

1 Modeling Memory: Exploring the Relationship Between Word Overlap and Single Word
2 Norms when Predicting Relatedness Judgments and Retrieval

3 Nicholas P. Maxwell¹ & Erin M. Buchanan¹

4 ¹ Missouri State University

5 Author Note

6 Nicholas P. Maxwell is a graduate student at Missouri State University. Erin M.
7 Buchanan is an Associate Professor of Psychology at Missouri State University

8 Correspondence concerning this article should be addressed to Nicholas P. Maxwell,
9 901 S. National Ave, Springfield, MO, 65897. E-mail: maxwell270@live.missouristate.edu

Abstract

10

11 Enter abstract here (note the indentation, if you start a new paragraph).

12 *Keywords:* judgments, memory, association, semantics, thematics

13 Word count: X

Modeling Memory: Exploring the Relationship Between Word Overlap and Single Word Norms when Predicting Relatedness Judgments and Retrieval

Previous research conducted on judgments of associative memory (JAM) has found that these judgments tend to be stable and highly generalizable across varying contexts (Maki, 2007a, 2007b; Valentine & Buchanan, 2013). The JAM task can be viewed as a manipulation of the traditional judgment of learning task (JOL). In a judgment of learning task, participants are presented with cue-target word pairs and are asked to make a judgment (typically on a scale of zero to 100) of how accurately they would be able to respond with the proper target word based on the presentation of a particular cue word (Dunlosky & Nelson, 1994; Nelson & Dunlosky, 1991). JAM tasks expand upon this concept by changing the focus of the judgments performed by participants. When presented with the item pair, such as cheese-mouse, participants are asked to judge the number of people out of 100 who would respond with the pair's target word if they were only shown the cue (Maki, 2007a).

This process mimics the creation of associative words norms (i.e., forward strength; D. L. Nelson, McEvoy, and Schreiber (2004)). As such, these judgments can be viewed as the participants' approximations of how associatively related they perceive the paired items to be. The JAM function can then be created by plotting participant judgments against the word's normed associative strength and calculating a line of best fit. This fit line typically displays a high intercept (bias) and a shallow slope (sensitivity), meaning that participants are biased towards overestimating the associative relatedness between word pairs, and show difficulties differentiating between different amounts of item relatedness (Maki, 2007a).

Building upon this research, we initially completed a pilot study in which we sought to examine recall accuracy within the context of item judgments, while also expanding the JAM task to incorporate judgments of semantic and thematic memory. In the pilot study, 63 word-pairs of varying associative, semantic, and thematic overlap were created and arranged into three blocks, consisting of 21 word-pairs each. Associative overlap was measured with forward strength (FSG; D. L. Nelson et al. (2004)), semantic overlap was measured with

cosine (COS; McRae, Cree, Seidenberg, and McNorgan (2005)), and thematic relatedness between pairs was measured with latent semantic analysis (LSA; Landauer and Dumais (1997); Landauer, Foltz, and Laham (1998)). Participants were randomly assigned to a condition in which they received a set of instructions explaining either an associative, semantic, or thematic relationship between words. Participants then judged the word-pairs in each block based on the instructions that they received. The order of block presentation and judgment instructions were counterbalanced so that each word-pair received each of the three types of judgments. After completing the judgment phase, participants then completed a cued recall task in which they were presented with the cue word from each of the previously presented word pairs and were asked to complete each pair with the missing target.

Multilevel modeling was then used to predict recall and judgment scores. This type of analysis was selected due to its ability to retain all data points while controlling for correlated error between participants. Significant three-way interactions were found between database norms when predicting judgments ($\beta = 3.324, p < .001$) and recall ($\beta = 24.571, p < .001$). Simple slopes analyses were then conducted to further examine these interactions. When semantic overlap was low, thematic and associative strength were competitive, with increases in thematic overlap decreasing the strength of associative overlap as a predictor. However, this trend saw a reversal when semantic overlap was high, with thematic and associative strength complimenting one another. This result was found when investigating the three-way interactions for both the judgment and recall tasks. Overall, our findings from this study indicated the degree to which the processing of associative, semantic, and thematic information impacts retrieval and judgment making, while also displaying the interactive relationship that exists between these three types of information.

