P\** is defined as follows:

$$P_{C_m}^* = \frac{H_{C,m}}{H_C}$$

where  $H_{C,m}$  is the same as above and  $H_C$  is the total of all hapax legomena in corpus  $C$ . This second and complementary measure approximates the likelihood that a new word contains the affix  $m$ . For example, *-itis* appears in 16 words (e.g., *meningitis*, *appendicitis*) but all of them have a frequency higher than 1, thus the  $P^*$  value for this suffix is 0. This means that the probability of a new word being created by adding *-itis* to a root is null, according to this calculation.

**Percentage of other words in the family that are more frequent (PFMF)** For each morpheme of each word, we computed their PFMF values by dividing the number of more frequent words in the family by the total number of members in the family, minus one. This last adjustment in the calculation was applied so that the variable would go from 0 to 100, where 0 means that no word in the family is more frequent and 100 means that all words in the family are more frequent. For example, *word*, *wordlessly*, and *wordiness* share a same root (i.e., *word*) and thus have an identical family size of 21. However, they all have a different PFMF: *word* does not have any more frequent competitor in the family and thus has a PFMF of 0 %, *wordiness* has 15 words that are more frequent in the family, which results in a 70 % PFMF, and *wordlessly* has 10 of those and has a PFMF of 45 %.

**Affix length** This variable indicates the number of letters of a particular affix. It is calculated for the canonic form of each affix and not for each of its allomorphs, as it is based on our final segmentation, which does not differentiate between allomorphs. Namely, *-ion* will always have a length of 3, even when it appears as *-tion* or *-ation*.

### The database

Each word in the database was tagged with a specific prefix-root-suffix (PRS) signature. This means that words that include one suffix and one root, but no prefix, share a 0-1-1 PRS signature (i.e., 0 for the number of prefixes, 1 for the number of roots and 1 for the number of suffixes), while words with two roots and a prefix will be tagged as 1-2-0 (i.e., 1 for the number of prefixes, 2 for the number of roots and 0 for the number of suffixes). The database is presented in an Excel file that is freely available at the following address: <https://github.com/hugomailhot/MorphoLex-en>.

Each PRS signature appears in different sheets that are titled with the corresponding PRS signature. This allows to directly access any specific subset of words depending on their morphological structure. The first page offers a list of all the variables and their corresponding headers, in order to facilitate the interpretation of the data. For each one of the morphemes on sheets 2 to 33, all the above mentioned morphological variables are provided in columns that are titled with the name of the variable, preceded by ROOT, PREF or

SUFF and a number (e.g., ROOT1, PREF2). That number indicates the situation of the morpheme in the word. For example, ROOT1 will be the first root in the word and PREF2 will be the second prefix in the word. In addition to the morphological variables for each morpheme of each word in the different PRS signatures, each word's ELP identification is provided, as well as its part of speech, PRS signature and number of morphemes. Sheets 34, 35, and 36 list all the prefixes, suffixes, and roots, respectively, organized by frequency. This will allow to obtain specific information about each morpheme, independent of the words it is in. This will be particularly useful when creating morphologically complex pseudowords.

### The influence of the morphological variables in Lexical Decision latencies

In order to exemplify how these new morphological variables (i.e. frequency, family size, P,  $P^*$ , PFMF, and affix length) can be used in a study on morphological processing, we extracted the LD RTs of morphologically complex nouns containing one suffix (i.e., a 0-1-1 PRS signature) from the ELP (Balota et al., 2007). For each of those words, we entered the values of their morphological variables and other relevant counts as predictors in a series of hierarchical regression models.

### Method

#### Material

Out of the 13,479 words of the ELP that had one root and one suffix (i.e., PRS 0-1-1), we selected only those that belonged to the noun category (from the variable POS, part of speech of the word, we chose only words with NN values). This left the database with 6,827 nouns. Then we eliminated all the words that had no information on LD latencies (RTs variable) in the ELP. This left the database with 5,678 nouns. Afterwards, we removed 237 words that had HAL frequency values of zero and 717 words for which no semantic values were available. The final database for the study was thus composed of 4,724 nouns. Table 1 shows the summary statistics for all the variables used in the LD study.

The psycholinguistic values for the 4,724 nouns for frequency, N-size and length in letters were obtained from the ELP online database (<http://ellexicon.wustl.edu>; Balota et al., 2007). We followed the recommendation of Balota et al. (2007) and preferred HAL frequency over other frequency estimates because it has been calculated from a sizeable corpus (~131 million words). While SUBTITLE frequency counts can be considered better estimates of frequency (Brysbaert & New, 2009), we chose HAL frequency due to the higher number of words from the ELP that included that

