

# **Evaluating the Replicability of Social Science Experiments in *Nature* and *Science***

**Colin F. Camerer<sup>1,†</sup>, Anna Dreber<sup>2,†</sup>, Felix Holzmeister<sup>3,†</sup>, Teck-Hua Ho<sup>4,†</sup>, Jürgen Huber<sup>3,†</sup>, Magnus Johannesson<sup>2,†</sup>, Michael Kirchler<sup>3,5,†</sup>, Gideon Nave<sup>6,†</sup>, Brian Nosek<sup>7,8,\*†</sup>, Thomas Pfeiffer<sup>9,†</sup>, Adam Altmejd<sup>2</sup>, Nick Buttrick<sup>7,8</sup>, Taizan Chan<sup>10</sup>, Yiling Chen<sup>11</sup>, Eskil Forsell<sup>12</sup>, Anup Gampa<sup>7,8</sup>, Emma Heikensten<sup>2</sup>, Lily Hummer<sup>8</sup>, Taisuke Imai<sup>13</sup>, Siri Isaksson<sup>2</sup>, Dylan Manfredi<sup>6</sup>, Julia Rose<sup>3</sup>, Eric-Jan Wagenmakers<sup>14</sup>, Hang Wu<sup>15</sup>**

<sup>1</sup> California Institute of Technology, 1200 E California Blvd, MC 228-77, Pasadena, CA 91125, USA

<sup>2</sup> Department of Economics, Stockholm School of Economics, Box 6501, SE-113 83 Stockholm, Sweden

<sup>3</sup> Department of Banking and Finance, University of Innsbruck, Universitätsstraße 15, 6020 Innsbruck, Austria

<sup>4</sup> NUS Business School, National University of Singapore, Singapore 119245

<sup>5</sup> Centre for Finance, Department of Economics, University of Göteborg, SE-40530 Göteborg, Sweden

<sup>6</sup> The Wharton School, University of Pennsylvania, 3730 Walnut Street, Philadelphia PA, 19104, USA

<sup>7</sup> Department of Psychology, University of Virginia, Charlottesville, VA 22904, USA

<sup>8</sup> Center for Open Science, Charlottesville, VA 22903, USA

<sup>9</sup> New Zealand Institute for Advanced Study, Private Bag 102904, North Shore Mail Centre, Auckland 0745, New Zealand

<sup>10</sup> Office of the Senior Deputy President and Provost, National University of Singapore, Singapore 119077

<sup>11</sup> John A. Paulsson School of Engineering and Applied Sciences, Harvard University, Cambridge MA 02138, USA

<sup>12</sup> Spotify Sweden AB, Birger Jarlsgatan 61, SE-113 56 Stockholm, Sweden

<sup>13</sup> LMU Munich, Department of Economics, Ludwigstraße 28, D-80539 Munich, Germany

<sup>14</sup> Department of Psychology, University of Amsterdam, Amsterdam 1018 VZ, The Netherlands

<sup>15</sup> School of Management, Harbin Institute of Technology, Harbin, 150001, China

\* To whom correspondence should be addressed. E-mail: [nosek@cos.io](mailto:nosek@cos.io).

† These first ten authors contributed equally to this work.

**Being able to replicate scientific findings is crucial for scientific progress<sup>1-15</sup>. We replicate 21 systematically selected experimental studies in the social sciences published in *Nature* and *Science* between 2010 and 2015. The replications follow analysis plans reviewed by the original authors and pre-registered prior to the replications. The replications are high powered with sample sizes on average about 5 times higher than in the original studies. We find a significant effect in the same direction as the original study for 13 (62%) studies, and the effect size of the replications is on average about 50% of the original effect size. Replicability varies between 12 (57%) and 14 (67%) for complementary replicability indicators. Consistent with these results, the estimated true positive rate is 67% in a Bayesian analysis. The relative effect size of true positives is estimated to be 71% suggesting that both false positives and inflated effect sizes of true positives contribute to imperfect reproducibility. We furthermore find that peer beliefs of replicability are strongly related to replicability, suggesting that the research community could predict which results would replicate and that failures to replicate were not the result of chance alone.**

To what extent can we trust scientific findings? The answer to this question is of fundamental importance<sup>1-3</sup>, and the reproducibility of published studies has been questioned in many fields<sup>4-10</sup>. Until recently, systematic evidence has been scarce<sup>11-15</sup>. The Reproducibility Project: Psychology<sup>12</sup> (*RPP*) put the question of scientific reproducibility at the forefront of scientific debate<sup>16-18</sup>. The *RPP* replicated 100 original studies in psychology, and found a significant effect in the same direction as the original for 36% of the 97 studies reporting “positive findings”<sup>12</sup>. The *RPP* was followed by the Experimental Economics Replication Project (*EERP*) which replicated

18 laboratory experiments in economics and found a significant effect in the same direction as the original studies for 61% of replications<sup>13</sup>. Both the *RPP* and the *EERP* had high statistical power to detect the effect sizes observed in the original studies. However, the effect sizes of published studies may be inflated even for true positive findings due to publication or reporting biases<sup>19-21</sup>. As a consequence, if replications were well-powered to detect effect sizes *smaller* than those observed in the original studies, replication rates might be higher than those estimated in the *RPP* and the *EERP*.

