

1 Examining the replicability of online experiments selected by 2 a decision market

3 Felix Holzmeister¹, Magnus Johannesson², Colin F. Camerer³, Yiling Chen⁴, Teck-Hua Ho⁵,
4 Suzanne Hoogeveen⁶, Juergen Huber⁷, Noriko Imai⁸, Taisuke Imai⁸, Lawrence Jin⁹, Michael
5 Kirchler⁷, Alexander Ly^{10,11}, Benjamin Mandl¹², Dylan Manfredi¹³, Gideon Nave¹³, Brian A.
6 Nosek^{14,15}, Thomas Pfeiffer¹⁶, Alexandra Sarafoglou¹⁰, Rene Schwaiger⁷, Eric-Jan
7 Wagenmakers¹⁰, Viking Waldén¹⁷, Anna Dreber^{1,2,*}

8 ¹Department of Economics, University of Innsbruck, Innsbruck, Austria. ²Department of Economics, Stockholm School
9 of Economics, Stockholm, Sweden. ³Division of the Humanities and Social Sciences, California Institute of
10 Technology, Pasadena, USA. ⁴John A. Paulson School of Engineering and Applied Sciences, Harvard University,
11 Boston, USA. ⁵Nanyang Technological University, Singapore. ⁶Faculty of Social and Behavioural Sciences, Utrecht
12 University, Utrecht, The Netherlands. ⁷Department of Banking and Finance, University of Innsbruck, Innsbruck,
13 Austria. ⁸Institute of Social and Economic Research, Osaka University, Osaka, Japan. ⁹Lee Kuan Yew School of
14 Public Policy, National University of Singapore, Singapore. ¹⁰Faculty of Social and Behavioural Sciences, University
15 of Amsterdam, Amsterdam, ¹¹Machine Learning, Centrum Wiskunde and Informatica, Amsterdam, The Netherlands.
16 ¹²Independent researcher, Vienna, Austria. ¹³Marketing Department, Wharton School, University of Pennsylvania,
17 Philadelphia, USA. ¹⁴Department of Psychology, University of Virginia, Charlottesville, USA. ¹⁵Center for Open
18 Science, Charlottesville, Virginia, USA. ¹⁶Institute for Advanced Study, Massey University, Auckland, New Zealand.
19 ¹⁷Sveriges Riksbank, Stockholm, Sweden.

20 * To whom correspondence should be addressed:

21 **Anna Dreber**
22 Department of Economics, Stockholm School of Economics
23 Box 6501, SE-113 83 Stockholm, Sweden.
24 Email: anna.dreber@hhs.se

25

26 **Abstract**

27 In this study, we test the feasibility of using decision markets to select studies for replication and
28 provide evidence about the replicability of online experiments. Social scientists ($n = 162$) traded
29 on the outcome of close replications of 41 systematically selected MTurk social science
30 experiments published in PNAS 2015–2018, knowing that the 12 studies with the lowest and the
31 12 with the highest final market prices would be selected for replication, along with two randomly
32 selected studies. The replication rate, based on the statistical significance indicator, was 83% for
33 the top-12 and 33% for the bottom-12 group. Overall, 54% of the studies were successfully
34 replicated, with replication effect size estimates averaging 45% of the original effect size
35 estimates. The replication rate varied between 54% and 62% for alternative replication indicators.
36 The observed replicability of MTurk experiments is comparable to that of previous systematic
37 replication projects involving laboratory experiments.

38 **Main**

39 Can published research findings be trusted? Unfortunately, the answer to this question is not
40 straightforward, and the credibility of scientific findings and methods has been questioned
41 repeatedly¹⁻⁹. A vital tool for evaluating and enhancing the reliability of published findings is to
42 carry out replications, which can be used to sort out likely true positive findings from likely false
43 positives. A replication essentially updates the probability of the hypothesis being true after
44 observing the replication outcome. A successful replication will move this probability towards
45 100%, while a failed replication will move it towards 0%^{10,11}. In recent years, several systematic
46 large-scale replication projects in the social sciences have been published¹²⁻¹⁷, reporting
47 replication rates of around 50% in terms of both the fraction of statistically significant replications
48 and the relative effect sizes of replications. Potential factors to explain these replication rates may
49 be low statistical power^{1,18,19} in the original studies, testing original hypotheses with low
50 priors^{1,10,20}, and questionable research practices^{1,21,22}. Systematic replication studies led to
51 discussions about improving research practices^{23,24} and have substantially increased the interest
52 in independent replications²⁵. However, as it is time-consuming and costly to conduct replications,
53 it has been argued that it is useful to have a principled mechanism to decide which replications
54 to prioritize to facilitate efficient and effective usage of resources²⁵⁻³⁷. In this study, we test the
55 feasibility of one potential method to select which studies to replicate. Building on previous work
56 using prediction markets³⁸⁻⁴⁰ to forecast replicability, we adapt the forecasting methodology to
57 what is referred to as decision markets⁴¹⁻⁴⁴.

58 The decisive distinction between prediction markets and decision markets is that prediction
59 markets elicit aggregate-level replicability forecasts on a predetermined set of studies, whereas
60 decision market forecasts determine which studies are going to be put to a replication test. While
61 previous studies provide evidence that prediction market forecasts are predictive of replication
62 outcomes^{10,16,17,45}, prediction efficiency might not generalize to decision markets, which involve
63 more complex procedures and incentives. The performance of decision markets as a tool for
64 selecting which empirical claims to replicate has not been systematically examined. Note that a
65 decision market in itself is not sufficient to provide a mechanism to select studies for replication,
66 but it has to be combined with an objective function of which studies to replicate (an example of
67 an objective function would be to replicate the studies with the lowest probability of replication).
68 For decision markets to be potentially useful for selecting studies for replication, it first has to be
69 established that the predictions of the decision markets are associated with the replication
70 outcomes. To provide such a “proof of concept” of using a decision market as a mechanism to
71 determine which studies to replicate, we first identified all social science experiments published
72 in the Proceedings of the National Academy of Sciences (PNAS) between 2015 and 2018 that
73 fulfilled our inclusion criteria for (i) the journal and period; (ii) the platform on which the experiment
74 was performed (Amazon Mechanical Turk; MTurk); (iii) the type of design (between-subjects or
75 within-subject treatment design); (iv) the equipment and materials needed to implement the
76 experiment (the experiment had to be logistically feasible for us to implement); and (v) the results

77 reported in the experiment (at least one main or interaction effect with $p < 0.05$ reported in the
78 main text). Based on our inclusion criteria, we identified 44 articles, three of which have been
79 excluded due to a lack of feasibility, leaving us with a final sample of 41 articles^{46–86} (see Methods
80 for details on the inclusion criteria). For each of these articles, we identified one critical finding
81 with $p < 0.05$ that we could potentially replicate (see Methods for details and Supplementary
82 Table 1 for the hypotheses selected for each of the 41 studies).

83 We then invited social science researchers to participate as forecasters in both a prediction survey
84 and an incentivized decision market on the 41 studies. In the survey, the forecasters
85 independently estimated the probability of replication for the 41 studies. In the decision market,
86 they could trade on whether the result of each of the 41 studies would replicate. Participants in
87 the decision market received an endowment of 100 tokens corresponding to USD 50, and 162
88 participants made a total of 4,412 trades. Traders in the market were informed about the
89 preregistered decision mechanism: The 12 studies with the highest and the 12 studies with the
90 lowest market prices were to be selected for close replication; additionally, two randomly chosen
91 studies (out of the remaining 17 studies) are replicated to ensure incentive compatibility, with
92 participant payoffs scaled up by the inverse of their probability in the decision rule (see Methods
93 for details). For incentive compatibility, all the 41 studies included need to have a strictly positive
94 probability of being selected for replication, which is ensured by having at least one randomly
95 selected study. Otherwise, traders would be incentivized to only trade on those studies that will
96 most likely be chosen according to the decision rule.

97 All replication experiments, just like all original studies, were conducted on Amazon Mechanical
98 Turk (MTurk), and the same sample restrictions and exclusion criteria as the original studies were
99 applied, which guards against concerns about the potential moderating effects of culture
100 differences in replications^{14,87}. Replication sample sizes were determined to have 90% power to
101 detect $\frac{2}{3}$ of the effect size reported in the original study at the 5% significance level in a two-sided
102 test (with the effect size estimates having been converted to Cohen's d to have a common
103 standardized effect size measure across the original studies and the replication studies; see
104 Methods for details). If sample size calculations led to replication sample sizes smaller than in the
105 original study, we targeted the same sample size as in the original study. The average sample
106 size in the replications ($n = 1,018$) was 3.5 times as large as the average sample size in the
107 original studies ($n = 292$).

108 The replication results for the 26 MTurk experiments selected by the decision market constitute
109 the second contribution of this project. Systematic evidence on the replicability of online
110 experiments in the social sciences is lacking, and concerns about the quality of online
111 experiments in general—and MTurk studies in particular—have been raised^{88–94}. Needless to say,
112 the replication results only pertain to the single focal result selected per paper, and the replication
113 outcome does not necessarily generalize to other results reported in the original articles^{95,96}. For
114 convenience, we refer to the replications as “replication of [study reference]” though. Also, our
115 assessment of the most central result may differ from that of the original authors.

116 Preregistering study protocols and analysis plans have been proposed as a means to reduce
117 questionable research practices. While empirical evidence is still limited, some recent studies
118 suggest that these practices enhance the credibility of published findings^{97–99}, although potential
119 issues with preregistration have also been raised^{100–102}. Prior to starting the survey data collection
120 (that preceded the decision market and replications), we preregistered^{103,104} an analysis plan
121 (“replication report”) for each of the 41 potential replications at OSF after obtaining feedback from
122 the original authors (<https://osf.io/sejyp>). After the replications had been conducted, the replication
123 reports of the 26 studies selected for replication were updated with the results of the replications
124 (and potential deviations from the protocol) and were posted to the same OSF repository. We
125 also preregistered an overall analysis plan at OSF before starting the data collection, detailing the
126 study's design and all planned analyses and tests (<https://osf.io/xsp6g>). Unless explicitly stated,
127 all analyses and tests reported in the manuscript have been preregistered and adhere exactly to
128 our pre-registered analysis plan. The Supplementary Notes details any deviations from the
129 planned design and analyses for the 26 replications.

130 We preregistered two primary replication indicators and two primary hypotheses. The two primary
131 replication indicators are the relative effect size of the replications and the statistical significance
132 indicator for replication (i.e., whether or not the replication results in a statistically significant effect
133 with $p < 0.05$ in the same direction as the original effect), which was the replication outcome
134 predicted by forecasters in the survey and the decision market.

135 The statistical significance indicator is a binary criterion of replication and is based on testing the
136 hypothesis for which the original study found support using standard null hypothesis significance
137 testing. The indicator crudely classifies replications as failed or successful depending on whether
138 the replication study yields evidence in support of the original hypothesis at a particular
139 significance threshold. (Critics of null hypothesis significance testing or privileging a p-value of
140 .05 will, justifiably so, object to this crude classification; that is why it is only one of several
141 indicators that we report.) A replication classified as failed based on this indicator, however, does
142 not imply that the estimated replication effect size is significantly different from the original
143 estimate (see more on this below). To keep the false negative risk at bay and to be informative,
144 the statistical significance indicator calls for well-powered replications (as in this study).^{105,106}
145 However, a limitation of this indicator for well-powered replication studies is that it may classify a
146 replication as successful even if the observed effect size is substantially smaller (or larger) than
147 in the original study. While the statistical significance indicator dichotomizes replication outcomes
148 into successful and failed, replicability may be perceived as a continuous matter of degree. This
149 is why we also consider the relative effect size—a continuous measure of replicability—as a
150 primary replication indicator. While the relative effect size constitutes an imprecise indicator for
151 an individual replication study, it arguably provides an informative measurement of the extent of
152 replicability for a group of studies as it quantifies the average degree of apparent inflation in the
153 original effect sizes.¹⁰⁷ As all replication indicators have limitations, we preregistered four
154 additional secondary replication indicators. In addition, we report the results for two non-

155 preregistered replication indicators, which were helpfully suggested during the review process of
156 the manuscript.

157 In our two primary hypotheses, we conjecture that (i) the decision market prices positively
158 correlate with the replication outcomes and (ii) the standardized effect sizes in the replications
159 are lower than in the original studies. All hypotheses are evaluated using two-tailed tests, and—
160 following Benjamin et al.¹⁰⁸—we interpret results with $p < 0.005$ as “statistically significant
161 evidence,” whereas results with $0.005 \leq p < 0.05$ are considered “suggestive evidence.” No
162 adjustments were made for multiple comparisons.

163

164 **Results**

165 **Replication outcomes and decision market performance.** Fig. 1 and Supplementary Table 2
166 show the results for the decision markets where the final market price can be interpreted as the
167 predicted replication probability. The predicted probabilities of replication range from 20.9% to
168 92.9% for the 41 studies, with a mean of 57.6% ($sd = 23.6\%$). The average predicted probability
169 for the 26 studies eventually selected for replication is 58.5%. Fig. 1 also delineates the replication
170 outcomes based on the statistical significance indicator, which allows for gauging the relationship
171 between the decision market prices and the replication outcomes. In Fig. 2 and Supplementary
172 Table 3, we show the replication results for the 26 studies selected for replication. Of the 26
173 claims, 14 (53.8%; 95% CI [33.4%, 73.4%]) replicated successfully according to the statistical
174 significance indicator. The point-biserial correlation between decision market prices and the
175 binary replication outcome, testing our first primary hypothesis, is $r = 0.505$ (95% CI [0.146,
176 0.712]; $t(24) = 2.867$, $p = 0.008$; $n = 26$). Thus, in support of our first primary hypothesis, we find
177 suggestive evidence of a positive association between decision market prices and replication
178 outcomes. As a related secondary hypothesis, we test if the replication rate is lower among the
179 12 studies with the lowest decision market prices than for the 12 studies with the highest decision
180 market prices. The replication rate is 33.3% (95% CI [9.9%, 65.1%]) for the studies in the “bottom-
181 12” group and 83.3% (95% CI [51.6%, 97.9%]) for the studies in the “top-12” group, yielding
182 suggestive evidence in support of our secondary hypothesis (Fisher’s exact test; $\chi^2(1) = 6.171$,
183 $p = 0.036$; $n = 24$). Note that Fisher’s test conditions on the margin totals; hence, it is only exact
184 for the conditional distribution and can be overly conservative if the margin totals are
185 unknown^{109,110}, as is the case in our analysis. Boschloo’s test¹¹¹, an exact unconditional procedure
186 uniformly more powerful than Fisher’s test, also yields suggestive evidence for the difference in
187 proportions between the “top-12” and the “bottom-12” group (not preregistered; 95% CI [0.089,
188 0.799], $p = 0.017$; $n = 24$).

189 **Relative effect sizes.** The mean estimated effect size of the 26 replication studies (in Cohen’s d
190 units) is 0.253 ($sd = 0.357$) compared to 0.563 ($sd = 0.426$) for the original studies, implying a
191 relative estimated average effect size, just dividing the two numbers, of 45.0%; the difference in
192 estimated effect sizes is statistically significant, supporting our second primary hypothesis of

193 systematically smaller estimated effect sizes in the replications (Wilcoxon signed-rank test;
194 $z = 4.203$, $p < 0.001$; $n = 26$). The relative effect size can also be estimated for each study
195 separately (reported in Supplementary Table 3) and varies between -17.0% and 136.2% , with a
196 mean estimate across studies of 41.1% (95% CI [24.5%, 57.7%]). For the 14 studies that
197 replicated according to the statistical significance indicator, the first and the second relative effect
198 size measures as defined above are 69.5% and 72.0% (95% CI [54.8%, 89.3%]), indicative of
199 some inflation in original effect sizes even for apparent true positives. The two estimated relative
200 effect size measures for the 12 studies that failed to replicate according to the statistical
201 significance indicator are 3.2% and 5.0% (95% CI [-2.6% , 12.5%]), respectively. Fig. 3 illustrates
202 the relationship between the estimated original and replication effect sizes.

