

Examining the replicability of online experiments selected by a decision market

Received: 29 November 2023

Accepted: 11 October 2024

Published online: 19 November 2024

 Check for updates

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

Here we test the feasibility of using decision markets to select studies for replication and provide evidence about the replicability of online experiments. Social scientists ($n = 162$) traded on the outcome of close replications of 41 systematically selected MTurk social science experiments published in PNAS 2015–2018, knowing that the 12 studies with the lowest and the 12 with the highest final market prices would be selected for replication, along with 2 randomly selected studies. The replication rate, based on the statistical significance indicator, was 83% for the top-12 and 33% for the bottom-12 group. Overall, 54% of the studies were successfully replicated, with replication effect size estimates averaging 45% of the original effect size estimates. The replication rate varied between 54% and 62% for alternative replication indicators. The observed replicability of MTurk experiments is comparable to that of previous systematic replication projects involving laboratory experiments.

Can published research findings be trusted? Unfortunately, the answer to this question is not straightforward, and the credibility of scientific findings and methods has been questioned repeatedly^{1–9}. A vital tool for evaluating and enhancing the reliability of published findings is to carry out replications, which can be used to sort out likely true positive findings from likely false positives. A replication essentially updates the probability of the hypothesis being true after observing the replication outcome. A successful replication will move this probability towards 100%, while a failed replication will move it towards 0% (refs. 10,11).

In recent years, several systematic large-scale replication projects in the social sciences have been published^{12–17}, reporting replication rates of around 50% in terms of both the fraction of statistically significant replications and the relative effect sizes of replications. Potential factors to explain these replication rates may be low statistical power^{1,18,19} in the original studies, testing original hypotheses with low priors^{1,10,20} and questionable research practices^{1,21,22}. Systematic replication studies led to discussions about improving research practices^{23,24} and have substantially increased the interest in independent replications²⁵.

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

✉ e-mail: anna.dreber@hhs.se

However, as it is time-consuming and costly to conduct replications, it has been argued that it is useful to have a principled mechanism to decide which replications to prioritize to facilitate efficient and effective usage of resources^{25–37}. Here we test the feasibility of one potential method to select which studies to replicate. Building on previous work using prediction markets^{38–40} to forecast replicability, we adapt the forecasting methodology to what is referred to as decision markets^{41–44}.

The decisive distinction between prediction markets and decision markets is that prediction markets elicit aggregate-level replicability forecasts on a predetermined set of studies, whereas decision market forecasts determine which studies are going to be put to a replication test. While previous studies provide evidence that prediction market forecasts are predictive of replication outcomes^{10,16,17,45}, prediction efficiency might not generalize to decision markets, which involve more complex procedures and incentives. The performance of decision markets as a tool for selecting which empirical claims to replicate has not been systematically examined. Note that a decision market in itself is not sufficient to provide a mechanism to select studies for replication, but it has to be combined with an objective function of which studies to replicate (an example of an objective function would be to replicate the studies with the lowest probability of replication). For decision markets to be potentially useful for selecting studies for replication, it first has to be established that the predictions of the decision markets are associated with the replication outcomes. To provide such a ‘proof of concept’ of using a decision market as a mechanism to determine which studies to replicate, we first identified all social science experiments published in the Proceedings of the National Academy of Sciences (PNAS) between 2015 and 2018 that fulfilled our inclusion criteria for (1) the journal and period; (2) the platform on which the experiment was performed (Amazon Mechanical Turk; MTurk); (3) the type of design (between-subjects or within-subject treatment design); (4) the equipment and materials needed to implement the experiment (the experiment had to be logistically feasible for us to implement); and (5) the results reported in the experiment (at least one main or interaction effect with $P < 0.05$ reported in the main text). On the basis of our inclusion criteria, we identified 44 articles, 3 of which have been excluded owing to a lack of feasibility, leaving us with a final sample of 41 articles^{46–86} (see Methods for details on the inclusion criteria). For each of these articles, we identified one critical finding with $P < 0.05$ that we could potentially replicate (see Methods for details and Supplementary Table 1 for the hypotheses selected for each of the 41 studies).

We then invited social science researchers to participate as forecasters in both a prediction survey and an incentivized decision market on the 41 studies. In the survey, the forecasters independently estimated the probability of replication for the 41 studies. In the decision market, they could trade on whether the result of each of the 41 studies would replicate. Participants in the decision market received an endowment of 100 tokens corresponding to USD 50, and 162 participants made a total of 4,412 trades. Traders in the market were informed about the preregistered decision mechanism: the 12 studies with the highest and the 12 studies with the lowest market prices were to be selected for close replication; in addition, 2 randomly chosen studies (out of the remaining 17 studies) are replicated to ensure incentive compatibility, with participant payoffs scaled up by the inverse of their probability in the decision rule (see Methods for details). For incentive compatibility, all the 41 studies included need to have a strictly positive probability of being selected for replication, which is ensured by having at least one randomly selected study. Otherwise, traders would be incentivized to only trade on those studies that will most likely be chosen according to the decision rule.

All replication experiments, just like all original studies, were conducted on Amazon Mechanical Turk (MTurk), and the same sample restrictions and exclusion criteria as the original studies were applied, which guards against concerns about the potential

moderating effects of culture differences in replications^{14,87}. Replication sample sizes were determined to have 90% power to detect 2/3 of the effect size reported in the original study at the 5% significance level in a two-sided test (with the effect size estimates having been converted to Cohen’s d to have a common standardized effect size measure across the original studies and the replication studies; see Methods for details). If sample size calculations led to replication sample sizes smaller than in the original study, we targeted the same sample size as in the original study. The average sample size in the replications ($n = 1,018$) was 3.5 times as large as the average sample size in the original studies ($n = 292$).

