

# Linguistic Obfuscation in Fraudulent Science

Journal of Language and Social Psychology  
2016, Vol. 35(4) 435–445  
© The Author(s) 2015  
DOI: 10.1177/0261927X15614605  
jls.sagepub.com



**David M. Markowitz<sup>1</sup> and Jeffrey T. Hancock<sup>1</sup>**

## **Abstract**

The rise of scientific fraud has drawn significant attention to research misconduct across disciplines. Documented cases of fraud provide an opportunity to examine whether scientists write differently when reporting on fraudulent research. In an analysis of over two million words, we evaluated 253 publications retracted for fraudulent data and compared the linguistic style of each paper to a corpus of 253 unretracted publications and 62 publications retracted for reasons other than fraud (e.g., ethics violations). Fraudulent papers were written with significantly higher levels of linguistic obfuscation, including lower readability and higher rates of jargon than unretracted and nonfraudulent papers. We also observed a positive association between obfuscation and the number of references per paper, suggesting that fraudulent authors obfuscate their reports to mask their deception by making them more costly to analyze and evaluate. This is the first large-scale analysis of fraudulent papers across authors and disciplines to reveal how changes in writing style are related to fraudulent data reporting.

## **Keywords**

text analysis, deception, scientific fraud, LIWC, Coh-Metrix

High-profile scandals and a rise in retractions due to scientific fraud (Fang, Steen, & Casadevall, 2012) have led to a focus on the scientific process, such as reconsidering peer review and reporting protocols (Casadevall & Fang, 2012; Simonsohn, 2013). Here we investigate scientific misconduct by asking a psychological question: Do scientists write differently when reporting fraudulent research?

Prior work has evaluated the writing style of a single fraudulent author, social psychologist Diederik Stapel, finding that his writing style differed across his fraudulent and genuine papers (Markowitz & Hancock, 2014). For example, compared with his writing in genuine papers, Stapel used fewer adjectives when describing false data but more words related to methods and procedures. This initial investigation suggests that

---

<sup>1</sup>Stanford University, Stanford, CA, USA

## **Corresponding Author:**

David M. Markowitz, Stanford University, McClatchy Hall, Building 120, Stanford, CA 94305, USA.  
Email: markowitz@stanford.edu

there may be linguistic differences signaling how fraudulent and genuine science reports are written. It is unclear, however, if linguistic patterns of fraud generalize when multiple authors from different countries and domains of science are involved. From a large database of science publications, we analyze the writing style of papers retracted for scientific fraud and compare them to matched unretracted papers and papers retracted for reasons other than fraud (e.g., ethics concerns, authorship issues).

The idea that deception can lead to changes in language use is consistent with a growing literature suggesting that psychological dynamics can be revealed in word patterns. For example, criminal psychopaths describe their murders differently than nonpsychopaths by using more causal terms (e.g., *because*, *result*) and fewer social words (e.g., *public*, *someone*) in an interview setting (Hancock, Woodworth, & Porter, 2013). Furthermore, individuals who report depressive symptoms are more self-focused and use more negative emotion terms (e.g., *hate*, *dislike*) than nondepressive individuals when writing about college life (Rude, Gortner, & Pennebaker, 2004).

The general finding that language can be used as a marker of psychological change (Pennebaker, 2011) has also been applied to deception with several studies using automated methods to analyze word patterns associated with false and truthful statements (Larcker & Zakolyukina, 2012; Pennebaker, 2011; Toma & Hancock, 2012). A recent meta-analysis of deception and language, for example, found that liars tend to express more negative emotion terms, use fewer first-person pronouns, and refer less often to cognitive processes than truth tellers (Hauch, Blandón-Gitlin, Masip, & Sporer, 2014). These effects, however, are substantially moderated by contextual factors, such as by the type of lie and the production mode of the communication.

One form of deception that resembles the context of scientific fraud is deceptive corporate financial reporting. Like scientific fraud, deceptive financial reporting involves writing about fraudulent data. Research on financial reports has revealed that deceptive reports have higher levels of linguistic obfuscation than accurate reports (Bloomfield, 2002; Courtis, 1998; Humpherys, Moffitt, Burns, Burgoon, & Felix, 2011; Li, 2008). Linguistic obfuscation involves reduced levels of readability and positive emotion terms (e.g., *improve*, *success*), but higher rates of causal terms (e.g., *depend*, *infer*), more abstraction (e.g., fewer articles, prepositions, and quantifiers), and more jargon than nonobfuscated language.