The proposed study seeks to expand upon this work by extending the original analysis to include multiple single word norms. These norms provide information about different “neighborhoods” of concept information. Broadly speaking, they can be separated into one of three categories. Base values refer to norms which capture information based on a word’s

structure. These include part of speech (PoS), word frequency, and the number of syllables, morphemes, and phonemes that comprise a word. Rated values refer to age of acquisition (AoA), concreteness, imageability, valence, and familiarity. Finally, we seek to examine norms that provide information about the connections a word shares with others based on context. These norms include orthographic neighborhood, phonographic neighborhood, cue and target set sizes, and feature set size.

First, we are interested in assessing the impact of base word norms. Chief amongst these is word frequency. Several sets of norms currently exist for measuring the frequency with which words occur in everyday language, and it is important to determine which of these offers the best representation of everyday language. One of the most commonly used collections of these norms is the Kucera and Francis (1967) frequency norms. This set consists of frequency values for words, which were generated by analyzing books, magazines, and newspapers. However, the validity of using these norms has been questioned on factors such as the properties of the sources analyzed, the size of the corpus analyzed, and the overall age of these norms. First, these norms were created from an analysis of written text. It is important to keep in mind that stylistically, writing tends to be more formal than everyday language and as a result, it may not be the best approximation of it (Brysbaert & New, 2009). Additionally, these norms were generated fifty years ago, meaning that these norms may not accurately reflect the current state of the English language. As such, the Kucera and Francis norms may not be the best choice for researchers interested in gauging the effects of word frequency.

Several viable alternatives to the KF frequency norms now exist. One popular method is to use frequency norms obtained from the HAL corpus, which consists of 131 million words (Burgess & Lund, 1997; Lund & Burgess, 1996). Other collections of frequency norms include CELEX (Baayen, Piepenbrock, & Gulikers, 1995) which is based on written text, the Zeno frequency norms (Zeno, Ivens, Millard, & Duvvuri, 1995) which were created from American children's textbooks, and Google Book's collection of word frequencies which is

95 derived from 131 billion words taken from books published in the United States. (See
96 Brysbaert, Keuleers, and New (2011) for an overview and comparison of these norms to
97 SUBLTEX). For the present study, we plan to use data taken from the both the SUBTLEX
98 project (Brysbaert & New, 2009), which is a collection of frequency norms derived from a
99 corpus of approximately 51 million words, which were generated from movie and television
100 subtitles, and the HAL corpus. SUBTLEX norms are thought to better approximate
101 everyday language, as lines from movies and television tend to be more reflective of everyday
102 speech than writing samples. Additionally, the larger corpus size of both SUBTLEX and
103 HAL contributes to the validity of these norms compared to KF frequency norms.

104 Next, we are interested in testing the effects of several measures of lexical information
105 related to the physical make-up of words. These measures include the numbers of phonemes,
106 morphemes, and syllables that comprise each word as well as its part of speech. The number
107 of phonemes refers to the number of individual sounds that comprise a word (i.e., the word
108 CAT has three phonemes, each of which correspond to the sounds its letters make), while
109 the term morpheme refers to the number of sound units that contain meaning. DRIVE
110 contains one morpheme, while DRIVER contains two. Morphemes typically consist of root
111 words and their affixes. We are also interested in word length (measured as the number of
112 individual characters a word consists of) and the number of syllables a word contains, as
113 previous research has suggested that the number of syllables may play a role in processing
114 time. In general, longer words require longer processing time (Kuperman,
115 Stadthagen-Gonzalez, & Brysbaert, 2012), and shorter words tend to be more easily
116 remembered (Cowan, Baddeley, Elliott, & Norris, 2003). Finally, we are interested in the
117 part of speech of each word. For the present study, part of speech will be coded as nouns,
118 verbs, adjectives, and other, and will be based on category size.