203 **Secondary replication indicators.** We also preregistered four secondary replication indicators:
204 the small-telescopes approach¹¹², the one-sided default Bayes factor¹¹³, the replication Bayes
205 factor¹¹⁴, and the fixed-effects weighted meta-analytic effect size (see Methods for details). When
206 relying on the small-telescopes approach, testing if the replication effect size is smaller than a
207 “small effect,”¹¹² 15 studies (57.7%; 95% CI [36.9%, 76.6%]) are considered successful
208 replications (Fig. 4 and Supplementary Table 4). The one-sided default Bayes factor (BF_{+0})
209 indicates the strength of evidence in favor of the alternative hypothesis as opposed to the null
210 hypothesis. BF_{+0} exceeds one for the 14 studies (53.8%; 95% CI [33.4%, 73.4%]) that replicated
211 according to the statistical significance indicator, with strong evidence ($BF_{+0} > 10$) for the tested
212 hypothesis for nine studies (34.6%; 95% CI [17.2%, 55.7%]); BF_{+0} is below 1 for the 12
213 replications (46.2%; 95% CI [26.6%, 66.6%]) that failed to replicate according to the statistical
214 significance indicator, with strong evidence ($BF_{+0} < 0.1$) for the null hypothesis for seven studies
215 (26.9%; 95% CI [11.6%, 47.8%]) based on the evidence categories proposed by Jeffreys¹¹⁵
216 (Fig. 5 and Supplementary Table 4). The one-sided replication Bayes factor (BF_{R0}) indicates the
217 strength of additional evidence in favor of the alternative hypothesis as opposed to the null
218 hypothesis, given the already acquired evidence based on the original data¹¹⁴. Replication Bayes
219 factors lead to similar conclusions as the one-sided default Bayes factors, with $BF_{R0} > 10$ for ten
220 studies (38.5%; 95% CI [20.2%, 59.4%]) and $BF_{R0} < 0.1$ for seven studies (26.9%; 95% CI
221 [11.6%, 47.8%]). One exception to this is the study by Cooney et al.⁵⁶, for which the default Bayes
222 factor exceeds one ($BF_{+0} = 8.01$) but the replication Bayes factor is below one ($BF_{R0} = 0.23$) due
223 to the replication effect size being only about a third of the original effect size and a larger sample
224 size in the replication compared to the original study (Fig. 5 and Supplementary Table 4). The
225 meta-analytic effect size is statistically significant at the 5% level for 16 studies (61.5%; 95% CI
226 [40.6%, 79.8%]) and significant at the 0.5% level for 14 studies [53.8%; 95% CI [33.4%, 73.4%]);
227 see Fig. 6 and Supplementary Table 4. The meta-analytic effect sizes should be interpreted
228 cautiously as original effect sizes reported as statistically significant are likely to be overestimated
229 on average due to insufficient sample sizes and, thereby, statistical power (and potentially due to
230 questionable research practices)^{18,19}. Overall, the primary and secondary replication indicators
231 yield the same binary conclusions for 23 of the 26 replications.

232 **Non-preregistered replication indicators.** Following the suggestion of a reviewer, we report the
233 results for two additional replication indicators. The first alternative replication indicator is a test
234 of whether the replication effect size is statistically significantly different from the original effect
235 size. This indicator is closely related to the prediction interval approach¹¹⁶ as the results can be
236 illustrated as prediction intervals that the replication effect sizes are evaluated against: if the
237 replication effect size falls outside the 95% prediction interval, the replication and original effect
238 sizes differ at the 5% significance level, and the replication is considered a failure. The prediction
239 interval approach, thus, yields a binary replication indicator which is complemented by a
240 continuous replicability measure defined as the p-value of the test of a significant difference
241 between the replication and original effect sizes. We illustrate the prediction interval results in Fig.
242 7 and report the z-statistics and p-values in Supplementary Table 5 (the p-values are also shown
243 in Fig. 7). According to the prediction interval indicator, 15 studies (57.7%, 95% CI [36.9%,
244 76.6%]) replicate. This replication rate is close to the result for the statistical significance indicator.
245 However, the classification of nine replication outcomes shifts: For four studies, the classification
246 changes from successful to failed, and for five studies, the classification changes from failed to
247 successful. These changes are due to the fact that low-powered original studies are more likely
248 to replicate, whereas high-powered original studies are less likely to replicate based on the
249 prediction interval indicator (as compared to evaluating replicability based on the statistical
250 significance indicator). According to the prediction interval indicator, six studies failed to replicate
251 among the top-12 studies in terms of decision market prices, whereas three studies failed to
252 replicate among the bottom-12 studies.

253 **Associations between indicators (not preregistered).** To examine the relationship between
254 the replication indicators, we estimated Kendall's rank correlations τ_b between all the replication
255 indicators used in the study (see Supplementary Table 6). All preregistered replication indicators
256 are strongly correlated with each other, with τ_b varying between 0.61 and 1.00 ($p < 0.005$ for all
257 correlations). However, they are more weakly to moderately correlated with the prediction interval
258 approach and p-values from z-tests comparing the replication and original effect sizes (with τ_b
259 varying between 0.12 and 0.56).

260 Each of the various replication indicators presented in this study has its strengths and
261 weaknesses. There is no general consensus about which indicator is most appropriate^{15-17,117-119}.
262 Therefore, we chose to report the results for a host of indicators and leave it to readers to judge
263 the suitability of the different indicators and their degree of consensus. The overall replication rate
264 is similar for all the binary replication indicators and varies between 14 (53.8%) and 16 (61.5%)
265 studies. The agreement about which results are classified as successfully replicated is large
266 between the indicators, with the exception of the prediction interval approach. The estimated
267 average relative effect size of around 45%, which can be interpreted in terms of a replicability
268 rate, yields an estimate in the same ballpark. The somewhat lower estimate for the relative effect
269 size is due to the fact that not only the false positive rate but also the inflation of true positive
270 effect sizes is factored in. Another advantage of the relative effect size indicator is that it is not
271 affected by replication power. The three other continuous replication indicators cannot be

272 aggregated across studies and are thus difficult to compare to the other indicators on an
273 aggregated level.

274 **Replicability forecasts and indicators (not preregistered).** In Supplementary Table 7, we also
275 provide Pearson correlations between our replicability forecasts (final decision market prices and
276 average prediction survey beliefs) and all the replication indicators. Both the decision market
277 prices and the survey beliefs are positively correlated with all the replication indicators except the
278 prediction interval approach (although positive, the correlations to the p-value of the test of a
279 significant difference between the replication and original effect sizes are also close to 0). Note
280 that these correlations should be interpreted cautiously as the forecasters predicted the replication
281 outcomes for the statistical significance indicator but not the other replication indicators.

282 **Survey forecasts vs. decision market predictions.** We tested three additional preregistered
283 secondary hypotheses based on the survey beliefs about replication (see Methods for details and
284 Supplementary Table 7 for the survey results). The point-biserial correlation between average
285 survey beliefs and the replication outcomes based on the statistical significance criterion is
286 $r = 0.476$ (95% CI [0.107, 0.694]; $t(24) = 2.650$, $p = 0.014$; $n = 26$). The survey beliefs and the
287 decision market prices are positively correlated with a Pearson correlation of 0.899 (95% CI
288 [0.814, 0.944]; $t(39) = 12.830$, $p < 0.001$; $n = 41$) (Fig. 8a). The final secondary hypothesis tests
289 if the prediction accuracy, measured in terms of the absolute prediction error and the Brier score
290 (i.e., the squared prediction error), is higher for the decision market than the survey forecasts (Fig.
291 8b). The mean absolute prediction error and the mean Brier score are 0.353 and 0.188 for the
292 decision market, and 0.421 and 0.202 for the survey, respectively, providing suggestive evidence
293 for higher accuracy for the market forecasts based on the absolute prediction error (Wilcoxon
294 signed-rank test: $z = 2.172$, $p = 0.030$; $n = 26$) but not the Brier score (Wilcoxon signed-rank test:
295 $z = 1.181$, $p = 0.238$; $n = 26$). The failure to reject the null hypothesis for the Brier score does not
296 imply that the null hypothesis is true. In the survey, we also elicited forecasters' self-rated
297 expertise for each study. The average self-rated expertise (of participants eventually active in the
298 markets, $n = 162$) for the 26 replicated studies was 2.31 (sd = 1.40; $n = 4,212$) on a scale from 1
299 ("no knowledge of the topic") to 7 ("very high knowledge of the topic"). Supplementary Figure 1
300 plots the absolute prediction error and the Brier score of the 26 survey and decision market
301 forecasts over the average self-rated expertise per study. We do not find evidence for the
302 prediction accuracy and the average self-rated expertise being significantly correlated (not
303 preregistered; see Supplementary Figure 1 for details).

304 **Beliefs about the Covid-19 pandemic and replicability.** A potential issue raised by some
305 original authors in giving feedback on the replication reports prior to the data collection was that
306 the replicability of some original results might be affected by the Covid-19 pandemic (as all the
307 original studies were conducted before the pandemic). We evaluate this possibility in a
308 preregistered exploratory analysis, relying on the forecasters' beliefs about the impact of the
309 pandemic on replicability. As part of the prediction survey, participants were asked to judge
310 whether the pandemic would have affected the likelihood of successful replication, measured on

311 a scale from -3 (“the pandemic has definitely decreased the probability of replication”) to 3 (“the
312 pandemic has definitely increased the probability of replication”). We test if the average response
313 to this question differs from zero using a one-sample t-test for each of the 26 replications, and we
314 test if the average response across all 26 studies differs from 0. We find a statistically significant
315 result for four and a suggestive result for two replications on beliefs that Covid-19 has affected
316 the replication probability (Supplementary Table 8). For the six studies with suggestive or
317 statistically significant evidence, the estimate is negative for two studies and positive for four; only
318 in two of the cases does the sign of the expectation match the eventual replication outcome. For
319 the average belief about the impact of the pandemic on replicability across the 26 studies of 162
320 forecasters (who were active in the decision markets), there is suggestive evidence that the mean
321 of 0.039 ($sd = 0.190$) differs from zero ($t(161) = 2.598$, $p = 0.010$; $n = 162$). Somewhat
322 surprisingly—and in contrast to the concerns raised by some of the original authors—there is thus
323 a tendency for forecasters to believe that the pandemic has *increased* the average likelihood that
324 the studies will replicate. However, the magnitude of the effect is small ($d = 0.204$; 95% CI [0.049 ,
325 0.360]).

326 In addition, we tested, estimating the point-biserial correlation, if the average belief (per study)
327 about the pandemic’s impact on replicability correlates with the replication outcomes based on
328 the statistical significance indicator; we do not find a statistically significant association ($r = 0.014$,
329 95% CI [-0.360 , 0.382]; $t(24) = 0.068$, $p = 0.946$; $n = 26$). Yet, we cannot rule out that Covid-19
330 has entailed effects on replicability not foreseen by scholars participating in the survey. Further
331 work is needed to gauge whether and to which extent experimental replications—and predictions
332 of replication success—might be sensitive to macro-historical secular change such as economic
333 upheaval, wars, pandemics, etc. Forecasters’ beliefs about the pandemic’s impact on replicability
334 are also neither statistically significantly correlated with the final decision market prices ($r = 0.387$,
335 95% CI [-0.008 , 0.669]; $t(24) = 2.055$, $p = 0.051$; $n = 26$) nor the average survey belief of
336 replication ($r = 0.347$, 95% CI [-0.053 , 0.644]; $t(24) = 1.815$, $p = 0.082$; $n = 26$), although the point
337 estimates of the correlations are quite sizeable.

338 **Original p-value and replication (not preregistered).** For comparison to previous systematic
339 large-scale replication projects, we also report the correlation between the original p-value and
340 the replication outcome for the statistical significance indicator. The point-biserial correlation
341 between the original p-value and the replication outcome for the statistical significance indicator
342 is -0.400 ($p = 0.043$; 95% CI [0.014 , 0.648]) and comparable in magnitude to correlations of -0.33
343 in the RPP¹⁵, -0.57 in the EERP¹⁶, and -0.40 in the SSRP¹⁷.

344

345 Discussion

346 We found suggestive evidence ($p < 0.05$) for our first primary hypothesis that final decision market
347 prices correlate with replication outcomes ($r = 0.505$). However, the estimated effect size is
348 somewhat smaller than the effect size of $r = 0.67$, as presumed in our a priori power calculations

349 (see Methods for details). The estimated correlation is within the range of previous prediction
350 markets on systematic replication projects with correlations of 0.42 in the Replication Project:
351 Psychology (RPP)^{10,15}, 0.30 in the Experimental Economics Replication Project (EERP)¹⁶, and
352 0.84 in the Social Sciences Replication Project (SSRP)¹⁷, but we expected a stronger correlation
353 because we selected studies with the highest and the lowest prices for replication. Consistent
354 with the primary hypothesis test, there is also suggestive evidence of a difference in the replication
355 rate between the “top-12” (10 of 12) and “bottom-12” (4 of 12) in our secondary hypothesis test.
356 The difference of 50 percentage points is also reflected in the difference between the forecasted
357 replication rates of 86.6% (“top-12”) vs. 29.6% (“bottom-12”) in the decision market. However, the
358 small sample size suggests caution against drawing firm conclusions about whether decision
359 markets are appropriate for selecting studies for replication.

360 The pooled evidence from previous prediction market studies on replication outcomes suggests
361 that markets are somewhat more accurate than surveys⁴⁵, although the difference tends to be
362 small. These indications are consistent with our results, yielding suggestive evidence of higher
363 accuracy in terms of the absolute prediction error but not in terms of the squared prediction error
364 (although, as noted above, failing to reject the null hypothesis for the squared prediction error
365 does not imply that the null hypothesis is true). The estimated correlation between the average
366 survey beliefs and the replication outcomes was almost as high for the survey as the prediction
367 market (0.476 vs. 0.505). The decision market prices and survey beliefs are also highly correlated
368 with each other ($r = 0.9$). Since surveys are less resource-intensive, simple polls can be an
369 expedient alternative to decision markets for selecting which studies to replicate, even if they
370 should be somewhat less accurate. Another potential method for selecting which studies to
371 replicate would be to rely on the original p -value for studies reporting statistically significant
372 results⁴⁵. Although the prediction accuracy appears to be somewhat lower for original p -values
373 than market and survey forecasts⁴⁵, relying on p -values may well be considered a practical
374 alternative as it does not involve any additional data collection. Another possibility would be to
375 use predicted replication probabilities from machine learning models to select studies for
376 replication. There has been some progress in developing such models^{120–123}, but evidence on
377 whether they outperform markets or surveys is yet missing. Other potential mechanisms for
378 selecting which studies to replicate include relying on general or study-specific characteristics
379 (e.g., connection to theory, surprise factor, sample size, effect size, importance)^{25–28}, relying on
380 cost-benefit considerations^{29,30}, employing Bayesian strategies^{31,32}, determining the “replication
381 value”³³, adopting empirical audit and review³⁴, selecting studies randomly³⁵, or using predictions
382 from laypeople^{36,37}.

383 Using decision markets to select the studies with the highest and lowest predicted probabilities
384 for replication is just one of many potential selection rules for this methodology. Our goal was to
385 test whether a decision market could distinguish findings that would replicate or not, and we aimed
386 to maximize the statistical power of detecting an association between market prices and
387 replication outcomes. For the practical application of decision markets, the choice of the selection
388 mechanism will largely depend on the objective function. One selection rule would be to choose

389 the studies with the highest predicted false positive likelihood, i.e., the studies with the smallest
390 market prices (in addition to at least one randomly selected study to ensure incentive
391 compatibility). This decision mechanism would align with the objective of identifying and correcting
392 false discoveries in the literature to facilitate an efficient allocation of resources for follow-up
393 investigations. Another selection rule would be to replicate the studies with market predictions
394 close to 50%, which reflects the highest possible uncertainty or disagreement regarding the
395 likelihood of the original finding being genuinely true. Providing additional evidence on these
396 claims could maximize the information value of replication studies, as well-powered replications
397 will move the probability that the tested hypothesis is genuinely true towards 0% or 100%.

