

# Grounding the Ungrounded: Estimating Locations of Unknown Place Names from Linguistic Associations and Grounded Representations

Gabriel Recchia (grecchia@memphis.edu)<sup>1</sup>  
Max M. Louwerse (M.M.Louwerse@uvt.nl)<sup>1,2</sup>

<sup>1</sup>Department of Psychology / Institute for Intelligent Systems, University of Memphis  
365 Innovation Drive, Suite 303, Memphis, TN 38152, USA

<sup>2</sup>Tilburg Center for Cognition and Communication, Tilburg University  
Warandelaan 2, 5037 AB Tilburg, The Netherlands

## Abstract

Spatial locations can be extracted from language statistics, based on the idea that nearby locations are mentioned in similar linguistic contexts, akin to Tobler's first law of geography. However, the performance of language-based estimates is inferior to human estimates, raising questions about whether human spatial representations can actually be informed by such (inferior) statistics. We show that alternative methods of computing co-occurrence statistics improve language-based estimates, illustrating that simple linguistic associations may in fact inform spatial representations. Most importantly, we show that by bootstrapping from grounded city locations, linguistic associations can be exploited to accurately estimate the locations of unknown cities, as well as human estimates of city locations. These results support the hypothesis that (ungrounded) linguistic associations can be productively combined with pre-existing spatial representations to yield new grounded representations, shedding light on the issue of symbol grounding in cognition.

**Keywords:** symbol grounding; geography; embodiment; symbolic cognition; embodied cognition; pointwise mutual association; latent semantic analysis

## Introduction

Accounts of meaning acquisition have long contrasted symbolic accounts with accounts that emphasize the importance of embodiment and perceptual grounding (de Vega, Glenberg, & Graesser, 2008). Historically, symbolic accounts treated semantics as being localized in an amodal system in which linguistic symbols<sup>1</sup> played a crucial role, whereas embodied accounts denied that the properties of a linguistic symbol (e.g. its orthography, length, or the other symbols with which it most frequently co-occurs) played any role in the symbol's meaning (Barsalou, 2009). Embodied accounts frequently emphasize the role of *perceptual simulation*, i.e., the construction of meaning via activation of the same perceptual and motor representations involved in first-hand experience. Today, the recognition that perceptual simulation plays an important role in semantic representation has shifted the focus of the debate

to a better understanding of the role of linguistic symbols in semantic representation. In some cases, the same data are used by different researchers to advance radically different hypotheses. For example, Louwerse and Jeuniaux (2010) and Barsalou, Santos, Simmons and Wilson (2008) agree that words are activated earlier than perceptual simulations, that words and perceptual simulations influence each other in the course of processing, and that linguistic associations allow some conceptual tasks to be accomplished solely by retrieving words and co-occurrence information, but their interpretations of these findings differ substantially. Louwerse and Jeuniaux (2010) interpret these findings as evidence for the Symbol Interdependency Hypothesis, which posits that humans rely heavily upon language-based statistics and invoke perceptual simulation only when cued by the task or by relevant sensorimotor cues (Louwerse & Jeuniaux, 2010; Louwerse, 2011). In contrast, Barsalou et al. (2008) use these points to argue that only situated simulations represent semantic information, but that we can use statistics of linguistic forms to achieve adequate performance on some laboratory tasks without accessing a word's actual meaning. That is, theories differ over whether the "surface" properties of linguistic symbols, such as their frequencies and the frequencies of words that they co-occur with, actually play a part in semantic representation.

What do these theories predict for words with spatial semantics? Consider a passage containing three familiar geographical terms (*San Francisco, Berkeley, Oakland*) and a novel one (*Hayward*). According to theories in which linguistic symbols play no role in semantic representations, even an incredibly weak representation for *Hayward*—say, knowing that it is a city, probably one in the San Francisco Bay area—is made possible by perceptually simulating these terms in the setting of the contents of the passage (e.g., envisaging a map of the area). In contrast, theories that admit a role for linguistic symbols argue that if the word *Hayward* is associated with the other linguistic symbols in the text, no perceptual simulation may be constructed until required by the task at hand (for example, if the reader were later asked, "Where is Hayward *exactly*?"). The burden on simulation-only theories is to show that perceptual simulation necessarily occurs in any condition in which semantic meaning (in this case, locations of place terms) is