**Table 1** Summary statistics for all the variables used in the LD study

	Mean	Standard deviation	Minimum	Maximum	Skewness
Length	8.30	1.82	2	15	.29
log Freq HAL	2.83	.91	.30	5.67	-.067
log Ortho N	.21	.28	0	1.18	1.23
log WN senses	.50	.18	0	1.26	.62
Root PFMF	41.49	34.64	0	100	.52
log Root Family Size	.84	.30	.30	2.33	.78
log Root Freq HAL	4.11	.86	.60	6.70	-.25
Suffix length	2.73	.86	1	7	.05
Suffix PFMF	27.99	22.96	0	100	.84
log Suffix Family Size	2.69	.72	.30	3.46	-1.18
log Suffix Freq HAL	5.98	.89	1.26	6.81	-1.54
Suffix P*	1.88	1.70	0	4.88	.57
RTs	779.70	118.51	535.47	1351.13	.74

*Length* length in letters, *log Freq HAL* log-transformed Hyperspace Analogue to Language (HAL) frequency, *log Ortho N* log-transformed orthographic neighborhood size, *log WN senses* log-transformed number of meanings, *Root PFMF* percentage of more frequent words in the morphological family of the root, *Suffix length* suffix length in letters, *Suffix PFMF* percentage of more frequent words in the morphological family of the suffix, *Suffix P\** suffix productivity, *RTs* reaction times, LD latencies

\* $p < .05$ , \*\* $p < .01$

variable (71,954) as compared to those for which SUBTITLE frequency counts are available (51,824). Future studies could run similar analyses using SUBTITLE frequency estimates to compare them to those of the present study. Additionally, we chose N-size as a measure of orthographic neighborhood. However, we acknowledge that Levenshtein Distance could be a better alternative for this type of analyses in future studies (Cortese & Schock, 2012; Yap & Balota, 2009). The semantic variables *WN\_senses* (i.e., log number of meanings) and *WM\_localsn* (i.e., the number of meanings of a target word in its different synsets or sets of semantically related words) are, to the best of our knowledge, the only ones available in English for a significant number of polymorphemic nouns. Cortese and Schock (2012) obtained values for imageability for 1,936 disyllabic words but these words do not necessarily overlap with the words used in the present study. *WN\_senses* was used in similar studies as a semantic measure (Baayen et al., 2006; Yap & Balota, 2009). Values for this variable were taken from the WordNet online database (<http://wordnet.princeton.edu>; Fellbaum, 1998). The variable *WM\_localsn* is a measure of the density of the semantic neighbourhood of a given word (Yap & Balota, 2009) and its values were taken from the Wordmine2 online database (<http://web2.uwindsor.ca/wordmine>; Durda & Buchanan, 2006).

We considered data to be skewed if they presented skewness values larger than  $\pm 2$  (Gravetter & Wallnau, 2014). Raw *Freq\_HAL* is generally skewed and this was the case for our data (skewness = 10.68). We thus followed current similar literature and used log transformed *Freq\_HAL* values (Baayen et al., 2007; Balota, Cortese, Sergent-Marshall,

Spieler, & Yap, 2004). Also, orthographic neighborhood (N-size), *WN\_senses*, and root family size were log-transformed due to skewness. Comparable variables for the morphological constituents (i.e., root and suffix *Freq\_HAL*, and suffix family size values) were also log-transformed. Suffix productivity *P* and its log-transformed equivalent were highly skewed (skewness = 12.74 and 9.33, respectively) and were consequently excluded from analysis. As can be seen in Table 1, all the variables kept for the regression analyses had skewness statistics smaller than  $\pm 2$ .

### Collinearity

Analyses of the correlation matrix of the variables revealed the presence of four coefficients greater than .60 (see Table 2). Such high correlation coefficients indicated a high level of multicollinearity (Balota et al., 2004; Cohen, Cohen, West, & Aiken, 2003). Log *Freq\_HAL* correlated highly with Suffix PFMF (i.e., the percentage of more frequent words in the morphological family of the suffix),  $r = -.813$ ,  $p < .01$ . Log Suffix *Freq\_HAL* correlated highly with log Suffix Family Size,  $r = .914$ ,  $p < .01$ . Log Root Family Size correlated highly with log Root *Freq\_HAL*,  $r = .645$ ,  $p < .01$ , and log Suffix Family Size correlated highly with Suffix *P\**,  $r = .743$ ,  $p < .01$ . To address this issue, we followed Balota et al. (2004) and ran seven additional regression models. In these models we excluded one of the correlated variables to determine if it influenced the remaining critical variables entered in the last step of the models. In other words, we ran additional models for: (1) Suffix PFMF entered in the last step when log *Freq\_HAL* was

**Table 2** Correlations between all the variables used as predictors (and the dependent variable RTs) in the lexical decision task

	Length	log Freq HAL	log Ortho N	log WN senses	Root PFMF	log Root Family Size	log Root Freq HAL	Suffix length	Suffix PFMF	log Suffix Family Size	log Suffix Freq HAL	Suffix P*	RTs
Length	1												
log Freq HAL	-.089**	1											
log Ortho N	-.562**	.023	1										
log WN senses	.036*	.435**	.035*	1									
Root PFMF	.111**	-.519**	-.034*	-.259**	1								
log Root Family Size	-.269**	.251**	.126**	.169**	-.400**	1							
log Root Freq HAL	-.232**	.538**	.158**	.267**	-.184**	.645**	1						
Suffix length	.489**	-.119**	-.422**	.014	.064**	.052**	-.008	1					
Suffix PFMF	.178**	-.813**	-.155**	-.322**	.467**	-.257**	-.536**	.141**	1				
log Suffix Family Size	-.080**	.068**	.349**	.168**	.017	-.101**	.003	-.372**	-.116**	1			
log Suffix Freq HAL	-.021	.147**	.270**	.208**	-.024	-.133**	-.017	-.379**	-.045**	.914**	1		
Suffix P*	-.360**	.021**	.583**	.068**	.018	-.038**	.061**	-.577**	-.233**	.743**	.573**	1	
RTs	.412**	-.571**	-.295**	-.305**	.323**	-.286**	-.454**	.223**	.572**	-.168**	-.150**	-.245**	1

