# Using Written Language as Indicator of Personality: a Meta-Analytic Study on Computational Models of Language

José David Moreno, José Ángel Martínez-Huertas,
Ricardo Olmos, Guillermo Jorge-Botana and Juan Botella

# Using language as indicator of personality: A meta-analytic study on computational models of language

José David Moreno[1], José Á. Martínez-Huertas[1], Ricardo Olmos[1], Guillermo Jorge-Botana[2], and Juan Botella[1]

[1] Universidad Autónoma de Madrid: Calle Iván Pavlov, 6, 28049 Madrid, Spain.

[2] UNED: Calle de Juan del Rosal, 10, 28040 Madrid, Spain.

**Author Note**

**Abstract**

Written language can be used to measure the Big Five personality traits using computational models of language. The aim of this study is to test the moderating role of different variables of computational models in a meta-analysis of 23 independent estimates. While the results showed significant combined estimates of the correlations for the five traits, these estimates were moderated by the type of information in the texts, the use of prediction mechanisms, and the source of publication of the primary studies. It is concluded that written language analyzed through computational methods could be used to extract relevant information of personality.

*Keywords:* language, computational models of language, meta-analysis, personality.

## Introduction

Over the last years, several authors have remarked on relevant relationships between personality and written language (e.g., Boyd & Pennebaker, 2015, 2017; Boyd & Schwartz, 2021; Chung & Pennebaker, 2018; Stachl et al., 2020). Keeping this in mind, we conducted a meta-analytic review of the automatic analysis methods of utterances and their correlation with the Big Five Personality questionnaire. The main goal of this meta-analysis was to obtain combined estimates of the relationships between personality traits and written language using computational models of language, to analyze the moderation role of different methodological characteristics of computational methods. This study has important theoretical and practical implications as the relevance of the different moderator variables will serve as a guide for researchers to design higher quality research. We expected higher relationships between personality traits and written language for the combination of semantic and syntactic information, and also when using prediction mechanisms. Also, we tested other relevant variables in the primary research as the source of text data from social networks, the language of texts, the publication source, the sex of the participants, or the text length of the materials of the studies. This conference paper aims to disseminate a summary of some of the main findings of the meta-analysis, but the full study can be found in Moreno et al. (2021).

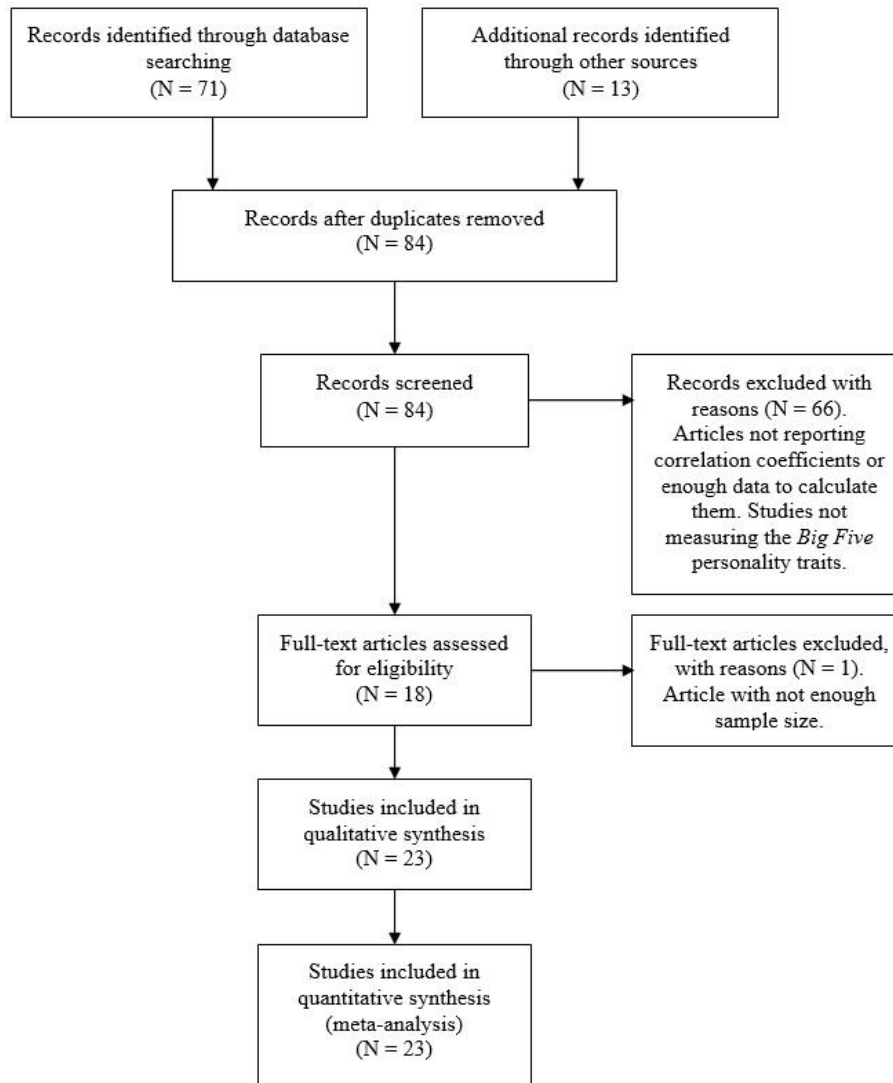## Method

### Selection of studies

The inclusion and exclusion criteria are shown in Figure 1. These criteria led to a final sample of 17 articles composed of 23 studies. The inter-rater reliability was high (93%; Cohen's Kappa = .86, $p<.001$).