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

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

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

The statistical significance indicator is a binary criterion of replication and is based on testing the hypothesis for which the original study found support using standard null hypothesis significance testing. The indicator crudely classifies replications as failed or successful depending on whether the replication study yields evidence in support of the original hypothesis at a particular significance threshold. (Critics of null hypothesis significance testing or privileging a P value of 0.05 will, justifiably so, object to this crude classification; that is why it is only one of the several indicators that we report.) A replication classified as failed based on this indicator, however, does not imply that the estimated replication effect size is significantly different from the original estimate (see more on this below). To keep the false negative risk at bay and to be informative, the statistical significance indicator calls for well-powered replications (as in this study)^{105,106}. However, a limitation of this indicator for well-powered replication studies is that it may classify a replication as successful even if the observed

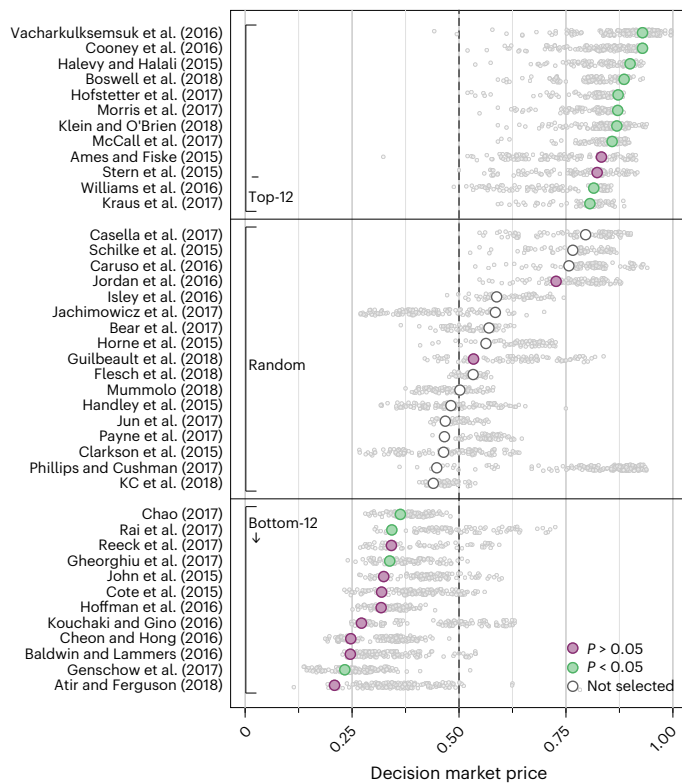


Fig. 1 | Decision market prices for the 41 included studies. Plotted are the decision market prices for the 41 MTurk social science experiments published in PNAS between 2015 and 2018. The small grey dots indicate the market prices after each market transaction; the larger dots indicate the final market price. The studies are ordered based on the final decision market prices, which can be interpreted as the market’s probability forecast of successful replication. The 12 studies with the highest decision market prices and the 12 studies with the lowest decision market prices were selected for replication; in addition, 2 of the remaining 17 studies were selected for replication at random to ensure that the decision market is incentive compatible. The replication outcomes for the statistical significance indicator are also illustrated for the 26 replicated studies. The point-biserial correlation between the decision market prices and the replication outcomes in primary hypothesis 1 is $r = 0.505$ (95% CI (0.146, 0.712), $t(24) = 2.867$, $P = 0.008$; $n = 26$, two-sided test).

effect size is substantially smaller (or larger) than in the original study. While the statistical significance indicator dichotomizes replication outcomes into successful and failed, replicability may be perceived as a continuous matter of degree. This is why we also consider the relative effect size—a continuous measure of replicability—as a primary replication indicator. While the relative effect size constitutes an imprecise indicator for an individual replication study, it arguably provides an informative measurement of the extent of replicability for a group of studies as it quantifies the average degree of apparent inflation in the original effect sizes¹⁰⁷. As all replication indicators have limitations, we preregistered four additional secondary replication indicators. In addition, we report the results for two non-preregistered replication indicators, which were helpfully suggested during the review process of the paper.

In our two primary hypotheses, we conjecture that (1) the decision market prices positively correlate with the replication outcomes and (2) the standardized effect sizes in the replications are lower than in the original studies. All hypotheses are evaluated using two-tailed tests, and—following Benjamin et al.¹⁰⁸—we interpret results with $P < 0.005$ as ‘statistically significant evidence’, whereas results with $0.005 \leq P < 0.05$ are considered ‘suggestive evidence’. No adjustments were made for multiple comparisons.

Results

Replication outcomes and decision market performance

Figure 1 and Supplementary Table 2 show the results for the decision markets where the final market price can be interpreted as the predicted replication probability. The predicted probabilities of replication range from 20.9% to 92.9% for the 41 studies, with a mean of 57.6% (s.d. = 23.6%). The average predicted probability for the 26 studies eventually selected for replication is 58.5%. Figure 1 also delineates the replication outcomes based on the statistical significance indicator, which allows for gauging the relationship between the decision market prices and the replication outcomes. In Fig. 2 and Supplementary Table 3, we show the replication results for the 26 studies selected for replication. Of the 26 claims, 14 (53.8%; 95% confidence interval (CI) (33.4%, 73.4%)) replicated successfully according to the statistical significance indicator. The point-biserial correlation between decision market prices and the binary replication outcome, testing our first primary hypothesis, is $r = 0.505$ (95% CI (0.146, 0.712); $t(24) = 2.867$, $P = 0.008$; $n = 26$). Thus, in support of our first primary hypothesis, we find suggestive evidence of a positive association between decision market prices and replication outcomes. As a related secondary hypothesis, we test if the replication rate is lower among the 12 studies with the lowest decision market prices than for the 12 studies with the highest decision market prices. The replication rate is 33.3% (95% CI (9.9%, 65.1%)) for the studies in the ‘bottom-12’ group and 83.3% (95% CI (51.6%, 97.9%)) for the studies in the ‘top-12’ group, yielding suggestive evidence in support of our secondary hypothesis (Fisher’s exact test; $\chi^2(1) = 6.171$, $P = 0.036$; $n = 24$). Note Fisher’s test conditions on the margin totals; hence, it is only exact for the conditional