*Length* length in letters, *log Freq HAL* log-transformed Hyperspace Analogue to Language (HAL) frequency, *log Ortho N* log-transformed orthographic neighborhood size, *log WN senses* log-transformed number of meanings, *Root PFMF* percentage of more frequent words in the morphological family of the root, *Suffix length* suffix length in letters, *Suffix PFMF* percentage of more frequent words in the morphological family of the suffix, *Suffix P\** suffix productivity, *RTs* reaction times, *LD* latencies

excluded; (2) log Suffix Freq\_HAL in the last step with the exclusion of log Suffix Family Size; (3) log Suffix Family Size in the last step without log Suffix Freq\_HAL; (4) log Root Family Size entered in the last step when log Root Freq\_HAL was excluded; (5) log Root Freq\_HAL entered in the last step when log Root Family Size was excluded; (6) log Suffix Family Size entered in the last step when Suffix P\* was excluded; and (7) Suffix P\* entered in the last step when log Suffix Family Size was excluded. Log Freq\_HAL highly correlated with both WM\_localsn,  $r = .764$ ,  $p < .01$ , and log WM\_localsn,  $r = .767$ ,  $p < .01$ . To avoid multiple additional regression models with control but not critical morphological variables (i.e., those entered in the last step of the regression models), we decided to keep log WN\_senses (i.e., the log number of meanings) as the only control semantic variable in our models.

## Data analysis

Following previous similar literature (Boukadi, Zouaidi, & Wilson, 2016; Cortese & Schock, 2012; Yap & Balota, 2009), we grouped and entered the variables in the regression

models in four different steps. Step 1 included three lexical variables: log Freq\_HAL; log N-size and length in letters. Step 2 included the semantic variable log WN\_senses (i.e., log number of meanings). Step 3 included all the new morphological variables for the root and suffix except one ( $n = 7$ ): root frequency, root family size, the percentage of more frequent words than the target word in its root morphological family, suffix length in letters, suffix frequency, suffix family size, suffix P\*, the percentage of more frequent words than the target word in its suffix morphological family. Step 4, the final step, included each one of the new 8 morphological variables separately. We ran thus eight different regression models in order to study the specific contribution of each morphological variable above and beyond that of the other variables. As was previously stated, seven additional regression models were run to control for collinearity.

## Results

Eight hierarchical regressions with four steps each were conducted with LD RTs as dependent variable, as well as seven



additional models that allowed to control for collinearity. Table 3 shows the results of these analyses. After controlling for the effect of lexical variables (Step 1), semantics (Step 2), and the other seven morphological variables (Step 3), log Root Freq\_HAL,  $\beta = -.084$ ,  $p < .001$ , Suffix length,  $\beta = -.059$ ,  $p < .001$ , and Suffix PFMF (percentage of more frequent words in the morphological family of the suffix),  $\beta = .158$ ,  $p < .001$ , were significant predictors of LD latencies. Log Root Freq\_HAL and Suffix length exerted a facilitatory effect, whereas Suffix PFMF exerted an inhibitory effect on RTs. None of the other morphological variables significantly predicted RTs in LD latencies, all  $ps > .05$ .

For the seven additional regression models conducted to control for collinearity, only two patterns of results remained unchanged. Suffix PFMF was still a significant predictor when log Freq\_HAL was excluded in Step 1 of the model in which Suffix PFMF was entered in the last step,  $\beta = .384$ ,  $p < .001$ . And log Root Freq\_HAL remained a significant predictor,  $\beta = -.087$ ,  $p < .001$ , when log Root Family Size was excluded in Step 3 of the model. Conversely, log Root Family Size reached significance,  $\beta = -.057$ ,  $p < .001$ , when log Root Freq\_HAL was excluded in Step 3. Log Suffix Family Size became a significant predictor when log Suffix Freq\_HAL was excluded,  $\beta = -.074$ ,  $p < .001$ , and when