119 Third, we are interested in exploring the effects of norms measuring word properties
120 that are rated by participants. The first of these is age of acquisition (AoA), which is a
121 measure of the age at which a word is learned. This norm is measured by presenting

participants with a word and having them enter the age (in years) in which they believe that they would have learned the word (Kuperman et al., 2012). AoA ratings have been found to be predictive of recall. For example, Dewhurst, Hitch, and Barry (1998) found recall to be higher for late acquired words. Also of interest are measures of a word's valence, which refers to its intrinsic pleasantness or perceived positiveness. Valence ratings are important across multiple psycholinguistic research settings. These include research on emotion, the impact of emotion of lexical processing and memory, estimating the sentiments of larger passages of text, and estimating the emotional value of new words based on valence ratings of semantically similar words (See Warriner, Kuperman, and Brysbaert (2013) for a review). The next of these rated measures is concreteness, which refers to the degree that a word relates to a perceptible object (Brysbaert, Warriner, & Kuperman, 2013). Similar to concreteness, imageability is described as being a measure of a word's ability to generate a mental image (Stadthagen-Gonzalez & Davis, 2006). Both imageability and concreteness have been linked to recall, as items rated higher in these areas tend to be more easily recalled (Nelson & Schreiber, 1992). Finally, familiarity norms can be described as an application of word frequency. These norms measure the frequency of exposure to a particular word (Stadthagen-Gonzalez & Davis, 2006).

The final group of norms that we are interested in examining are those which provide information based on connections with neighboring words. Phonographic neighborhood refers to the number of words that can be created by changing one sound in a word (i.e., CAT to KITE). Similarly, orthographic neighborhood refers to the number of words created by changing a single letter in word (i.e., CAT to BAT, Adelman and Brown (2007); Peereman and Content (1997)). Previous findings have suggested that the frequency of a target word relative to that of its orthographic neighbors has an effect on recall, increasing the likelihood of recall for that word (Carreiras, Perea, & Grainger, 1997). Additionally, both of measures have been found to effect processing speed for items (Adelman & Brown, 2007; Buchanan, Holmes, Teasley, & Hutchison, 2013; Coltheart, Davelaar, Jonasson, &

Besner, 1977). Next, we are interested in examining two single word norms that are directly related to item associations. These norms measure the number of associates a word shares connections with. Cue set size (QSS) refers to the number of cue words that a target word is connected to, while target set size (TSS) is a count of the number of target words a cue word is connected to (Schreiber and Nelson (1998)). Previous research has shown evidence for a cue set size effect in which cue words that are linked to a larger number of associates (target words) are less likely to be recalled than cue words linked to fewer target words (D. L. Nelson, Schreiber, & Xu, 1999). As such, we will also calculate set size values for the semantic feature overlap and thematic overlap word norms. Finally, feature list sizes will be calculated for each word overlap norm from the Buchanan et al. 2013 semantic feature norm set.

In summary, this study seeks to expand upon previous work by examining how single word norms belonging to these three neighborhoods of item information impact the accuracy of item judgments and recall. These findings will be assessed within the context of associative, semantic, and thematic memory systems. Specifically, we utilize a three-tiered view of the interconnections between these systems as it relates to processing concept information. First, semantic information is processed, which provides a means for categorizing concepts based on feature similarity. Next, processing moves into the associative memory network, where contextual information pertaining to the items is added. Finally, the thematic network incorporates information from both the associative and semantic networks to generate a mental representation of the concept containing both the items meaning and its place in the world.

As such, the present study has two aims. First, we seek to replicate the interaction results from the pilot study using a new set of stimuli. These three-way interactions occurred between the associative, semantic, and thematic database norms when predicting participant judgments and recall. Second, we wish to expand upon these findings by extending the analysis to include neighborhood information for the item pairs. The extended analysis will be run by introducing the different types single word norms through a series of steps based

on the type of neighborhood they belong to. First, base word norms will be analyzed. Next, measures of word ratings will be analyzed. Third, single word norms measuring connections between concepts will be analyzed. Finally, network norms and their interactions will be reanalyzed. The end goal is to determine both which neighborhood of norms have the greatest overall impact on recall and judgment ability, and to further assess the impact of network connections after controlling for the various neighborhoods of single word information.

Methods

Participants

A power analysis was conducted using the *SIMR* package in *R* (Green & MacLeod, 2016), which uses simulations to calculate power for mixed linear models created from the *LME4* and *nlme* packages (D. Bates, Machler, Bolker, & Walker, 2015; Pinheiro, Bates, Debroy, Sarkar, & R Core Team, 2017). The results of this analyses suggested a minimum of 35 participants was required to find an effect at 80% power. However, because power often is underestimated (Brysbaert & Stevens, 2018), we plan to extend the analysis to include 200 participants, a number determined by the amount of available funding. Participants will be recruited from Amazon's Mechanical Turk, which is a website where individuals can host projects and be connected with a large respondent pool who complete tasks for small amounts of money (Buhrmester, Kwang, & Gosling, 2011). Participants will be paid \$2.00 for their participation. Participant responses will be screened for a basic understanding of study instructions.

