1    Preregistration: Replication Studies in Linguistic Journals

2    Kristina Kobrock[1] & Timo B. Roettger[2]

3    [1] University of Osnabrück

4    [2] University of Oslo

5    Author Note

<sub>10</sub> Preregistration: Replication Studies in Linguistic Journals

<sub>11</sub> **Introduction**

<sub>12</sub> Coordinated efforts to replicate published findings have uncovered surprisingly low

<sub>13</sub> rates of successful replications across the psychological sciences (Open Science

<sub>14</sub> Collaboration, 2015), economics (Camerer et al., 2016), and social sciences (Camerer et al.,

<sub>15</sub> 2018). Experimental linguistics shares research practices that have been shown to decrease

<sub>16</sub> the replicability of findings. Thus, there are raising concerns about a similarly low number

<sub>17</sub> of replication studies conducted and published in this field (e.g. Marsden, Morgan-Short,

<sub>18</sub> Thompson, & Abugaber, 2018; Roettger & Baer-Henney, 2019). A number of failed

<sub>19</sub> replication attempts in various subfields of linguistics indicate that these concerns warrant

<sub>20</sub> attention (e.g. in language comprehension: Papesh, 2015; predictive processing: Nieuwland

<sub>21</sub> et al., 2018; among others: e.g. Chen, 2007; Stack, James, & Watson, 2018; Westbury,

<sub>22</sub> 2018). One driving factor for this phenomenon is an asymmetric incentive system that

<sub>23</sub> rewards novel confirmatory findings more than direct replications and null results. This

<sub>24</sub> leads to an abundance of positive findings in the absence of possible conflicting negative

<sub>25</sub> evidence. In order to thoroughly understand and be able to address this problem, it is

<sub>26</sub> important to assess the number of replication attempts and their contributing factors.

<sub>27</sub> Other fields such as psychology (Makel, Plucker, & Hegarty, 2012), eudcation science

<sub>28</sub> (Makel & Plucker, 2014), and special education research (Makel et al., 2016) have assessed

<sub>29</sub> the amount of direct replications in their respective field and report alarmingly low

<sub>30</sub> replication rates (0.13% - 1.07%).

<sub>31</sub> In order to evaluate the replication rate in linguistics, the present study aims at

<sub>32</sub> assessing the frequency and typology of replication studies that are published in a

<sub>33</sub> representative sample of linguistic journals. The study consists of two parts: First, we will

<sub>34</sub> assess the frequency of self-reported replication attempts across 100 linguistic journals and

<sub>35</sub> relate the replication rate to factors related to journal policy, impact factor and publication

36  type. Second, we will assess the type of replication studies (direct, partial, conceptual)

37  published in a subset of 20 journals and relate their frequency to factors like the year of

38  publication, and the citation and publication year of the original study.

### Overview Analysis: Rate of Replication Mention

40  The key dependent variable of the first part of this study is the rate of replication

41  mention for journals relevant to the field of experimental linguistics. In order to determine

42  these rates for the individual journals, we will draw on a method introduced by Makel et

43  al. (2012). We will use the search engine "Web of Science" (https://webofknowledge.com)

44  for journal articles that contain the search term "replicat*" in title, abstract or keywords

45  and compute the rate of replication mention.

46  **Research Questions.**  We intend to answer the following research questions: How

47  many replication studies have been published in journals representative for experimental

48  linguistic research? How did the rate change over time and how does it relate to journal

49  policy, impact factor, and publication type?

50  **Sample.**  To obtain a representative sample of journals relevant for the field of

51  experimental linguistics, we follow the procedure presented here: First, using the Web of

52  Science advanced search on the "Web of Science Core Collection" database, we filter for the

53  category "Linguistics" (WC=(Linguistics)) that lists the articles of every journal covered

54  by Web of Science that was assigned to the subject category of Linguistics (see a list of

55  categories here). All English language articles from the full available range of complete

56  years (1945-2020) are taken into account. From the resulting set (159002) only those

57  articles are selected which contain the search term "experiment*" in their title, abstract or

58  keywords using TS="experiment*" in order to filter for experimental linguistic studies.