<sup>1</sup> In computational models, the linguistic symbols under consideration are typically single words, but could also include morphemes or grammatical constituents.

accessed. The burden on the Symbol Interdependency Hypothesis is to show (a) that it is actually possible to infer semantic representations (in this case, locations) by combining information from language with pre-existing grounded representations, and (b) that the human semantic system actually does so. The present work offers support for claim (a); claim (b) is beyond the scope of this paper (but see Louwerse & Benesh, 2012; Louwerse & Connell, 2011).

Estimates of the locations to which place names refer can be made on the basis of information implicit in patterns of place name usage in text (Davies, 2013; Louwerse & Zwaan, 2009; Louwerse & Benesh, 2012). This is possible because place names referring to nearby spatial locations tend to co-occur in language (Hecht & Moxley, 2009), consistent with Tobler's first law of geography (Tobler, 1970). Returning to the Hayward example in the San Francisco Bay area, the Symbol Interdependency Hypothesis predicts that if readers are asked for the location of Hayward after having read the passage containing Hayward, San Francisco, Oakland, and Berkeley, they should be able to reconstruct its probable location on the fly, even if they did not simulate its location when they first encountered the term. If so, it must be possible to infer the approximate location of Hayward by virtue of the fact that Hayward co-occurred with several cities for which the reader *does* have perceptually grounded representations. It has never been demonstrated that this is actually possible, nor is it self-evident. One concern is that the results of Louwerse and Zwaan (2009) suggest that the cognitive maps of individuals who relied solely upon language statistics would be fairly unreliable. Low accuracy of language statistics is not necessarily problematic for the Symbol Interdependency Hypothesis, which suggests that humans do not rely on language statistics alone, but combine it with perceptual information about the locations of known places (when necessary). But no mechanism has ever been proposed by which this might take place, nor has it been demonstrated that doing so would reduce the error in location estimates to reasonable levels of accuracy.

The current study aims to fill this gap. Study 1 aims to replicate Louwerse and Zwaan's (2009) study estimating geographical locations in the United States using newspaper articles. In addition, it tries to optimize various parameters known to influence the performance of co-occurrence-based methods in cognitive science. Study 2 extends Study 1 by investigating whether location estimates can be improved when the algorithm is given access to grounded information (actual locations of some cities) in order to demonstrate how an "ungrounded" algorithm with a long history of use in cognitive science (i.e., Latent Semantic Analysis) might be combined with grounded information to achieve more reliable estimates. Finally, Study 3 illustrates how large-scale co-occurrence frequencies can be combined with grounded information to achieve even higher levels of accuracy.

## Study 1

Study 1 aimed to replicate and extend Louwerse and Zwaan (2009), who illustrated how language-based statistics obtained from Latent Semantic Analysis could estimate U.S. city locations across three different corpora with low-to-moderate bidimensional correlations ( $r = .53, .28, \text{ and } .43$  for the *Wall Street Journal*, *New York Times*, and *Los Angeles Times*, respectively). Latent Semantic Analysis is a commonly used method of estimating the degree to which two words are associated in language, the mathematics of which are described in detail in Landauer and Dumais (1997). The current study investigates several parameters that are known to influence the performance of computational models of lexical semantics, such as document length, amount of data analyzed, and other properties of the semantic space (Bullinaria & Levy, 2007; Quesada, 2006), to determine whether alternative parameterizations than those used by Louwerse and Zwaan (2009) result in more accurate location estimates. If so, language statistics may encode geographical information more robustly than is currently acknowledged.