398 For our second primary hypothesis, we found strong evidence that original effect sizes are inflated
399 on average compared to replication effect sizes, with a relative estimated average effect size of
400 45%. This is comparable to previous systematic replication studies, with relative average effect
401 size estimates of 49% in the RPP¹⁵, 59% in the EERP¹⁴, and 54% in the SSRP¹⁷. The replication
402 rate of 54% based on the statistical significance indicator is also similar to previous replication
403 studies, with 36% in RPP¹⁵, 61% in the EERP¹⁴, and 62% in the SSRP¹⁷. Caution should be
404 exercised when comparing the replication results across these studies: the number of replications
405 in each of the projects is small, only one focal result per paper has been selected for replication,
406 and the particular journals and time periods considered differ. However, the results of all these
407 studies are consistent with a replication rate of about 50% for both the binary statistical
408 significance indicator and the continuous relative effect size indicator; compatible replication
409 results have also been observed in the Many Labs replication projects¹²⁻¹⁴.

410 The ability of the statistical significance indicator to discriminate between true positives and false
411 positives depends on replication power, and the relative average effect size of the studies that
412 failed to replicate should be close to zero if the systematic replication study successfully separates
413 false positives from true positives. The relative average estimated effect size of the 12 studies
414 that failed to replicate according to the statistical significance indicator was 3.2%, which is close
415 to zero and consistent with a successful separation between true positives and false positives.
416 But also true positive findings can be expected to have exaggerated effect sizes in the published
417 literature due to a lack of statistical power^{18,19}. In line with this, we found an estimated average
418 effect size of 69.5% for the 14 studies that were successfully replicated based on the statistical
419 significance indicator. These findings are consistent with similar analyses in the SSRP¹⁷ in which
420 the estimated mean relative effect size among the studies that failed to replicate according to the
421 statistical significance indicator was 0.3%, and the estimated mean relative effect size among the
422 studies that replicated successfully was 73.1%. This illustrates how the combination of statistical
423 significance and relative effect size can contribute to revealing possible false positives and true
424 positives with exaggerated effect sizes.

425 Previous systematic replication studies have focused on laboratory experiments rather than
426 online experiments. Concerns have been raised over data quality in online data collections using
427 “crowd workers,” as via MTurk⁸⁸⁻⁹⁴, and part of the rationale for zeroing in on experiments

428 conducted via MTurk was that we tend to share these concerns. However, the results of this study
429 do not suggest that replicability is substantively lower for experiments conducted via MTurk
430 compared to experiments conducted in physical laboratories for studies published in top journals;
431 more evidence is needed to draw strong conclusions. Relatedly, the predicted average
432 replicability rate of 57.6% in the decision market—despite widespread concerns about data quality
433 on MTurk—is within the range of replication rate forecasts in previous prediction markets of 56%
434 in the RPP¹⁰, 75% in the EERP¹⁶, and 63% in the SSRP¹⁷. We used IP quality checks^{90,124} to
435 minimize the chances of low-quality participant data (see Methods for details), screening out
436 participants before the random assignment into treatments. In total, across all 26 replications,
437 29% of the participants who accepted a “human intelligence task” (HIT) failed the IP check and
438 were excluded (this descriptive result was not preregistered; see Methods for further details). The
439 replication results from our study should thus not be extrapolated to MTurk experiments not using
440 a comparable screening procedure. An important caveat is that although our IP quality checks
441 seem to have been effective in filtering out bots, this may not be the case for artificial responses
442 generated by large language models like ChatGPT, which could pose a challenge for collecting
443 data online via platforms such as MTurk¹²⁵.

444 There are several important limitations to our study. A successful replication, on its own, does not
445 provide valid evidence for the tested hypothesis. It goes without saying that inference in replication
446 studies is subject to type-I and type-II errors, just as in original studies. Moreover, a finding can
447 be replicable while being based on an invalid experimental design, leading to biased results. An
448 example of this would be an experimental design that systematically results in more attrition in
449 one experimental treatment, causing selection bias in favor of the tested conjecture.⁸⁸ Likewise,
450 a failed replication, on its own, does not provide direct evidence against the tested hypothesis. A
451 finding can be unreplicable and based on an invalid experimental design, leaving the hypothesis
452 untested. Although the replication rate for online experiments in our study appears to be similar
453 to previous laboratory evidence, it does not necessarily imply that online and laboratory
454 experiments provide equally valid evidence of the tested hypotheses.

455 Another limitation is that we only replicate a single focal result per paper, and the replication
456 outcome does not necessarily generalize to other results reported in the original articles.
457 Furthermore, we only gathered data from one online population using the same experimental
458 design as in the original study. It cannot be ruled out that the difference in timing between the
459 replication studies and the original studies has affected the replication results as a consequence
460 of changes in the composition of the MTurk subject pool or the tested phenomenon having
461 changed over time. Large-scale, multi-site replication studies that collect data across various
462 populations and settings, similar to the Many Labs replication projects^{12–14}, qualify as a promising
463 method to shed light on the heterogeneity of replication effect sizes across populations and
464 designs^{126–128} in future replication work, potentially increasing the strength of evidence for whether
465 the hypothesis supported in the original study is likely true or not. Collaborative networks such as
466 the Psychological Science Accelerator¹³³ facilitate multi-site replication studies and can be a door
467 opener to large and diverse samples.

468 Another caveat in interpreting our results is the lack of agreement about how to define and
469 measure replicability. We chose to report the results for a broad set of replication indicators
470 proposed in the literature and leave it to readers to gauge the strengths and weaknesses of the
471 various measures. Decision markets come with the limitation of being a relatively resource-
472 intensive tool, rendering simple polls an appealing alternative.

473 In our proof-of-concept investigation of using decision markets to assess replicability, decision
474 markets show potential as a tool for selecting studies for replications, but further work is needed
475 to draw strong conclusions. The observed replication rate of social science experiments based
476 on data collections via MTurk published in PNAS is comparable to previous systematic replication
477 projects of experimental studies in the social sciences, primarily based on lab experiments.
478 However, the sample size of 26 replication studies is small, implying substantial uncertainty about
479 both the estimated replication rate and the estimated association between the decision market
480 prices and the replication rate. Our study is also limited to one scientific journal, and may not be
481 representative of social science results based on MTurk samples published in other journals, or
482 studies using other online platforms for the data collection. Thus, prudence should be exercised
483 in generalizing our findings and comparing replication results across studies.

484

485 **Methods**

486 We preregistered an analysis plan for the project at OSF on October 7, 2021, prior to starting the
487 survey data collection (that preceded the decision market and replications), which detailed the
488 design of the study and the exact analyses for all planned analyses and tests
489 (<https://osf.io/xsp6g>). Unless explicitly mentioned in the main text, we adhere exactly to our pre-
490 analysis plan. The information in this Methods section follows the pre-analysis plan (PAP; with
491 some of the information from the pre-analysis plan reported in the Supplementary Notes). Note
492 that previous systematic replication projects like the RPP¹⁵, the EERP¹⁶, and the SSRP¹⁷ did not
493 file preregistrations of the overall study protocol and planned analyses.

494 Prior to starting the survey data collection, we also preregistered an analysis plan (“replication
495 report”) for each of the 41 potential replications included in the decision market at OSF after
496 obtaining feedback from the original authors (<https://osf.io/sejyp>). After the replications had been
497 conducted, the 26 replication reports of the replications selected for replication by the decision
498 market were updated with the results of the replications and posted in the same OSF repository.
499 Any deviations from the preregistered analysis plans for the 26 replications are detailed in the 26
500 “post-replication reports” and listed in Supplementary Notes. We provided all original authors the
501 opportunity to comment on the replication results (without a particular due date) and make the
502 comments publicly available as we receive them alongside the post-replication reports on OSF
503 (<https://osf.io/sejyp>).

504 Below, we provide further details on the inclusion criteria, the decision market setup and the
505 survey, the replications, and the replication rate indicators included in the study. The preregistered

506 analyses and tests were divided into descriptive results of the replication rate among the 26
507 replicated studies and hypothesis tests. The preregistered descriptive results were furthermore
508 divided into primary replication indicators and secondary replication indicators, and the pre-
509 registered hypothesis tests were divided into (i) primary hypotheses, (ii) secondary hypotheses,
510 and (iii) exploratory analyses. See Supplementary Notes for more details about the preregistered
511 hypothesis tests and exploratory analyses.

512

513 **Inclusion criteria for studies**

514 We reviewed all *PNAS* articles from 2015–2018 and searched for the terms Amazon Mechanical
515 Turk, MTurk, and Turk. When we began planning our study at the start of 2019, we started
516 reviewing the most recent articles published in *PNAS*. We then continued to look back in time,
517 year by year, until we reached a sufficiently large number of studies to run a decision market.
518 However, data collection was delayed after some original authors expressed concerns that the
519 Covid-19 pandemic could affect the replication outcomes. We included all social sciences articles
520 that fulfilled our inclusion criteria for (i) the journal and time period, (ii) the platform on which the
521 experiment was performed (MTurk), (iii) the type of design (between-subjects or within-subject
522 treatment design), (iv) the equipment and materials needed to implement the experiment (the
523 experiment had to be logistically feasible for us to implement), and (v) the results reported in the
524 experiment (that there is at least one statistically significant $p < 0.05$ main or interaction effect in
525 the main text). Based on the inclusion criteria, we identified 44 articles. After contacting the
526 original authors, we ended up with 41 articles (the three excluded articles^{129–131} involved either
527 software or platforms that no longer existed or methods we were unfamiliar with). In these 41
528 articles, we identified at least one critical finding that we could replicate. In cases where several
529 studies in the same article fit the inclusion criteria, we randomly picked one of the studies; this
530 was the case for 17 of the 26 replicated studies (Ames and Fiske⁴⁶, Atir and Ferguson⁴⁷, Baldwin
531 and Lammers⁴⁸, Boswell et al.⁵⁰, Cooney et al.⁵⁶, Genschow et al.⁵⁹, Gheorghiu et al.⁶⁰, Halevy
532 and Halali⁶², Hofstetter et al.⁶⁵, John et al.⁶⁹, Jordan et al.⁷⁰, Klein and O'Brien⁷³, Kouchaki and
533 Gino⁷⁴, McCall et al.⁷⁶, Rai et al.⁸¹, Stern et al.⁸⁴, and Williams et al.⁸⁶). In cases where the
534 (randomly picked) study contained several conditions, we randomly picked which to compare to
535 the control condition. After that, we looked for the central result with $p < 0.05$ for that particular
536 study. If there were several statistically significant results, one was selected at random. The
537 replication results thus only pertain to the single central result selected per paper, and the
538 replication outcome does not necessarily generalize to other results reported in the original
539 articles. For convenience, we refer to the replications as “replication of [study reference],” though.
540 For Cheon and Hong⁵⁴, the result chosen for replication is reported as part of a 2x2 ANOVA in
541 the original article; since the paper does not report the main effect, the original authors kindly
542 provided us with the corresponding estimates. For Gheorghiu et al.⁶⁰, the result to be replicated
543 is only reported with its p-value in the paper; a precise estimate of the test statistic has been
544 obtained from a re-analysis of the original data, which the original authors kindly provided. For

545 the study by Kraus et al.⁷⁵, we could not reproduce the result reported in the original article using
546 the original data. The original authors acknowledged that there had been a reporting error in the
547 original article. For the replication, we use the analysis described in the paper; the effect size and
548 the test statistic reported in the original paper were replaced by the re-estimated result. For the
549 study by Williams et al.⁸⁶, the focal hypothesis test in the replication is based on a composite
550 score of five suites of behavior (which are tested separately in the original article) to have a single
551 test. The original authors also report tests on composite measures in the Supporting Information
552 of their article, and they approved the choice to investigate the replicability of the focal hypothesis
553 using a composite score. These changes are transparently reported in the replication reports for
554 each study (see <https://osf.io/sejyp> for details).

555

556 **Decision market and prediction survey**

557 We invited researchers to voluntarily participate in the decision market through public mailing lists
558 (ESA and JDM lists) and social media (e.g., Twitter/X); we also emailed colleagues asking them
559 to distribute the call to participants within their professional networks. Participants were required
560 to hold a Ph.D. degree or to be a Ph.D. student currently. In the decision market, participants bet
561 on whether or not the specific result chosen for each study would replicate based on the statistical
562 significance indicator ($p < 0.05$ in the replication and an effect in the same direction as in the
563 original study) as a criterion for replication (thus a binary outcome, as discussed below). Prior to
564 the decision market, participants filled out a survey where we asked them to assign a probability
565 of successful replication to each of the 41 results. The survey is available at <https://osf.io/a24zq>.
566 Completing the survey was a prerequisite for participating in the markets. We started the
567 recruitment of participants for the decision market on October 4 (2021), and we started sending
568 out the prediction survey on October 8 to those who had signed up for the study (participants who
569 signed up after October 8 received the survey invitation a few days after their registration). The
570 deadline for registering as a participant was October 29, and the deadline for completing the
571 survey was November 5. Overall, 289 participants signed up to participate and were forwarded
572 the link to the survey; 193 participants started the survey, and 162 completed it by the due date.
573 The forecasters were from the following fields of research: psychology (37.7%), economics
574 (34.6%), management (7.4%), political science (4.9%), sociology (1.9%), and other fields
575 (13.6%). No additional demographic information was collected.

576 In the survey, we asked participants to assess, for each replication study, (i) the likelihood that
577 the hypothesis will successfully replicate (on a scale from 0% to 100%); (ii) their stated expertise
578 for the study/the hypothesis (on a scale from 1 to 7); and (iii) whether they believe the pandemic
579 has affected the likelihood of replication. The question about the pandemic was measured on a
580 scale from -3 (“the pandemic has definitely decreased the probability of replication”) to 3 (“the
581 pandemic has definitely increased the probability of replication”); the 0 midpoint implies that they
582 do not think that the pandemic has affected the probability of replication. The survey was not
583 incentivized.

584 The decision market opened on November 8 (2021) and closed after two weeks on November 22
585 (and before the decision market opened, participants had at least one week to complete the
586 prediction survey). In the decision market, participants could trade (bet) on whether they expected
587 the 41 studies to replicate. While participants had the opportunity to bet on the replication outcome
588 of the 41 studies, we did not carry out replications for all 41 studies, but the final decision market
589 prices determined which studies to replicate. We replicated the 12 studies that had the highest
590 and the 12 studies that had the lowest market prices when the market closed. Additionally, two
591 out of the remaining 17 studies were randomly selected for replication to ensure a non-zero
592 probability for each study to be replicated (i.e., we replicated $12 + 12 + 2 = 26$ studies in total).
593 Since payoffs are only determined based on forecasts of studies that were eventually replicated,
594 payoffs were scaled up by the inverse of their probability of being selected for replication in the
595 decision rule (see below for details). This incentive scheme encourages trading based on traders'
596 true beliefs, even though some studies will not be replicated. Consequently, participants have the
597 incentive to buy shares of a particular study whenever they believe that the likelihood of replication
598 is higher than the current market price; likewise, participants have an incentive to (short) sell
599 shares whenever they believe that the likelihood of replication is lower than the current market
600 price. Thus, as long as the market price differs from the predicted likelihood of replication for a
601 participant, the participant has an incentive to buy or (short) sell shares of a particular study and
602 realizing a trade according to a trader's belief will move the market price in the direction of the
603 trader's belief. The decision rule for which studies to replicate was based on final market prices
604 and was common knowledge to the market participants; the instructions (provided to participants
605 who completed the prediction survey by the due date) are available at <https://osf.io/a24zq/>.