Material

First, mimicking the design of the original pilot study, sixty-three word pairs of varying associative, semantic, and thematic overlap were created to use as stimuli. These word pairs were created using the Buchanan et al. (2013) word norm database. Next, neighborhood

information for all cue and target items was collected. Word frequency was collected from the SUBTLEX project (Brysbaert & New, 2009) and the HAL corpus (Burgess & Lund, 1997). Part of speech (POS), word length, and the number of morphemes, phonemes, and syllables of each item was derived from the Buchanan et al. (2013) word norms. For items with multiple parts of speech (for example, Drink can refer to both a beverage and the act of drinking a beverage), the most commonly used form was used. Following the design of Buchanan et al. (2013), this was determined using Google's "Define" feature. Concreteness, cue set size (QSS), and target set size (TSS) were taken from the South Florida Free Association Norms (D. L. Nelson et al., 2004). Imageability and familiarity norms were taken from the (Toglia, 2009; Toglia & Battig, 1978) semantic word norms. Age of acquisition ratings (AoA) were pulled from the (Kuperman et al., 2012) database. Finally, valence ratings for all items were obtained from the (Warriner et al., 2013) norms. After gathering neighborhood information, network norms measuring associative, semantic, and thematic overlap were generated for each pair. Forward strength (FSG) was used as a measure of associative overlap. FSG is a value ranging from zero to one which measures of the probability that a cue word will elicit a particular target word in response to it (D. L. Nelson et al., 2004). Cosine (COS) strength was used to measure semantic overlap between concepts (Buchanan et al. (2013); McRae et al. (2005); Vinson and Vigliocco (2008)). As with FSG, this value ranges from zero to one, with higher values indicating more shared features between concepts. Finally, thematic overlap was measured with Latent Semantic Analysis (LSA), which is a measure generated based upon the co-occurrences of words within a document (Landauer & Dumais, 1997; Landauer et al., 1998). Like the measures of associative and semantic overlap, LSA values range from zero to one, with higher values indicating higher co-occurrence between items. As such, the selected stimuli contained a range of values across both the network and neighborhood norms. As with the pilot study, stimuli will be arranged into three blocks, with each block consisting of 21 word pairs. The blocks will be structured to have seven words of low COS (0 - .33), medium COS (.34 - .66),

and high COS (.67 - 1). COS was chosen due to both limitations with the size of the available dataset and the desire to recreate the selection process used for the pilot study. The result of this selection process is that values for the remaining network norms (FSG and LSA) and information neighborhood norms will be contingent upon the COS strengths of the selected stimuli. To counter this, we selected stimuli at random based on the different COS groupings so as to cover a broad range of FSG, LSA, and information neighborhood values.

The stimuli will be presented to the participants online via Qualtrics surveys. Three different surveys will be created, which will counter-balance the order in which stimuli blocks are presented. Judgment conditions will be counter-balanced across blocks, so that each word pair receives a judgment for each type of memory. Finally, word pairs will be randomized within blocks.

Procedure

This study will be divided into three sections. First, participants will be presented with word pairs and will be asked to judge how related the items are to one another. This section will comprise three blocks, with each block containing 21 word pairs. Each item block will be preceded by a set of instructions explaining one of the three types of relationships. Participants will also be provided with examples illustrating the type of relationship to be judged. The associative instructions explain associative relationships between concepts, how these relationships can be strong or weak, and the role of free association tasks in determining the magnitude of these relationships. The semantic instructions will provide participants with a brief overview of how words can be related by meaning and will give participants examples of item pairs with high and low levels of semantic overlap. Finally, the thematic instructions will explain how concepts can be connected by overarching themes. These instruction sets are modeled after Buchanan (2010) and Valentine and Buchanan (2013).