59  This search results in 11093 articles. The relevant journals are selected based on the

60  obtained article counts. From all journals that include at least one experimental linguistic

61  study according to our criterium (418), journals with less than 100 published articles are

62 excluded, yielding 259 remaining journals. Because we are interested in journals with a

63 high proportion of experimental studies, we calculate the ratio of studies that contained the

64 search term "experiment*" by the total amount of articles per journal and sort the results

65 in descending order. Our sample constitutes the first 100 journals of that list. Counts have

66 been obtained on the 21st February 2021. See here for more details: https://osf.io/q2e9k/

67 **Procedure.** The total number of articles containing the search term "replicat*" in

68 title, abstract or keywords is obtained via Web of Science search for the 100 sampled

69 journals. This number and and the total number of experimental studies described above

70 serve as a baseline for calculating the rates of replication mention, following the method

71 used by Makel et al. (2012). The rates of replication mention are calculated by dividing the

72 number of articles containing the term "replicat*" by the number of articles that contain

73 the term "experiment*" for each journal, respectively.

74 In order to relate the rate of replication mention to journal policies, we further

75 examine the journals' submission guidelines adopting the procedure used by Martin and

76 Clarke (2017). They grouped psychology journals into four classes determined by what was

77 stated in the "instructions to authors" and "aims and scope" sections on the websites of

78 the respective journals: (1) Journals which stated that they accepted replications; (2)

79 Journals which did not state they accepted replications but did not discourage replications

80 either; (3) Journals which implicitly discouraged replications through the use of emphasis

81 on the scientific originality of submissions, (4) Journals which actively discouraged

82 replications by stating explicitly that they did not accept replications for publication

83 (Martin & Clarke, 2017, p. 3).

84 Journal impact factors are extracted via Journal Citation Reports

85 (https://jcr.clarivate.com). The 2019 journal impact factors are calculated by dividing the

86 citations in 2019 to items published in 2017 and 2018 by the total number of citable items

87 in 2017 and 2018. The open access category of journals is assessed via Web of Science. We

88 distinguish between three categories: journals which are listed on the Directory of Open

Access Journals (DOAJ) ("DOAJ gold"), journals with some articles being published as

open access articles ("partial") and journals with no openly accessible articles ("no").

**Data Analysis.** We will use Bayesian parameter estimation based on generalized

linear regression models with a binomial link function in order to estimate the rate of

replication mention relative to the following predictors: journal impact factors (continuous),

open access (binary: open access journal or not), and replication policies (binary: either

explicitly encourage or not). The model will be fitted to the proportion of replication

mentions per journal using the R package brms (Bürkner, 2016). We will use (weakly)

informative normal priors centered on -2.1973 (corresponding to a 10% base rate, sd = 2.5)

for the intercept since we expect very low base rates (e.g. Makel et al., 2012). We will use

weakly informative Cauchy priors centered on zero (scale = 2.5) for all population-level

regression coefficients. These priors are what is referred to as regularizing (Gelman,

Jakulin, Pittau, Su, & others, 2008), i.e. our prior assumption is agnostic as to whether the

predictors affect the dependent variable, thus making our model conservative with regards

to the predictors under investigation. Four sampling chains with 2000 iterations each will

be run for each model, with a warm-up period of 1000 iterations. For relevant predictor

levels and contrasts between predictor levels, we will report the posterior probability for the

rate of replication mention. We summarize these distributions by reporting the posterior

mean and the 95% credible intervals (calculated as the highest posterior density interval).