For instance, Li (2008) observed that financial companies with poor earnings write their yearly reports to the U.S. Securities and Exchange Commission differently than companies with positive earnings. Such reports are obfuscated with less readable writing (e.g., the documents are longer and sentence structure is more complex), more causal language, and a lower rate of positive emotion terms relative to thriving companies. Another study offers that fraudulent reports contain more complex words and extended sentences compared with nonfraudulent reports (Humpherys et al., 2011). This research suggests that companies with information to hide obfuscate in their reports to make it more difficult for readers to assess their company's performance or intentions, and this obfuscation effect can be observed in writing style patterns.

We test the linguistic obfuscation hypothesis on science papers retracted for data fraud. The hypothesis predicts that scientific papers retracted for fraud will have higher levels of linguistic obfuscation than unretracted papers and papers retracted for

reasons other than fraud (e.g., ethics concerns, authorship issues). Given that science papers typically include separate sections with distinct objectives and discourse requirements that affect writing style (Biber, Connor, & Upton, 2007), we also measure obfuscation by section.

## Method

Publications retracted for scientific misconduct were identified from the PubMed archive from 1973 through 2013. Using each paper's retraction notice, two independent coders determined whether the retraction was due to data fraud versus other scientific misconduct (e.g., ethics violations or authorship issues) and had good agreement ( $\kappa = .64$ ,  $p < .001$ ). Discrepancies were resolved by consulting additional sources where available (see Fang et al., 2012). Of the 315 retracted papers identified, primarily from biomedical journals, 253 papers were retracted for data fraud (e.g., faking data, manipulating data) with the remaining papers ( $N = 62$ ) retracted for ethics concerns or authorship issues.

For each retracted paper, an unretracted control paper was identified from the same journal, during the same year of publication, and matched where possible on keywords. There were some exceptions to this matching process. For 19% (48/253) of the retracted publications, there was no match from the same year, so the match was selected from an adjacent year. If a fraudulent paper did not contain keywords, terms from the abstract section of that paper were used as a substitute in the PubMed search. If this method did not yield any results from the PubMed archive, a keyword was not used and a paper was selected at random from the same journal and year. This occurred in 9% of the matches (24/253 publications). Thirteen *Science* publications were not included in the section-by-section analysis, as they did not contain distinct Introduction, Method, Results, and Discussion sections.

It is important to note that papers in the unretracted corpus are presumed genuine, but at some rate, there may be undetected fraud. Fang et al. (2012) estimate that approximately .007% of published research papers are retracted because of fraud or suspected fraud, suggesting that the papers in the unretracted corpus are unlikely to be fraudulent.

The final corpus of 253 fraudulent papers contained 1,033,400 words, whereas the matched unretracted corpus of 253 papers contained 1,005,929 words. The 62 publications retracted for reasons other than fraud contained 191,748 words.

## Database Preprocessing

Each text file was preprocessed according to the following method. First, a Python script converted words from British English (e.g., *tumour*, *analysed*) to American English (e.g., *tumor*, *analyzed*). This script is available from the authors.<sup>1</sup> Second, brackets, parentheses, and percent signs were removed to more accurately capture words per sentence, which is a component of readability statistics. Third, periods were removed from certain words to avoid influencing words per sentence (e.g., *Dr.*, *Inc.*, *Figs.*). Only main body text, excluding section titles, figures, tables, legends, and supplementary materials, was included in the analysis.

## The Obfuscation Index

Linguistic obfuscation was calculated as a single index by summing the standardized rates of causal terms (Li, 2008), the abstraction index (Larrimore, Jiang, Larrimore, Markowitz, & Gorski, 2011), and jargon, and subtracting the rate of positive emotion terms (Li, 2008) and Flesch Reading Ease readability (Flesch, 1948). A higher score on this index indicates that the text is more obfuscated than a lower score.