606 We chose 12 studies with the lowest predicted probability and 12 studies with the highest
607 predicted probability based on a power calculation using the pooled data from our previous
608 prediction market studies⁴⁵. The power calculations were conducted by randomly sampling 41
609 studies from the dataset described in Gordon et al.⁴⁵ in a simulation with 10,000 iterations and
610 then selecting the forecasts and outcomes from the 12 studies with the lowest predicted
611 probability, the 12 studies with the highest predicted probability, and two random studies. We
612 failed to set a random seed for the simulation when the study was conducted, implying that the
613 pre-registered power estimates could not be numerically reproduced when we wrote up the study
614 results. For full transparency, we report the power estimates included in the PAP in parentheses
615 below for full transparency. The median point-biserial correlation coefficient across the 10,000
616 runs is 0.671 (reported as 0.66 in the PAP), and we have 91.0% power (reported as >90% in the
617 PAP) to detect a statistically significant correlation ($n = 26$) between the decision market prices
618 and the replication outcomes at the 0.5% level and 99.4% power (reported as >95% in the PAP)
619 to detect a statistically significant correlation at the 5% level, which is our first primary hypothesis
620 test. As a secondary hypothesis test, we test if the fraction of studies that successfully replicate
621 differs between the 12 studies with the highest and the 12 studies with the lowest predicted
622 replication probabilities using Fisher's exact test. Applying the same sampling approach as for
623 primary hypothesis 1, the median difference in replication rates between the 12 studies with the

624 highest and the 12 studies with the lowest market prices is 0.663; the secondary test (n = 24) has
625 66.5% power at the 0.5% level (reported as 66% in the PAP) and 94.9% power at the 5% level
626 (reported as 95% in the PAP). The code for the power simulations is available at
627 <https://osf.io/47drs>.

628 **Implementation of the decision market.** We used a web-based trading platform, similar to the
629 ones used in Camerer et al.^{16,17} and identical to the one used in Botvinnik-Nezer et al.⁶. The trading
630 platform involves two main views: (i) the market overview and (ii) the trading page. The market
631 overview listed the 41 assets (i.e., one corresponding to each study) in tabular format, including
632 information on the current price for buying a share and the number of shares held (separated for
633 long and short positions). Via the trading page, which was shown after clicking on any of the
634 assets, participants could make investment decisions (i.e., buy or sell shares) and view price
635 developments in graphical format for the particular asset.

636 **Trading and incentivization.** Decision market participants received an endowment of 100 tokens
637 corresponding to USD 50. Once the markets opened, market participants could use the tokens to
638 trade shares of the assets available in the market. An automated market maker, implementing a
639 logarithmic market scoring rule¹³², determined the assets' prices. At the beginning of the markets,
640 all assets were valued at 0.50 tokens per share. The market maker calculated the share price for
641 each infinitesimal transaction and updated the price based on the scoring rule. With this
642 mechanism, participants had incentives to invest according to their beliefs^{43,44}. With the
643 logarithmic scoring rule, the price p for an infinitesimal trade is determined as $p = e^{s/b} \div (e^{s/b} + 1)$,
644 where s denotes the net sales (shares held – shares borrowed) that the market maker has done
645 so far in a market; the liquidity parameter b determines how strongly the market price is affected
646 by trade and was set to $b = 100$, implying that by investing ten tokens, traders could move the
647 price of a single asset from 0.50 to about 0.55. We opted for the same value of b as the one
648 employed in the prediction markets in SSRP¹⁷, which appears intuitively sensible in terms of
649 striking a good balance between price sensitivity and liquidity. Notwithstanding, it is worth noting
650 that it is unclear whether or not our results are sensitive to the choice of this parameter. Decision
651 market participants were paid only for studies chosen for replication (based on their final
652 holdings). Participants received one token per correct share for the replications with the 12 lowest
653 and 12 highest final market prices. For the two randomly selected replications, participants
654 received $17 \div 2 = 8.5$ tokens for each share; for replications that were not chosen for replication,
655 participants received no compensation for their holdings. We followed this procedure to keep
656 information revelation in the decision market incentive-compatible, with the increased payouts for
657 the randomly selected studies compensating for the “voided” shares in studies not chosen for
658 replication. Participants were paid after all 26 replications had been completed.

659 **Participation.** A total of 193 participants completed the prediction survey (a prerequisite to
660 participate in market trading) after providing consent to participate and were subsequently invited
661 to trade on the decision market. Of these 193 participants, 162 (83.9%) traded in the market at
662 least once. During the two-week trading period, a total of 4,412 transactions were recorded. On

663 average, each trader prompted 27.2 transactions (sd = 30.7; min = 1, max = 185). The average
664 number of traders per hypothesis was 65.1 (sd = 15.3; min = 35, max = 98); the average number
665 of transactions recorded per hypothesis was 107.6 (sd = 35.2; min = 56, max = 213). See
666 Supplementary Table 2 for descriptive statistics on the trading activity for each market.

667

668 **Replications**

669 We carry out close replications¹⁰⁷ as closely as possible following the experimental design,
670 sample restrictions, exclusion criteria, and analysis as used in the original studies and carried out
671 in the same population (Amazon Mechanical Turk) as the original studies. The replications started
672 in January 2022 and were completed in October 2023. The replications were planned and pre-
673 registered by five replication teams: a team at CalTech, LMU, and Wharton; a team at the
674 Stockholm School of Economics; a team at the National University of Singapore; a team at the
675 University of Amsterdam; and a team at the University of Innsbruck.

676 **Participants in replication studies.** All replications were carried out at Amazon Mechanical Turk
677 as in the original studies. We ensured that participants could only participate once using the same
678 account in a specific study. If the original study had not specified a HIT approval rate, we recruited
679 participants with a HIT approval rate of at least 95%; if the original study had specified a higher
680 approval rate, we applied the same requirement as used in the original study.

681 To ward off concerns about impaired data quality due to low-attention participants and bots^{88–94},
682 we implemented several “quality filters.” Particularly, before redirecting participants to each study,
683 we forwarded the IP addresses to <https://www.ipqualityscore.com/> for a quality check to minimize
684 the chances of low-quality participant data (we initially planned to use this filter ex-post but during
685 the data collection of the first two replication studies of Klein & O’Brien⁷³ and Halevy & Halali⁶² we
686 decided to set it up so that the IP address quality check happened before participants got
687 redirected to the study). Participants for whom one or more of the following was true could not
688 proceed with participating in the study: fraud score ≥ 85 ; TOR = True; VPN = True; Bot = True;
689 abuse velocity = high. This means that, for example, participants were not allowed to use a virtual
690 private network (VPN) or Tor connections or participate if they had IP addresses that had recently
691 engaged in automated bot activity (the VPN exclusions were made ex-ante, i.e., before
692 participants were redirected to the study, for four studies and ex-post for 22 studies). After that,
693 in all replications, participants were first shown a Captcha and then provided informed consent.
694 After this, we included an attention check that participants had to pass to proceed to the study
695 (with the exception of Reeck et al.⁸²; see Section 4 in the Supplementary Information for details).
696 The attention check was implemented in addition to any other potential attention check(s) used
697 in the original study. All these exclusions based on the “quality filters” were preregistered, but the
698 pre-analysis plan did not specify if participants would be excluded before or after participating in
699 the study.

700 The individual replication studies sometimes also used additional exclusion criteria that are
701 detailed in the preregistered replication report for each replication (we tried to use the same

702 exclusion criteria for the replications as used in the original studies as much as possible). The
703 replication sample sizes defined below are the sample sizes after any exclusions of participants.

704 **Replication sample sizes.** The replications were carried out with high statistical power.
705 Replication sample sizes were based on having 90% power to detect $\frac{2}{3}$ of the effect size reported
706 in the original study (with the effect size converted to Cohen's d to have a common standardized
707 effect size measure across the original studies and the replication studies). See Supplementary
708 Notes for more details about the power calculations and replication sample sizes. The criteria for
709 replication were an effect in the same direction as the original study and a p -value < 0.05 (in a
710 two-sided test). In cases where this power estimation led to a sample size smaller than the original
711 one, we used the same sample size as in the original study. The average replication sample
712 ($\bar{n} = 1,018$) size was 3.5 times as large as the average sample size of the original studies
713 ($\bar{n} = 292$). We continued the data collection for each replication until we reached at least the
714 preregistered sample size after exclusions for that replication, and this led to slightly larger
715 replication sample sizes than preregistered in all replications except one (as it is not possible with
716 exclusion criteria to get an exact sample size as the number of exclusions is not known ex-ante).

717 **Conversion of effect sizes to Cohen's d .** We converted the effect sizes of all the original studies
718 and all the replication studies to Cohen's d to have a standardized effect size (the effect size in
719 the original study was always assigned a positive sign; the effect size in the replication study was
720 assigned a positive sign if the effect was in the same direction as in the original study and a
721 negative sign if the effect was in the opposite direction of the original study). See Supplementary
722 Notes for details about the conversion of effect sizes to Cohen's d .

723 **Replication reports.** For each of the 41 studies, we prepared a pre-replication plan/report stating
724 the hypothesis we had chosen from each paper and how we planned to proceed with the
725 replication study. These reports were shared with the original authors for feedback, and at least
726 one original author from each paper replied. These pre-replication reports were posted at OSF
727 (<https://osf.io/sejyp>) at the same time as the pre-analysis plan and prior to the start of the prediction
728 survey (that preceded the decision markets and the replication data collections). For those studies
729 that were selected for replication, we have updated the replication reports with the replication
730 results after the replications were completed. After sharing them with the original authors for
731 feedback, we have posted the updated replication reports at OSF as well (<https://osf.io/sejyp>).
732 Additionally, we reached out to the original authors for their comments on the replication reports
733 and results. We promised to make their comments available along with the replication reports,
734 and any comments received can be found at <https://osf.io/sejyp>.

735 **Incentivization in the replication experiments.** We standardized payments across all
736 replications such that studies had a certain show-up fee depending on the expected length of the
737 study. In particular, we paid an hourly fee of USD 8.00 for all studies, and we calculated the show-
738 up fee for each study based on the expected length of the study. For all studies, we implemented
739 a minimum payoff of USD 1.00. For studies with incentive payments, we used the same incentive

740 payment as in the original study, paid in addition to the show-up fee. If we faced problems in
741 recruiting participants, we increased the show-up fee, which happened for two studies^{61,65}.

742

743 **Replication indicators**

744 **Statistical significance criterion (primary indicator).** The first primary replication indicator was
745 the statistical significance criterion – i.e., whether the replication resulted in an effect size in the
746 same direction as the original study and a two-sided p-value less than 0.05. Unless otherwise
747 stated above, we used the same statistical test as in the original study. We report the replication
748 rate (i.e., the fraction of the 26 studies that replicated according to this criterion) and the 95%
749 Clopper-Pearson CI of this fraction in the Results section. We also report the 95% CI of the
750 replication effect size for each of the 26 replication studies in Fig. 2 and Supplementary Table 3.

751 **Relative effect sizes (primary indicator).** As a second primary replication indicator, we used
752 relative effect sizes. Relative effect sizes were estimated in two different ways. We report the
753 mean effect size of all 26 replications and compare it to the mean effect size of the 26 original
754 studies (see also primary hypothesis test 2 below). We furthermore estimate the relative effect
755 size of each replication (the replication effect size divided by the original effect size) and estimate
756 the mean of this variable for the 26 replication studies and the 95% CI of this mean (based on a
757 one-sample t-test). We report both of these measures of the relative effect size separately for the
758 replications that replicate and those that do not. These results are reported in the Results section,
759 Fig. 3, and Supplementary Table 3.

760 **Small-telescopes approach (secondary indicator).** We also used the small-telescopes
761 approach¹¹². For this indicator, we estimated whether the replication effect size was significantly
762 smaller (using a one-sided test at the 5% level) than a “small effect,” defined as the effect size
763 the original study would have had 33% power to detect. For studies using t-tests (or F-tests
764 converted to a t-test statistic), we based “the small effect size” on the effect size that a t-test had
765 33% power to detect (at the 5% level in a two-sided test); for studies using z-test statistics (or chi-
766 square tests converted to a z-test statistic), we based “the small effect size” on the effect size that
767 a z-test had 33% power to detect (at the 5% level in a two-sided test). To test whether the
768 replication effect size was significantly smaller than “the small effect size” in a one-sided test at
769 the 5% level, we estimated a 90% CI of the replication effect size. We tested if the 90% CI
770 overlapped the small effect size with CIs constructed as described in Supplementary Notes. If the
771 effect size in the replication was significantly smaller than this “small effect size,” the result was
772 considered a failed replication; otherwise, it was considered successful. We report the fraction of
773 studies that replicate according to this criterion and the 95% Clopper-Pearson CI of this fraction.
774 The small-telescopes results are reported in Fig. 4 and Supplementary Table 4.

775 **Bayes factors (secondary indicators).** We also compute the one-sided default Bayes factors
776 on the replication data, allowing us to obtain the strength of evidence in favor of the hypothesis
777 that stipulates an effect in the direction of the original experiment (where a default prior in terms
778 of a truncated Cauchy distribution with scale 0.707 was assigned to the size of the effect) versus

779 the null hypothesis that stipulates the effect to be absent¹¹³. In addition, we also computed (one-
780 sided) replication Bayes factors, which quantifies the additional evidence for the hypothesis given
781 the evidence already provided by the original study¹¹⁴. (We are counting the one-sided default
782 and replication as Bayes factors as two separate indicators, which they are.) These results are
783 reported in Fig. 5 and Supplementary Table 4. We use the evidence categories proposed by
784 Jeffreys¹¹⁵ to interpret the Bayes Factors. A detailed report on the estimation of the Bayes factors
785 is available at <https://osf.io/47drs/>.

786 **Meta-analytic effect sizes (secondary indicator).** We estimated the meta-analytic estimate of
787 the effect size by combining the original result and the replication result in a fixed-effect meta-
788 analysis. We report the fraction of the 26 studies that replicated according to the 0.05 and the
789 0.005 significance threshold and the 95% Clopper-Pearson CI of these fractions. We also use the
790 stricter 0.005 significance threshold as a replication indicator for the meta-analytic effect sizes
791 because this is similar to observing two studies (an original study and a replication study) that are
792 significant at the 0.05 level. We report these results in the Results section, Fig. 6, and
793 Supplementary Table 4.

794

795

796

797 **Data Availability**

798 The data reported in this paper is tabulated in Supplementary Tables 1–8. The replication reports
799 (both the pre-replication and the post-replication versions), the pre-analysis plan, the data from
800 the survey and the decision market, and the data for each of the 26 replications are available at
801 the project’s OSF repository (<https://osf.io/sk82q>).

802

803 **Code Availability**

804 The analysis scripts, generating all results, figures, and tables reported in the main text and the
805 Supplementary Information, are available at the project’s OSF repository (<https://osf.io/sk82q>).

806

807 **Acknowledgments**

808 We thank Alexander Andevall for helping with the data collection and programming of experiments
809 and Robb Willer for helpful advice on defining the IP address check and exclusion criteria used
810 to exclude individuals from participating to minimize low-quality participant data. For financial
811 support, we thank the Austrian Science FWF (grant SFB F63 to J.H. and M.K.), Jan Wallander
812 and Tom Hedelius Foundation (grants P21-0091 and P23-0098 to A.D.), Knut and Alice
813 Wallenberg Foundation and Marianne and Marcus Wallenberg Foundation (Wallenberg Scholar
814 grant to A.D.), and Riksbankens Jubileumsfond (grant P21-0168 to M.J.). The funders had no role
815 in study design, data collection and analysis, decision to publish or preparation of the manuscript.
816 One author (V.W.) is currently employed by Sveriges Riksbank but did this work before being
817 employed by Sveriges Riksbank; the opinions expressed in this article are the sole responsibility
818 of the authors and should not be interpreted as reflecting the views of Sveriges Riksbank. We
819 sought ethical approval from the Swedish Ethical Review Authority who had no ethical objections
820 to the decision market part of the project and judged the replication part of the project to not be
821 covered by the Swedish ethical review law (Dnr 2019-06501).

822

823 **Author Contributions**

824 A.D., F.H., J.H., M.J., M.K., B.A.N., and T.P. designed the study; A.D., F.H., and M.J. managed
825 the study; Y.C., A.D., F.H., M.J., and T.P. designed and implemented the decision market; A.D.,
826 F.H., M.J., B.M., and V.W. selected articles and critical findings for (potential) replication; A.D.,
827 C.F.C., F.H., T.-H.H., S.H., J.H., N.I., T.I., L.J., M.J., M.K., B.M., D.M., G.N., A.S., R.S., E.-J.W.,
828 and V.W., designed the replications and collected replication data; F.H., S.H., T.I., L.J., A.S., R.S.,
829 and V.W. conducted the preregistered statistical tests on the individual replications; A.L.