## Method

The Los Angeles Times corpus, the New York Times News Syndicate corpus, and the Wall Street Journal corpus were obtained from the North American News Text Corpus (Graff, 1995), and documents not containing at least one of the 50 U.S. cities whose locations were estimated by Louwerse and Zwaan (2009) were removed from the corpus. Location estimates were obtained as in Louwerse & Zwaan (2009) by applying a multidimensional scaling (MDS) algorithm to a 50 x 50 matrix of text-based similarities populated with LSA cosines between city names. The matrix was converted to a matrix of dissimilarities with a Euclidean transformation (SPSS 20), and the standard MDS algorithm included with SPSS 20 (ALSCAL) was applied. Results were evaluated by computing differences between actual and predicted coordinates via bidimensional regression. Affine bidimensional regressions (Friedman & Kohler, 2003) were computed with the BiDimRegression package in the R software environment.

**Document length.** We varied minimum document length, maximum document length, and both at once. Tails were cropped with respect to the absolute deviation around the median (Leys et al., 2013).

**Corpus size.** Investigating precisely how performance varies with corpus size can help us understand exactly how much text is required to obtain acceptable estimates. Nineteen randomly generated subsets were extracted from each corpus, nine ranging from 2,000-10,000 documents and ten ranging from 10,000-100,000 documents.

**Dimensionality.** LSA proceeds by computing a low-rank approximation of a term-by-document matrix in which rows represent terms and columns represent documents. The

mathematical rank of the resulting matrix is known as its dimensionality. Generally speaking, the lower the rank, the coarser the representation of each word. For example, a dimensionality of 100 means that each word can be represented by a vector of only 100 numbers. Selecting an appropriate number of dimensions for the task at hand is an important component of an LSA analysis, although it is common to select 300, a number that has been empirically shown to work well on a wide variety of tasks (Landauer, Laham, & Derr, 2004). However, as Davies (2013) notes, “the oft-quoted 300 dimensions may well not be the optimum number to choose for this kind of LSA application... It seems unlikely that the same number of domains will be optimal for [geographic estimation]” (p. 333). There is no well-accepted method for selecting the optimum number of dimensions for a particular task other than exhaustively testing several choices (Quesada, 2006). We computed LSA vectors for varying numbers of dimensions from 50 to 300 to determine whether a peak existed at a lower dimensionality.

## Results and Discussion

Due to the great number of possible combinations of parameter values, parameters are discussed independently, using the default settings (all document lengths, all documents, 300 dimensions) as the baseline.

**Document length.** The effect of imposing a maximum document length was inconsistent and did not lead to significant increments or decrements in performance. In contrast, imposing a minimum document length improved performance. LSA's mean performance across corpora peaked when documents shorter than 1.8 median deviations were excluded from the space (mean  $r = .66$ ; Figure 1).

**Corpus size.** LSA cosines performed relatively poorly on small corpora (Figure 2). However, results also suggested that respectable performance can be achieved on corpora as small as twenty thousand documents.

**Dimensionality.** Average performance revealed a U-shaped curve with a peak at 200 dimensions. Building the LSA space with a dimensionality of 2 resulted in severely suboptimal performance, suggesting that it is necessary to first build the space with a higher dimensionality and to transform the resulting matrix to a two-dimensional approximation via MDS. Each of the three corpora exhibited a generally U-shaped curve with downturn in performance at 300 dimensions, suggesting that this number is suboptimal for this particular task.

**Conclusion.** Study 1 demonstrated that spatial information can be extracted from language statistics more effectively than previously reported (Louwerse & Zwaan, 2009). We obtained  $r$ s of  $> .6$  for all three corpora, vs. previous bests of  $r = .43$ ,  $r = .28$ , and  $r = .53$  for the Los Angeles Times, New York Times, and Wall Street Journal, respectively.

However, humans also represent perceptual information that could theoretically be combined with linguistic information to generate more informed location estimates for unknown place names. Study 2 explored a novel way to extend the method employed in Study 1 to accomplish this task.