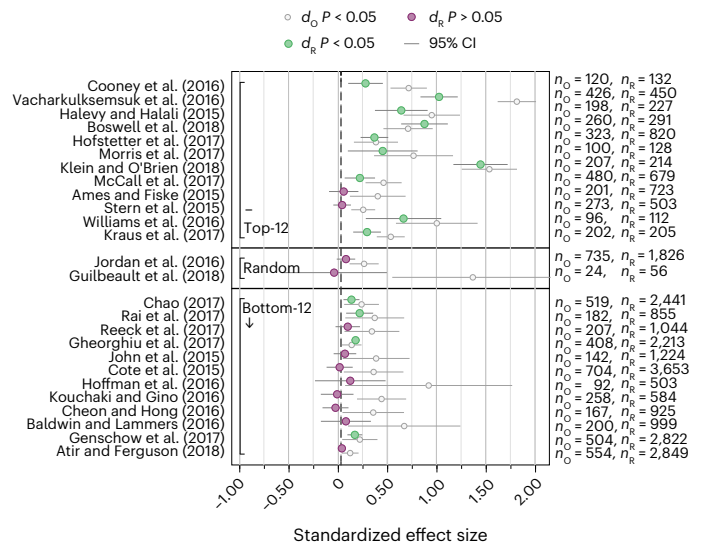


Fig. 2 | Replication results. Plotted are the point estimates and the 95% CIs (standardized to Cohen’s d units) of the 26 replications (d_r) and original studies. Studies within each of the three panels (top-12, random, bottom-12) are sorted based on the decision market prices as in Fig. 1. There is a statistically significant effect ($P < 0.05$) in the same direction as the original study for 14 out of 26 replications (53.8%; 95% CI (33.4%, 73.4%)). For the 12 studies with the highest decision market prices, there is a statistically significant effect ($P < 0.05$) in the same direction as the original study for 10 out of 12 replications (83.3%; 95% CI (51.6%, 97.9%)). For the 12 studies with the lowest decision market prices, there is a statistically significant effect ($P < 0.05$) in the same direction as the original study for 4 out of 12 replications (33.3%; 95% CI (9.9%, 65.1%)). Our secondary hypothesis test provides suggestive evidence that the difference in replication rates between the top-12 and the bottom-12 group is different from zero (Fisher’s exact test; $\chi^2(1) = 6.171$, $P = 0.036$; $n = 24$, two-sided test). The error bars denote the 95% CIs of the original and the replication effect size estimates. The numbers of observations used to estimate the 95% CIs are the original and replication sample sizes noted on the right as n_o and n_r .

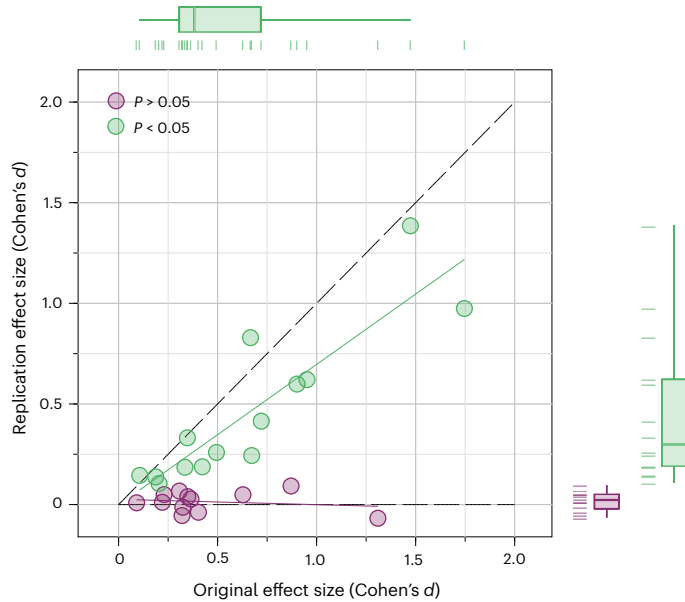


Fig. 3 | Relationship between estimated original and replication effect sizes. Plotted are the estimated original and replication effect sizes for each of the 26 replication studies (the estimated effect sizes of both the original and replication studies are standardized to Cohen's *d* units). The 95% CIs for the original and replication effect size estimates are illustrated in Fig. 2 and tabulated in Supplementary Table 3. The mean estimated effect size of the 26 replication studies is 0.253 (s.d. = 0.357) compared with 0.563 (s.d. = 0.426) for the original studies, resulting in a relative estimated average effect size of 45.0%, confirming our second primary hypothesis (Wilcoxon signed-rank test, $z = 4.203$, $P < 0.001$; $n = 26$, two-sided test). The estimated relative effect size of the 13 replications that have been successfully replicated according to the statistical significance indicator is 69.5%, and the estimated relative effect size of the 13 studies that did not replicate is 3.2%. The box plots show the median, the interquartile range, and the 5th and 95th percentile of the effect size estimates in the 26 original studies and the 26 replication studies.

distribution and can be overly conservative if the margin totals are unknown^{109,110}, as is the case in our analysis. Boschloo's test¹¹¹, an exact unconditional procedure uniformly more powerful than Fisher's test, also yields suggestive evidence for the difference in proportions between the 'top-12' and the 'bottom-12' group (not preregistered; 95% CI (0.089, 0.799), $P = 0.017$; $n = 24$).

Relative effect sizes

The mean estimated effect size of the 26 replication studies (in Cohen's *d* units) is 0.253 (s.d. = 0.357) compared with 0.563 (s.d. = 0.426) for the original studies, implying a relative estimated average effect size, just dividing the two numbers, of 45.0%; the difference in estimated effect sizes is statistically significant, supporting our second primary hypothesis of systematically smaller estimated effect sizes in the replications (Wilcoxon signed-rank test, $z = 4.203$, $P < 0.001$; $n = 26$). The relative effect size can also be estimated for each study separately (reported in Supplementary Table 3) and varies between -17.0% and 136.2%, with a mean estimate across studies of 41.1% (95% CI (24.5%, 57.7%)). For the 14 studies that replicated according to the statistical significance indicator, the first and the second relative effect size measures as defined above are 69.5% and 72.0% (95% CI (54.8%, 89.3%)), indicative of some inflation in original effect sizes even for apparent true positives. The two estimated relative effect size measures for the 12 studies that failed to replicate according to the statistical significance indicator are 3.2% and 5.0% (95% CI (-2.6%, 12.5%)), respectively. Figure 3 illustrates the relationship between the estimated original and replication effect sizes.