We used Linguistic Inquiry and Word Count (LIWC; Pennebaker, Booth, & Francis, 2007) and Coh-Metrix (McNamara, Louwrese, Cai, & Graesser, 2013) to capture the writing style elements that comprised the linguistic obfuscation index. Both LIWC and Coh-Metrix are well-validated computerized text analysis programs that calculate relative word frequencies across psychological, semantic, and part of speech categories (McNamara, Graesser, McCarthy, & Cai, 2014; Tausczik & Pennebaker, 2010) and have frequently been used to analyze the language of deception (for a review, see Hauch et al., 2014).

**Jargon.** We operationalized jargon by calculating the percentage of words not identified by the LIWC dictionary, which is an overall measure of common words in English (Pennebaker et al., 2007). LIWC contains a large dictionary ranging from part of speech (e.g., conjunctions, prepositions) to content categories (e.g., affect words, tentative words, causal terms). Words outside of the dictionary are specialized terms that are uncommon in everyday communication (Pennebaker, 2011; Tausczik & Pennebaker, 2010). Therefore, we calculated jargon for each paper by using the formula  $(100 - \text{Dictionary})$  and then standardizing the values.

**Abstraction.** An abstraction index was constructed by taking the inverse of the sum of the standardized LIWC scores for articles, prepositions, and quantifiers (Larrimore et al., 2011). Articles (e.g., *a*, *an*, *the*) make references to nouns, prepositions (e.g., *after*, *unless*, *except*) specify relationships between objects and people, and quantifiers (e.g., *more*, *less*, *significant*) express degrees of difference between objects. A high abstraction score suggests that the language is less descriptive and less concrete than a low abstraction score.

**Positive Emotion Terms and Causal Terms.** Both of these language categories were drawn from the standard LIWC dictionary. Positive emotion terms are words such as *support*, *worthwhile*, and *inspired*, whereas causal terms are words such as *depend*, *induce*, and *manipulated*.

**Flesch Reading Ease.** A single Flesch Reading Ease score was computed for each full paper and individual section using Coh-Metrix (McNamara et al., 2013). A lower Flesch Reading Ease score suggests that the text is less readable than text with a higher score (Flesch, 1948).

Correlations for the five obfuscation features and components of the abstraction index are described in Tables 1 and 2, respectively.

**Table 1.** Correlations Between Variables in the Obfuscation Index ( $N = 506$ ).

	Abstraction	Jargon	PE terms	Causal terms	FRE readability
Abstraction	—	.698**	-.231**	.256**	-.095*
Jargon	.698**	—	-.352**	.159**	-.061
PE terms	-.231**	-.352**	—	-.049	-.033
Causal terms	.256**	.159**	-.049	—	-.049
FRE readability	-.095*	-.061	-.033	-.049	—

Note. PE = Positive Emotion, FRE = Flesch Reading Ease.

\* $p < .05$ . \*\*  $p < .01$  (two-tailed).

**Table 2.** Correlations Between Variables in the Abstraction Index ( $N = 506$ ).

	Articles	Prepositions	Quantifiers
Articles	—	.206***	.263***
Prepositions	.206***	—	.176***
Quantifiers	.263***	.176***	—

\*\*\* $p < .001$ .

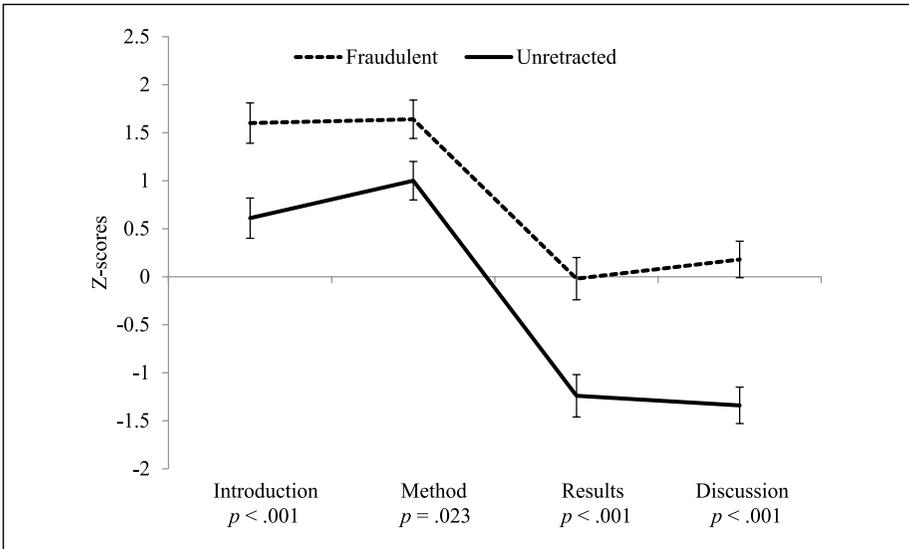
## Results

These data were analyzed using linear mixed models with paper type (fraudulent vs. unretracted) as a between-subjects factor. The number of words per paper,  $F(1, 504) = 0.57, p > .25$ , and authors per paper,  $F(1, 504) = 0.34, p > .25$ , were not statistically different across the corpora.