## Study 2

In previous work (Davies, 2013; Louwerse & Zwaan, 2009; Louwerse & Benesh, 2012), distances have been approximated from language statistics alone. However, the Symbol Interdependency Hypothesis predicts that, when necessary, individuals integrate language statistics with perceptually grounded information (Louwerse & Jeuniaux, 2010; Louwerse, 2011). One way to amend Study 1 to include perceptually grounded information about known places is to replace estimated distances for all cities except for one—the “unknown place”—with actual distances among those known places. In classical MDS, there is no way to specify fixed latitudes and longitudes of known locations. These deficiencies are addressed with PROXSCAL (Busing, Commandeur, and Heiser, 1999), an implementation of a common majorization algorithm. PROXSCAL permits fixed vectors to be specified for some variables, while other variables are permitted to vary freely. In Study 2, we exploit this property to perform a replication of study 1 that integrates grounded locations. We hypothesize that this method will achieve more reliable estimates, demonstrating that it is possible to infer more accurate spatial representations by combining language statistics with grounded representations.

## Method

Text-based estimates for each city location were obtained via the PROXSCAL algorithm. For each run of the algorithm, one of the 50 U.S. cities from the previous study was treated as an unknown place, while the other 49 were treated as known places. A dissimilarity matrix was constructed in which cells were populated with actual distances between known places. As in Study 1, each row and column corresponded to one of the 50 city names. Cells for which the row or column corresponded to the unknown place were populated with co-occurrence-based similarity estimates (LSA cosines) as in the previous study. LSA cosines in the row and column corresponding to the unknown place were converted to distances by subtracting each from 1. The row and column corresponding to the unknown place were transformed so as to have a mean and standard deviation comparable to the other vectors in the matrix (mean = the average mean of the vectors of distances among known places; SD = the average SD of the vectors of distances among known places). Finally, we applied PROXSCAL, specifying that the known places should be fixed, and providing their latitudes and longitudes.

This process was run 50 times for each parameter setting (with a different city treated as the “unknown place”) each time, yielding a set of 50 estimated city locations, and all parameters were investigated independently as before.

## Results and Discussion

Performance of the PROXSCAL-based method was higher than the ALSICAL-based method (Study 1) across the full range of parameter settings (Figures 1 and 2), although mean  $r$ s after parameter optimization were similar (Study 1,  $r = .66$ ; Study 2,  $r = .68$ )<sup>2</sup>. Accuracy of PROXSCAL was more stable as parameters varied, with mean  $r$ s ranging from .63-.68 for all choices of minimum document length (vs. .24-.66 for Study 1), .63-.65 for corpus sizes ranging from 20,000 to 100,000 documents (vs. .46-.61 for Study 1), and .58-.68 for all choices of dimensionality of 50 or greater (vs. .49-.66 for Study 1). In all, estimates from the PROXSCAL-based method that included information about the locations of known places were consistently more accurate than those in Study 1.

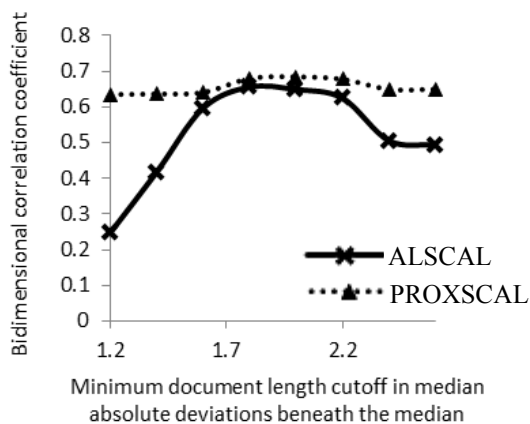


Figure 1: Performance as a function of minimum document length, averaged across corpora. Higher values on the x-axis correspond to more documents included (e.g., at the far right, only documents with lengths 2.6 deviations below the median were excluded).

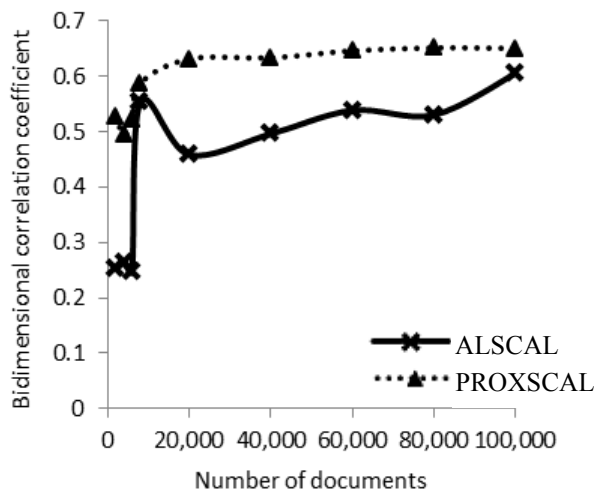


Figure 2: Performance as a function of corpus size, averaged across corpora.