We provide evidence about the replicability of experimental studies in the social sciences published in the two most prestigious general science journals, *Nature* and *Science* (the Social Sciences Replication Project; *SSRP*). Articles published in these journals are considered exciting, innovative, and important. We include all experimental studies published between 2010 and 2015 that (i) test for an experimental treatment effect between or within subjects, (ii) test at least one clear hypothesis with a statistically significant finding, and (iii) were performed on students or other accessible subject pools. Twenty-one studies were identified to meet these criteria. We used the following three criteria in descending order to determine which treatment effect to replicate within each of these 21 papers: (a) select the first study reporting a significant treatment effect for papers reporting more than one study, (b) from that study, select the statistically significant result identified in the original study as the most important result among all within and between subjects treatment comparisons, and (c) if there was more than one equally central result, randomly select one of them for replication. The interpretation of which was the most central and important statistically significant result within a study in criteria (b) above was made by us, and not by the original authors. (See Supplementary Information, section 1 and Tables S1-S2 for details.)

To address the possibility of inflated effect sizes in the original studies, we used a high-powered design and a two-stage procedure for conducting the replications. In Stage 1 we had 90% power to detect 75% of the original effect size at the 5% significance level in a two-sided test. If the original result replicated in Stage 1 (a two-sided  $p$ -value  $< 0.05$  and an effect in the same direction as in the original study), no further data collection was carried out. If the original result did not replicate in Stage 1, we carried out a second data collection in Stage 2 to have 90% power to detect 50% of the original effect size for the first and second data collections pooled.

The motivation for having 90% power to detect 50% of the original effect size was based on the replication effect sizes in the *RPP* being on average about 50% of the original effect sizes (12) (see Supplementary Information, section 1, for details; the average relative effect size of the replications in the *EERP* was 66%<sup>13</sup>). On average, replication sample sizes in Stage 1 were about three times as large as the original sample sizes and replication sample sizes in Stage 2 were about six times as large as the original sample sizes. All of the replication and analysis plans were made publicly known on the project website, pre-registered at the OSF and sent to the original authors for feedback and verification prior to data collection (see Supplementary Information, section 1, for details and see individual replication reports for methodological details and reporting of any deviations in the protocols from the original studies).

There is no universally agreed upon criterion for replication<sup>12, 22-25</sup>, but our power analysis strategy is based on detecting a significant effect in the same direction as the original study using the same statistical test. As such, we treat this as the primary indicator of replication and refer to it as the statistical significance criterion. This approach is appealing for its simplicity as a binary measure of replication, but does not

fully represent evidence of reproducibility. We also provide results for the relative effect size of the replication as a continuous measure of the degree of replication. To complement these indicators, we present results for: (i) a meta-analytic estimate of the original and replication results combined<sup>12</sup>, (ii) 95% prediction intervals<sup>26</sup>, (iii) the “Small Telescopes” approach<sup>25</sup>, (iv) the one-sided default Bayes factor<sup>27</sup>, (v) a Bayesian mixture model<sup>28</sup>, (vi) and peer beliefs about replicability<sup>29</sup>. See Supplementary Information, section 2, and Fig. S1–S3 for additional robustness tests of the replication results.

In Stage 1 we find a significant effect in the same direction as the original study for 12 replications (57.1%) (Fig. 1a and Table S3). When we increase the statistical power further in Stage 2 (Fig. 1b and Table S4), 2 additional studies replicate based on this criterion. By mistake, a second data collection was carried out for one study replicating in Stage 1, and we therefore also include this study in the Stage 2 results to base our results on all the data collected. This study does not replicate in Stage 2. This may suggest that replication studies should routinely be powered to detect at least 50% of the original effect size, or that one should use a lower  $p$ -value threshold than 0.05 for not continuing to Stage 2 in our two-stage testing procedure. Based on all data collected, 13 (61.9%) studies replicate after Stage 2 using the statistical significance criterion.

The mean standardized effect size (correlation coefficient,  $r$ ) of the replications is 0.249, compared to 0.460 in the original studies (Fig. S4). This difference is significant (Wilcoxon signed-ranks test,  $z = 3.667$ ,  $p < 0.001$ ,  $n = 21$ ), and the mean *relative* effect size of the replications is 46.2%. For the 13 studies that replicated, the mean relative effect size is 74.5%, and for the 8 studies that did not replicate, the mean

relative effect size is 0.3%. It is not surprising that the mean relative effect size is smaller for the non-replicating effects than replicating effects as these are correlated indicators. However, it is notable that, even among the replicating effects, effect sizes for the replications were weaker than the original findings and, for the non-replicating effects, the mean effect sizes were approximately zero.

We also combined the original result and the replication in a meta-analytic estimate of the effect size. As seen in Fig. 1c, 16 studies (76.2%) have a significant effect in the same direction as the original study in the meta-analysis. However, the meta-analysis assume that the results of the original studies are not influenced by publication or reporting biases making the meta-analytic results an overly optimistic indicator compared to criteria focused on the replication evidence<sup>12</sup>. A team recently suggested that the  $p$ -value threshold for statistically significant findings should be lowered from 0.05 to 0.005 for new discoveries<sup>30</sup>. In a replication context it would be relevant to apply this stricter threshold to meta-analytic results. In this case, the meta-analysis leads to the same conclusions about replication as our primary replication indicator (i.e., 13 studies or 61.9% have a  $p$ -value  $< 0.005$  in the meta-analysis). It is obvious that the 13 successful replications would achieve  $p < 0.005$  when the original and replication results were pooled, but this criterion could have also included replications that did not achieve  $p < 0.05$  but were in the right direction and were combined with an original study with particularly strong evidence.

A complementary replication criterion is to count how many replicated effects lie in a 95% prediction interval (26), which takes into account the variability in both the original study and the replication study. Using this method, 14 effects replicate (66.7%; see Fig. 2a and Supplementary Information, section 2, for details). This method yields

the same replication outcome as the statistical significance criterion for 20 of the 21 studies.