Consistent with the obfuscation hypothesis, papers retracted for data fraud ( $M = 1.17, SE = 0.23$ ) had higher levels of linguistic obfuscation than unretracted papers ( $M = -0.46, SE = 0.23$ ),  $F(1, 504) = 24.51, p < .001$ . As seen in Figure 1, fraudulent papers had higher levels of linguistic obfuscation across paper sections, and paper section did not interact with paper type,  $F(3, 1912) = 1.67, p = .17$ .

The linguistic obfuscation effect was observed for each variable that comprised the obfuscation index, suggesting that the effect was robust across these dimensions. Fraudulent papers contained more jargon,  $F(1, 504) = 11.37, p < .001$ ; more causal terms,  $F(1, 504) = 5.36, p = .021$ ; and were written more abstractly than unretracted papers,  $F(1, 504) = 14.92, p < .001$ . Fraudulent papers were less readable,  $F(1, 504) = 5.26, p = .022$ , and included fewer positive emotion terms than unretracted papers,  $F(1, 504) = 8.13, p = .005$ . Descriptive statistics and section results for each obfuscation dimension can be found in the Supplementary Material (available at <http://jls.sagepub.com/content/by/supplemental-data>).

To test the possibility that retracted papers are different from unretracted papers for reasons other than fraud, we examined whether linguistic obfuscation was higher for papers retracted for fraud ( $N = 253$ ) or for other misconduct (e.g., ethics or authorship issues;  $N = 62$ ). Papers retracted for fraud ( $M = 1.17, SE = 0.21$ ) had more obfuscation



**Figure 1.** Standardized obfuscation differences between paper type and across sections. Note. Dashed and solid lines represent fraudulent and unretracted papers, respectively. Error bars represent one standard error above and below the sample mean.

**Table 3.** Prediction Accuracy Rates Using a Cross-Validated Model.

	Hit	Miss	Hit Rate (%)	Accuracy (%)
Fraudulent ( $N = 253$ )	153	100	60.5	57.2
Unretracted ( $N = 253$ )	136	117	53.8	

than other retractions ( $M = -0.86$ ,  $SE = 0.43$ ),  $F(1, 313) = 17.96$ ,  $p < .001$ , suggesting that the linguistic obfuscation effect cannot be explained by retraction status only.

### Statistical Classification

Using a standard leave-one-out cross-validation technique with the dependent variable as paper type (fraudulent or unretracted), we examined the ability to detect deception statistically with the five obfuscation features (see Larcker & Zakolyukina, 2012). The model had good fit ( $\chi^2 = 25.18$ ,  $p < .001$ ) with a classification accuracy of 57.2%, a level consistent with human performance on deception detection (Bond & DePaulo, 2006; see Table 3).

Although this represents a statistically significant improvement over chance, it is clear that our limited model is not feasible for detecting fraudulent science with an especially problematic false-positive rate (46%). To improve the classification accuracy, more computationally sophisticated methods to analyze language patterns

(e.g., machine learning, natural language processing) will be required. These steps, in addition to widening the feature set beyond the theoretically derived obfuscation dimensions, should improve deception detection accuracy.

### *Psychological Mechanism*

What is the mechanism that leads to obfuscation in science writing when researchers publish fake data? The obfuscation literature suggests that fraudulent financial companies dissimulate their reports to make them more costly to analyze (Humpherys et al., 2011). If the observed effect reflects this goal, then obfuscation should be correlated with other cues in scientific reporting that are difficult to evaluate. One such cue is citations, which can serve as credibility markers that are costly to assess because they require the reader to obtain and appraise claims from an external source. Indeed, the obfuscation index was positively correlated with the number of references per paper ( $r = 0.31, p < .001$ ) with fraudulent papers ( $M = 42.47, SE = 0.99$ ) containing approximately 3.5 more references than unretracted papers ( $M = 38.92, SE = 0.99$ ),  $F(504) = 6.50, p = .011$ . These data suggest that fraudulent scientists obfuscate by increasing the cost of evaluating a paper to mask their deception. Given that our analysis controlled for journal and keywords, this effect cannot be explained by reporting conventions alone.