Participants will then rate the relatedness of the word pairs based on the set of

instructions they receive at the start of each judgment block. These judgments will be made using a scale of zero (no relatedness between pairs) to one hundred (a perfect relationships). Judgments were recorded by the participant typing it into the survey. Participants will complete each of the three judgment blocks in this manner, with judgment instructions changing with each block. Three versions of the study will be created to counter balance the order in which judgment blocks appear. Participants will be randomly assigned to survey conditions. After completing the judgment blocks, participants will be presented with a short distractor task to account for recency effects. This section will be timed to last two minutes, and will task participants with alphabetizing a scrambled list of the fifty U.S. states. Once two minutes elapses, participants will automatically progress to a cued recall task, in which they will be presented with each of the 63 cues that had previously been judged as cue-target pairs. Participants will be asked to complete each word pair with the appropriate target word, based on the available cue word. Presentation of these pairs will be randomized, and participants will be informed that there is no penalty for guessing.

Results

First, the results from the recall section will be coded as zero for incorrect responses and one for correct responses. NA will be used to denote missing responses from participants who did not complete the recall section. Responses that are words instead of numbers in the judgment phase will be deleted and treated as missing data. Data will then be screened for out of range judgment responses (i.e., responses greater than 100), recall and judgment scores will be screened for outliers using Mahalanobis distance at $p < .001$, and multicollinearity between predictor variables will be measured with Pearson correlations. Mean judgment and recall scores will also be reported for each judgment condition.

Multilevel modeling will then be used to analyze the data. First, network norms and neighborhood norms will be mean centered, so as to control for multicollinearity. Next, two maximum likelihood multilevel models will be created. These models will be both use the

280 network norms as predictors and will examine their effects on recall and judgments. The goal
281 of these models is to replicate three-way interaction findings from the pilot study. If
282 significant three-way interactions are found between the network norms, these interactions
283 will be broken down with moderation analyses. Finally, neighborhood norms will be added
284 introduced into each model in steps. Initially, base word norms will be added, followed by
285 lexical information, rated properties, and norms measuring neighborhood connections.

References

- Adelman, J. S., & Brown, G. D. A. (2007). Phonographic neighbors, not orthographic neighbors, determine word naming latencies. *Psychonomic Bulletin & Review*, 14, 455–459.
- Baayen, R. H., Piepenbrock, R., & Gulikers, L. (1995). The CELEX lexical database (CD-ROM). Philadelphia.
- Bates, D., Machler, M., Bolker, B., & Walker, S. (2015). Fitting Linear Mixed-Effects Models Using lme4. *Journal of Statistical Software*, 67(1), 1–48.
- Brysbaert, M., & New, B. (2009). Moving beyond Kučera and Francis: A critical evaluation of current word frequency norms and the introduction of a new and improved word frequency measure for American English. *Behavior Research Methods*, 41(4), 977–990. doi:[10.3758/BRM.41.4.977](https://doi.org/10.3758/BRM.41.4.977)
- Brysbaert, M., & Stevens, M. (2018). Power Analysis and Effect Size in Mixed Effects Models: A Tutorial. *Journal of Cognition*, 1(1), 1–20. doi:[10.5334/joc.10](https://doi.org/10.5334/joc.10)
- Brysbaert, M., Keuleers, E., & New, B. (2011). Assessing the usefulness of Google Books' word frequencies for psycholinguistic research on word processing. *Frontiers in Psychology*, 2(MAR), 1–8. doi:[10.3389/fpsyg.2011.00027](https://doi.org/10.3389/fpsyg.2011.00027)
- Brysbaert, M., Warriner, A. B., & Kuperman, V. (2013). Concreteness ratings for 40 thousand generally known English word lemmas. *Behavior Research Methods*, 41, 977–990.
- Buchanan, E. M. (2010). Access into Memory: Differences in Judgments and Priming for Semantic and Associative Memory. *Journal of Scientific Psychology*, (March), 1–8. Retrieved from <http://www.psychencelab.com/images/Access%7B%7Dinto%7B%7DMemory%7B%7D%7B%7DDifferences%7B%7Din%7B%7D>
- Buchanan, E. M., Holmes, J. L., Teasley, M. L., & Hutchison, K. A. (2013). English semantic word-pair norms and a searchable Web portal for experimental stimulus

creation. *Behavior Research Methods*, 45(3), 746–757. doi:10.3758/s13428-012-0284-z

Buhrmester, M., Kwang, T., & Gosling, S. D. (2011). Amazon’s Mechanical Turk.