The Small Telescopes approach estimates whether the replication effect size is significantly smaller than a “small effect” in the original study with a one-sided test at the 5% level. A small effect is defined as the effect size the original study would have had 33% power to detect. Following the Small Telescopes approach<sup>25</sup> 12 studies (57.1%) replicate (see Fig. 2b and Supplementary Information, section 2, for details). One replication has a significant effect in the same direction as the original study, but the effect size is significantly smaller than a small effect as defined by the Small Telescopes approach. This is the only difference compared to the statistical significance criterion.

Another way to represent the strength of evidence in favor of the original result versus the null hypothesis of no effect is to estimate the Bayes factor<sup>24, 27, 31-32</sup>. The Bayes factor compares the predictive performance of the null hypothesis against that of an alternative hypothesis in which the uncertainty about the true effect size is quantified by a prior distribution. The prior distributions were first set to their generic defaults; they were then folded across the test value so that all prior mass is consistent with the direction of the effect from the original study, thereby implementing a Bayesian one-sided test (see the Supplementary Information, section 2, for details). For example, the replication of Pyc and Rawson yielded a one-sided default Bayes factor of  $BF_{+0} = 6.8$ , meaning that the one-sided alternative hypothesis outpredicted the null hypothesis of no effect by a factor of almost 7.

The one-sided default Bayes factor exceeds 1, providing evidence in favor of an effect in the direction of the original study for the 13 (61.9%) studies that replicated

according to our primary replication indicator (Fig. 3). This evidence is strong to extreme for 9 (42.9%) studies. The default Bayes factor is below 1 for 8 (38.1%) studies providing evidence in support of the null hypothesis; this evidence is strong to extreme for 4 (19.0%) studies.

In additional Bayesian analyses, we use an errors-in-variables mixture model (28) to estimate the true positive rate in the total sample (see the Supplementary Information, section 2 and Fig. S5 for details). The estimated true positive rate is 67% (Fig. S5), which is close to the other replicability estimates. The mixture model also estimates that the average relative effect size of true positives is 71% (Fig. S5) suggesting that the original studies overestimated the effect sizes of true positives.

We also estimate peer beliefs about replicability using surveys and prediction markets<sup>29,33</sup> (see Supplementary Information, section 3 and Table S5 and Fig. S6 for details). The prediction markets produce a collective peer estimate of the probability of replication that can be interpreted as a reproducibility indicator<sup>29</sup>. The average prediction market belief of replicating after Stage 2 is a replication rate of 63.4% and the average survey belief is 60.6%, both close to the observed replication rate of 61.9% (Fig. 4; see Supplementary Information, section 4, Fig. S7–S8 and Tables S5–S6 for more details). The prediction market beliefs and the survey beliefs are highly correlated, and both are highly correlated with a successful replication (Fig. 4 and Fig. S7). That is, in the aggregate, peers were very effective at predicting future replication success.

In the *RPP*<sup>12</sup> and the *EERP*<sup>13</sup>, replication success was negatively correlated with the *p*-value of the original study, suggesting that original study *p*-values might be a predictor of replicability. We also find a negative correlation between the *p*-value of the original study and replication success, although it is not significant (Spearman



correlation coefficient  $-0.405$ ,  $p = 0.069$ ,  $n = 21$ ); the estimate is in between the correlations found in the *RPP* ( $-0.327$ ) and the *EERP* ( $-0.572$ ) (Table S7). That peers are to some extent able to predict which studies are most likely to replicate suggests that there are features of the original studies that journals or researchers can use in determining *ex ante* whether a study is likely to replicate. The results from the *RPP*, *EERP*, and *SSRP* taken together suggest that the *p*-value of the original study is one such important determinant of replication. The *SSRP* with  $n = 21$  studies is too small to reliably test determinants of replications, but pooling the results of all large scale replication projects may offer a higher powered opportunity to explore moderators of replication.

To summarize, we successfully replicated 13 of 21 findings from experimental social and behavioral science studies published in *Science* or *Nature* between 2010 and 2015 based on the statistical significance criterion with very high-powered studies compared to the *RPP*<sup>12</sup> and the *EERP*<sup>13</sup>. This number is larger than the *RPP*'s replication rate and similar to the *EERP*'s replication rate (Fig. S9). However, the small sample of studies and different selection criteria makes it difficult to draw any interpretation confidently in comparison with those studies. We can conclude, however, that increasing power substantially is not sufficient to reproduce all published studies. Also, we observe that the conclusions across binary replication criteria converge with increased statistical power. The Small Telescopes and the 95% prediction interval indicators drew different conclusions on only one of the replications compared to the statistical significance criterion.

Considering statistical significance and effect sizes simultaneously, we observe two major outcomes. First, even among successful replications, estimated effect sizes

were smaller than the original study. For the 13 studies that replicated according to the statistical significance criterion the replication effect sizes were about 75% of the original effect size. This provides an estimate of the overestimation of effect sizes of true positives in original studies. The Bayesian mixture model corroborates this result yielding an estimate of the relative effect size of true positives of 71%. This implies that meta-analyses of true positive findings will overestimate effect sizes on average. This finding bolsters evidence that the existing literature contains exaggerated effect sizes because of pervasive low powered research coupled with bias selecting significant results for publication.<sup>8,12</sup> Also, if this finding generalizes to the literatures investigated by the *RPP* and the *EERP*, it suggests that the statistical power of these two projects, where the sample sizes were determined to obtain 90% power to detect the original effect size, was de-facto smaller than intended. This would imply that the replication rates, based on the statistical significance criterion, were underestimated in these studies consistent with those authors' speculation.