<sup>2</sup> Reported here are optimal  $r$ s (mean) when the same parameter settings are used across all three corpora. Higher  $r$ s are achieved if parameters are fit separately for each individual news corpus.

## Study 3

The mean estimates in Studies 1 and 2 correlated more closely with actual locations than the best correlations previously reported from language statistics alone. In particular, including grounded information made the method much less sensitive to the specific parameter settings employed. However, it remains uncertain whether estimates based on raw co-occurrences can take advantage of grounded information in an analogous manner, and whether language-based statistics combined with perceptual grounding can produce spatial representations that are accurate enough to be useful. Study 3 investigates the degree of accuracy that can be achieved from by combining grounded representations with linguistic associations.

### Method

We obtained a set of U.S. place—all U.S. cities, towns, and Hawaiian census-designated places<sup>3</sup> with a population of at least 20,000 in 2010—from the *2012 National Population Projections* of the U.S. Census Bureau. Of these, all that shared a name with another U.S. place of population greater than 20,000 were omitted, as were all place names that also happened to be English words appearing in the official Scrabble® dictionary. This yielded a set of 1,283 U.S. place names that we linked to their geographic coordinates in the Geographic Names Information System (U.S. Geological Survey, 2012). This constituted the set of *localized cities*. In contrast, the 50 place names from Louwerse & Zwaan (2009) were treated as *unlocalized cities*, the locations of which had to be inferred on the basis of their co-occurrences with localized cities.

We identified all 4-grams and 5-grams appearing in the Google Web 1T 5-gram corpus in which a localized and an unlocalized city co-occurred. Pointwise mutual information (PMI) scores (Church & Hanks, 1990) were computed between every pair of localized and unlocalized cities as in Manning & Schütze (1999) via the standard formula

$$I(x, y) = \log_2 \frac{P(x, y)}{P(x)P(y)}$$

Here,  $P(x)$  and  $P(y)$  can each be calculated as the frequency of  $x$  and  $y$  (respectively) divided by the total number of tokens in the corpus.  $P(x, y)$  is computed by dividing the number of times that  $x$  and  $y$  co-occur (in the same  $n$ -gram) by the total number of tokens in the corpus. PMI essentially normalizes the probability with which  $x$  and  $y$  co-occur by their overall frequencies. Thus, even though the terms “horse” and “the” co-occur more frequently than “horse” and “saddle” do, “horse” and “saddle” may still have a higher PMI, because “saddle” is a much less frequent word than “the.” To estimate the location of an unlocalized city name  $u$ , we calculated the city  $c$  that was closest to the  $k$  nearest neighbors of  $u$  (that is, the  $k$  cities with the highest PMIs to  $u$ ). After excluding outliers that occurred more than

<sup>3</sup> Officially, Honolulu is the only incorporated city in Hawaii.

600 miles away from  $c$ , we calculated the mean latitude and longitude of the remaining neighbors, which served as the estimate. Bidimensional regressions were computed between the actual locations of the 50 unlocalized cities and the coordinates estimated using this method.  $k$  was treated as a free parameter ranging from 1 to 20.

### Results and Discussion

After bidimensional regression, correlation coefficients between actual and estimated locations ranged from .78-.91, depending on the value of  $k$  used. Mean distance between actual and estimated locations was minimized at  $k = 10$ .