### *Alternative Explanations*

One possible explanation for our results is that different science domains have conventions that shape writing style and our findings were a reflection of genre differences instead of deception. To address this concern, we analyzed the influence of science domain on the overall obfuscation index. Four domains were identified from prior literature (Lu, Jin, Uzzi, & Jones, 2013): Biology and Medicine ( $N = 210$ ), Multidisciplinary Sciences ( $N = 29$ ), Other ( $N = 10$ ), and Social Sciences ( $N = 4$ ). The interaction of paper type (fraudulent vs. unretracted) and science domain for the obfuscation index was not significant,  $F(3, 498) = 0.18, p > .25$ , suggesting that the rate of obfuscation was not a function of the type of science reported on.

Another possible explanation for our results may be the fact that the papers were written by authors from different countries, where English may be a second language. To address whether writing style changes resulted from differences in geographical location, we organized each paper across four continents by the first author's home institution at the time of publication: Asia (fraudulent  $N = 79$ , unretracted  $N = 42$ ), Europe (fraudulent  $N = 61$ , unretracted  $N = 83$ ), North America (fraudulent  $N = 109$ , unretracted  $N = 116$ ), and Other (fraudulent  $N = 4$ , unretracted  $N = 12$ ). The interaction of paper type (fraudulent vs. unretracted) and continent for the obfuscation index was not significant,  $F(3, 498) = 0.89, p > .25$ , suggesting that geographic location was also not a factor in driving the obfuscation effect.

### Comparison to Related Work

Markowitz and Hancock (2014) evaluated the writing style patterns of Diederik Stapel, a prominent social psychologist found guilty of scientific fraud after extensive investigations of his research papers. Supporting the obfuscation hypothesis, Stapel's fraudulent writing was more obfuscated ( $M = 0.21$ ,  $SE = 0.65$ ) than his genuine writing ( $M = -0.21$ ,  $SE = 0.64$ ), although this trend did not reach significance. Given that the Stapel analysis (49 papers) was less than 10% of the current corpus (506 papers), this lack of power is not surprising.

### Discussion

The prevailing discussion on research misconduct has focused on scientific practice and policy such as the need for replication studies and avoiding problematic reporting practices, including *p*-hacking (e.g., altering data collection, conditions, or analyses to reach statistical significance; Simonsohn, Nelson, & Simmons, 2014). The present research suggests that linguistic analyses of scientific fraud can also advance our understanding of how deception affects communication. Scientists reporting fraudulent data wrote their reports with a significantly more obfuscated writing style than unretracted papers and papers retracted for reasons other than fraud (e.g., ethics violations, authorship issues). Furthermore, we found that linguistic obfuscation was correlated with the number of references per paper, suggesting that fraudulent scientists were using obfuscation to make claims in their papers more difficult and costly to assess (Humpherys et al., 2011).

While the classification accuracy for detecting fraudulent papers based on linguistic obfuscation was low, the obfuscation effect is consistent with research examining deception in other domains, such as annual financial reporting and deceptive conference calls from corporate officers (Burgoon et al., 2015). In line with the results from fraudulent financial reporting, our data suggest that fraudulent scientists obfuscate by making their writing less comprehensible with higher rates of technical terminology (e.g., jargon) and less readable text, compared with scientists not engaged in fraud.

It is also important to position this study within the broader deception research. Consistent with theoretical approaches to deception that emphasize the goal-oriented and strategic nature of a liar's communication patterns (e.g., Buller & Burgoon, 1996), we found that scientists wrote fraudulent papers with more obfuscation to make them more difficult to assess. This observation suggests that the obfuscated writing style was not simply a reflection of emotional distress from the deception or cognitive load from making up data, but rather a strategic and purposeful tactic consistent with their goals.