**Table 3** Standardized  $\beta$ s, R2s, and  $\Delta$ R2s for the regression analyses of lexical decision (LD)

Step		Root PFMF	log Root Family Size	log Root Freq HAL	Suffix length	Suffix PFMF	log Suffix Family Size	log Suffix Freq HAL	Suffix P*
Step 1	Length	.300***							
	log Freq HAL	-.541***							
	log Ortho N	-.112***							
	R2	.466							
	$\Delta$ R2	.466***							
Step 2	Semantic variable	-.096***							
	R2	.473							
	$\Delta$ R2	.007***							
Step 3	All variables but one								
	Root PFMF	n/a	.012	-.021	.009	.029*	.008	.010	.010
	log Root Family Size	-.010	n/a	-.056***	.008	.012	-.007	-.004	-.003
	log Root Freq HAL	-.078***	-.087***	n/a	-.089***	-.107***	-.083***	-.084***	-.084***
	Suffix length	-.056***	-.059***	-.063***	n/a	-.065***	-.067***	-.059***	-.051***
	Suffix PFMF	.160***	.157***	.178***	.163***	n/a	.170***	.157***	.164***
	log Suffix Family Size	-.069	-.072	-.066	-.114**	-.146***	n/a	-.074***	-.097**
	log Suffix Freq HAL	-.004	-.001	-.007	.036	.097**	-.059***	n/a	.010
	Suffix P*	-.026	-.026	-.029	.013	-.064**	-.051**	-.026	n/a
	R2	.497	.497	.495	.495	.492	.497	.496	.497
	$\Delta$ R2	.024***	.024***	.021***	.022***	.019***	.023***	.024***	.024***
Step 4	Critical variable								
	Variable	.010	-.004	-.084***	-.059***	.158***	-.072	-.002	-.027
	R2	.497	.497	.497	.497	.497	.497	.497	.497
	$\Delta$ R2	.00005	.000007	.002***	.002***	.005***	.0004	.0000004	.0001

Length length in letters, log Freq HAL log-transformed Hyperspace Analogue to Language (HAL) frequency, log Ortho N log-transformed orthographic neighborhood size, log WN senses log-transformed number of meanings, Root PFMF percentage of more frequent words in the morphological family of the root, Suffix length suffix length in letters, Suffix PFMF percentage of more frequent words in the morphological family of the suffix, Suffix P\* suffix productivity

Columns refer to the different regression models and their title shows the variables entered in the last step of the models. Critical variable refers to the variable entered in the last step of the regression models

$\Delta$ R2 is the incremental increase in the model R2 that results from the addition of a predictor or set of predictors in a new step of the model

\* $p < .05$ , \*\*  $p < .01$ , \*\*\*  $p < .001$

Suffix P\* was excluded,  $\beta = -.097$ ,  $p < .01$ , both in Step 3. Suffix P\* became a significant predictor,  $\beta = -.051$ ,  $p < .01$ , when log Suffix Family Size was excluded in Step 3. Also, in the model in which log Suffix Family Size was excluded in Step 3 of the model and log Suffix Freq\_HAL was entered in the last step, this latter reached significance,  $\beta = -.059$ ,  $p < .001$ . All the variables exerted a facilitatory effect on RTs.

## Discussion

The present study aimed to present MorphoLex, a database that includes six morphological variables for affixes and three for roots for each of the 68,624 words in the ELP that were already segmented into morphemes. The novelty of this database is that it includes calculations of some of the most studied variables in the morphological processing literature and applies the same computations to a highly representative number of different words. This represents a significant improvement when compared to the current situation where researchers in morphology have to manually calculate these variables and no specific guidelines are shared on how exactly to proceed when doing so. In this context, MorphoLex presents identical calculations of these variables for almost 70,000 English words and aims to facilitate comparisons across studies that use these variables. Additionally, MorphoLex includes lists of all the roots, prefixes and suffixes encountered in the words from the ELP, along with their corresponding morphological variables. This offers a helpful tool for researchers interested in creating pseudowords that share specific morphological characteristics. For instance, if researchers want to assess the effect of suffix family size on naming RTs using pseudowords, they could easily select suffixes from MorphoLex that have a small family size and compare them to those with a large family size. The database is freely accessible at the following URL: <https://github.com/hugomailhot/MorphoLex-en>

In order to illustrate how MorphoLex can be used in psycholinguistic research we entered those variables as predictors of LD RTs for 4,724 suffixed nouns in a series of regression models. Results of the regression study indicate that root frequency and suffix length have a facilitative effect. This means that the higher the root frequency and suffix length, the shorter the LD latencies. Conversely, the percentage of more frequent words in the morphological family (PFMF) of the suffix had an inhibitory effect. Thus, the higher the percentage of more frequent words than the target sharing the same suffix, the longer the LD latencies. Several variables showed a high degree of collinearity and new analyses excluding one of the highly correlated variables were conducted. When root frequency was excluded from the regression, root family size exerted a significant facilitatory effect on latencies. When suffix frequency or suffix productivity P\* were eliminated from the regression models, suffix family size became a significant

predictor of RTs, with a facilitative effect. When suffix family size was eliminated from the model, suffix frequency then became a significant predictor of RTs with a facilitative effect. Additionally, suffix productivity P\* became a significant facilitatory predictor of latencies when suffix family size was excluded. This indicates that suffix family size shares an important part of its variance with suffix frequency and suffix productivity P\*. It is noteworthy that both suffix family size and productivity P\* have been used as measures of suffix productivity (Baayen & Renouf, 1996; Kuperman et al., 2010). It is not surprising, then, that they share a common variance.