830 computed the Bayes factors; F.H. conducted all analyses reported in the manuscript; A.D., F.H.,
831 and M.J. wrote the paper; all authors reviewed and approved the final manuscript.

832

833 **Competing Interest Statement**

834 The authors declare no competing interests.

835

836 **Figure legends**

837 **Fig. 1. Decision market prices for the 41 included studies.** Plotted are the decision market
838 prices for the 41 MTurk social science experiments published in PNAS between 2015 and 2018.
839 The small gray dots indicate the market prices after each market transaction; the larger dots
840 indicate the final market price. The studies are ordered based on the final decision market prices,
841 which can be interpreted as the market's probability forecast of successful replication. The 12
842 studies with the highest decision market prices and the 12 studies with the lowest decision market
843 prices were selected for replication; in addition, two of the remaining 17 studies were selected for
844 replication at random to ensure that the decision market is incentive compatible. The replication
845 outcomes for the statistical significance indicator are also illustrated for the 26 replicated studies.
846 The point-biserial correlation between the decision market prices and the replication outcomes in
847 primary hypothesis 1 is $r = 0.505$ (95% CI [0.146, 0.712], $t(24) = 2.867$, $p = 0.008$; $n = 26$, two-
848 sided test).

849 **Fig. 2. Replication results.** Plotted are the point estimates and the 95% CIs (standardized to
850 Cohen's d units) of the 26 replications (d_R) and original studies. Studies within each of the three
851 panels (top-12, random, bottom-12) are sorted based on the decision market prices as in Fig. 1.
852 There is a statistically significant effect ($p < 0.05$) in the same direction as the original study for
853 14 out of 26 replications (53.8%; 95% CI [33.4%, 73.4%]). For the 12 studies with the highest
854 decision market prices, there is a statistically significant effect ($p < 0.05$) in the same direction as
855 the original study for 10 out of 12 replications (83.3%; 95% CI [51.6%, 97.9%]). For the 12 studies
856 with the lowest decision market prices, there is a statistically significant effect ($p < 0.05$) in the
857 same direction as the original study for 4 out of 12 replications (33.3%; 95% CI [9.9%, 65.1%]).
858 Our secondary hypothesis test provides suggestive evidence that the difference in replication
859 rates between the top-12 and the bottom-12 group is different from zero (Fisher's exact test;
860 $\chi^2(1) = 6.171$, $p = 0.036$; $n = 24$, two-sided test). The error bars denote the 95% confidence
861 intervals (CIs) of the original and the replication effect size estimates. The numbers of
862 observations used to estimate the 95% CIs are the original and replication sample sizes noted on
863 the right as n_O and n_R .

864 **Fig. 3. Relationship between estimated original and replication effect sizes.** Plotted are the
865 estimated original and replication effect sizes for each of the 26 replication studies (the estimated
866 effect sizes of both the original and replication studies are standardized to Cohen's d units). The
867 95% confidence intervals for the original and replication effect size estimates are illustrated in Fig.

868 2 and tabulated in Supplementary Table 3. The mean estimated effect size of the 26 replication
869 studies is 0.253 (sd = 0.357) compared to 0.563 (sd = 0.426) for the original studies, resulting in
870 a relative estimated average effect size of 45.0%, confirming our second primary hypothesis
871 (Wilcoxon signed-rank test; $z = 4.203$, $p < 0.001$; $n = 26$, two-sided test). The estimated relative
872 effect size of the 13 replications that have been successfully replicated according to the statistical
873 significance indicator is 69.5%, and the estimated relative effect size of the 13 studies that did not
874 replicate is 3.2%. The box plots show the median, the interquartile range, and the 5th and 95th
875 percentile of the effect size estimates in the 26 original studies and the 26 replication studies.

876 **Fig. 4. Replication results based on the small-telescopes approach (a secondary**
877 **replication indicator).** Plotted are the 90% CIs of replication effect sizes in relation to small-
878 effect sizes as defined by the small-telescopes approach¹¹² (the effect size that the original study
879 would have had 33% power to detect). Studies within the three panels (top-12, random, bottom-
880 12) are sorted based on the decision market prices as in Fig. 1. A study is defined as failing to
881 replicate if the 90% CI is below the small effect. According to the small-telescopes approach, 15
882 out of 26 studies (57.7%; 95% CI [36.9%, 76.6%]) replicate. The error bars denote the 90%
883 confidence intervals (CIs) of the estimated replication effect sizes. The numbers of observations
884 used to estimate the 90% CIs are the replication sample sizes noted on the right as n_R .

885 **Fig. 5. Replication results based on Bayes factors (secondary replication indicators).** The
886 figure plots the one-sided default Bayes factor (BF_{+0}) and the replication Bayes factor (BF_{R0}) for
887 the 26 replications¹¹³. $BF_{+0} > 1$ favors the hypothesis of an effect in the direction of the original
888 paper, whereas $BF_{+0} < 1$ favors the null hypothesis of no effect. BF_{R0} quantifies the additional
889 evidence provided by the replication results on top of the original evidence. $BF_{R0} > 1$ indicates
890 additional evidence in favor of the alternative over the null, whereas $BF_{R0} < 1$ indicates additional
891 evidence for the null instead. The evidence categories proposed by Jeffreys¹¹⁵ are also shown
892 (from extreme support for the null hypothesis to extreme support for the original hypothesis).
893 Studies within the three panels (top-12, random, bottom-12) are sorted based on the decision
894 market prices as in Fig. 1. The BF_{+0} is above one for all 14 replication studies that successfully
895 replicated according to the statistical significance indicator and below one for all 12 replication
896 studies that failed to replicate according to the statistical significance indicator. The BF_{R0} is above
897 one for 13 of the 14 replication studies that replicated according to the statistical significance
898 indicator and below one for Cooney et al.⁵⁶ whose estimated relative effect size of 0.36 is the
899 lowest among these 14 studies; the BF_{R0} is below one for all of the 12 replication studies that
900 failed to replicate according to the statistical significance indicator. The numbers of observations
901 used to estimate BF_{+0} and BF_{R0} are the original and replication sample sizes noted on the right
902 as n_O and n_R .

903 **Fig. 6. Meta-analytic estimated effect sizes combining the original and the replication**
904 **estimated effect sizes (a secondary replication indicator).** The figure plots the point estimates
905 and 95% and 99.5% CIs of the fixed-effects weighted meta-analytic effect sizes, combining the
906 original and the replication studies (standardized to Cohen's d units). Studies within the three

907 panels (top-12, random, bottom-12) are sorted based on the decision market prices as in Fig. 1.
908 As preregistered, we report the significance of the estimated meta-analytic effect sizes for both
909 the 0.05 significance threshold and the 0.005 significance threshold (based on a two-sided z-test).
910 Sixteen out of 26 (61.5%; 95% CI [40.6%, 79.8%]) studies replicated according to the statistical
911 significance indicator using the 0.05 significance threshold, and 14 out of 26 (53.8%; 95% CI
912 [33.4%, 73.4%]) studies replicated using the 0.005 significance threshold. The error bars denote
913 the 95% confidence intervals (CIs) of the estimated meta-analytic effect sizes. The number of
914 observations used to estimate the 95% CIs are the sums of the original and replication sample
915 sizes noted next to the study identifier on the Y-axis as n_{O+R} .

916 **Fig. 7. Replication results based on prediction intervals (not preregistered).** Plotted are the
917 95% prediction intervals¹¹⁶ for the standardized original effect sizes (Cohen's d). Studies within
918 the three panels (top-12, random, bottom-12) are sorted based on the decision market prices as
919 in Fig. 1. Fifteen replications out of 26 (57.7%; CI [36.9%, 76.6%]) are within the 95% prediction
920 interval and replicate according to this indicator. The p-values reported on the right are based on
921 two-sample z-tests for a difference between the replication effect size and the original effect size.
922 The gray lines denote the 95% prediction intervals, and the small circles denote the mean
923 replication effect sizes. All tests are two-sided. The numbers of observations used to estimate the
924 95% prediction intervals are the original and replication sample sizes noted next to the study
925 identifier on the y-axis as n_O and n_R .

926 **Fig. 8. Relationship between decision market prices and mean survey beliefs and**
927 **forecasting accuracy. a,** Plotted are the decision market prices and the mean survey beliefs
928 about replication for the 41 studies included in the decision market and the survey; the color
929 coding highlights the replication outcomes for the 26 replicated studies. The decision market
930 prices and the mean survey beliefs about replication are highly correlated with a Pearson
931 correlation of $r = 0.899$ (95% CI [0.814, 0.944]; $t(39) = 12.830$, $p = 1.4e^{-15}$; $n = 41$, two-sided test).
932 **b,** Plotted are the absolute prediction errors and the Brier scores (the squared prediction errors)
933 for the decision market and the prediction survey for the 26 replicated studies. There is suggestive
934 evidence of higher prediction accuracy for the decision market in terms of the absolute prediction
935 error (0.353 for the decision markets and 0.421 for the survey; Wilcoxon signed-rank test:
936 $z = 2.172$, $p = 0.030$; $n = 26$, two-sided test), but not in terms of the Brier score (0.188 for the
937 decision markets and 0.202 for the survey; Wilcoxon signed-rank test: $z = 1.181$, $p = 0.238$;
938 $n = 26$, two-sided test). The box plots show the median, the interquartile range, and the 5th and
939 95th percentile of the absolute prediction errors and Brier scores for the survey and decision
940 market predictions of the 26 replication studies.

941 **References**

- 942 1. Ioannidis, J. P. A. Why most published research findings are false. *PLoS Med.* **2**, e124
943 (2005).
- 944 2. Leamer, E. E. Let's take the con out of econometrics. *Am. Econ. Rev.* **73**, 31–43 (1983).
- 945 3. Begley, C. G. & Ellis, L. M. Raise standards for preclinical cancer research. *Nature* **483**,
946 531–533 (2012).
- 947 4. McNutt, M. Reproducibility. *Science* **343**, 229–229 (2014).
- 948 5. Gertler, P., Galiani, S. & Romero, M. How to make replication the norm. *Nature* **554**, 417–
949 419 (2018).
- 950 6. Botvinik-Nezer, R. *et al.* Variability in the analysis of a single neuroimaging dataset by
951 many teams. *Nature* **582**, 84–88 (2020).
- 952 7. Breznau, N. *et al.* Observing many researchers using the same data and hypothesis
953 reveals a hidden universe of uncertainty. *Proc. Natl. Acad. Sci.* **119**, e2203150119 (2022).
- 954 8. Delios, A. *et al.* Examining the generalizability of research findings from archival data.
955 *Proc. Natl. Acad. Sci.* **119**, e2120377119 (2022).
- 956 9. Huber, C. *et al.* Competition and moral behavior: A meta-analysis of forty-five crowd-
957 sourced experimental designs. *Proc. Natl. Acad. Sci.* **120**, e2215572120 (2023).
- 958 10. Dreber, A. *et al.* Using prediction markets to estimate the reproducibility of scientific
959 research. *Proc. Natl. Acad. Sci.* **112**, 15343–15347 (2015).
- 960 11. Maniadis, Z., Tufano, F. & List, J. A. To replicate or not to replicate? Exploring
961 reproducibility in economics through the lens of a model and a pilot study. *Econ. J.* **127**,
962 F209–F235 (2017).
- 963 12. Klein, R. A. *et al.* Investigating variation in replicability: A “many labs” replication project.
964 *Soc. Psychol.* **45**, 142–152 (2014).
- 965 13. Ebersole, C. R. *et al.* Many Labs 3: Evaluating participant pool quality across the academic
966 semester via replication. *J. Exp. Soc. Psychol.* **67**, 68–82 (2016).
- 967 14. Klein, R. A. *et al.* Many Labs 2: Investigating variation in replicability across samples and
968 settings. *Adv. Methods Pract. Psychol. Sci.* **1**, 443–490 (2018).
- 969 15. Open Science Collaboration. Estimating the reproducibility of psychological science.
970 *Science* **349**, aac4716 (2015).
- 971 16. Camerer, C. F. *et al.* Evaluating replicability of laboratory experiments in economics.
972 *Science* **351**, 1433–1436 (2016).
- 973 17. Camerer, C. F. *et al.* Evaluating the replicability of social science experiments in Nature
974 and Science between 2010 and 2015. *Nat. Hum. Behav.* **2**, 637–644 (2018).

- 975 18. Ioannidis, J. P. A. Why most discovered true associations are inflated. *Epidemiology* **19**,
976 640 (2008).
- 977 19. Button, K. S. *et al.* Power failure: Why small sample size undermines the reliability of
978 neuroscience. *Nat. Rev. Neurosci.* **14**, 365–376 (2013).
- 979 20. Wegener, D. T., Fabrigar, L. R., Pek, J. & Hoisington-Shaw, K. Evaluating research in
980 personality and social psychology: Considerations of statistical power and concerns about
981 false findings. *Pers. Soc. Psychol. Bull.* **48**, 1105–1117 (2022).
- 982 21. Simmons, J. P., Nelson, L. D. & Simonsohn, U. False-positive psychology: Undisclosed
983 flexibility in data collection and analysis allows presenting anything as significant. *Psychol.*
984 *Sci.* **22**, 1359–1366 (2011).
- 985 22. John, L. K., Loewenstein, G. & Prelec, D. Measuring the Prevalence of Questionable
986 Research Practices With Incentives for Truth Telling. *Psychol. Sci.* **23**, 524–532 (2012).
- 987 23. Finkel, E. J., Eastwick, P. W. & Reis, H. T. Replicability and other features of a high-quality
988 science: Toward a balanced and empirical approach. *J. Pers. Soc. Psychol.* **113**, (2017).
- 989 24. Flake, J. K., Davidson, I. J., Wong, O. & Pek, J. Construct validity and the validity of
990 replication studies: A systematic review. *Am. Psychol.* **77**, 576–588 (2022).
- 991 25. Pittelkow, M.-M. *et al.* The process of replication target selection in psychology: What to
992 consider? *R. Soc. Open Sci.* **10**, 210586 (2023).
- 993 26. Makel, M. C., Plucker, J. A. & Hegarty, B. Replications in psychology research: How often
994 do they really occur? *Perspect. Psychol. Sci.* **7**, 537–542 (2012).
- 995 27. Lindsay, D. S. Replication in psychological science. *Psychol. Sci.* **26**, 1827–1832 (2015).
- 996 28. Block, J. & Kuckertz, A. Seven principles of effective replication studies: Strengthening
997 the evidence base of management research. *Manag. Rev. Q.* **68**, 355–359 (2018).
- 998 29. Coles, N. A., Tiokhin, L., Scheel, A. M., Isager, P. M. & Lakens, D. The costs and benefits
999 of replication studies. *Behav. Brain Sci.* **41**, e124 (2018).
- 1000 30. Alipourfard, N. *et al.* Systematizing Confidence in Open Research and Evidence
1001 (SCORE). Preprint at <https://doi.org/10/hn4g> (2021).
- 1002 31. Hardwicke, T. E., Tessler, M. H., Peloquin, B. N. & Frank, M. C. A Bayesian decision-
1003 making framework for replication. *Behav. Brain Sci.* **41**, e132 (2018).
- 1004 32. Field, S. M., Hoekstra, R., Bringmann, L. & van Ravenzwaaij, D. When and why to
1005 replicate: As easy as 1, 2, 3? *Collabra Psychol.* **5**, 46 (2019).
- 1006 33. Isager, P. M. *et al.* Deciding what to replicate: A decision model for replication study
1007 selection under resource and knowledge constraints. *Psychol. Methods* **28**, 438–451
1008 (2023).
- 1009 34. O'Donnell, M. *et al.* Empirical audit and review and an assessment of evidentiary value in