## Detailed Analysis: Types and Contributing Factors

**Methods**

The second part of the analysis aims at obtaining a better understanding of the

underlying mechanisms of replication attempts published in the field of experimental

linguistics. Because the term "replication" is commonly used in ambiguous ways, the

articles that contain the search term "replicat*" require further analysis to determine

whether the articles in question indeed report a replication study or use the term in a different way.

**Research Questions.** We are interested in which kinds of replication studies are published and which factors contribute to their publication. We aim at investigating what types of replication studies are prevalent in the field. We are further interested in the relationship of direct replications and whether the paper was published as open access or not, the number of citations of the initial study and the years between publication of the initial study and the replication attempt.

**Sample.** From the superset of 100 journals obtained above, the first 20 journals (i.e. those journals with the highest proportion of experimental studies) are selected for a more detailed analysis. We exclude those journals for which less than 2 hits (TS=(replicat*)) can be obtained. This method yields a total number of 274 articles (see here for a list of article counts per journal: https://osf.io/f3yp8/). Because of the skewed distribution of our sample (114 hits for Journal of Memory and Language, and less than 40 for all other journals), we randomly select 50 out of the 114 articles for the Journal of Memory and Language to achieve a more balanced distribution of papers across journals by drawing from a uniform distribution in R without replacement (see here for details).

**Procedure.** The sampling procedure above results in 210 possible self-labeled replication studies. In a first step, we will identify whether the article indeed presents a replication study or not. By reading title and abstract of the paper a first intuition of what the article is about can be obtained. The main task is to assess whether the authors claim that their underlying aim was to replicate or reproduce findings or methods of another study (henceforth initial study). A search for occurrences of the search term "replicat" in the text and an assessment of the paragraph before the Methods section as well as the first paragraph of the Discussion section (following the procedure specified by Makel et al. (2016)) helps to further identify the intention communicated by the authors regarding their use of the term replication. If the authors communicate that (one of) the underlying

141 aim(s) was to replicate an original study, this article can be treated as a replication. It

142 then qualifies for further analysis after the coding scheme that can be viewed here:

143 https://osf.io/ct2xj/.

144      Assuming that the authors did not make any drastic changes to the initial study

145 *without* reporting them, number and type of changes made by the replication study are

146 extracted. The replication studies are classified according to three types: direct replication

147 (0 changes), partial replication (1 change) and conceptual replication (2 or more changes),

148 following Marsden et al. (2018). We note the nature of the change as one of the following

149 categories (yes/no): experimental paradigm, sample, materials/experimental set-up,

150 dependent variable, independent variable, and control. Coding the articles also involves

151 examining the factors open access of that article (yes/no), years between initial study and

152 replication attempt, author overlap of initial study and replication attempt (yes/no),

153 citation counts of both studies and the language under investigation. The information on

154 whether the article is open access as well as citation counts and years of publication for

155 both studies can be obtained from Web of Science. An author overlap is attested when one

156 of the authors is a (co-)author on both articles.

157      **Data Analysis.**   We will use Bayesian parameter estimation based on generalized

158 linear regression models with a logit link function in order to estimate the rate of direct

159 replications relative to the following predictors: year of publication (continuous), open

160 access (binary: open access article or not), time lag between publication of initial study

161 and replication attempt (continuous) and number of citations of initial study (continuous).

162 The model will be fitted to whether the replication mention was a direct replication or not

163 using the R package brms (Bürkner, 2016). The model further includes random intercepts

164 for individual journals to account for varying rates of direct replications across journals.

165 We will use (weakly) informative normal priors centered on -2.1973 (corresponding to 10%

166 base rate, sd = 2.5) for the intercept (since we expect very low base rates) and weakly

167 informative Cauchy priors centered on zero (scale = 2.5) for all population-level regression

168 coefficients. Four sampling chains with 2000 iterations each will be run for each model,

169 with a warm-up period of 1000 iterations. For relevant predictor levels and contrasts

170 between predictor levels, we will report the posterior probability for the rate of direct

171 replication. We summarize these distributions by reporting the posterior mean and the

172 95% credible intervals (calculated as the highest posterior density interval).