Secondary replication indicators

We also preregistered four secondary replication indicators: the small-telescopes approach¹¹², the one-sided default Bayes factor¹¹³, the replication Bayes factor¹¹⁴ and the fixed-effects weighted meta-analytic effect size (see Methods for details). When relying on the small-telescopes approach, testing if the replication effect size is smaller than a 'small effect'¹¹², 15 studies (57.7%; 95% CI (36.9%, 76.6%)) are considered successful replications (Fig. 4 and Supplementary Table 4). The one-sided default Bayes factor (BF_{+0}) indicates the strength of evidence in favour of the alternative hypothesis as opposed to the null hypothesis. BF_{+0} exceeds 1 for the 14 studies (53.8%; 95% CI (33.4%, 73.4%)) that replicated according to the statistical significance indicator, with strong evidence ($BF_{+0} > 10$) for the tested hypothesis for 9 studies (34.6%; 95% CI (17.2%, 55.7%)); BF_{+0} is below 1 for the 12 replications (46.2%; 95% CI (26.6%, 66.6%)) that failed to replicate according to the statistical significance indicator, with strong evidence ($BF_{+0} < 0.1$) for the null hypothesis for 7 studies (26.9%; 95% CI (11.6%, 47.8%)) based on the evidence categories proposed by Jeffreys¹¹⁵ (Fig. 5 and Supplementary Table 4). The one-sided replication Bayes factor (BF_{R0}) indicates the strength of additional evidence in favour of the alternative hypothesis as opposed to the null hypothesis, given the already acquired evidence based on the original data¹¹⁴. Replication Bayes factors lead to similar conclusions as the one-sided default Bayes factors, with $BF_{R0} > 10$ for 10 studies (38.5%; 95% CI (20.2%, 59.4%)) and $BF_{R0} < 0.1$ for 7 studies (26.9%; 95% CI (11.6%, 47.8%)). One exception to this is the study by Cooney et al.⁵⁶, for which the default Bayes factor exceeds one ($BF_{+0} = 8.01$) but the replication Bayes factor is below one ($BF_{R0} = 0.23$) owing to the replication effect size being only about a third of the original effect size and a larger sample size in the replication compared with the original study (Fig. 5 and Supplementary Table 4). The meta-analytic effect size is statistically significant at the 5% level for 16 studies (61.5%; 95% CI (40.6%, 79.8%)) and significant at the 0.5% level for 14 studies (53.8%; 95% CI (33.4%, 73.4%)) (Fig. 6 and Supplementary Table 4). The meta-analytic effect sizes should be interpreted cautiously as original effect sizes reported as statistically

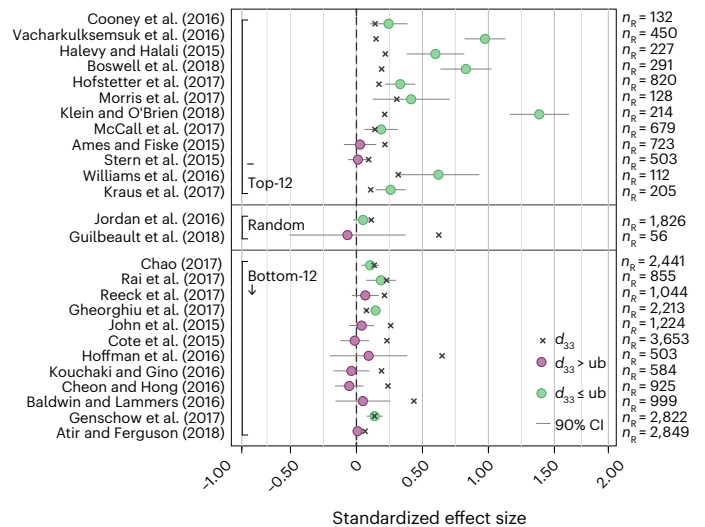


Fig. 4 | Replication results based on the small-telescopes approach (a secondary replication indicator). Plotted are the 90% CIs of replication effect sizes in relation to small-effect sizes as defined by the small-telescopes approach¹¹² (the effect size that the original study would have had 33% power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

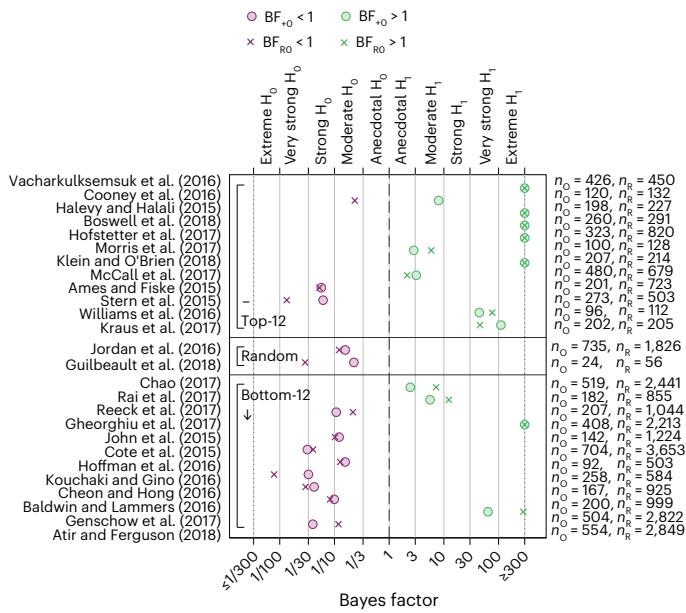


Fig. 5 | Replication results based on Bayes factors (secondary replication indicators). The figure plots the one-sided default Bayes factor (BF_{0}) and the replication Bayes factor (BF_{R0}) for the 26 replications¹¹³. $BF_{0} > 1$ favours the hypothesis of an effect in the direction of the original paper, whereas $BF_{0} < 1$ favours the null hypothesis of no effect. BF_{R0} quantifies the additional evidence provided by the replication results on top of the original evidence. $BF_{R0} > 1$ indicates additional evidence in favour of the alternative over the null, whereas $BF_{R0} < 1$ indicates additional evidence for the null instead. The evidence categories proposed by Jeffreys¹¹⁵ are also shown (from extreme support for the null hypothesis to extreme support for the original hypothesis). Studies within the three panels (top-12, random, bottom-12) are sorted based on the decision market prices as in Fig. 1. The BF_{0} is above 1 for all 14 replication studies that successfully replicated according to the statistical significance indicator and below 1 for all 12 replication studies that failed to replicate according to the statistical significance indicator. The BF_{R0} is above 1 for 13 of the 14 replication studies that replicated according to the statistical significance indicator and below 1 for Cooney et al.⁵⁶ whose estimated relative effect size of 0.36 is the lowest among these 14 studies; the BF_{R0} is below 1 for all of the 12 replication studies that failed to replicate according to the statistical significance indicator. The numbers of observations used to estimate BF_{0} and BF_{R0} are the original and replication sample sizes noted on the right as n_0 and n_R .

significant are likely to be overestimated on average owing to insufficient sample sizes and, thereby, statistical power (and potentially owing to questionable research practices)^{18,19}. Overall, the primary and secondary replication indicators yield the same binary conclusions for 23 of the 26 replications.