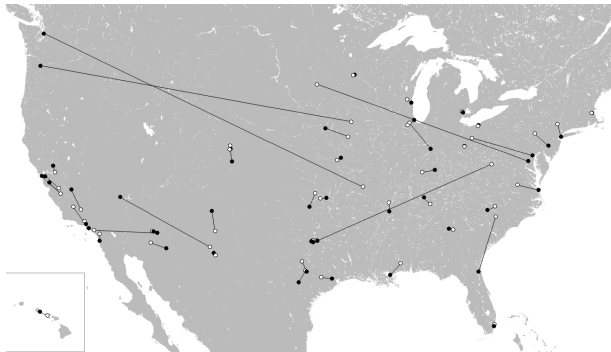


Figure 4. Actual locations of the fifty U.S. cities in the unlocalized dataset from Study 3 (black circles) are connected with lines to estimated locations of the same cities when  $k = 10$  (white circles).

Although these estimates are more accurate than those achieved in Study 2, they come with some caveats. First, Study 3 used much more linguistic data: an  $n$ -gram corpus based on roughly a terabyte of text, in contrast to the 100,000-document corpora used in Studies 1 and 2. In addition, a much greater amount of grounded information was employed. As such, the improvements of Study 3 should not be taken as evidence that the method employed is necessarily more effective than LSA, but solely as an existence proof that if one has grounded representations of a large number of cities, it is possible to accurately infer the locations of a place  $u$  from only linguistic co-occurrences, linguistic frequencies, and the locations of cities other than  $u$ . In our final study, we considered the extent to which language-based estimates predict human estimates.

### Study 4

Study 1 demonstrated that language-based estimates predict the actual locations of cities, while Study 2 and 3 showed that the predictive ability of language increases substantially when it is combined with grounded information (i.e., when the locations of the cities that a word co-occurs with in text are known). Study 4 considered whether inclusion of grounded information also improves the ability of language to predict human estimates of city locations.

### Method

Human estimates of city locations were obtained from Louwerse & Zwaan (2009). Participants in this study had estimated the location of the 50 unlocalized cities from Studies 1-3 on a blank sheet of paper. Affine bidimensional regression correlations were computed between each participant's estimates, the best-performing language-based estimates from Study 1, and the best-performing language/perception-based estimates from Studies 2 and 3.

### Results and Discussion

Median correlations between participant estimates and computational measures are reported in Table 1. For each corpus, correlations between human estimates and the language/perception-based estimates (Study 2, 3) were higher than or equal to the correlations between human estimates and the language-based estimates alone. Mean correlations exhibited the same pattern. These results suggest that human geographical estimates might be based in part on information implicitly coded in language and part on explicitly grounded spatial information.

	Correlation language estimate (Study 1) and human estimates	Correlation language/perception estimates (Study 2, 3) and human estimates
LA Times	.541***	.574***
NY Times	.546***	.592***
Wall St. J.	.569***	.602***
Google	--	.769***

Table 1. Median correlations between computational and human estimates. All correlations are significant at  $p < .001$ .

### Conclusion

According to the Symbolic Interdependency Hypothesis, it should be possible to infer much of the content of semantic representations—e.g., the locations to which place names refer—by combining information from language statistics with pre-existing grounded representations. In Study 1, we demonstrated that information from language statistics alone can estimate place names more accurately than had been previously recognized. Study 2 illustrated that by combining language statistics with pre-existing grounded representations, locations to which place names refer can indeed be estimated more accurately than from linguistic information alone. Study 3 showed that with enough data, even extremely simple co-occurrence-based measures can be combined with grounded representations to yield accurate city locations, showing that it is possible to bootstrap spatial semantics from associations. Finally, Study 4 showed that the computational estimates from Studies 1, 2, and 3 also predict human judgments of city locations, with the grounded language-based estimates having correlations greater than or equal to the estimates that relied on language statistics alone. Explorations of other methods

of quantifying relatedness between places may further refine our understanding of the relationship between place name co-occurrences and cognitive measures of place relatedness.

### Acknowledgements

This project was supported by a grant from the Intelligence Community Postdoctoral Research Fellowship Program through funding from the Office of the Director of National Intelligence. All statements of fact, opinion, or analysis expressed are those of the author and do not reflect the official positions or views of the Intelligence Community or any other U.S. Government agency. Nothing in the contents should be construed as asserting or implying U.S. Government authentication of information or Intelligence Community endorsement of the author's views.