Perspectives on Psychological Science, 6(1), 3–5. doi:10.1177/1745691610393980

Burgess, C., & Lund, K. (1997). Representing abstract words and emotional connotation in

a high-dimensional memory space. *Proceedings of the Cognitive Science Society*,

61–66. Retrieved from

http://books.google.com/books?hl=en&lr=&id=sQyJiDk45HEC&oi=fnd&pg=PA61&dq=Jbgs6O8i27OKajqGo{_}ADoko{\\%}5Cnpapers3://publication/uuid/FE5168D9-C7C7-4C0F

[Jbgs6O8i27OKajqGo{_}ADoko{\\%}5Cnpapers3://publication/uuid/FE5168D9-](http://books.google.com/books?hl=en&lr=&id=sQyJiDk45HEC&oi=fnd&pg=PA61&dq=Jbgs6O8i27OKajqGo{_}ADoko{\\%}5Cnpapers3://publication/uuid/FE5168D9-C7C7-4C0F)

[C7C7-4C0F](http://books.google.com/books?hl=en&lr=&id=sQyJiDk45HEC&oi=fnd&pg=PA61&dq=Jbgs6O8i27OKajqGo{_}ADoko{\\%}5Cnpapers3://publication/uuid/FE5168D9-C7C7-4C0F)

Carreiras, M., Perea, M., & Grainger, J. (1997). Effects of orthographic neighborhood in

visual word recognition: cross-task comparisons. *Journal of Experimental Psychology.*

Learning, Memory, and Cognition, 23(4), 857–871. doi:10.1037/0278-7393.23.4.857

Coltheart, M., Davelaar, E., Jonasson, T., & Besner, D. (1977). Access to the internal

lexicon. In S. Dornic (Ed.), *Attention and performance vi* (pp. 535–555). Hillsdale,

NJ: Earlbaum.

Cowan, N., Baddeley, A. D., Elliott, E. M., & Norris, J. (2003). List composition and the

word length effect in immediate recall: A comparison of localist and globalist

assumptions. *Psychonomic Bulletin and Review*, 10(1), 74–79.

doi:10.3758/BF03196469

Dewhurst, S. a., Hitch, G. J., & Barry, C. (1998). Separate effects of word frequency and age

of acquisition in recognition and recall. *Journal of Experimental Psychology:*

Learning, Memory, and Cognition, 24(2), 284–298. doi:10.1037/0278-7393.24.2.284

Dunlosky, J., & Nelson, T. O. (1994). Does the sensitivity of judgments of learning (JOLs)

to the effects of various study activities depend on when the JOLs occur?

doi:10.1006/jmla.1994.1026

Green, P., & MacLeod, C. J. (2016). SIMR: An R Package for Power Analysis of Generalized

Linear Mixed Models by Simulation. *Methods in Ecology and Evolution*, 7(4), 493–498.

Kucera, H., & Francis, W. N. (1967). *Computational analysis of present-day English*. Providence, RI: Brown University Press.

Kuperman, V., Stadthagen-Gonzalez, H., & Brysbaert, M. (2012). Age-of-acquisition ratings for 30,000 English words. *Behavior Research Methods*, 44(4), 978–990. doi:[10.3758/s13428-012-0210-4](https://doi.org/10.3758/s13428-012-0210-4)

Landauer, T. K., & Dumais, S. T. (1997). A solution to Plato’s problem: The latent semantic analysis theory of acquisition, induction, and representation of knowledge. *Psychological Review*, 104(2), 211–240. doi:[10.1037//0033-295X.104.2.211](https://doi.org/10.1037//0033-295X.104.2.211)

Landauer, T. K., Foltz, P. W., & Laham, D. (1998). An introduction to latent semantic analysis. *Discourse Processes*, 25(2), 259–284. doi:[10.1080/01638539809545028](https://doi.org/10.1080/01638539809545028)

Lund, K., & Burgess, C. (1996). Producing high-dimensional semantic spaces from lexical co-occurrence. *Behavior Research Methods, Instruments, & Computers*, 28(2), 203–208. doi:[10.3758/BF03204766](https://doi.org/10.3758/BF03204766)