**Figure 1**

*Flowchart (inclusion and exclusion criteria).*



**Personality measurement**

      We studied the personality models used in all the primary studies analyzed in the present meta-analysis. The vast majority of the studies measured the personality traits using the Big Five Personality model (*Big Five*; McCrae & Costa, 2008).

**Effect size calculation and statistical analysis**

Given the correlational design of the studies analyzed and the statistics provided, our choice for the effect size measure was the Pearson correlation coefficient *r*, as all primary studies reported it directly. For statistical analysis, the coefficients were previously transformed through Fisher's formula, to have a more symmetrical distribution (e.g., Borenstein et al., 2009). The final results were back-transformed from Fisher's value to its corresponding correlation coefficient to facilitate the interpretation of the results.

All the statistical analyses were performed with R Statistical Software in the 4.1.0 version, using the *metafor* package (Viechtbauer, 2010) for the combined estimates, the $Q$ statistic, and the $I^2$ statistic estimates. In addition, we used the SPSS macros of Lipsey and Wilson (2001) in order to analyze the categorical moderator variables. Separated meta-analysis were performed for each one of the five personality traits.

**Moderator variables**

In order to analyze the heterogeneity between the results of the studies, we conducted several moderator analyses for each of the five personality traits. Eight different moderator variables were selected based on their potential explanatory role in the results of the analyzed studies in the present meta-analysis and grouped on the basis of their methodological or theoretical nature.

On the one hand, six categorical and one numerical moderator variables were included in the analyses due to their methodological nature. Firstly, we studied the *type of information* of the input (semantic vs. syntactic vs. a combination of both). We consider "semantic" the kind of information that is tagged in a semantic category no matter if the category is produced by pattern detection, vector space models, or a predictive model layer. We consider

"syntactic" the kind of information that comes from non-semantic cues, as verbal persons, verbal tenses, discourse markers, n-grams and even sub-lexical cues as punctuation, spelling, capitalization, number of letters, syllables, etc. In addition, we used *prediction mechanisms* (no vs. yes) as a moderator variable. That is, whether a predictive model was trained with the semantic or syntactic information. Some studies use the raw semantic or syntactic information (i.e., no complex model was used) while other studies use complex predictive models with that information as input (e.g., artificial neural networks or support vector machine).

Taking into account the high amount of written materials analyzed from social networks in the primary research, we decided to include the use of *social networks* as the input in the primary studies (no vs. yes) as other moderator variable. In addition, to test the potential influence of the source of the results of the primary studies, we also analyzed the *language* of the texts analyzed in the primary studies (Chinese vs. English), and the *publication source* of the studies analyzed in the present meta-analysis (conference paper vs. journal paper). As a control measure, we also tested the potential influence of the instruments used to measure personality in the primary studies (*personality instrument*; BFI vs. others). The most of the primary studies analyzed in the present meta-analysis ($k = 18$) used a 44-items version of the BFI, but other studies ($k = 5$) used different instruments as reduced 10-items versions of the BFI, TIPI, Goldberg's 100-adjectives questionnaire, or BFI completed by indirect expert judges. Due to their low number, they were grouped on the same moderator category. Finally, the text length of the written materials analyzed in the primary studies (number of words) was included as a numerical moderator variable to test the influence of the quantity of linguistic information used as input. Additionally, sex (proportion of woman in the sample) was used as continuous moderator variable due to previous studies have already found sex differences that show different patterns in written language, which could also be reflected on personality.

<div align="center">**Results**</div>

**Combined effect size estimates**

All the combined effect size estimates were statistically significant (*r* ranged from .26 to .30) for each personality trait, although small to moderate combined effect size estimates according to Cohen's conventions. Significance tests were performed with the Hartung-Knapp-Sidik-Jonkman method.

**Moderator analyses**

Focusing on their methodological nature, the effects of six categorical variables and one numerical were analyzed for each personality trait. Regarding to the categorical variables, firstly, for the *type of information* (only syntactic information vs. only semantic information vs. a combination of both), a statistically significant effect was observed for Openness, Extraversion, Agreeableness and Neuroticism ($Q_B$ ranged from 6.51 to 18.12; $p \leq .05$) (the same tendency can be observed for Conscientiousness). This result can be explained by the higher performance of the combination of syntactic and semantic information versus using just one type of information. Secondly, for the use of *prediction mechanisms*, a statistically significant effect was observed favoring their use in all the personality traits ($Q_B$ ranged from 8.82 to 16.63; $p \leq .05$). Thirdly, for the use of *social networks'* text data, no statistically significant differences were observed in any personality trait. Fourthly, for the *language* of the text analyzed in the primary study, a significantly higher prediction of Conscientiousness was observed in favor of Chinese ($Q_B = 6.99$; $p \leq .05$) (a similar tendency can be observed in the rest of the personality traits). Fifthly, for the *publication source*, a statistically significant effect was observed in favor of conference publications as compared with journal publications in all the personality traits ($Q_B$ ranged from 4.09 to 7.11; $p \leq .05$). Finally, no statistically significant differences were found in any personality trait due to the different *personality*

*instruments* used in the primary studies. Regarding to the numerical variable, a meta-regression model was fit for text length (number of words), but it did not show significant effects on the estimated effect size.