Second, among the unsuccessful replications, there was essentially no evidence for the original finding. The average relative effect size was very close to zero for the eight findings that failed to replicate according to the statistical significance criterion. The expected relative effect size for a sample of false positives is zero, but this observation does not demand the conclusion that the eight original findings were false positives. Another possibility is that the replication studies failed to implement necessary features of the protocol to detect the effect<sup>17</sup>. We cannot rule out this alternative, but we also do not have evidence for necessary features missing from the replications that would reduce the observed effect sizes to zero. Indeed, it would be surprising but interesting to identify an unintended difference that completely

eliminated the effect rather than just reduce the effect size. One suggested indicator for whether differences between studies are a likely cause for bias, is the endorsement of the original authors.<sup>17</sup> In the current project, we took extensive efforts to ensure that the replications would be as close as possible to the originals. All of the replications but one were designed with the collaboration of the original authors (for one replication the original authors did not respond to our queries). And, all of the reviewed replications but one were approved by the original authors. However, none of this implies that original authors agree with the final outcomes or interpretation. For example, changes in planned implementation or insights after observing the results could lead to different interpretations of the replication outcome and ideas for subsequent research to clarify the understanding of the phenomenon. See Supplementary Information, section 1 and the posted replication reports for each study for more details including follow-up comments from original authors if provided.

Another hypothesis that could account for replication failures, at least partly, is the result of chance, such as a large degree of heterogeneity in treatment effects in different samples<sup>17</sup>. However, such heterogeneity would not affect the average relative effect size of replications, as replications would be as likely to over- as underestimate original effect sizes. It cannot therefore explain why the average effect sizes of our replications is only about 50% of the original effect sizes. Furthermore, the strong correlation between the peer predictions and the observed replicability is discordant with the possibility that replication failures occurred by chance alone. That is, researchers appear to have identified a priori systematic differences between the studies that replicated and those that did not. This capacity to predict the replicability of effects is a reason for optimism that methods will emerge to anticipate reproducibility

challenges and guide efficient use of replication resources toward exciting but uncertain findings.

The observed replication rate of 62%, based on the statistical significance criterion, adds to a growing pool of replicability rates from a variety of systematic replication efforts with distinct selection and design criteria: the *RPP*<sup>12</sup> (36%,  $n = 100$  studies), the *EERP*<sup>13</sup> (61%,  $n = 18$  studies), Many Labs 1<sup>11</sup> (77%;  $n = 13$  studies), Many Labs 2<sup>15</sup> (50%,  $n = 28$  studies), and Many Labs 3<sup>14</sup> (30%,  $n = 10$  studies). It is too early to draw a specific conclusion about the reproducibility rates of experimental studies in the social and behavioral sciences. Each investigation has a relatively small sample of studies with idiosyncratic inclusion criteria and unknown generalizability. But, the diversity in approaches provides some confidence that considering them in the aggregate may provide more general insight about reproducibility in the social-behavioral sciences. As a descriptive and speculative interpretation of these findings in the aggregate, we believe that reasonable lower and upper bound estimates are 35% and 75% for an average reproducibility rate of published findings in social and behavioral sciences. Accumulating additional evidence will reveal if there are systematic biases in these reproducibility estimates themselves.

When assessing reproducibility we are interested in both the systematic bias in the estimated effect sizes of original studies and the fraction of original hypotheses that are directionally true. The average relative effect size of 50% in the *SSRP* is a direct estimate of the systematic bias in the published findings of the 21 studies, as it should be 100% if original studies provide unbiased estimates of true effect sizes. This estimate assumes that there is no systematic difference in the effectiveness of implementing the study procedures or the appropriateness of testing circumstances between original and

replication studies. If both of those assumptions are true, then our data suggests that the systematic bias is partly due to false positives and partly due to the overestimated effect sizes of true positives. These systematic biases can be reduced by implementation of pre-registration of analysis plans to reduce the likelihood of false positives, and registration and reporting of all study results to reduce the effects of publication bias inflating effect sizes<sup>34</sup>. With notable progress on these practices, particularly in the social and behavioral sciences<sup>35</sup>, we predict that replicability will improve over time.

### **Limitations**

The *SSRP* is a small sample of studies with specific selection criteria for experimental studies from two high-profile journals. Work that is published in *Nature* and *Science* may be atypical to the field as a whole, and may have a stronger focus on novelty, which may also lead to greater – or lesser – editorial scrutiny. The small sample and selective criteria significantly reduce confidence in generalizing these findings to the social science literature more generally. Indeed, like all other research, replications require an accumulation of evidence across multiple efforts to identify and address sampling biases and to obtain increasingly precise estimates of replicability. This study adds to this accumulating literature with a focused, high-powered investigation of high-profile studies published in *Nature* and *Science*. Notably, with replication sample sizes about five times larger as the original studies, we get relatively precise estimates of the individual effects of these single replications and average relative effect sizes that are very similar to what was observed in *RPP*.

Another important limitation is that for papers reporting a series of studies we only replicate one of those studies, and for studies testing more than one hypothesis we only replicate one hypothesis. Like prior large-scale replication projects, this study does

not provide definitive insight on any of the original papers from which we designed replication studies. An alternative methodology would be to replicate all results within the selected study or all results within all studies in a paper reporting a series of studies. This would give more information from each replication and a more precise estimate of reproducibility of each study and paper. All investigations involve tradeoffs. The advantage of an in-depth examination of a paper is greater insight and precision of the reproducibility of its findings. The disadvantage is that many fewer findings can be investigated to learn about reproducibility of findings more generally. Some other findings reported in the original papers can be tested with the data available in our study's replications. We did not consider those secondary findings in this paper or in deciding the statistical power plans for the design. However, all of our data and materials are publicly posted as part of open science and will be available to other researchers who may want to pursue this issue further in follow-up work.