We can also consider how the present data fit with other research on linguistic cues associated with deception. One cue of interest across a variety of deception experiments is first-person singular pronouns (e.g., *I*, *me*, *my*), as these words can suggest a person's psychological attachment to his or her lie (Hauch et al., 2014; Newman, Pennebaker, Berry, & Richards, 2003; Pennebaker, 2011). The science genre, however, constrains the writing style to language dimensions that are normative for the

science community (Biber et al., 2007; Markowitz & Hancock, 2014), and these conventions do not allow for frequent use of first-person singular. Consistent with the approach of considering deceptive content in context (Blair, Levine, & Shaw, 2010), we examined linguistic dimensions that are appropriate for deception in the science genre. That is, lies about science data should produce linguistic patterns that are different from lies about one's online dating profile (Toma & Hancock, 2012) or fake reviews about a hotel (Ott, Choi, Cardie, & Hancock, 2011), given the different psychological aspects of the lie and the radical differences in genre conventions. This context-contingent approach to hypothesis testing in deception research should improve predictions about how deception affects language use in future studies.

Finally, the present study had several important limitations. Chief among them was that more than 80% of our corpora were limited to the biomedical sciences. Future studies should consider research fraud beyond the biomedical domain to evaluate linguistic patterns of deceptive science more broadly. Also, as noted, we used a relatively narrow set of theoretically derived linguistic cues related to obfuscation to detect differences between fraudulent and unretracted papers. A more bottom-up, natural language processing approach to linguistic patterns in fraudulent science papers is likely to generate not only a more robust classification of papers but is also likely to uncover unexpected linguistic differences in fraudulent science writing.

## Conclusion

The highly edited, constrained, and collaborative nature of science writing suggests that any effect of fraud on writing style may be difficult to uncover and that scientists should be able to conceal their deception linguistically. This is not the case. Participating in scientific fraud altered how researchers wrote their reports, highlighting the role of language as a marker of psychological change in deceptive communication. Our test of the obfuscation hypothesis contributes to a larger body of work supporting how language can reveal social and psychological dynamics, such as deception.

## Authors' Note

This work was presented at the 48th Annual Hawaii International Conference on System Sciences (January 2015). Data collection and analyses were performed at the authors' prior institution (Cornell University).

## Declaration of Conflicting Interests

The author(s) declared no potential conflicts of interest with respect to the research, authorship, and/or publication of this article.

## Funding

The author(s) disclosed receipt of the following financial support for the research, authorship, and/or publication of this article: This research was supported by the NSF SATC Grant TWC SBES- 1228857.

## Note

1. The list of papers (authors and journals) and data files from the fraudulent and unretracted corpora are also available from the authors on request.

## References

- Biber, D., Connor, U., & Upton, T. A. (2007). *Discourse on the move: Using corpus analysis to describe discourse structure*. Philadelphia, PA: John Benjamins.
- Blair, J. P., Levine, T. R., & Shaw, A. S. (2010). Content in context improves deception detection accuracy. *Human Communication Research, 36*, 423-442.
- Bloomfield, R. J. (2002). The "incomplete revelation hypothesis" and financial reporting. *Accounting Horizons, 16*, 233-243.
- Bond, C. F., & DePaulo, B. M. (2006). Accuracy of deception judgments. *Personality and Social Psychology Review, 10*, 214-234.
- Buller, D. B., & Burgoon, J. K. (1996). Interpersonal deception theory. *Communication Theory, 6*, 203-242.
- Burgoon, J. K., Mayew, W. J., Giboney, J. S., Elkins, A. C., Moffitt, K., Dorn, B., . . . Spitzley, L. (2015). Which spoken language markers identify deception in high-stakes settings? Evidence from earnings conference calls. *Journal of Language and Social Psychology*. Advance online publication. doi:10.1177/0261927X15586792
- Casadevall, A., & Fang, F. C. (2012). Reforming science: Methodological and cultural reforms. *Infection and Immunity, 80*, 891-896.
- Courtis, J. K. (1998). Annual report readability variability: Tests of the obfuscation hypothesis. *Accounting, Auditing and Accountability Journal, 11*, 459.
- Fang, F. C., Steen, R. G., & Casadevall, A. (2012). Misconduct accounts for the majority of retracted scientific publications. *Proceedings of the National Academy of Sciences of the United States of America, 109*, 17028-17033.
- Flesch, R. (1948). A new readability yardstick. *Journal of Applied Psychology, 32*, 221-233.
- Hancock, J. T., Woodworth, M. T., & Porter, S. (2013). Hungry like the wolf: A word-pattern analysis of the language of psychopaths. *Legal and Criminological Psychology, 18*, 102-114.
- Hauch, V., Blandón-Gitlin, I., Masip, J., & Sporer, S. L. (2014). Are computers effective lie detectors? A meta-analysis of linguistic cues to deception. *Personality and Social Psychology Review, 19*, 307-342.
- Humpherys, S. L., Moffitt, K. C., Burns, M. B., Burgoon, J. K., & Felix, W. F. (2011). Identification of fraudulent financial statements using linguistic credibility analysis. *Decision Support Systems, 50*, 585-594.
- Larcker, D. F., & Zakolyukina, A. A. (2012). Detecting deceptive discussions in conference calls. *Journal of Accounting Research, 50*, 495-540.
- Larrimore, L., Jiang, L., Larrimore, J., Markowitz, D. M., & Gorski, S. (2011). Peer to peer lending: The relationship between language features, trustworthiness, and persuasion success. *Journal of Applied Communication Research, 39*, 19-37.
- Li, F. (2008). Annual report readability, current earnings, and earnings persistence. *Journal of Accounting and Economics, 45*, 221-247.
- Lu, S. F., Jin, G. Z., Uzzi, B., & Jones, B. (2013). The retraction penalty: Evidence from the web of science. *Scientific Reports, 3*, 1-4.
- Markowitz, D. M., & Hancock, J. T. (2014). Linguistic traces of a scientific fraud: The case of Diederik Stapel. *PLoS ONE, 9*, e105937.