Maki, W. S. (2007a). Judgments of associative memory. *Cognitive Psychology*, 54(4), 319–353. doi:[10.1016/j.cogpsych.2006.08.002](https://doi.org/10.1016/j.cogpsych.2006.08.002)

Maki, W. S. (2007b). Separating bias and sensitivity in judgments of associative memory. *Journal of Experimental Psychology. Learning, Memory, and Cognition*, 33(1), 231–7. doi:[10.1037/0278-7393.33.1.231](https://doi.org/10.1037/0278-7393.33.1.231)

McRae, K., Cree, G. S., Seidenberg, M. S., & McNorgan, C. (2005). Semantic feature production norms for a large set of living and nonliving things. *Behavior Research Methods*, 37(4), 547–559. doi:[10.3758/BRM.40.1.183](https://doi.org/10.3758/BRM.40.1.183)

Nelson, D. L., McEvoy, C. L., & Schreiber, T. A. (2004). The University of South Florida free association, rhyme, and word fragment norms. *Behavior Research Methods, Instruments, & Computers*, 36(3), 402–407. doi:[10.3758/BF03195588](https://doi.org/10.3758/BF03195588)

Nelson, D. L., Schreiber, T. A., & Xu, J. (1999). Cue set size effects: sampling activated

associates or cross-target interference? *Memory & Cognition*, 27(3), 465–477.

doi:[10.3758/BF03211541](https://doi.org/10.3758/BF03211541)

Nelson, T. O., & Dunlosky, J. (1991). When people’s judgments of learning (JOLs) are extremely accurate at predicting subsequent recall: The delayed-JOL effect.

Psychological Science, 2(4), 267–270. doi:[10.1111/j.1467-9280.1991.tb00147.x](https://doi.org/10.1111/j.1467-9280.1991.tb00147.x)

Nelson, T. O., & Schreiber, T. A. (1992). Word concreteness and word structure as independent determinants of recall. *Journal of Memory and Language*, 31, 237–260.

Peereman, R., & Content, A. (1997). Orthographic and phonological neighborhoods in naming: Not all neighbors are equally influential in orthographic space. *Journal of Memory and Language*, 37, 382–410.

Pinheiro, J., Bates, D., Debroy, S., Sarkar, D., & R Core Team. (2017). nlme: Linear and Nonlinear Mixed Effects Models. Retrieved from

<https://cran.r-project.org/package=nlme>

Schreiber, T. A., & Nelson, D. L. (1998). The relation between feelings of knowing and the number of neighboring concepts linked to the test cue. *Memory & Cognition*, 26(5), 869–883. doi:[10.3758/BF03201170](https://doi.org/10.3758/BF03201170)

Stadthagen-Gonzalez, H., & Davis, C. J. (2006). The Bristol norms for age of acquisition, imageability, and familiarity. *Behavior Research Methods*, 38, 598–605.

Toglia, M. P. (2009). Withstanding the test of time: The 1978 semantic word norms.

Behavior Research Methods, 41(2), 531–533. doi:[10.3758/BRM.41.2.531](https://doi.org/10.3758/BRM.41.2.531)

Toglia, M. P., & Battig, W. F. (1978). *Handbook of semantic word norms*. Hillsdale, NJ: Earlbaum.

Valentine, K. D., & Buchanan, E. M. (2013). JAM-boree: An application of observation oriented modelling to judgements of associative memory. *Journal of Cognitive*

Psychology, 25(4), 400–422. doi:[10.1080/20445911.2013.775120](https://doi.org/10.1080/20445911.2013.775120)

Vinson, D. P., & Vigliocco, G. (2008). Semantic feature production norms for a large set of objects and events. *Behavior Research Methods*, 40(1), 183–190.

doi:[10.3758/BRM.40.1.183](https://doi.org/10.3758/BRM.40.1.183)

- Warriner, A. B., Kuperman, V., & Brysbaert, M. (2013). Norms of Valence, Arousal, and Dominance for 13,915 English Lemmas. *Behavior Research Methods*, 45(4), 1191–1207.
- Zeno, S. M., Ivens, S. H., Millard, R. T., & Duvvuri, R. (1995). *The educators's word frequency guide*. Brewster, NY: Touchstone Applied Science.