The facilitative effect of cumulative root frequency had been previously found in several studies (Burani & Thornton, 2003; Caramazza et al., 1988; Luke & Christianson, 2011; Taft & Forster, 1975, 1976), while others had failed to find compelling evidence of it (Schreuder & Baayen, 1997; Sereno & Jongman, 1997) or had even found root frequency to exert an inhibitory effect (Baayen et al., 2002, 2007). Colé et al. (1989) argued that these inconsistent findings were due to the differential effect of root frequency on distinct types of complex words, as they only found a facilitative effect of root frequency for suffixed, but not prefixed, words. Our study, thus, confirms that root frequency plays a central role in the processing of English suffixed words. Future studies could use words from different PRS signatures in MorphoLex to further explore root frequency effects on different types of morphologically complex words.

Similar to other studies that have looked into root family size effects (Balling & Baayen, 2008; Bertram, Baayen, et al., 2000; De Jong et al., 2000; Ford et al., 2010; Moscoso del Prado Martín et al., 2004), we also found that words whose roots are part of a numerous family are recognized faster than those whose roots pertain to a smaller family. However, we only found this pattern of results after the exclusion of root frequency from the predictor variables, which indicates that root family size and root frequency share part of their variance.

While root frequency and family size effects have been widely investigated, morphological variables related to affixes have generally been kept out of the spotlight, with only a few studies focusing on their role in complex word processing. However, studying the characteristics of all the morphemes in complex words is the only way to obtain a thorough picture of the morphological variables that affect their processing. Our study confirms the importance of introducing several affix variables in the same model, as all the variables linked to the suffix (i.e., suffix length, PFMF, frequency, family size, and P\*) affected LD latencies.

When Burani and Thornton (2003) studied suffix frequency effects, they suggested that these might be due to the number of words that shared a suffix (i.e., suffix family size) and not solely to their summed frequency. They argued that highly

frequent suffixes are also part of big morphological families and, thus, suffix frequency effects might actually emerge from the family size of the suffix instead of its actual cumulative frequency. Our study confirms this hypothesis, as both variables showed a concurrent facilitative effect on LD latencies. Thus, words that have a higher number of suffix family members and, hence, a higher suffix frequency, are recognized faster than words with smaller suffix frequencies and family sizes.

Contrary to the facilitative effect of the number of suffix family members and their frequency, the percentage of more frequent words in the morphological family of the suffix showed an inhibitory effect on LD latencies. Indeed, being part of a numerous suffixal family facilitates the processing of a word but only if this word has fewer more frequent suffix family members. In other words, a word that has a considerable number of competitors (i.e., more frequent words) in its suffix morphological family will be processed more slowly than a word that is amongst the most frequent ones in its suffixal family. This finding corroborates the results from previous studies, such as those by Meunier and Segui (1999) and Colé et al. (1989).

The contribution of affixal length to faster processing through morphological decomposition (Kuperman et al., 2010) also seems to be confirmed by our data, as words with longer suffixes were processed faster than words with shorter suffixes. This could be interpreted in line with Laudanna and Burani's (1995) idea that affixes that are more salient trigger faster processing through the decompositional route. Indeed, all variables that increased suffix salience (i.e., length, frequency, family size) in our study presented a facilitative effect on latencies, and this effect of affix salience was only moderated by the number of competitors that a specific word had in the suffixal family (i.e., PFMF).

Finally, our study confirms the relevance of productivity variables for LD latencies. To the best of our knowledge, the P value had only been entered as a predictor in a regression study in Baayen et al. (2007), where it exerted a small facilitative effect on LD latencies, whereas the P\* had never been used before in a study of these characteristics. Our results show for the first time that suffix productivity P\* also has a facilitative effect on LD latencies when suffix family size is eliminated from the model.

In conclusion, as exemplified in this series of regression models, the data from MorphoLex offer the opportunity to investigate the effects of each of these new morphological variables with a plethora of designs, tasks and combinations of morphological complexity in English. We argue that the use of a single database as the one we have created and rendered freely available here will facilitate the comparison of future experimental studies on the effect of morphology in visual word recognition. Hopefully, this will have a positive impact

on the future development of models of morphological processing.

**Author note** This research was supported by an Insight Development Grant awarded to M.A.W., S.H.D., and C.S.G. by the Social Sciences and Humanities Research Council (CRSH) of Canada, grant number: 430-2015-00699.