The original authors in reviewing our paper and replication results have noted some limitations on the replications of their individual studies. These are discussed more in the Supporting Information (Section 1.2); and several of the original authors have also posted comments on the replications at OSF alongside our Replication Reports. For example, previously unidentified or inadvertent changes to the protocol may have affected replication success for some studies. Also, for papers reporting a series of studies we replicated the first study reporting a significant treatment effect. In some cases the original authors argue that other studies in their papers report more important results or use stronger research designs.<sup>46,54</sup> If the replicability of the first study systematically differ from the replicability of subsequent studies in a paper our

criteria for deciding which study to replicate will systematically over- or under-estimate replicability.

Inspired by our replication, the original authors of Shah et al.<sup>54</sup> decided to carry out a replication study of their own on all of their five studies (with results posted at <https://osf.io/vzm23/>). They did replicate what they consider to be their most important finding, that scarcity itself leads to over-borrowing. They also failed to replicate study 1 in their paper consistent with our findings. Their approach of conducting replications of their own studies is admirable and provides additional insight and precision for understanding those effects.

Five of our replications were carried out on Amazon Mechanical Turk (AMT) and for one of those (Rand et al.<sup>53</sup>), the original authors argue that increasing familiarity with economic game paradigms among AMT samples may have decreased the replicability of their result. It cannot be ruled out that changes in the AMT subject pool over time have affected our results, but we also note that the two other studies based on economic game paradigms and AMT data replicated successfully<sup>43,50</sup>. It would be interesting in future work to test if replicability differ for older versus newer studies or depends on the time that has elapsed between the original study and the replication.

For the Sparrow et al.<sup>55</sup> replication, the original authors did not provide us with materials for the replication or feedback on our inquiries. This made it more difficult to replicate the experimental design of the original study. After the replication had been completed the original authors noted some design differences compared to the original study. These design differences are discussed further in the Supplementary Information,

and we cannot rule out that they impacted the replication result. This illustrates the importance of open access to all the materials of published studies for conducting direct replications and accumulating scientific knowledge.

## References

1. McNutt, M. Reproducibility. *Science* **343**, 229–229 (2014). DOI: [10.1126/science.1250475](https://doi.org/10.1126/science.1250475)
2. Baker, M. Is there a reproducibility crisis? *Nature* **533**, 452–454 (2016).
3. Munafò, M. R. et al. A manifesto for reproducible science. *Nat. Hum. Behav.* **1**, 0021 (2017). DOI: [10.1038/s41562-016-0021](https://doi.org/10.1038/s41562-016-0021)
4. Ioannidis, J. P. A. Why most published research findings are false. *PLOS Med.* **2**, e124 (2005). DOI: [10.1371/journal.pmed.0020124](https://doi.org/10.1371/journal.pmed.0020124)
5. Prinz, F., Schlange, T. & Asadullah, K. Believe it or not: how much can we rely on published data on potential drug targets? *Nat. Rev. Drug Discov.* **10**, 712–712 (2011). DOI: [10.1038/nrd3439-c1](https://doi.org/10.1038/nrd3439-c1)
6. Begley, C. G. & Ellis, L.M. Drug development: raise standards for preclinical cancer research. *Nature* **483**, 531–33 (2012). DOI: [10.1038/483531a](https://doi.org/10.1038/483531a)
7. Lawrence, M. S. et al. Mutational heterogeneity in cancer and the search for new cancer-associated genes. *Nature* **499**, 214–18 (2013). DOI: [10.1038/nature12213](https://doi.org/10.1038/nature12213)
8. Button, K. S. et al. Power failure: why small sample size undermines the reliability of neuroscience.” *Nat. Rev. Neurosci.* **14**, 365–76 (2013). DOI: [10.1038/nrn3475](https://doi.org/10.1038/nrn3475)
9. Maniadis, Z., Tufano, F. & List, J. A. One swallow doesn’t make a summer: new evidence on anchoring effects. *Am. Econ. Rev.* **104**, 277–290 (2014). DOI: [10.1257/aer.104.1.277](https://doi.org/10.1257/aer.104.1.277)
10. Freedman, L. P., Cockburn, I. M. & Simcoe, T. S. The economics of reproducibility in preclinical research. *PLOS Biol.* **13**, e1002165 (2015). DOI: [10.1371/journal.pbio.1002165](https://doi.org/10.1371/journal.pbio.1002165)