Regarding to the variables selected due to their theoretical nature, finally only a meta-regression model was fit for sex (proportion of women in the sample), showing significant effects on the estimated effect size ($t = 2.58$; $p \leq .05$). Agreeableness was significantly more predictable when the proportion of women in the sample was higher. That is, when the proportion of women in the sample was higher, the estimated effect size was significantly higher. The same tendency can be observed for the five personality traits, but no statistically significant effects were obtained.

### Discussion

This conference paper aims to disseminate a summary of some of the main findings of the meta-analysis, but the full study can be found in Moreno et al. (2021). We conducted a meta-analysis from 23 primary studies, providing a synthesis of the combined effect size estimations of the predictive validity of the Big Five personality traits through computational models of language. We found that written language shows significant relationships with the basic five personality dimensions so that it can be used as a predictor of the personality profile of the individual. These results reinforce the relevance of personality and language relationships (Boyd & Pennebaker, 2015, 2017). Also, we found statistically significant moderating effects about the type of information used in the text materials, the use of prediction mechanisms, and the publication source of the primary studies, and some interesting differences when analyzing language and sex as moderator variables. These results raise the potential of written language as a consistent and reliable indirect personality predictor, and also that computational methods are equally reliable to measure written

language information to predict the Big Five personality traits (McCrae & Costa, 2008). These estimates are very informative to provide an effect size reference for the predictions of personality traits using computational models of language.

The present meta-analysis shows that written language can be considered a fruitful personality indirect indicator, and also that it is possible to accurately extract text information using computational models of language. But it is worthy to note, that the effect sizes of the present meta-analysis are small to moderate, which highlights that there is room for improvement in this research field. Thus, we found the information of this meta-analysis relevant to improve future research proposals. Building more connections between these methods and strong psychological theories could be a key issue for future developments in the field. Finally, we would like to encourage authors to conduct future research on this promising field, taking into account some of the main points highlighted in the present meta-analysis**,** as data about the moderator variables analyzed in this study can be useful to design higher quality research for future research.

## References

Borenstein, M., Hedges, L.V., Higgins, J.P.T., & Rothstein, H.R. (2009). *Introduction to Meta-Analysis.* John Wiley and Sons. https://doi.org/10.1002/9780470743386

Boyd, R.L., & Pennebaker, J.W. (2015). Did Shakespeare write Double Falsehood? Identifying individuals by creating psychological signatures with text analysis. *Psychological Science, 26*(5), 570-582. https://doi.org/10.1177/0956797614566658

Boyd, R.L., & Pennebaker, J.W. (2017). Language-based personality: a new approach to personality in a digital world. *Current Opinion in Behavioral Sciences, 18*, 63-68. https://doi.org/10.1016/j.cobeha.2017.07.017

Boyd, R.L., & Schwartz, H.A. (2021). Natural Language Analysis and the Psychology of Verbal Behavior: The Past, Present, and Future States of the Field. *Journal of*

*Language and Social Psychology, 40*(1), 21-41.
https://doi.org/10.1177/0261927X20967028

Chung, C.K., & Pennebaker, J.W. (2018). What do we know when we LIWC a person? Text analysis as an assessment tool for traits, personal concerns and life stories. In V. Zeigler-Hill & T. K. Shackelford (Eds), *The SAGE Handbook of Personality and Individual Differences: The Science of Personality and Individual Differences* (pp. 341–360). SAGE. https://doi.org/10.4135/9781526451163.n16

Lipsey, M.W., & Wilson, D.B. (2001). *Practical Meta-Analysis.* SAGE.

McCrae, R.R., & Costa, P.T., Jr. (2008). The Five-Factor Theory of personality. In O.P. John, R.W. Robins, & L.A. Pervin (Eds.), *Handbook of personality: Theory and research* (3$^{rd}$ ed., pp. 159-181). Guilford.

Moreno, J.D., Martínez-Huertas, J.A., Olmos, R., Jorge-Botana, G., & Botella, J. (2021). Can personality traits be measured by analyzing written language? A meta-analytic study on computational methods. *Personality and Individual Differences*. https://doi.org/10.1016/j.paid.2021.110818

Stachl, C., Pargent, F., Hilbert, S., Harari, G.M., Schoedel, R., Vaid, S., Gosling, S.D., & Bühner, M. (2020). Personality research and assessment in the era of machine learning. *European Journal of Personality*. https://doi.org/10.1002/per.2257

Viechtbauer, W. (2010). Metafor: Meta-Analysis Package for R. R package version 1.4-0, URL http://CRAN.R-project.org/package=metafor