**References**

Bürkner, P.-C. (2016). Brms: An R package for Bayesian multilevel models using Stan. *Journal of Statistical Software*, *80*(1), 1–28.

Camerer, C. F., Dreber, A., Forsell, E., Ho, T.-H., Huber, J., Johannesson, M., . . . Wu, H. (2016). Evaluating replicability of laboratory experiments in economics. *Science*, *351*(6280), 1433–1436. https://doi.org/10.1126/science.aaf0918

Camerer, C. F., Dreber, A., Holzmeister, F., Ho, T.-H., Huber, J., Johanesson, M., . . . Wu, H. (2018). Evaluating the replicability of social science experiments in nature and science between 2010 and 2015. *Nature*, *2*, 637–644. https://doi.org/10.1038/s41562-018-0399-z

Chen, J.-Y. (2007). Do Chinese and English speakers think about time differently? Failure of replicating Boroditsky (2001). *Cognition*, *104*(2), 427–436.

Gelman, A., Jakulin, A., Pittau, M. G., Su, Y.-S., & others. (2008). A weakly informative default prior distribution for logistic and other regression models. *The Annals of Applied Statistics*, *2*(4), 1360–1383.

Makel, M. C., & Plucker, J. A. (2014). Facts are more important than novelty: Replication in the education sciences. *Educational Researcher*, *43*(6), 304–316.

Makel, M. C., Plucker, J. A., Freeman, J., Lombardi, A., Simonsen, B., & Coyne, M. (2016). Replication of special education research: Necessary but far too rare. *Remedial and Special Education*, *37*(4), 205–212.

Makel, M. C., Plucker, J. A., & Hegarty, B. (2012). Replications in psychology research: How often do they really occur? *Perspectives on Psychological Science*, *7*(6), 537–542.

Marsden, E., Morgan-Short, K., Thompson, S., & Abugaber, D. (2018). Replication in Second Language Research: Narrative and Systematic Reviews and

198  Recommendations for the Field. *Language Learning*, *68*(2), 321–391.

199  https://doi.org/10/gc3h3b

200  Martin, G. N., & Clarke, R. M. (2017). Are psychology journals anti-replication? A

201  snapshot of editorial practices. *Frontiers in Psychology*, *8*.

202  https://doi.org/10.3389/fpsyg.2017.00523

203  Nieuwland, M. S., Politzer-Ahles, S., Heyselaar, E., Segaert, K., Darley, E., Kazanina, N.,

204  . . . Huettig, F. (2018). Large-scale replication study reveals a limit on probabilistic

205  prediction in language comprehension. *eLife*, *7*, e33468.

206  https://doi.org/10.7554/eLife.33468.001

207  Open Science Collaboration. (2015). Estimating the reproducibility of psychological

208  science. *Science*, *349*(6251). https://doi.org/10.1126/science.aac4716

209  Papesh, M. H. (2015). Just out of reach: On the reliability of the action-sentence

210  compatibility effect. *Journal of Experimental Psychology: General*, *144*(6),

211  e116–e141. https://doi.org/10.1037/xge0000125

212  Roettger, T. B., & Baer-Henney, D. (2019). Toward a replication culture: Speech

213  production research in the classroom. *Phonological Data and Analysis*, *1*(4), 1–23.

214  Stack, C. M. H., James, A. N., & Watson, D. G. (2018). A failure to replicate rapid

215  syntactic adaptation in comprehension. *Memory & Cognition*, *46*(6), 864–877.

216  https://doi.org/10.3758/s13421-018-0808-6

217  Westbury, C. (2018). Implicit sound symbolism effect in lexical access, revisited: A

218  requiem for the interference task paradigm. *Journal of Articles in Support of the*

219  *Null Hypothesis*, *15*(1), 1–12.