## References

- Alegre, M., & Gordon, P. (1999). Frequency effects and the representational status of regular inflections. *Journal of Memory and Language*, 40(1), 41–61. doi:<https://doi.org/10.1006/jmla.1998.2607>
- Amenta, S., & Crepaldi, D. (2012). Morphological processing as we know it: An analytical review of morphological effects in visual word identification. *Frontiers in Psychology*, 3(Article 232), 1–12. doi:<https://doi.org/10.3389/fpsyg.2012.00232>
- Andrews, S., Miller, B., & Rayner, K. (2004). Eye movements and morphological segmentation of compound words: There is a mouse in mousetrap. *European Journal of Cognitive Psychology*, 16(1–2), 285–311. doi:<https://doi.org/10.1080/09541440340000123>
- Angelelli, P., Marinelli, C. V., & Burani, C. (2014). The effect of morphology on spelling and reading accuracy: A study on Italian children. *Frontiers in Psychology*, 5(Article 1373), 1–10. doi:<https://doi.org/10.3389/fpsyg.2014.01373>
- Baayen, R. H. (2009). Corpus linguistics in morphology: Morphological productivity. In A. Lüdeling & M. Kyto (Eds.), *Corpus Linguistics. An international handbook* (pp. 900–919). Berlin: Mouton De Gruyter.
- Baayen, R. H., Dijkstra, T., & Schreuder, R. (1997). Singulars and plurals in Dutch: Evidence for a parallel dual-route model. *Journal of Memory and Language*, 37(1), 94–117. doi:<https://doi.org/10.1006/jmla.1997.2509>
- Baayen, R. H., Feldman, L. B., & Schreuder, R. (2006). Morphological influences on the recognition of monosyllabic monomorphemic words. *Journal of Memory and Language*, 55, 290–313. doi:<https://doi.org/10.1016/j.jml.2006.03.008>
- Baayen, R. H., & Lieber, R. (1991). Productivity and English derivation: a corpus-based study. *Linguistics*, 29(5), 801–844. doi:<https://doi.org/10.1515/ling.1991.29.5.801>
- Baayen, R. H., Piepenbrock, R., & Gullikers, L. (1995). *The CELEX lexical database [CD-ROM]*. Philadelphia: Linguistic Data Consortium, University of Pennsylvania.
- Baayen, R. H., & Renouf, A. (1996). Chronically the times: Productive lexical innovations in an English newspaper. *Language*, 72, 69–96. doi:<https://doi.org/10.2307/416794>
- Baayen, R. H., Tweedie, F. J., & Schreuder, R. (2002). The subjects as a simple random effect fallacy: Subject variability and morphological family effects in the mental lexicon. *Brain and Language*, 81(1), 55–65. doi:<https://doi.org/10.1006/brln.2001.2506>
- Baayen, R. H., Wurm, L. H., & Aycocock, J. (2007). Lexical dynamics for low-frequency complex words. A regression study across tasks and modalities. *The Mental Lexicon*, 2(3), 419–463. doi:<https://doi.org/10.1075/ml.2.3.06baa>
- Balling, L., & Baayen, R. H. (2008). Morphological effects in auditory word recognition: Evidence from Danish. *Language and Cognitive Processes*, 23(7–8), 1159–1190. doi:<https://doi.org/10.1080/01690960802201010>