### References

- Barsalou, L. W. (2009). Simulation, situated conceptualization, and prediction. *Philosophical Transactions of the Royal Society of London: Biological Sciences*, 364, 1281-1289.
- Barsalou, L. W., Santos, A., Simmons, W. K., & Wilson, C. D. (2008). Language and simulation in conceptual processing. In M. D. Vega, A. M. Glenberg, & A. C. Graesser, *Symbols, embodiment, and meaning* (pp. 245-283). Oxford: Oxford University Press.
- Bullinaria, J., & Levy, J. (2007). Extracting semantic representations from word co-occurrence statistics: A computational study. *Behavior Research Methods*, 39(3), 510-526.
- Busing, F. M., Commandeur, J. J., & Heiser, W. J. (1997). PROXSCAL: A multidimensional scaling program for individual differences scaling with constraints. In W. Bandilla, & F. Faulbaum, *Softstat '97: Advances in statistical software* (pp. 237-258). Stuttgart, Germany: Lucius.
- Davies, C. (2013). Reading geography between the lines: Extracting local place knowledge from text. In T. Tenbrink, J. Stell, A. Galton, & Z. Wood (Ed.), *Conference on Spatial Information Theory (COSIT)* (pp. 320-337). Scarborough, UK: Springer.
- de Vega, M., Glenberg, A. M., & Graesser, A. C. (2008). *Symbols, embodiment, and meaning*. Oxford: Oxford University Press.
- Friedman, A., & Kohler, B. (2003). Bidimensional regression: Assessing the configural similarity and accuracy of cognitive maps and other two-dimensional data sets. *Psychological Methods*, 8(4), 468-491.
- Graff, D. (1995). North American News Text Corpus. [Software resource]. Philadelphia: Linguistic Data Consortium.
- Hecht, B., & Moxley, E. (2009). Terabytes of Tobler: Evaluating the first law in a massive, domain-neutral representation of world knowledge. In K. Stewart Hornsby, C. Claramunt, M. Denis, & G. Ligozat (Ed.), *Conference on Spatial Information Theory (COSIT)* (pp. 88-105). Aber Wrac'h, France: Springer.
- 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.
- Landauer, T. K., Laham, D., & Derr, M. (2004). From paragraph to graph: Latent semantic analysis for information visualization. *Proceedings of the National Academy of Sciences of the United States of America*, 101, 5214-5219.
- Lays, C., Ley, C., Klein, O., Bernard, P., & Licata, L. (2013). Detecting outliers: Do not use standard deviation around the mean, use absolute deviation around the median. *Journal of Experimental Social Psychology*.
- Louwerse, M. M. (2011). Symbol interdependency in symbolic and embodied cognition. *TopiCS in Cognitive Science*, 3, 273-302.
- Louwerse, M. M., & Benesh, N. (2012). Representing spatial structure through maps and language: Lord of the Rings encodes the spatial structure of Middle Earth. *Cognitive Science*, 36(8), 1556-1569.
- Louwerse, M. M., & Connell, L. C. (2011). A taste of words: Linguistic context and perceptual simulation predict the modality of words. *Cognitive Science*, 35, 381-398.
- Louwerse, M. M., & Jeuniaux, P. (2010). The linguistic and embodied nature of conceptual processing. *Cognition*, 114, 96-104.
- Louwerse, M. M., & Zwaan, R. A. (2009). Language encodes geographical information. *Cognitive Science*, 33(1), 51-73.
- Manning, C. D., & Schütze, H. (1999). *Foundations of statistical natural language processing*. Cambridge, MA: MIT Press.
- Quesada, J. (2006). Creating your own LSA spaces. In T. Landauer, D. McNamara, S. Dennis, & W. Kintsch, *Latent Semantic analysis: A Road to Meaning* (pp. 71-85). Mahwah: Erlbaum.
- Tobler, W. R. (1970). A computer movie simulating urban growth in the Detroit region. *Economic Geography*, 46, 234-240.
- U.S. Geological Survey. (2012). *BGN: Domestic Names*. Retrieved from Geonames: <http://geonames.usgs.gov/domestic/index.html>