Non-preregistered replication indicators

Following the suggestion of a reviewer, we report the results for two additional replication indicators. The first alternative replication indicator is a test of whether the replication effect size is statistically significantly different from the original effect size. This indicator is closely related to the prediction interval approach¹¹⁶ as the results can be illustrated as prediction intervals that the replication effect sizes are evaluated against: if the replication effect size falls outside the 95% prediction interval, the replication and original effect sizes differ at the 5% significance level, and the replication is considered a failure. The prediction interval approach, thus, yields a binary replication indicator, which is complemented by a continuous replicability measure defined as the P value of the test of a significant difference between the replication and original effect sizes. We illustrate the prediction interval results in Fig. 7 and report the z -statistics and

P values in Supplementary Table 5 (the P values are also shown in Fig. 7). According to the prediction interval indicator, 15 studies (57.7%, 95% CI (36.9%, 76.6%)) replicate. This replication rate is close to the result for the statistical significance indicator. However, the classification of nine replication outcomes shifts: for four studies, the classification changes from successful to failed, and for five studies, the classification changes from failed to successful. These changes are due to the fact that low-powered original studies are more likely to replicate, whereas high-powered original studies are less likely to replicate based on the prediction interval indicator (compared with evaluating replicability based on the statistical significance indicator). According to the prediction interval indicator, six studies failed to replicate among the top-12 studies in terms of decision market prices, whereas three studies failed to replicate among the bottom-12 studies.

Associations between indicators (not preregistered)

To examine the relationship between the replication indicators, we estimated Kendall’s rank correlations τ_b between all the replication indicators used in the study (Supplementary Table 6). All preregistered replication indicators are strongly correlated with each other, with τ_b varying between 0.61 and 1.00 ($P < 0.005$ for all correlations). However, they are more weakly to moderately correlated with the prediction interval approach and P values from z -tests comparing the replication and original effect sizes (with τ_b varying between 0.12 and 0.56).

Each of the various replication indicators presented in this study has its strengths and weaknesses. There is no general consensus about which indicator is most appropriate^{15–17,117–119}. Therefore, we chose to report the results for a host of indicators and leave it to readers to judge the suitability of the different indicators and their degree of consensus. The overall replication rate is similar for all the binary replication indicators and varies between 14 (53.8%) and 16 (61.5%) studies.

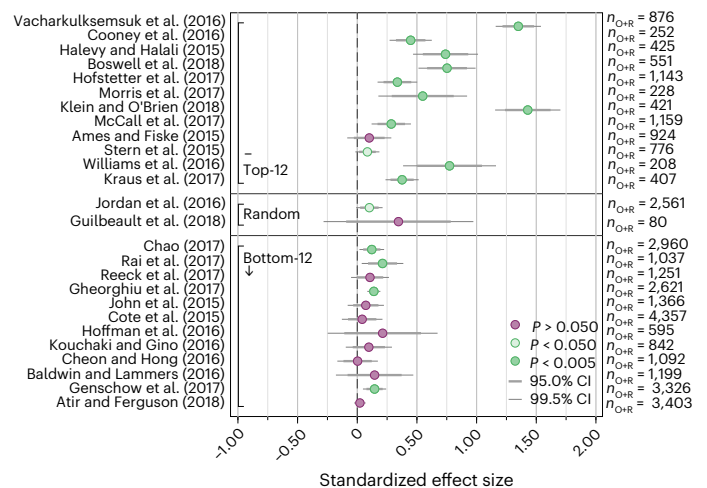


Fig. 6 | Meta-analytic estimated effect sizes combining the original and the replication estimated effect sizes (a secondary replication indicator).

The figure plots the point estimates and 95% and 99.5% CIs of the fixed-effects weighted meta-analytic effect sizes, combining the original and the replication studies (standardized to Cohen’s d units). Studies within the three panels (top-12, random, bottom-12) are sorted based on the decision market prices as in Fig. 1. As preregistered, we report the significance of the estimated meta-analytic effect sizes for both the 0.05 significance threshold and the 0.005 significance threshold (based on a two-sided z -test). Sixteen out of 26 (61.5%; 95% CI (40.6%, 79.8%)) studies replicated according to the statistical significance indicator using the 0.05 significance threshold, and 14 out of 26 (53.8%; 95% CI (33.4%, 73.4%)) studies replicated using the 0.005 significance threshold. The error bars denote the 95% CIs of the estimated meta-analytic effect sizes. The number of observations used to estimate the 95% CIs are the sums of the original and replication sample sizes noted next to the study identifier on the y -axis as n_{0+R} .

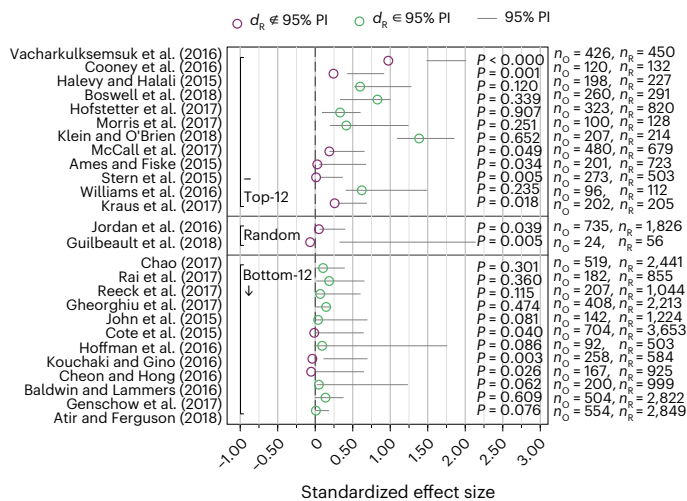


Fig. 7 | Replication results based on prediction intervals (not preregistered).

Plotted are the 95% prediction intervals¹¹⁶ (PIs) for the standardized original effect sizes (Cohen's d). Studies within the three panels (top-12, random, bottom-12) are sorted based on the decision market prices as in Fig. 1. Fifteen replications out of 26 (57.7%; CI (36.9%, 76.6%)) are within the 95% prediction interval and replicate according to this indicator. The P values reported on the right are based on two-sample z -tests for a difference between the replication effect size and the original effect size. The grey lines denote the 95% prediction intervals, and the small circles denote the mean replication effect sizes. All tests are two-sided. The numbers of observations used to estimate the 95% prediction intervals are the original and replication sample sizes noted next to the study identifier on the y -axis as n_o and n_r .