- Balota, D. A., Cortese, M. J., Sergent-Marshall, S. D., Spieler, D. H., & Yap, M. J. (2004). Visual word recognition of single-syllable words. *Journal of Experimental Psychology: General*, *133*(2), 283–316. doi:<https://doi.org/10.1037/0096-3445.133.2.283>
- Balota, D. A., Yap, M. J., Cortese, M. J., Hutchison, K. A., Kessler, B., Loftis, B., ... Treiman, R. (2007). The English Lexicon project. *Behavior Research Methods*, *39*(3), 445–459. doi:<https://doi.org/10.3758/BF03193014>
- Bauer, L., & Nation, P. (1993). Word families. *International Journal of Lexicography*, *6*(4), 253–279. doi:<https://doi.org/10.1093/ijl/6.4.253>
- Beauvillain, C. (1996). The integration of morphological and whole-word form information during eye fixations on prefixed and suffixed words. *Journal of Memory and Language*, *35*(6), 801–820. doi:<https://doi.org/10.1006/jmla.1996.0041>
- Bertram, R., Baayen, R. H., & Schreuder, R. (2000). Effects of family size for complex words. *Journal of Memory and Language*, *42*(3), 390–405. doi:<https://doi.org/10.1006/jmla.1999.2681>
- Bertram, R., Schreuder, R., & Baayen, R. H. (2000). The balance of storage and computation in morphological processing: The role of word formation type, affixal homonymy, and productivity. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *26*(2), 489–511. doi:<https://doi.org/10.1037/0278-7393.26.2.489>
- Boukadi, M., Zouaidi, C., & Wilson, M. A. (2016). Norms for name agreement, familiarity, subjective frequency, and imageability for 348 object names in Tunisian Arabic. *Behavior Research Methods*, *48*(2), 585–599. doi:<https://doi.org/10.3758/s13428-015-0602-3>
- Bradley, D. C. (1979). Lexical representation of derivational relations. In M. Aronoff & M. L. Kean (Eds.), *Juncture* (pp. 37–55). Cambridge, MA: MIT Press.
- Brysbart, M., and 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*, 977–990. doi:<https://doi.org/10.3758/BRM.41.4.977>
- Burani, C., Dovetto, F. M., Thornton, A. M., & Laudanna, A. (1997). Accessing and naming suffixed pseudo-words. In G. E. Booij & J. van Marle (Eds.), *Yearbook of Morphology 1996* (pp. 55–72). Dordrecht: Kluwer.
- Burani, C., & Laudanna, A. (1992). Units of representation for derived words in the lexicon. *Advances in Psychology*, *94*, 361–376. doi:[https://doi.org/10.1016/S0166-4115\(08\)62803-4](https://doi.org/10.1016/S0166-4115(08)62803-4)
- Burani, C., & Thornton, A. M. (2003). The interplay of root, suffix and whole-word frequency in processing derived words. In R. H. Baayen & R. Schreuder (Eds.), *Morphological Structure in Language Processing* (pp. 157–208). Berlin-New York: Mouton de Gruyter.
- Burgess, C., & Livesay, K. (1998). The effect of corpus size in predicting reaction time in a basic word recognition task: Moving on from Kučera and Francis. *Behavior Research Methods, Instruments, & Computers*, *30*(2), 272–277. doi:<https://doi.org/10.3758/BF03200655>
- Caramazza, A., Laudanna, A., & Romani, C. (1988). Lexical access and inflectional morphology. *Cognition*, *28*(3), 297–332. doi:[https://doi.org/10.1016/0010-0277\(88\)90017-0](https://doi.org/10.1016/0010-0277(88)90017-0)
- Cohen, J., Cohen, P., West, S. G., & Aiken, L. S. (2003). *Applied multiple regression/correlation analysis for the behavioral sciences* (3rd ed.). New York, NY: Routledge.
- Colé, P., Beauvillain, C., & Segui, J. (1989). On the representation and processing of prefixed and suffixed derived words: A differential frequency effect. *Journal of Memory and Language*, *28*(1), 1–13. doi:[https://doi.org/10.1016/0749-596X\(89\)90025-9](https://doi.org/10.1016/0749-596X(89)90025-9)
- Cortese, M. J., & Schock, J. (2012). Imageability and age of acquisition effects in disyllabic word recognition. *The Quarterly Journal of Experimental Psychology*, *66*(5), 946–972. doi:<https://doi.org/10.1080/17470218.2012.722660>
- Davies, S. K., Izura, C., Socas, R., & Dominguez, A. (2016). Age of acquisition and imageability norms for base and morphologically complex words in English and in Spanish. *Behavior Research Methods*, *48*(1), 349–365. doi:<https://doi.org/10.3758/s13428-015-0579-y>
- De Jong, N. H., Schreuder, R., & Baayen, R. H. (2000). The morphological family size effect and morphology. *Language and Cognitive Processes*, *15*(4–5), 329–365. doi:<https://doi.org/10.1080/01690960050119625>
- Durda, K., & Buchanan, L. (2006). WordMine2. Retrieved from <http://web2.uwindsor.ca/wordmine>
- Fellbaum, C. (1998). *WordNet: An electronic database*. Cambridge, MA: The MIT Press.
- Ford, M. A., Davis, M. H., & Marslen-Wilson, W. D. (2010). Derivational morphology and base morpheme frequency. *Journal of Memory and Language*, *63*, 117–130. doi:<https://doi.org/10.1016/j.jml.2009.01.003>
- Gravetter, F. J., & Wallnau, L. B. (2014). *Essentials of statistics for the behavioral sciences* (8th ed.). Belmont, CA: Wadsworth.
- Hay, J. (2001). Lexical frequency in morphology: Is everything relative? *Linguistics*, *39*(6), 1041–1070. doi:<https://doi.org/10.1515/ling.2001.041>
- Kuperman, V., Bertram, R., & Baayen, H. (2010). Processing trade-offs in the reading of Dutch derived words. *Journal of Memory and Language*, *62*, 83–97. doi:<https://doi.org/10.1016/j.jml.2009.10.001>
- Laudanna, A., & Burani, C. (1995). Distributional properties of derivational affixes: Implications for processing. In L. B. Feldman (Ed.), *Morphological aspects of language processing: Cross-linguistic perspectives* (pp. 345–364). Hillsdale: Lawrence Erlbaum Associates.
- Lehtonen, M., Niska, H., Wande, E., Niemi, J., & Laine, M. (2006). Recognition of inflected words in a morphologically limited language: Frequency effects in monolinguals and bilinguals. *Journal of Psycholinguistic Research*, *35*(2), 121–146. doi:<https://doi.org/10.1007/s10936-005-9008-1>
- Longtin, C. M., & Meunier, F. (2005). Morphological decomposition in early visual word processing. *Journal of Memory and Language*, *53*(1), 26–41. doi:<https://doi.org/10.1016/j.jml.2005.02.008>
- Luke, S. G., & Christianson, K. (2011). Stem and whole-word frequency effects in the processing of inflected verbs in and out of a sentence context. *Language and Cognitive Processes*, *26*(8), 1173–1192. doi:<https://doi.org/10.1080/01690965.2010.510359>
- Marcolini, S., Traficante, D., Zoccolotti, P., & Burani, C. (2011). Word frequency modulates morpheme-based reading in poor and skilled Italian readers. *Applied Psycholinguistics*, *32*(3), 513–532. doi:<https://doi.org/10.1017/S0142716411000191>
- McCormick, S. F., Brysbart, M., & Rastle, K. (2009). Is morphological decomposition limited to low-frequency words? *The Quarterly Journal of Experimental Psychology*, *62*(9), 1706–1715. doi:<https://doi.org/10.1080/17470210902849991>
- McCormick, S. F., Rastle, K., & Davis, M. H. (2009). Adorable not adorable? Orthographic underspecification studied with masked repetition priming. *European Journal of Cognitive Psychology*, *21*(6), 813–836. doi:<https://doi.org/10.1080/09541440802366919>
- Meunier, F., & Segui, J. (1999). Frequency effects in auditory word recognition: The case of suffixed words. *Journal of Memory and Language*, *41*(3), 327–344. doi:<https://doi.org/10.1006/jmla.1999.2642>
- Moscato del Prado Martín, F., Bertram, R., Häikiö, T., Schreuder, R., & Baayen, R. H. (2004). Morphological family size in a morphologically rich language: The case of Finnish compared with Dutch and Hebrew. *Journal of Experimental Psychology: Learning, Memory and Cognition*, *30*(6), 1271–1278. doi:<https://doi.org/10.1037/0278-7393.30.6.1271>
- New, B., Brysbart, M., Segui, J., Ferrand, L., & Rastle, K. (2004). The processing of singular and plural nouns in French and English.