11. Klein, R. A. et al. Investigating variation in replicability: a ‘many labs’ replication project.” *Soc. Psychol.* **45**, 142–152 (2014). DOI: [10.1027/1864-9335/a000178](https://doi.org/10.1027/1864-9335/a000178)
12. Open Science Collaboration. Estimating the reproducibility of psychological science. *Science* **349**, aac4716 (2015). DOI: [10.1126/science.aac4716](https://doi.org/10.1126/science.aac4716)
13. Camerer, C. F. et al. Evaluating replicability of laboratory experiments in economics. *Science* **351**, 1433–1436 (2016). DOI: [10.1126/science.aaf0918](https://doi.org/10.1126/science.aaf0918)
14. Ebersole, C. R. et al. Many Labs 3: Evaluating participant pool quality across the academic semester via replication. *J. Exp. Soc. Psychol.* **67**, 68–82 (2016). DOI: [10.1016/j.jesp.2015.10.012](https://doi.org/10.1016/j.jesp.2015.10.012)
15. Klein, R. A. et al. Many Labs 2: Investigating variation in replicability across sample and setting. *Adv. in Meth. and Prac. in Psychol. Sci.*, in-principle accepted (2018).
16. Bohannon, J. Replication effort provokes praise—and ‘bullying’ charges. *Science* **344**, 788–89 (2014). DOI: [10.1126/science.344.6186.788](https://doi.org/10.1126/science.344.6186.788)
17. Gilbert, D. T., King, G., Pettigrew, S. & Wilson, T. D. Comment on ‘Estimating the reproducibility of psychological science.’ *Science* **351**, 1037 (2016). DOI: [10.1126/science.aad7243](https://doi.org/10.1126/science.aad7243)
18. Anderson, C. J. et al. Response to comment on ‘Estimating the reproducibility of psychological science’. *Science* **351**, 1037 (2016). DOI: [10.1126/science.aad9163](https://doi.org/10.1126/science.aad9163)
19. Ioannidis, J. P. A. Why most discovered true associations are inflated.” *Epidemiology* **19**, 640–648 (2008). DOI: [10.1097/EDE.0b013e31818131e7](https://doi.org/10.1097/EDE.0b013e31818131e7)
20. Simmons, J. P., Nelson, L. D. & Simonsohn, U. False-positive psychology undisclosed flexibility in data collection and analysis allows presenting anything as significant. *Psychol. Sci.* **22**, 1359–1366 (2011). DOI: [10.1177/0956797611417632](https://doi.org/10.1177/0956797611417632)
21. Etz, A. & Vandekerckhove, J. A Bayesian perspective on the Reproducibility Project: Psychology. *PLOS ONE* **11**, e0149794 (2016). DOI: [10.1371/journal.pone.0149794](https://doi.org/10.1371/journal.pone.0149794)
22. Gelman, A. & Stern, H. The difference between “significant” and “not significant” is not itself statistically significant. *Am. Stat.* **60**, 328–331 (2006). DOI: [10.1198/00036810600000000](https://doi.org/10.1198/00036810600000000)

[10.1198/000313006X152649](https://doi.org/10.1198/000313006X152649)

23. Cumming, G. Replication and p intervals: p values predict the future only vaguely, but confidence intervals do much better. *Psychol. Sci.* **3**, 286–300 (2008). DOI: [10.1111/j.1745-6924.2008.00079.x](https://doi.org/10.1111/j.1745-6924.2008.00079.x)
24. Verhagen, J. & Wagenmakers, E.-J. Bayesian tests to quantify the result of a replication attempt. *J. Exp. Psychol. Gen.* **143**, 1457–75 (2014). DOI: [10.1037/a0036731](https://doi.org/10.1037/a0036731)
25. Simonsohn, U. Small Telescopes: detectability and the evaluation of replication results. *Psychol. Sci.* **26**, 559–569 (2015). DOI: [10.1177/0956797614567341](https://doi.org/10.1177/0956797614567341)
26. Patil, P., Peng, R. D. & Leek, J. T. What should researchers expect when they replicate studies? A statistical view of replicability in psychological science. *Perspect. Psychol. Sci.* **11**, 539–544 (2016). DOI: [10.1177/1745691616646366](https://doi.org/10.1177/1745691616646366)
27. Wagenmakers, E.-J. et al. Bayesian inference for psychology. Part II: Example applications with JASP. *Psychon. Bull. Rev.*, in press (2017). DOI: [10.3758/s13423-017-1323-7](https://doi.org/10.3758/s13423-017-1323-7)
28. Lee, M. D. & Wagenmakers, E.-J. *Bayesian cognitive modeling: a practical course*. (Cambridge University press, Cambridge, UK, 2013). ISBN: [9781107603578](https://doi.org/9781107603578)
29. Dreber, A, et al. Using prediction markets to estimate the reproducibility of scientific research. *Proc. Natl Acad. Sci. U.S.A.* **112**, 15343–15347 (2015). DOI: [10.1073/pnas.1516179112](https://doi.org/10.1073/pnas.1516179112)
30. Benjamin, D, et al. Redefine statistical significance. *Nat. Hum. Behav.* **1** (2017). DOI: [10.1038/s41562-017-0189-z](https://doi.org/10.1038/s41562-017-0189-z)
31. Jeffreys, H. *Theory of probability*, (Oxford University Press, Oxford, UK, ed. 3, 1961). ISBN: [9780198503682](https://doi.org/9780198503682)
32. Kass, R. E. & Raftery, A. E. Bayes factors. *J. Am. Stat. Assoc.* **90**, 773–795 (1995). DOI: [10.2307/2291091](https://doi.org/10.2307/2291091)
33. Arrow, K. J. et al. R. The promise of prediction markets. *Science* **320**, 877 (2008). DOI: [10.1126/science.1157679](https://doi.org/10.1126/science.1157679)
34. Nosek, B. A., Ebersole, C. R., DeHaven, A. & Mellor, D. M. The preregistration