Table 1

Summary Statistics for Network Norms

Variable	Citation	Mean	SD	Min	Max
FSG	Nelson, McEvoy, and Schrieber, 2004	0.13	0.19	0.01	0.83
COS	Maki, McKinley, and Thompson, 2004	0.42	0.29	0	0.84
LSA	Landauer and Dumais, 1997	0.38	0.2	0.05	0.88

Note. COS: Cosine, FSG: Forward Strength, LSA: Latent Semantic Analysis.

Table 2

Summary Statistics of Single Word Norms for Cue Items

Variable	Citation	Mean	SD	Min	Max
QSS	Nelson et al., 2004	14.76	4.45	4	24
TSS	Nelson et al., 2004	14.59	4.54	4	24
Concreteness	Nelson et al., 2004	5.35	1	1.98	7
HAL Frequency	Lund and Burgess, 1996	9.34	1.67	6.26	13.39
SUBTLEX Frequency	Brysbaert and New, 2009	3.15	0.74	1.76	5.2
Length	Buchanan et al., 2013	4.9	1.5	3	10
Ortho N	Buchanan et al., 2013	7.44	5.91	0	19
Phono N	Buchanan et al., 2013	19	15.11	0	51
Phonemes	Buchanan et al., 2013	3.94	1.39	2	9
Syllables	Buchanan et al., 2013	1.35	0.6	1	3
Morphemes	Buchanan et al., 2013	1.1	0.3	1	2
AOA	Kuperman et al., 2012	5.15	1.53	2.47	8.5
Valence	Warriner et al., 2013	5.77	1.23	1.91	7.72
Imageability	Toglia and Battig, 1978	5.52	0.68	3.22	6.61
Familiarity	Toglia and Battig, 1978	6.17	0.28	5.58	6.75

Note.

Table 3

Summary Statistics of Single Word Norms for Target Items

Variable	Citation	Mean	SD	Min	Max
QSS	Nelson et al., 2004	15.44	4.86	5	26
TSS	Nelson et al., 2004	15.44	4.86	5	26
Concreteness	Nelson et al., 2004	5.4	1.01	1.28	7
HAL Frequency	Lund and Burgess, 1996	9.78	1.52	6.05	13.03
SUBTLEX Frequency	Brysbaert and New, 2009	3.34	0.64	1.59	4.74
Length	Buchanan et al., 2013	4.62	1.67	3	10
Ortho N	Buchanan et al., 2013	9.02	7.77	0	29
Phono N	Buchanan et al., 2013	21.51	16.71	0	59
Phonemes	Buchanan et al., 2013	3.7	1.5	1	10
Syllables	Buchanan et al., 2013	1.25	0.54	1	3
Morphemes	Buchanan et al., 2013	1.05	0.21	1	2
AOA	Kuperman et al., 2012	4.87	1.56	2.5	9.16
Valence	Warriner et al., 2013	5.84	1.27	1.95	7.89
Imageability	Toglia and Battig, 1978	5.5	0.71	2.95	6.43
Familiarity	Toglia and Battig, 1978	6.28	0.32	5.19	6.85

Note.

Table 4

Summary Statistics of Single Word Norms for All Items

Variable	Citation	Mean	SD	Min	Max
QSS	Nelson et al., 2004	15.1	4.65	4	26
TSS	Nelson et al., 2004	15.02	4.7	4	26
Concreteness	Nelson et al., 2004	5.38	1	1.28	7
HAL Frequency	Lund and Burgess, 1996	9.56	1.6	6.05	13.39
SUBTLEX Frequency	Brysbaert and New, 2009	3.25	0.7	1.59	5.2
Length	Buchanan et al., 2013	4.76	1.59	3	10
Ortho N	Buchanan et al., 2013	8.23	6.92	0	29
Phono N	Buchanan et al., 2013	20.26	15.92	0	59
Phonemes	Buchanan et al., 2013	3.82	1.44	1	10
Syllables	Buchanan et al., 2013	1.3	0.57	1	3
Morphemes	Buchanan et al., 2013	1.08	0.26	1	2
AOA	Kuperman et al., 2012	5.01	1.55	2.47	9.16
Valence	Warriner et al., 2013	5.8	1.24	1.91	7.89
Imageability	Toglia and Battig, 1978	5.51	0.69	2.95	6.61
Familiarity	Toglia and Battig, 1978	6.22	0.3	5.19	6.85

Note.