- 1010 research on the psychological consequences of scarcity. *Proc. Natl. Acad. Sci.* **118**,
1011 e2103313118 (2021).
- 1012 35. Kuehberger, A. & Schulte-Mecklenbeck, M. Selecting target papers for replication. *Behav.*
1013 *Brain Sci.* **41**, e139 (2018).
- 1014 36. Hoogeveen, S., Sarafoglou, A. & Wagenmakers, E.-J. Laypeople can predict which social-
1015 science studies will be replicated successfully. *Adv. Methods Pract. Psychol. Sci.* **3**, 267–
1016 285 (2020).
- 1017 37. Marcoci, A. *et al.* Predicting the replicability of social and behavioural science claims from
1018 the COVID-19 Preprint Replication Project with structured expert and novice groups. *Nat.*
1019 *Hum. Behav.* (forthcoming).
- 1020 38. Wolfers, J. & Zitzewitz, E. Prediction markets. *J. Econ. Perspect.* **18**, 107–126 (2004).
- 1021 39. Arrow, K. J. *et al.* The promise of prediction markets. *Science* **320**, 877–878 (2008).
- 1022 40. Tziralis, G. & Tatsiopoulou, I. Prediction markets: An extended literature review. *J. Predict.*
1023 *Mark.* **1**, 75–91 (2012).
- 1024 41. Hanson, R. Decision markets. *IEEE Intell. Syst.* **14**, 16–19 (1999).
- 1025 42. Hanson, R. Combinatorial information market design. *Inf. Syst. Front.* **5**, 107–119 (2003).
- 1026 43. Chen, Y., Kash, I., Ruberry, M. & Shnayder, V. Decision markets with good incentives. in
1027 *Internet and Network Economics* (eds. Chen, N., Elkind, E. & Koutsoupias, E.) 72–83
1028 (Springer, Berlin, 2011). doi:10/b4vtjj.
- 1029 44. Wang, W. & Pfeiffer, T. Securities based decision markets. in *Distributed Artificial*
1030 *Intelligence* (eds. Chen, J., Lang, J., Amato, C. & Zhao, D.) vol. 13170 79–92 (Springer,
1031 Shanghai, China, 2022).
- 1032 45. Gordon, M., Viganola, D., Dreber, A., Johannesson, M. & Pfeiffer, T. Predicting
1033 replicability—Analysis of survey and prediction market data from large-scale forecasting
1034 projects. *PLoS One* **16**, e0248780 (2021).
- 1035 46. Ames, D. L. & Fiske, S. T. Perceived intent motivates people to magnify observed harms.
1036 *Proc. Natl. Acad. Sci.* **112**, 3599–3605 (2015).
- 1037 47. Atir, S. & Ferguson, M. J. How gender determines the way we speak about professionals.
1038 *Proc. Natl. Acad. Sci.* **115**, 7278–7283 (2018).
- 1039 48. Baldwin, M. & Lammers, J. Past-focused environmental comparisons promote
1040 proenvironmental outcomes for conservatives. *Proc. Natl. Acad. Sci.* **113**, 14953–14957
1041 (2016).
- 1042 49. Bear, A., Fortgang, R. G., Bronstein, M. V. & Cannon, T. D. Mistiming of thought and
1043 perception predicts delusionality. *Proc. Natl. Acad. Sci.* **114**, 10791–10796 (2017).
- 1044 50. Boswell, R. G., Sun, W., Suzuki, S. & Kober, H. Training in cognitive strategies reduces

- 1045 eating and improves food choice. *Proc. Natl. Acad. Sci.* **115**, E11238–E11247 (2018).
- 1046 51. Caruso, E. M., Burns, Z. C. & Converse, B. A. Slow motion increases perceived intent.
1047 *Proc. Natl. Acad. Sci.* **113**, 9250–9255 (2016).
- 1048 52. Casella, A., Kartik, N., Sanchez, L. & Turban, S. Communication in context: Interpreting
1049 promises in an experiment on competition and trust. *Proc. Natl. Acad. Sci.* **115**, 933–938
1050 (2018).
- 1051 53. Chao, M. Demotivating incentives and motivation crowding out in charitable giving. *Proc.*
1052 *Natl. Acad. Sci.* **114**, 7301–7306 (2017).
- 1053 54. Cheon, B. K. & Hong, Y.-Y. Mere experience of low subjective socioeconomic status
1054 stimulates appetite and food intake. *Proc. Natl. Acad. Sci.* **114**, 72–77 (2017).
- 1055 55. Clarkson, J. J. *et al.* The self-control consequences of political ideology. *Proc. Natl. Acad.*
1056 *Sci.* **112**, 8250–8253 (2015).
- 1057 56. Cooney, G., Gilbert, D. T. & Wilson, T. D. When fairness matters less than we expect.
1058 *Proc. Natl. Acad. Sci.* **113**, 11168–11171 (2016).
- 1059 57. Côté, S., House, J. & Willer, R. High economic inequality leads higher-income individuals
1060 to be less generous. *Proc. Natl. Acad. Sci.* **112**, 15838–15843 (2015).
- 1061 58. Flesch, T., Balaguer, J., Dekker, R., Nili, H. & Summerfield, C. Comparing continual task
1062 learning in minds and machines. *Proc. Natl. Acad. Sci.* **115**, E10313–E10322 (2018).
- 1063 59. Genschow, O., Rigoni, D. & Brass, M. Belief in free will affects causal attributions when
1064 judging others' behavior. *Proc. Natl. Acad. Sci.* **114**, 10071–10076 (2017).
- 1065 60. Gheorghiu, A. I., Callan, M. J. & Skylark, W. J. Facial appearance affects science
1066 communication. *Proc. Natl. Acad. Sci.* **114**, 5970–5975 (2017).
- 1067 61. Guilbeault, D., Becker, J. & Centola, D. Social learning and partisan bias in the
1068 interpretation of climate trends. *Proc. Natl. Acad. Sci.* **115**, 9714–9719 (2018).
- 1069 62. Halevy, N. & Halali, E. Selfish third parties act as peacemakers by transforming conflicts
1070 and promoting cooperation. *Proc. Natl. Acad. Sci.* **112**, 6937–6942 (2015).
- 1071 63. Handley, I. M., Brown, E. R., Moss-Racusin, C. A. & Smith, J. L. Quality of evidence
1072 revealing subtle gender biases in science is in the eye of the beholder. *Proc. Natl. Acad.*
1073 *Sci.* **112**, 13201–13206 (2015).
- 1074 64. Hoffman, K. M., Trawalter, S., Axt, J. R. & Oliver, M. N. Racial bias in pain assessment
1075 and treatment recommendations, and false beliefs about biological differences between
1076 blacks and whites. *Proc. Natl. Acad. Sci.* **113**, 4296–4301 (2016).
- 1077 65. Hofstetter, R., Ruppell, R. & John, L. K. Temporary sharing prompts unrestrained
1078 disclosures that leave lasting negative impressions. *Proc. Natl. Acad. Sci.* **114**, 11902–
1079 11907 (2017).

- 1080 66. Horne, Z., Powell, D., Hummel, J. E. & Holyoak, K. J. Countering antivaccination attitudes.
1081 *Proc. Natl. Acad. Sci.* **112**, 10321–10324 (2015).
- 1082 67. Isley, S. C., Stern, P. C., Carmichael, S. P., Joseph, K. M. & Arent, D. J. Online purchasing
1083 creates opportunities to lower the life cycle carbon footprints of consumer products. *Proc.*
1084 *Natl. Acad. Sci.* **113**, 9780–9785 (2016).
- 1085 68. Jachimowicz, J. M., Chafik, S., Munrat, S., Prabhu, J. C. & Weber, E. U. Community trust
1086 reduces myopic decisions of low-income individuals. *Proc. Natl. Acad. Sci.* **114**, 5401–
1087 5406 (2017).
- 1088 69. John, L. K., Barasz, K. & Norton, M. I. Hiding personal information reveals the worst. *Proc.*
1089 *Natl. Acad. Sci.* **113**, 954–959 (2016).
- 1090 70. Jordan, J. J., Hoffman, M., Nowak, M. A. & Rand, D. G. Uncalculating cooperation is used
1091 to signal trustworthiness. *Proc. Natl. Acad. Sci.* **113**, 8658–8663 (2016).
- 1092 71. Jun, Y., Meng, R. & Johar, G. V. Perceived social presence reduces fact-checking. *Proc.*
1093 *Natl. Acad. Sci.* **114**, 5976–5981 (2017).
- 1094 72. KC, R. P., Kunter, M. & Mak, V. The influence of a competition on noncompetitors. *Proc.*
1095 *Natl. Acad. Sci.* **115**, 2716–2721 (2018).
- 1096 73. Klein, N. & O'Brien, E. People use less information than they think to make up their minds.
1097 *Proc. Natl. Acad. Sci.* **115**, 13222–13227 (2018).
- 1098 74. Kouchaki, M. & Gino, F. Memories of unethical actions become obfuscated over time.
1099 *Proc. Natl. Acad. Sci.* **113**, 6166–6171 (2016).
- 1100 75. Kraus, M. W., Rucker, J. M. & Richeson, J. A. Americans misperceive racial economic
1101 equality. *Proc. Natl. Acad. Sci.* **114**, 10324–10331 (2017).
- 1102 76. McCall, L., Burk, D., Laperrière, M. & Richeson, J. A. Exposure to rising inequality shapes
1103 Americans' opportunity beliefs and policy support. *Proc. Natl. Acad. Sci.* **114**, 9593–9598
1104 (2017).
- 1105 77. Morris, A., MacGlashan, J., Littman, M. L. & Cushman, F. Evolution of flexibility and rigidity
1106 in retaliatory punishment. *Proc. Natl. Acad. Sci.* **114**, 10396–10401 (2017).
- 1107 78. Mummolo, J. Militarization fails to enhance police safety or reduce crime but may harm
1108 police reputation. *Proc. Natl. Acad. Sci.* **115**, 9181–9186 (2018).
- 1109 79. Payne, B. K., Brown-Iannuzzi, J. L. & Hannay, J. W. Economic inequality increases risk
1110 taking. *Proc. Natl. Acad. Sci.* **114**, 4643–4648 (2017).
- 1111 80. Phillips, J. & Cushman, F. Morality constrains the default representation of what is
1112 possible. *Proc. Natl. Acad. Sci.* **114**, 4649–4654 (2017).
- 1113 81. Rai, T. S., Valdesolo, P. & Graham, J. Dehumanization increases instrumental violence,
1114 but not moral violence. *Proc. Natl. Acad. Sci.* **114**, 8511–8516 (2017).

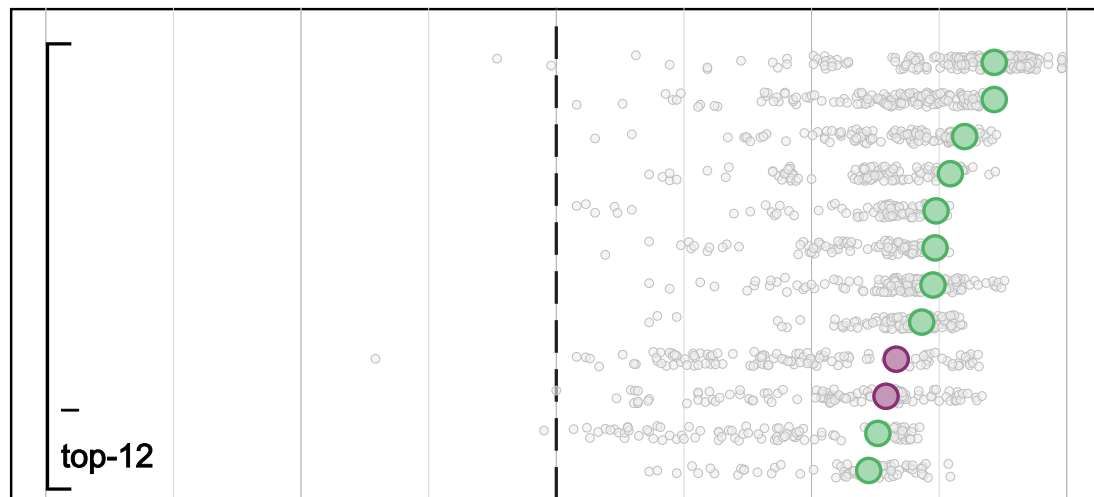
- 1115 82. Reeck, C., Wall, D. & Johnson, E. J. Search predicts and changes patience in
1116 intertemporal choice. *Proc. Natl. Acad. Sci.* **114**, 11890–11895 (2017).
- 1117 83. Schilke, O., Reimann, M. & Cook, K. S. Power decreases trust in social exchange. *Proc.*
1118 *Natl. Acad. Sci.* **112**, 12950–12955 (2015).
- 1119 84. Stern, C., West, T. V. & Rule, N. O. Conservatives negatively evaluate
1120 counterstereotypical people to maintain a sense of certainty. *Proc. Natl. Acad. Sci.* **112**,
1121 15337–15342 (2015).
- 1122 85. Vacharkulksemsuk, T. *et al.* Dominant, open nonverbal displays are attractive at zero-
1123 acquaintance. *Proc. Natl. Acad. Sci.* **113**, 4009–4014 (2016).
- 1124 86. Williams, K. E. G., Sng, O. & Neuberg, S. L. Ecology-driven stereotypes override race
1125 stereotypes. *Proc. Natl. Acad. Sci.* **113**, 310–315 (2016).
- 1126 87. Schimmelpfennig, R. *et al.* The moderating role of culture in the generalizability of
1127 psychological phenomena. *Adv. Methods Pract. Psychol. Sci.* **7**, 25152459231225163
1128 (2024).
- 1129 88. Zhou, H. & Fishbach, A. The pitfall of experimenting on the web: How unattended selective
1130 attrition leads to surprising (yet false) research conclusions. *J. Pers. Soc. Psychol.* **111**,
1131 493–504 (2016).
- 1132 89. Chmielewski, M. & Kucker, S. C. An MTurk crisis? Shifts in data quality and the impact on
1133 study results. *Soc. Psychol. Personal. Sci.* **11**, 464–473 (2020).
- 1134 90. Aguinis, H., Villamor, I. & Ramani, R. S. MTurk research: Review and recommendations.
1135 *J. Manag.* **47**, 823–837 (2021).
- 1136 91. Brodeur, A., Cook, N. & Heyes, A. We need to talk about Mechanical Turk: What 22,989
1137 hypothesis tests tell us about publication bias and p-hacking in online experiments.
1138 Preprint at <https://doi.org/10/nd9h> (2022).
- 1139 92. Peer, E., Rothschild, D., Gordon, A., Evernden, Z. & Damer, E. Data quality of platforms
1140 and panels for online behavioral research. *Behav. Res. Methods* **54**, 1643–1662 (2022).
- 1141 93. Webb, M. A. & Tangney, J. P. Too good to be true: Bots and bad data from Mechanical
1142 Turk. *Perspect. Psychol. Sci.* 174569162211200 (2022) doi:10/gq7nw4.
- 1143 94. Douglas, B. D., Ewell, P. J. & Brauer, M. Data quality in online human-subjects research:
1144 Comparisons between MTurk, Prolific, CloudResearch, Qualtrics, and SONA. *PLoS One*
1145 **18**, e0279720 (2023).
- 1146 95. Abelson, R. P. *Statistics as Principled Argument*. (Psychology Press, New York, 1995).
1147 doi:10/gf5svk.
- 1148 96. Macdonald, R. R. Statistical inference and Aristotle's Rhetoric. *Br. J. Math. Stat. Psychol.*
1149 **57**, 193–203 (2004).