The agreement about which results are classified as successfully replicated is large between the indicators, with the exception of the prediction interval approach. The estimated average relative effect size of around 45%, which can be interpreted in terms of a replicability rate, yields an estimate in the same ballpark. The somewhat lower estimate for the relative effect size is due to the fact that not only the false positive rate but also the inflation of true positive effect sizes is factored in. Another advantage of the relative effect size indicator is that it is not affected by replication power. The three other continuous replication indicators cannot be aggregated across studies and are thus difficult to compare to the other indicators on an aggregated level.

Replicability forecasts and indicators (not preregistered)

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

Survey forecasts versus decision market predictions

We tested three additional preregistered secondary hypotheses based on the survey beliefs about replication (see Methods for details and Supplementary Table 7 for the survey results). The point-biserial correlation between average survey beliefs and the replication outcomes based on the statistical significance criterion is $r = 0.476$ (95% CI (0.107, 0.694); $t(24) = 2.650$, $P = 0.014$; $n = 26$). The survey beliefs and the decision market prices are positively correlated with a Pearson correlation

of 0.899 (95% CI (0.814, 0.944); $t(39) = 12.830$, $P < 0.001$; $n = 41$) (Fig. 8a).

The final secondary hypothesis tests if the prediction accuracy, measured in terms of the absolute prediction error and the Brier score (that is, the squared prediction error), is higher for the decision market than the survey forecasts (Fig. 8b). The mean absolute prediction error and the mean Brier score are 0.353 and 0.188 for the decision market, and 0.421 and 0.202 for the survey, respectively, providing suggestive evidence for higher accuracy for the market forecasts based on the absolute prediction error (Wilcoxon signed-rank test, $z = 2.172$, $P = 0.030$; $n = 26$) but not the Brier score (Wilcoxon signed-rank test, $z = 1.181$, $P = 0.238$; $n = 26$). The failure to reject the null hypothesis for the Brier score does not imply that the null hypothesis is true. In the survey, we also elicited forecasters' self-rated expertise for each study. The average self-rated expertise (of participants eventually active in the markets, $n = 162$) for the 26 replicated studies was 2.31 (s.d. = 1.40; $n = 4,212$) on a scale from 1 ('no knowledge of the topic') to 7 ('very high knowledge of the topic'). Supplementary Fig. 1 plots the absolute prediction error and the Brier score of the 26 survey and decision market forecasts over the average self-rated expertise per study. We do not find evidence for the prediction accuracy and the average self-rated expertise being significantly correlated (not preregistered; see Supplementary Fig. 1 for details).

Beliefs about the Covid-19 pandemic and replicability

A potential issue raised by some original authors in giving feedback on the replication reports before the data collection was that the replicability of some original results might be affected by the Covid-19 pandemic (as all the original studies were conducted before the pandemic). We evaluate this possibility in a preregistered exploratory analysis, relying on the forecasters' beliefs about the impact of the pandemic on replicability. As part of the prediction survey, participants were asked to judge whether the pandemic would have affected the likelihood of successful replication, measured on a scale from -3 ('the pandemic has definitely decreased the probability of replication') to 3 ('the pandemic has definitely increased the probability of replication'). We test if the average response to this question differs from zero using a one-sample t -test for each of the 26 replications, and we test if the average response across all 26 studies differs from 0. We find a statistically significant result for four and a suggestive result for two replications on beliefs that Covid-19 has affected the replication probability (Supplementary Table 8). For the six studies with suggestive or statistically significant evidence, the estimate is negative for two studies and positive for four; only in two of the cases does the sign of the expectation match the eventual replication outcome. For the average belief about the impact of the pandemic on replicability across the 26 studies of 162 forecasters (who were active in the decision markets), there is suggestive evidence that the mean of 0.039 (s.d. = 0.190) differs from zero ($t(161) = 2.598$, $P = 0.010$; $n = 162$). Somewhat surprisingly—and in contrast to the concerns raised by some of the original authors—there is thus a tendency for forecasters to believe that the pandemic has increased the average likelihood that the studies will replicate. However, the magnitude of the effect is small ($d = 0.204$; 95% CI (0.049, 0.360)).

In addition, we tested, estimating the point-biserial correlation, if the average belief (per study) about the pandemic's impact on replicability correlates with the replication outcomes based on the statistical significance indicator; we do not find a statistically significant association ($r = 0.014$, 95% CI (−0.360, 0.382); $t(24) = 0.068$, $P = 0.946$; $n = 26$). Yet, we cannot rule out that Covid-19 has entailed effects on replicability not foreseen by scholars participating in the survey. Further work is needed to gauge whether and to which extent experimental replications—and predictions of replication success—might be sensitive to macro-historical secular change such as economic upheaval, wars, pandemics and so on. Forecasters' beliefs about the pandemic's impact on replicability are also neither statistically significantly correlated with the final decision market prices ($r = 0.387$, 95% CI (−0.008, 0.669);

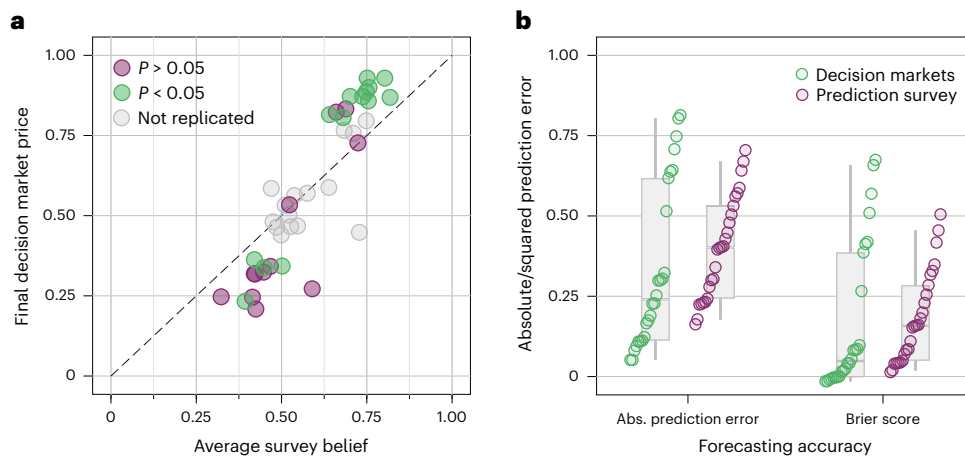


Fig. 8 | Relationship between decision market prices and mean survey beliefs and forecasting accuracy. **a**, Plotted are the decision market prices and the mean survey beliefs about replication for the 41 studies included in the decision market and the survey; the colour coding highlights the replication outcomes for the 26 replicated studies. The decision market prices and the mean survey beliefs about replication are highly correlated with a Pearson correlation of $r = 0.899$ (95% CI (0.814, 0.944); $t(39) = 12.830$, $P = 1.4 \times 10^{-15}$; $n = 41$, two-sided test).