- McNamara, D. S., Graesser, A. C., McCarthy, P. M., & Cai, Z. (2014). *Automated evaluation of text and discourse with Coh-Metrix*. New York, NY: Cambridge University Press.
- McNamara, D. S., Louwerse, M. M., Cai, Z., & Graesser, A. C. (2013). Coh-Metrix version 3.0. Retrieved from <http://cohmetrix.com>.
- Newman, M. L., Pennebaker, J. W., Berry, D. S., & Richards, J. M. (2003). Lying words: Predicting deception from linguistic styles. *Personality and Social Psychology Bulletin, 29*, 665-675.
- Ott, M., Choi, Y., Cardie, C., & Hancock, J. T. (2011, June). *Finding deceptive online spam by any stretch of the imagination*. Paper presented at the 49th Annual Meeting of the Association for Computational Linguistics, Portland, OR.
- Pennebaker, J. W. (2011). *The secret life of pronouns: What our words say about us*. London, UK: Bloomsbury Press.
- Pennebaker, J. W., Booth, R. J., & Francis, M. E. (2007). *Linguistic Inquiry and Word Count*. Austin, TX: Pennebaker Conglomerates. Retrieved from [www.liwc.net](http://www.liwc.net).
- Rude, S., Gortner, E.-M., & Pennebaker, J. (2004). Language use of depressed and depression-vulnerable college students. *Cognition and Emotion, 18*, 1121-1133.
- Simonsohn, U. (2013). Just post it: The lesson from two cases of fabricated data detected by statistics alone. *Psychological Science, 24*, 1875-1888.
- Simonsohn, U., Nelson, L. D., & Simmons, J. P. (2014). p-Curve and effect size: Correcting for publication bias using only significant results. *Perspectives on Psychological Science, 9*, 666-681.
- Tausczik, Y. R., & Pennebaker, J. W. (2010). The psychological meaning of words: LIWC and computerized text analysis methods. *Journal of Language and Social Psychology, 29*, 24-54.
- Toma, C., & Hancock, J. T. (2012). What lies beneath: The linguistic traces of deception in online dating profiles. *Journal of Communication, 62*, 78-97.

## Author Biographies

**David M. Markowitz** (MSc, Cornell University) is a PhD candidate in the Department of Communication at Stanford University. His research uses computational methods to analyze how language is affected by social and psychological dynamics, including deception and persuasion.

**Jeffrey T. Hancock** (PhD, Dalhousie University, 2002) is a professor in the Department of Communication at Stanford University. He works on understanding psychological and interpersonal processes in social media by using computational linguistics and behavioral experiments to examine deception and trust, emotional dynamics, intimacy and relationships, and social support.