- 1150 97. Scheel, A. M., Schijen, M. R. M. J. & Lakens, D. An excess of positive results: Comparing
1151 the standard psychology literature with registered reports. *Adv. Methods Pract. Psychol.*
1152 *Sci.* **4**, 25152459211007467 (2021).
- 1153 98. Soderberg, C. K. *et al.* Initial evidence of research quality of registered reports compared
1154 with the standard publishing model. *Nat. Hum. Behav.* **5**, 990–997 (2021).
- 1155 99. Brodeur, A., Cook, N., Hartley, J. & Heyes, A. Do pre-registration and pre-analysis plans
1156 reduce p-hacking and publication bias? Evidence from 15,992 test statistics and
1157 suggestions for improvement. *J. Polit. Econ. Microecon.* **2**, (2024).
- 1158 100. Yamada, Y. How to crack pre-registration: Toward transparent and open science. *Front.*
1159 *Psychol.* **9**, (2018).
- 1160 101. Flis, I. The function of literature in psychological science. *Rev. Gen. Psychol.* **26**, 146–156
1161 (2022).
- 1162 102. Rubin, M. Questionable metascience practices. *J. Trial Error* **4**, (2024).
- 1163 103. Wagenmakers, E.-J., Wetzels, R., Borsboom, D., van der Maas, H. L. J. & Kievit, R. A. An
1164 agenda for purely confirmatory research. *Perspect. Psychol. Sci.* **7**, 632–638 (2012).
- 1165 104. Nosek, B. A., Ebersole, C. R., DeHaven, A. C. & Mellor, D. T. The preregistration
1166 revolution. *Proc. Natl. Acad. Sci.* **115**, 2600–2606 (2018).
- 1167 105. Maxwell, S. E., Lau, M. Y. & Howard, G. S. Is psychology suffering from a replication
1168 crisis? What does “failure to replicate” really mean? *Am. Psychol.* **70**, 487–498 (2015).
- 1169 106. Shrout, P. E. & Rodgers, J. L. Psychology, science, and knowledge construction:
1170 Broadening perspectives from the replication Crisis. *Annu. Rev. Psychol.* **69**, 487–510
1171 (2018).
- 1172 107. Dreber, A. & Johannesson, M. A framework for evaluating reproducibility and replicability
1173 in economics. *Econ. Inq.* 1–19 (2024) doi:10/gt3vmw.
- 1174 108. Benjamin, D. J. *et al.* Redefine statistical significance. *Nat. Hum. Behav.* **2**, 6–10 (2018).
- 1175 109. Barnard, G. A. Significance tests for 2x2 tables. *Biometrika* **34**, 123–138 (1947).
- 1176 110. Mehrotra, D. V., Chan, I. S. F. & Berger, R. L. A cautionary note on exact unconditional
1177 inference for a difference between two independent binomial proportions. *Biometrics* **59**,
1178 441–450 (2003).
- 1179 111. Boschloo, R. D. Raised conditional level of significance for the 2x2-table when testing the
1180 equality of two probabilities. *Stat. Neerlandica* **24**, 1–9 (1970).
- 1181 112. Simonsohn, U. Small Telescopes: Detectability and the Evaluation of Replication Results.
1182 *Psychol. Sci.* **26**, 559–569 (2015).
- 1183 113. Ly, A., Verhagen, J. & Wagenmakers, E.-J. Harold Jeffreys’s default Bayes factor
1184 hypothesis tests: Explanation, extension, and application in psychology. *J. Math. Psychol.*

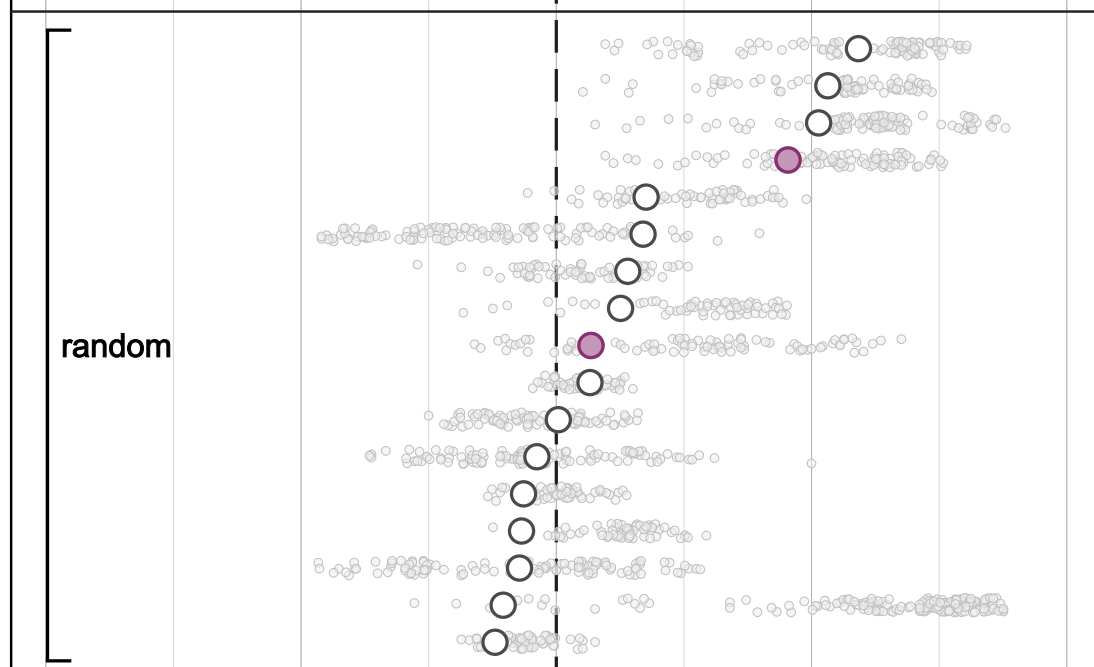
- 1185 **72**, 19–32 (2016).
- 1186 114. Ly, A., Etz, A., Marsman, M. & Wagenmakers, E.-J. Replication Bayes factors from
1187 evidence updating. *Behav. Res. Methods* **51**, 2498–2508 (2019).
- 1188 115. Jeffreys, H. *The Theory of Probability*. (Oxford University Press, Oxford, 1961).
- 1189 116. Patil, P., Peng, R. D. & Leek, J. T. What should researchers expect when they replicate
1190 studies? A statistical view of replicability in psychological science. *Perspect. Psychol. Sci.*
1191 **11**, 539–544 (2016).
- 1192 117. Gelman, A. & Stern, H. The difference between “significant” and “not significant” is not
1193 itself statistically significant. *Am. Stat.* **60**, 328–331 (2006).
- 1194 118. Cumming, G. Replication and p intervals: p values predict the future only vaguely, but
1195 confidence intervals do much better. *Perspect. Psychol. Sci.* **3**, 286–300 (2008).
- 1196 119. Muradchanian, J., Hoekstra, R., Kiers, H. & van Ravenzwaaij, D. How best to quantify
1197 replication success? A simulation study on the comparison of replication success metrics.
1198 *R. Soc. Open Sci.* **8**, 201697 (2021).
- 1199 120. Altmejd, A. *et al.* Predicting the replicability of social science lab experiments. *PLoS One*
1200 **14**, e0225826 (2019).
- 1201 121. Yang, Y., Youyou, W. & Uzzi, B. Estimating the deep replicability of scientific findings using
1202 human and artificial intelligence. *Proc. Natl. Acad. Sci.* **117**, 10762–10768 (2020).
- 1203 122. Rajtmajer, S. *et al.* A synthetic prediction market for estimating confidence in published
1204 work. *Proc. AAAI Conf. Artif. Intell.* **36**, 13218–13220 (2022).
- 1205 123. Youyou, W., Yang, Y. & Uzzi, B. A discipline-wide investigation of the replicability of
1206 psychology papers over the past two decades. *Proc. Natl. Acad. Sci.* **120**, e2208863120
1207 (2023).
- 1208 124. Agle, J., Xiao, Y., Nolan, R. & Golzarri-Arroyo, L. Quality control questions on Amazon’s
1209 Mechanical Turk (MTurk): A randomized trial of impact on the USAUDIT, PHQ-9, and
1210 GAD-7. *Behav. Res. Methods* **54**, 885–897 (2022).
- 1211 125. Veselovsky, V., Ribeiro, M. H. & West, R. Artificial artificial artificial intelligence: Crowd
1212 workers widely use large language models for text production tasks. Preprint at
1213 <https://doi.org/10/gsc3gr> (2023).
- 1214 126. Olsson-Collentine, A., Wicherts, J. M. & van Assen, M. A. L. M. Heterogeneity in direct
1215 replications in psychology and its association with effect size. *Psychol. Bull.* **146**, 922–940
1216 (2020).
- 1217 127. Linden, A. H. & Hönekopp, J. Heterogeneity of research results: A new perspective from
1218 which to assess and promote progress in psychological science. *Perspect. Psychol. Sci.*
1219 **16**, 358–376 (2021).

- 1220 128. Holzmeister, F. *et al.* Heterogeneity in effect size estimates. *Proc. Natl. Acad. Sci.* **121**,
1221 e2403490121 (2024).
- 1222 129. Epstein, R. & Robertson, R. E. The search engine manipulation effect (SEME) and its
1223 possible impact on the outcomes of elections. *Proc. Natl. Acad. Sci.* **112**, E4512–E4521
1224 (2015).
- 1225 130. Gallo, E. & Yan, C. The effects of reputational and social knowledge on cooperation. *Proc.*
1226 *Natl. Acad. Sci.* **112**, 3647–3652 (2015).
- 1227 131. Li, V., Michael, E., Balaguer, J., Herce Castañón, S. & Summerfield, C. Gain control
1228 explains the effect of distraction in human perceptual, cognitive, and economic decision
1229 making. *Proc. Natl. Acad. Sci.* **115**, E8825–E8834 (2018).
- 1230 132. Hanson, R. Logarithmic market scoring rules for modular combinatorial information
1231 aggregation. *J. Predict. Mark.* **1**, 3–15 (2007).
- 1232 133. Moshontz, H. *et al.* The Psychological Science Accelerator: Advancing psychology
1233 through a distributed collaborative network. *Adv. Methods Pract. Psychol. Sci.* **1**, 501–515
1234 (2018).
- 1235

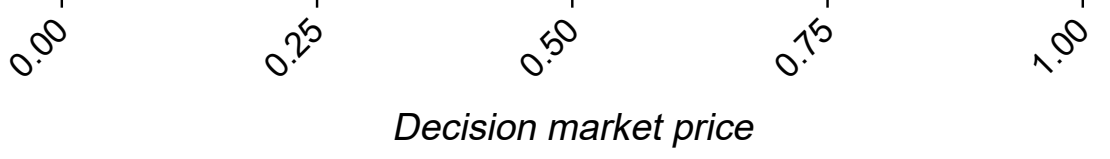
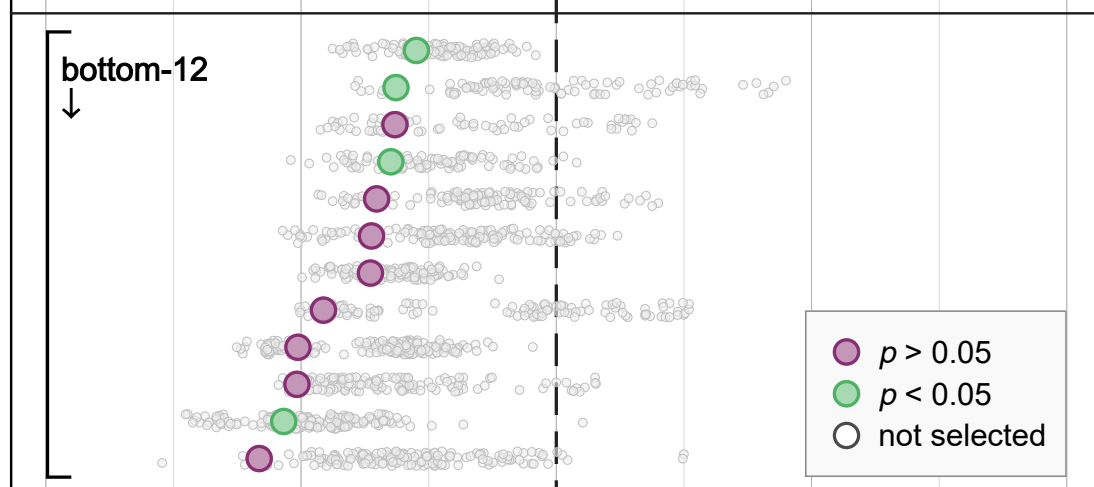
Vacharkulksemsuk et al. (2016)
 Cooney et al. (2016)
 Halevy and Halali (2015)
 Boswell et al. (2018)
 Hofstetter et al. (2017)
 Morris et al. (2017)
 Klein and O'Brien (2018)
 McCall et al. (2017)
 Ames and Fiske (2015)
 Stern et al. (2015)
 Williams et al. (2016)
 Kraus et al. (2017)



Casella et al. (2017)
 Schilke et al. (2015)
 Caruso et al. (2016)
 Jordan et al. (2016)
 Isley et al. (2016)
 Jachimowicz et al. (2017)
 Bear et al. (2017)
 Horne et al. (2015)
 Guilbeault et al. (2018)
 Flesch et al. (2018)
 Mummolo (2018)
 Handley et al. (2015)
 Jun et al. (2017)
 Payne et al. (2017)
 Clarkson et al. (2015)
 Phillips and Cushman (2017)
 KC et al. (2018)



Chao (2017)
 Rai et al. (2017)
 Reeck et al. (2017)
 Gheorghiu et al. (2017)
 John et al. (2015)
 Cote et al. (2015)
 Hoffman et al. (2016)
 Kouchaki and Gino (2016)
 Cheon and Hong (2016)
 Baldwin and Lammers (2016)
 Genschow et al. (2017)
 Atir and Ferguson (2018)



○ $d_O: p < 0.05$ ● $d_R: p > 0.05$ ● $d_R: p < 0.05$ — 95% CI

Cooney et al. (2016)
 Vacharkulksemsuk et al. (2016)
 Halevy and Halali (2015)
 Boswell et al. (2018)
 Hofstetter et al. (2017)
 Morris et al. (2017)
 Klein and O'Brien (2018)
 McCall et al. (2017)
 Ames and Fiske (2015)
 Stern et al. (2015)
 Williams et al. (2016)
 Kraus et al. (2017)

$n_O = 120, n_R = 132$
 $n_O = 426, n_R = 450$
 $n_O = 198, n_R = 227$
 $n_O = 260, n_R = 291$
 $n_O = 323, n_R = 820$
 $n_O = 100, n_R = 128$
 $n_O = 207, n_R = 214$
 $n_O = 480, n_R = 679$
 $n_O = 201, n_R = 723$
 $n_O = 273, n_R = 503$
 $n_O = 96, n_R = 112$
 $n_O = 202, n_R = 205$

top-12

Jordan et al. (2016)
 Guilbeault et al. (2018)