- Journal of Memory and Language*, 51(4), 568-585. doi:<https://doi.org/10.1016/j.jml.2004.06.010>
- Niswander-Klement, E., & Pollatsek, A. (2006). The effects of root frequency, word frequency, and length on the processing of prefixed English words during reading. *Memory & Cognition*, 34(3), 685-702. doi:<https://doi.org/10.3758/BF03193588>
- Oldfield, R. C., & Wingfield, A. (1965). Response latencies in naming objects. *Quarterly Journal of Experimental Psychology*, 17(4), 273-281. doi:<https://doi.org/10.1080/17470216508416445>
- Plag, I. (2006). Productivity. In K. Brown (Ed.), *Encyclopedia of language & linguistics* (2nd ed., pp. 121-128). Oxford: Elsevier.
- Rastle, K., & Davis, M. H. (2008). Morphological decomposition based on the analysis of orthography. *Language and Cognitive Processes*, 23(7-8), 942-971. doi:<https://doi.org/10.1080/01690960802069730>
- Rastle, K., Davis, M. H., & New, B. (2004). The broth in my brother's brothel: Morpho-orthographic segmentation in visual word recognition. *Psychonomic Bulletin & Review*, 11(6), 1090-1098. doi:<https://doi.org/10.3758/BF03196742>
- Scarborough, D. L., Cortese, C., & Scarborough, H. S. (1977). Frequency and repetition effects in lexical memory. *Journal of Experimental Psychology: Human Perception and Performance*, 3(1), 1-17. doi:<https://doi.org/10.1037/0096-1523.3.1.1>
- Schreuder, R., & Baayen, R. H. (1997). How simplex complex words can be. *Journal of Memory and Language*, 37, 118-139. doi:<https://doi.org/10.1006/jmla.1997.2510>
- Sereno, J. A., & Jongman, A. (1997). Processing of English inflectional morphology. *Memory and Cognition*, 25(4), 425-437. doi:<https://doi.org/10.3758/BF03201119>
- Stemberger, J. P., & MacWhinney, B. (1986). Frequency and the lexical storage of regularly inflected forms. *Memory and Cognition*, 14(1), 17-26. doi:<https://doi.org/10.3758/BF03209225>
- Taft, M. (1979). Recognition of affixed words and the word frequency effect. *Memory and Cognition*, 7(4), 263-272. doi:<https://doi.org/10.3758/BF03197599>
- Taft, M., & Forster, K. I. (1975). Lexical storage and retrieval of prefixed words. *Journal of Verbal Learning and Verbal Behavior*, 14(6), 638-647. doi:[https://doi.org/10.1016/S0022-5371\(75\)80051-X](https://doi.org/10.1016/S0022-5371(75)80051-X)
- Taft, M., & Forster, K. I. (1976). Lexical storage and retrieval of polymorphemic and polysyllabic words. *Journal of Verbal Learning and Verbal Behavior*, 15(6), 607-620. doi:[https://doi.org/10.1016/0022-5371\(76\)90054-2](https://doi.org/10.1016/0022-5371(76)90054-2)
- Vannest, J., Polk, T. A., & Lewis, R. L. (2005). Dual-route processing of complex words: New fMRI evidence from derivational suffixation. *Cognitive, Affective, & Behavioral Neuroscience*, 5(1), 67-76. doi:<https://doi.org/10.3758/CABN.5.1.67>
- Whaley, C. P. (1978). Word-nonword classification time. *Journal of Verbal Learning and Verbal Behavior*, 17(2), 143-154. doi:[https://doi.org/10.1016/S0022-5371\(78\)90110-X](https://doi.org/10.1016/S0022-5371(78)90110-X)
- Yap, M. J., & Balota, D. A. (2009). Visual word recognition of multisyllabic words. *Journal of Memory and Language*, 60, 502-529. doi:<https://doi.org/10.1016/j.jml.2009.02.001>