b, Plotted are the absolute prediction errors and the Brier scores (the squared prediction errors) for the decision market and the prediction survey for the

26 replicated studies. There is suggestive evidence of higher prediction accuracy for the decision market in terms of the absolute prediction error (0.353 for the decision markets and 0.421 for the survey; Wilcoxon signed-rank test, $z = 2.172$, $P = 0.030$; $n = 26$, two-sided test), but not in terms of the Brier score (0.188 for the decision markets and 0.202 for the survey; Wilcoxon signed-rank test, $z = 1.181$, $P = 0.238$; $n = 26$, two-sided test). The box plots show the median, the interquartile range, and the 5th and 95th percentile of the absolute prediction errors and Brier scores for the survey and decision market predictions of the 26 replication studies.

$t(24) = 2.055$, $P = 0.051$; $n = 26$) nor the average survey belief of replication ($r = 0.347$, 95% CI (-0.053, 0.644); $t(24) = 1.815$, $P = 0.082$; $n = 26$), although the point estimates of the correlations are quite sizeable.

Original P value and replication (not preregistered)

For comparison to previous systematic large-scale replication projects, we also report the correlation between the original P value and the replication outcome for the statistical significance indicator. The point-biserial correlation between the original P value and the replication outcome for the statistical significance indicator is -0.400 ($P = 0.043$; 95% CI (0.014, 0.648)) and comparable in magnitude to correlations of -0.33 in the Replication Project: Psychology (RPP)¹⁵, -0.57 in the Experimental Economics Replication Project (EERP)¹⁶ and -0.40 in the Social Sciences Replication Project (SSRP)¹⁷.

Discussion

We found suggestive evidence ($P < 0.05$) for our first primary hypothesis that final decision market prices correlate with replication outcomes ($r = 0.505$). However, the estimated effect size is somewhat smaller than the effect size of $r = 0.67$, as presumed in our a priori power calculations (see Methods for details). The estimated correlation is within the range of previous prediction markets on systematic replication projects with correlations of 0.42 in the RPP^{10,15}, 0.30 in the EERP¹⁶ and 0.84 in the SSRP¹⁷, but we expected a stronger correlation because we selected studies with the highest and the lowest prices for replication. Consistent with the primary hypothesis test, there is also suggestive evidence of a difference in the replication rate between the ‘top-12’ (10 of 12) and ‘bottom-12’ (4 of 12) in our secondary hypothesis test. The difference of 50 percentage points is also reflected in the difference between the forecasted replication rates of 86.6% (‘top-12’) versus 29.6% (‘bottom-12’) in the decision market. However, the small sample size suggests caution against drawing firm conclusions about whether decision markets are appropriate for selecting studies for replication.

The pooled evidence from previous prediction market studies on replication outcomes suggests that markets are somewhat more accurate than surveys⁴⁵, although the difference tends to be small. These indications are consistent with our results, yielding suggestive

evidence of higher accuracy in terms of the absolute prediction error but not in terms of the squared prediction error (although, as noted above, failing to reject the null hypothesis for the squared prediction error does not imply that the null hypothesis is true). The estimated correlation between the average survey beliefs and the replication outcomes was almost as high for the survey as the prediction market (0.476 versus 0.505). The decision market prices and survey beliefs are also highly correlated with each other ($r = 0.9$). Since surveys are less resource intensive, simple polls can be an expedient alternative to decision markets for selecting which studies to replicate, even if they should be somewhat less accurate. Another potential method for selecting which studies to replicate would be to rely on the original P value for studies reporting statistically significant results⁴⁵. Although the prediction accuracy appears to be somewhat lower for original P values than market and survey forecasts⁴⁵, relying on P values may well be considered a practical alternative as it does not involve any additional data collection. Another possibility would be to use predicted replication probabilities from machine learning models to select studies for replication. There has been some progress in developing such models^{120–123}, but evidence on whether they outperform markets or surveys is yet missing. Other potential mechanisms for selecting which studies to replicate include relying on general or study-specific characteristics (for example, connection to theory, surprise factor, sample size, effect size and importance)^{25–28}, relying on cost–benefit considerations^{29,30}, using Bayesian strategies^{31,32}, determining the ‘replication value’³³, adopting empirical audit and review³⁴, selecting studies randomly³⁵ or using predictions from laypeople^{36,37}.

Using decision markets to select the studies with the highest and lowest predicted probabilities for replication is just one of the many potential selection rules for this methodology. Our goal was to test whether a decision market could distinguish findings that would replicate or not, and we aimed to maximize the statistical power of detecting an association between market prices and replication outcomes. For the practical application of decision markets, the choice of the selection mechanism will largely depend on the objective function. One selection rule would be to choose the studies with the highest predicted false positive likelihood, that is, the studies with the smallest market prices (in addition to at least one randomly selected study to

ensure incentive compatibility). This decision mechanism would align with the objective of identifying and correcting false discoveries in the literature to facilitate an efficient allocation of resources for follow-up investigations. Another selection rule would be to replicate the studies with market predictions close to 50%, which reflects the highest possible uncertainty or disagreement regarding the likelihood of the original finding being genuinely true. Providing additional evidence on these claims could maximize the information value of replication studies, as well-powered replications will move the probability that the tested hypothesis is genuinely true towards 0% or 100%.