$n_O = 735, n_R = 1826$
 $n_O = 24, n_R = 56$

random

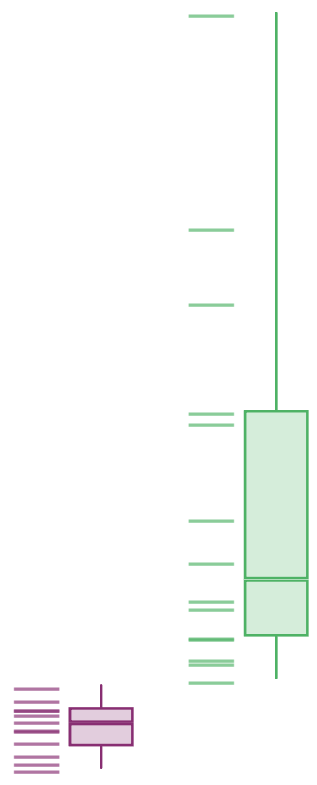
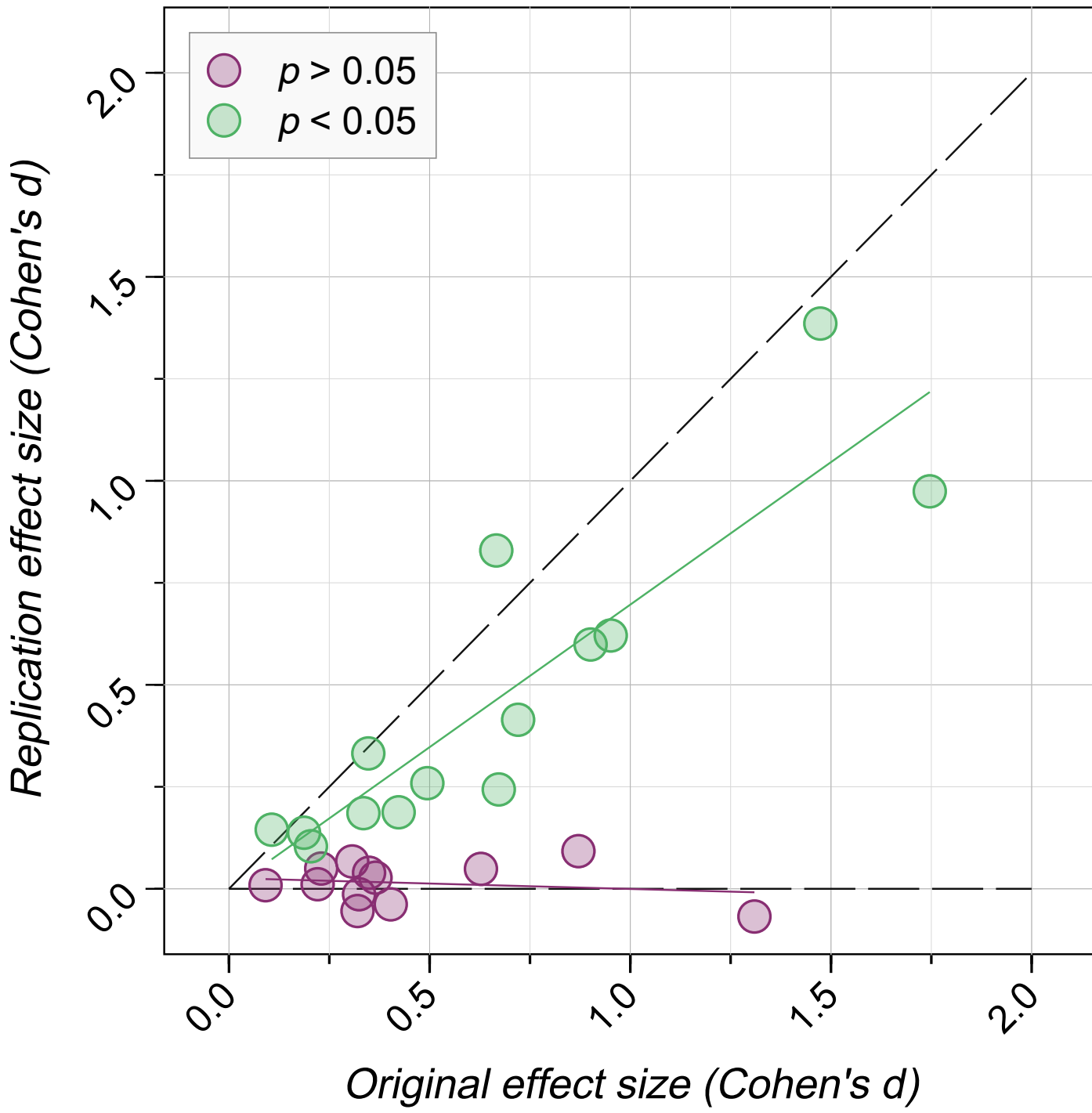
Chao (2017)
 Rai et al. (2017)
 Reeck et al. (2017)
 Gheorghiu et al. (2017)
 John et al. (2015)
 Cote et al. (2015)
 Hoffman et al. (2016)
 Kouchaki and Gino (2016)
 Cheon and Hong (2016)
 Baldwin and Lammers (2016)
 Genschow et al. (2017)
 Atir and Ferguson (2018)

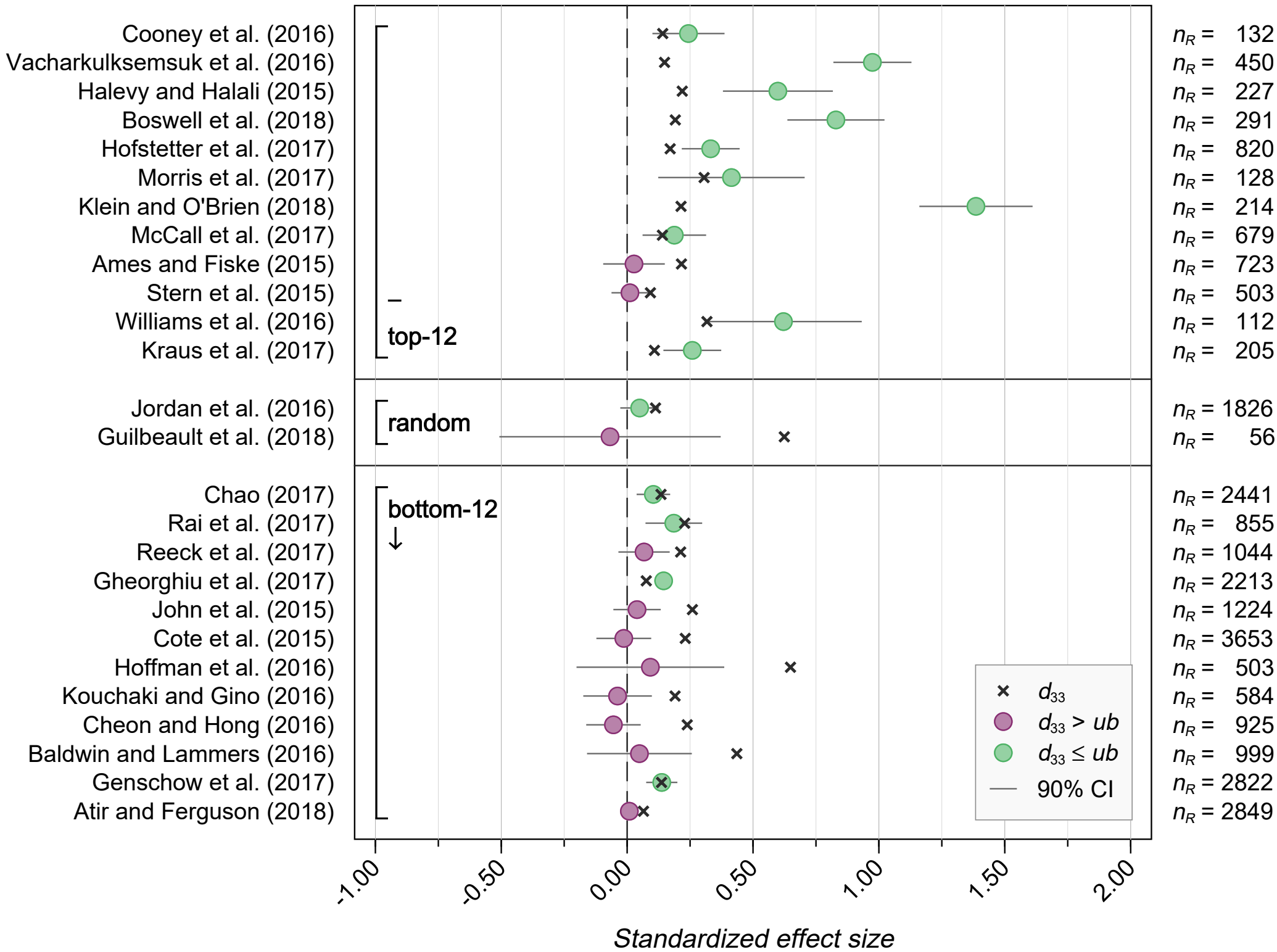
$n_O = 519, n_R = 2441$
 $n_O = 182, n_R = 855$
 $n_O = 207, n_R = 1044$
 $n_O = 408, n_R = 2213$
 $n_O = 142, n_R = 1224$
 $n_O = 704, n_R = 3653$
 $n_O = 92, n_R = 503$
 $n_O = 258, n_R = 584$
 $n_O = 167, n_R = 925$
 $n_O = 200, n_R = 999$
 $n_O = 504, n_R = 2822$
 $n_O = 554, n_R = 2849$

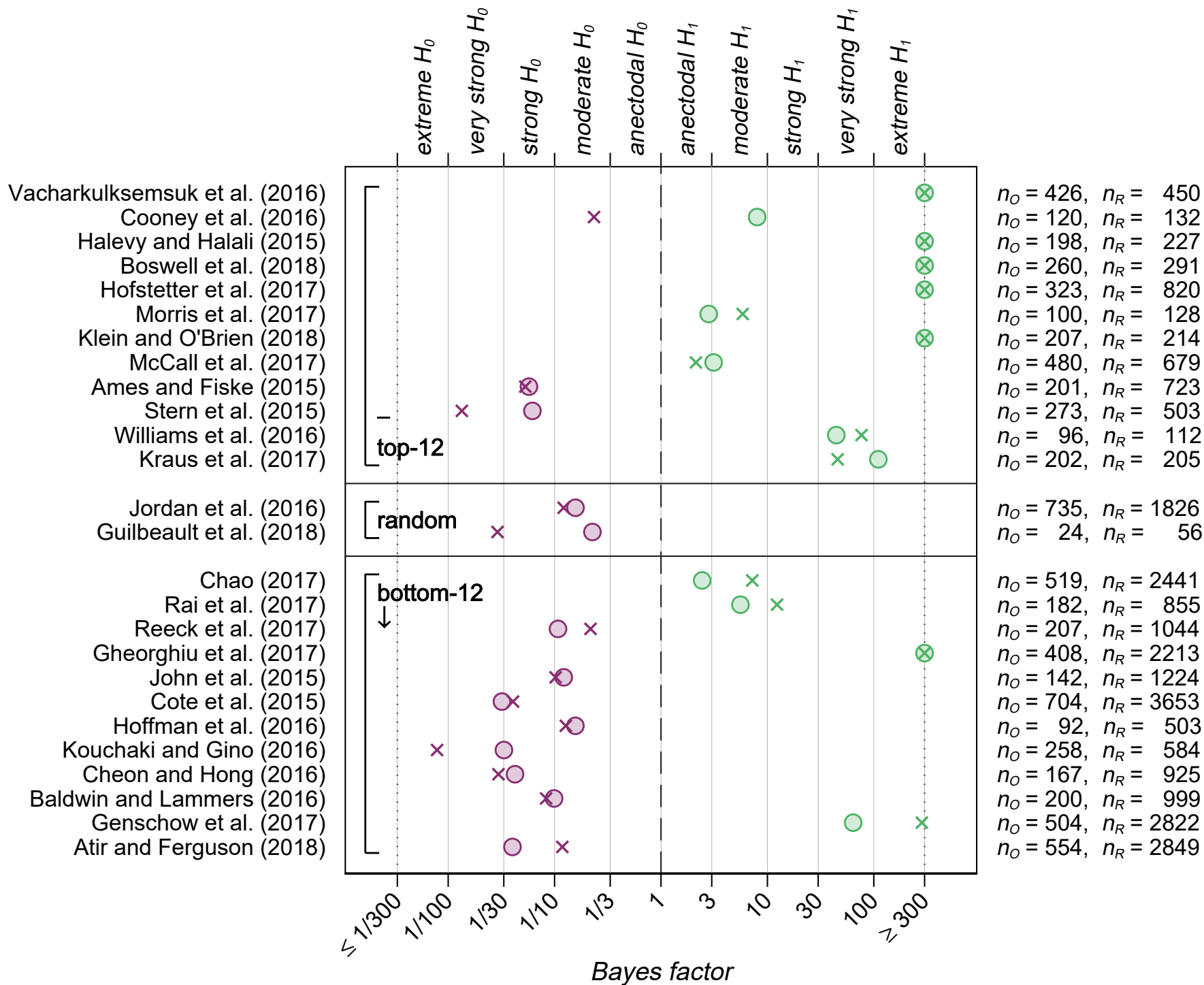
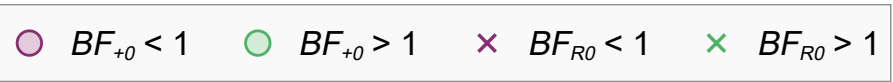
bottom-12

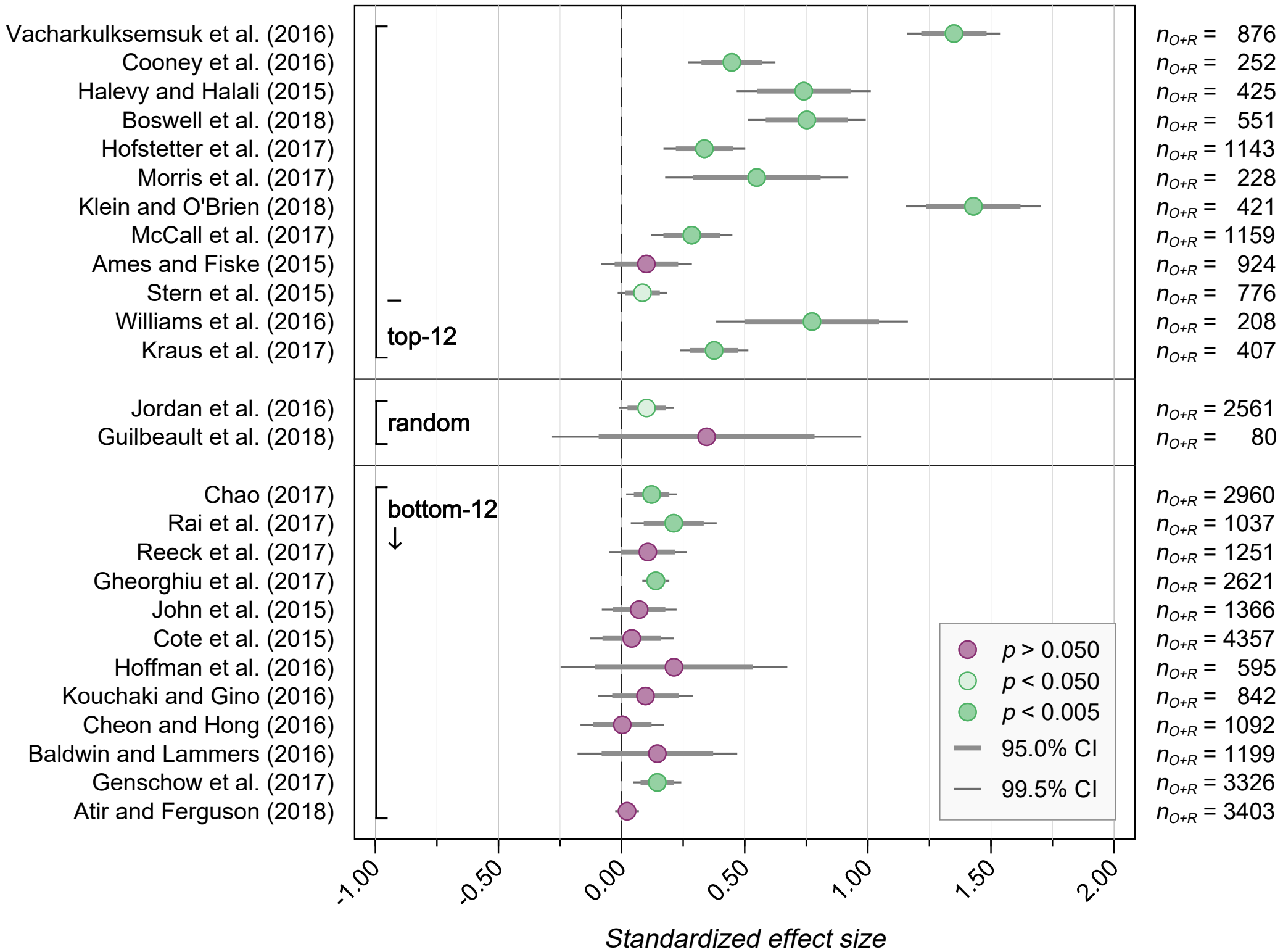
-1.00 -0.50 0.00 0.50 1.00 1.50 2.00

Standardized effect size









○ $d_R \notin 95\% \text{ PI}$ ○ $d_R \in 95\% \text{ PI}$ — 95% PI