- revolution. *Proc. Natl. Acad. Sci.*, in press (2018).
35. Nosek, B. A. et al. Promoting an open research culture: author guidelines for journals could help to promote transparency, openness, and reproducibility. *Science* **348**, 1422 (2015). DOI: [10.1126/science.aab2374](https://doi.org/10.1126/science.aab2374)
  36. Ackerman, J. M., Nocera, C. C. & Bargh, J. A. Incidental haptic sensations influence social judgments and decisions. *Science* **328**, 1712–1715 (2010). DOI: [10.1126/science.1189993](https://doi.org/10.1126/science.1189993)
  37. Aviezer, H., Trope, Y & Todorov, A. Body cues, not facial expressions, discriminate between intense positive and negative emotions. *Science* **338**, 1225–1229 (2012). DOI: [10.1126/science.1224313](https://doi.org/10.1126/science.1224313)
  38. Balafoutas, L. & Sutter, M. Affirmative action policies promote women and do not harm efficiency in the laboratory. *Science* **335**, 579–582 (2012). DOI: [10.1126/science.1211180](https://doi.org/10.1126/science.1211180)
  39. Derex, M., Beugin, M.-P., Godelle, B. & Raymond, M. Experimental evidence for the influence of group size on cultural complexity. *Nature* **503**, 389–391 (2013). DOI: [10.1038/nature12774](https://doi.org/10.1038/nature12774)
  40. Duncan, K., Sadanand, A. & Davachi, L. Memory’s penumbra: episodic memory decisions induce lingering mnemonic biases. *Science* **337**, 485–487 (2012). DOI: [10.1126/science.1221936](https://doi.org/10.1126/science.1221936)
  41. Gervais, W. M. & Norenzayan, A. Analytic thinking promotes religious disbelief. *Science* **336**, 493–496 (2012). DOI: [10.1126/science.1215647](https://doi.org/10.1126/science.1215647)
  42. Gneezy, U., Keenan, E. A. & Gneezy, A. Avoiding overhead aversion in charity. *Science* **346**, 632–635 (2014). DOI: [10.1126/science.1253932](https://doi.org/10.1126/science.1253932)
  43. Hauser, O. P., Rand, D. G., Peysakhovich, A. & Nowak, M. A. Cooperating with the future. *Nature* **511**, 220–223 (2014). DOI: [10.1038/nature13530](https://doi.org/10.1038/nature13530)
  44. Janssen, M. A., Holahan, R., Lee, A. & Ostrom, E. Lab experiments for the study of social-ecological systems. *Science* **328**, 613–617 (2010). DOI: [10.1126/science.1183532](https://doi.org/10.1126/science.1183532)

45. Karpicke, J. D. & Blunt, J. R. Retrieval practice produces more learning than elaborative studying with concept mapping. *Science* **331**, 772–775 (2011). DOI: [10.1126/science.1199327](https://doi.org/10.1126/science.1199327)
46. Kidd, D. C. & Castano, E. Reading literary fiction improves theory of mind. *Science* **342**, 377–380 (2013). DOI: [10.1126/science.1239918](https://doi.org/10.1126/science.1239918)
47. Kovacs, Á. M., Téglás, E. & Endress, A. D. The social sense: susceptibility to others' beliefs in human infants and adults. *Science* **330**, 1830–1834 (2010). DOI: [10.1126/science.1190792](https://doi.org/10.1126/science.1190792)
48. Lee, S. W. S. & Schwarz, N. Washing away postdecisional dissonance. *Science* **328**, 709 (2010). DOI: [10.1126/science.1186799](https://doi.org/10.1126/science.1186799)
49. Morewedge, C. K., Huh, Y. E. & Vosgerau, J. Thought for food: imagined consumption reduces actual consumption. *Science* **330**, 1530–1533 (2010). DOI: [10.1126/science.1195701](https://doi.org/10.1126/science.1195701)
50. Nishi, A., Shirado, H., Rand, D. G. & Christakis, N. A. Inequality and visibility of wealth in experimental social networks. *Nature* **526**, 426–429 (2015). DOI: [10.1038/nature15392](https://doi.org/10.1038/nature15392)
51. Pyc, M. A. & Rawson, K. A. Why testing improves memory: mediator effectiveness hypothesis. *Science* **330**, 335 (2010). DOI: [10.1126/science.1191465](https://doi.org/10.1126/science.1191465)
52. Ramirez, G. & Beilock, S. L. Writing about testing worries boosts exam performance in the classroom. *Science* **331**, 211–213 (2011). DOI: [10.1126/science.1199427](https://doi.org/10.1126/science.1199427)
53. Rand, D. G., Greene, J. D. & Nowak, M. A. Spontaneous giving and calculated greed. *Nature* **489**, 427–430 (2012). DOI: [10.1038/nature11467](https://doi.org/10.1038/nature11467)
54. Shah, A. K., Mullainathan, S. & Shafir, E. Some consequences of having too little. *Science* **338**, 682–685 (2012). DOI: [10.1126/science.1222426](https://doi.org/10.1126/science.1222426)
55. Sparrow, B., Liu, J. & Wegner, D. M. Google effects on memory: cognitive consequences of having information at our fingertips. *Science* **333**, 776–778 (2011). DOI: [10.1126/science.1207745](https://doi.org/10.1126/science.1207745)

56. Wilson, T. D. et al. Just think: the challenges of the disengaged mind. *Science* **345**, 75–77 (2014). DOI: [10.1126/science.1250830](https://doi.org/10.1126/science.1250830)

**Acknowledgements:** Neither Nature Human Behaviour nor the publisher had any involvement with the conduct of this study prior to its submission to the journal. For financial support we thank: Austrian Science Fund FWF (SFB F63, START-grant Y617-G11), Austrian National Bank (grant OeNB 14953), Behavioral and Neuroeconomics Discovery Fund (CFC), Jan Wallander and Tom Hedelius Foundation (P2015-0001:1 and P2013-0156:1), Knut and Alice Wallenberg Foundation (Wallenberg Academy Fellows grant to A. Dreber), Swedish Foundation for Humanities and Social Sciences (NHS14-1719:1), Netherlands Organisation for Scientific Research (Vici grant to E.J. Wagenmakers NWO; 016.Vici.170.083), Sloan Foundation (G-2015-13929), and the Singapore National Research Foundation's Returning Singaporean Scientists Scheme (grant to T.H.Ho; NRF-RSS2014-001). We thank the following persons for assistance with the experiments and analyses: Don van den Bergh, Parampret-Chrisopher Bindra, Johnny van Doorn, Christoph Huber, Alexander Ly, Maarten Marsman, and Jayendra Zambre. The data reported in this paper are tabulated in Tables S2–S4 and the Replication Reports, analyses code, and the data from the replications are available at [www.socialsciencesreplicationproject.com](http://www.socialsciencesreplicationproject.com) and at OSF (<https://osf.io/pfdyw/>).

**Author Contributions:** C.C., A.D., F.H., J.H., T.H., M.J., M.K., G.N., B.N., and T.P. designed research; C.C., A.D., F.H., T.H., J.H., M.J., M.K., D.M., G.N., B.N., T.P., E.J.W. wrote the paper; T.C., A.D., E.F., F.H., T.H., M.J., T.P., and Y.C. helped design

the prediction market part; F.H. and E.J.W. analyzed data; A.A., N.B., A.G., E.H., F.H., L.H., T.I., S.I., D.M., J.R. and H.W. carried out the replications (including re-estimating the original estimate with the replication data); all authors approved the final manuscript.

**Competing interests:** The authors report no potential conflicts of interest. No MTAs, patents or patent applications apply to methods or data in the paper.

**Additional Information:** Supplementary information is available for this paper at doi:

**Fig. 1. Replication results after Stage 1 and Stage 2.** (a) Plotted are 95% CIs of replication effect sizes (standardized to correlation coefficients  $r$ ) after Stage 1. The standardized effect sizes are normalized so that 1 equals the original effect size. There is a significant effect in the same direction as in the original study for 12 replications [57.1%; 95% CI = (34.1%, 80.2%)]. (b) Plotted are 95% CIs of replication effect sizes (standardized to correlation coefficients  $r$ ) after Stage 2 (replications not proceeding to Stage 2 are included with their Stage 1 results). The standardized effect sizes are normalized so that 1 equals the original effect size. There is a significant effect in the same direction as in the original study for 13 replications [61.9%; 95% CI = (39.3%, 84.6%)]. (c) Meta-analytic estimates of effect sizes combining the original and replication studies. 95% CIs of standardized effect sizes (correlation coefficient  $r$ ). The standardized effect sizes are normalized so that 1 equals the original effect size. 16 studies have a significant effect in the same direction as the original study in the meta-analysis [76.2%; 95% CI = (56.3%, 96.1%)].

**Fig. 2. Replication results for two complementary replication indicators; 95% prediction intervals<sup>26</sup> in panel a and the Small Telescopes approach<sup>25</sup> in panel b.**

(a) Plotted are 95% prediction intervals for the standardized original effect sizes (correlation coefficient  $r$ ). The standardized effect sizes are normalized so that 1 equals the original effect size. 14 replications [66.7%; 95% CI = (44.7%, 88.7%)] are within the 95% prediction interval and replicate according to this indicator. (b) Plotted are 90% CIs of replication effect sizes in relation to small effect sizes as defined by the Small Telescopes approach (the effect size the original study would have had 33% power to detect). Effect sizes are standardized to correlation coefficients  $r$  and normalized so that 1 equals the original effect size. A study is defined as failing to replicate if the 90% confidence interval is below the small effect. According to the Small Telescopes approach 12 [57.1%; 95% CI = (34.1%, 80.2%)] studies replicate.



**Fig. 3. Default Bayes factors (one-sided)<sup>27</sup> for the 21 replications.** A default Bayes factor above 1 favors the hypothesis of an effect in the direction of the original paper and a default Bayes factor below 1 favors the null hypothesis of no effect. The evidence categories proposed by Jeffreys<sup>31</sup> are also shown in the Figure (from extreme support for the null hypothesis to extreme support for the original hypothesis). The default Bayes factor is above 1 and provide evidence in favor of an effect in the direction of the original study for the 13 (61.9%) studies that replicated according to the statistical significance criterion. This evidence is strong to extreme for 9 (42.9%) studies. The default Bayes factor is below 1 for 8 (38.1%) studies providing evidence in support of the null hypothesis; this evidence is strong to extreme for 4 (19.0%) studies.

**Fig. 4. Prediction market and survey beliefs.** The Figure shows the prediction market beliefs and the survey beliefs of replicating (from Treatment 2 for measuring beliefs; see the supplementary materials, section 3 for details and Fig. S6 for results from Treatment 1). The replication studies are ranked in terms of prediction market beliefs on the y-axis. The mean prediction market belief of replication is 63.4% [range of 23.1% to 95.5%, 95% CI = (53.7%, 73.0%)], and the mean survey belief is 60.6% [range of 27.8% to 81.5%, 95% CI = (53.0%, 68.2%)]. This is similar to the actual replication rate of 61.9%. The prediction market beliefs and survey beliefs are highly correlated, but imprecisely estimated (Spearman correlation coefficient 0.845, 95% CI = (0.652, 0.936),  $p < 0.001$ ,  $n = 21$ ). Both the prediction market beliefs (Spearman correlation coefficient 0.842, 95% CI = (0.645, 0.934),  $p < 0.001$ ,  $n = 21$ ), and the survey beliefs (Spearman correlation coefficient 0.761, 95% CI = (0.491, 0.898),  $p < 0.001$ ,  $n = 21$ ) are also highly correlated with a successful replication.