For our second primary hypothesis, we found strong evidence that original effect sizes are inflated on average compared with replication effect sizes, with a relative estimated average effect size of 45%. This is comparable to previous systematic replication studies, with relative average effect size estimates of 49% in the RPP¹⁵, 59% in the EERP¹⁴ and 54% in the SSRP¹⁷. The replication rate of 54% based on the statistical significance indicator is also similar to previous replication studies, with 36% in the RPP¹⁵, 61% in the EERP¹⁴ and 62% in the SSRP¹⁷. Caution should be exercised when comparing the replication results across these studies: the number of replications in each of the projects is small, only one focal result per paper has been selected for replication, and the particular journals and time periods considered differ. However, the results of all these studies are consistent with a replication rate of about 50% for both the binary statistical significance indicator and the continuous relative effect size indicator; compatible replication results have also been observed in the Many Labs replication projects^{12–14}.

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

Previous systematic replication studies have focused on laboratory experiments rather than online experiments. Concerns have been raised over data quality in online data collections using ‘crowd workers’, as via MTurk^{88–94}, and part of the rationale for zeroing in on experiments conducted via MTurk was that we tend to share these concerns. However, the results of this study do not suggest that replicability is substantively lower for experiments conducted via MTurk compared with experiments conducted in physical laboratories for studies published in top journals; more evidence is needed to draw strong conclusions. Relatedly, the predicted average replicability rate of 57.6% in the decision market—despite widespread concerns about data quality on MTurk—is within the range of replication rate forecasts in previous prediction markets of 56% in the RPP¹⁰, 75% in the EERP¹⁶ and 63% in the SSRP¹⁷. We used IP quality checks^{90,124} to minimize the chances of low-quality participant data (see Methods for details), screening out participants before the random assignment into treatments. In total, across all 26 replications, 29% of the participants who accepted

a ‘human intelligence task’ (HIT) failed the IP check and were excluded (this descriptive result was not preregistered; see Methods for further details). The replication results from our study should thus not be extrapolated to MTurk experiments not using a comparable screening procedure. An important caveat is that although our IP quality checks seem to have been effective in filtering out bots, this may not be the case for artificial responses generated by large language models such as ChatGPT, which could pose a challenge for collecting data online via platforms such as MTurk¹²⁵.

There are several important limitations to our study. A successful replication, on its own, does not provide valid evidence for the tested hypothesis. It goes without saying that inference in replication studies is subject to type-I and type-II errors, just as in original studies. Moreover, a finding can be replicable while being based on an invalid experimental design, leading to biased results. An example of this would be an experimental design that systematically results in more attrition in one experimental treatment, causing selection bias in favour of the tested conjecture⁸⁸. Likewise, a failed replication, on its own, does not provide direct evidence against the tested hypothesis. A finding can be unreplicable and based on an invalid experimental design, leaving the hypothesis untested. Although the replication rate for online experiments in our study appears to be similar to previous laboratory evidence, it does not necessarily imply that online and laboratory experiments provide equally valid evidence of the tested hypotheses.

Another limitation is that we only replicate a single focal result per paper, and the replication outcome does not necessarily generalize to other results reported in the original articles. Furthermore, we only gathered data from one online population using the same experimental design as in the original study. It cannot be ruled out that the difference in timing between the replication studies and the original studies has affected the replication results as a consequence of changes in the composition of the MTurk subject pool or the tested phenomenon having changed over time. Large-scale, multi-site replication studies that collect data across various populations and settings, similar to the Many Labs replication projects^{12–14}, qualify as a promising method to shed light on the heterogeneity of replication effect sizes across populations and designs^{126–128} in future replication work, potentially increasing the strength of evidence for whether the hypothesis supported in the original study is likely true or not. Collaborative networks such as the Psychological Science Accelerator¹²⁹ facilitate multi-site replication studies and can be a door opener to large and diverse samples.

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

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

Methods

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

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

Below, we provide further details on the inclusion criteria, the decision market set-up and the survey, the replications and the replication rate indicators included in the study. The preregistered analyses and tests were divided into descriptive results of the replication rate among the 26 replicated studies and hypothesis tests. The preregistered descriptive results were furthermore divided into primary replication indicators and secondary replication indicators, and the preregistered hypothesis tests were divided into (1) primary hypotheses, (2) secondary hypotheses and (3) exploratory analyses. See Supplementary Notes for more details about the preregistered hypothesis tests and exploratory analyses. We sought ethical approval from the Swedish Ethical Review Authority who had no ethical objections to the decision market part of the project and judged the replication part of the project to not be covered by the Swedish ethical review law (Dnr 2019-06501).

Inclusion criteria for studies

We reviewed all PNAS articles from 2015 to 2018 and searched for the terms Amazon Mechanical Turk, MTurk and Turk. When we began planning our study at the start of 2019, we started reviewing the most recent articles published in PNAS. We then continued to look back in time, year by year, until we reached a sufficiently large number of studies to run a decision market. However, data collection was delayed after some original authors expressed concerns that the Covid-19 pandemic could affect the replication outcomes. We included all social sciences articles that fulfilled our inclusion criteria for (1) the journal and time period, (2) the platform on which the experiment was performed (MTurk), (3) the type of design (between-subjects or within-subject treatment design), (4) the equipment and materials needed to implement the experiment (the experiment had to be logistically feasible for us to implement) and (5) the results reported in the experiment (that there is at least one statistically significant $P < 0.05$ main or interaction effect in the main text). On the basis of the inclusion criteria, we identified 44 articles. After contacting the original authors, we ended up with 41 articles (the three excluded articles^{130–132} involved either software or platforms that no longer existed or methods we were unfamiliar with). In these 41 articles, we identified at least one critical finding that we could replicate. In cases where several studies in the same article fit the inclusion criteria, we randomly picked one of the studies; this was the case for 17 of the 26 replicated studies (Ames and Fiske⁴⁶, Atir and Ferguson⁴⁷, Baldwin and Lammers⁴⁸, Boswell et al.⁵⁰,

Cooney et al.⁵⁶, Genschow et al.⁵⁹, Gheorghiu et al.⁶⁰, Halevy and Halali⁶², Hofstetter et al.⁶⁵, John et al.⁶⁹, Jordan et al.⁷⁰, Klein and O'Brien⁷³, Kouchaki and Gino⁷⁴, McCall et al.⁷⁶, Rai et al.⁸¹, Stern et al.⁸⁴ and Williams et al.⁸⁶). In cases where the (randomly picked) study contained several conditions, we randomly picked which to compare to the control condition. After that, we looked for the central result with $P < 0.05$ for that particular study. If there were several statistically significant results, one was selected at random. The replication results thus only pertain to the single central result selected per paper, and the replication outcome does not necessarily generalize to other results reported in the original articles. For convenience, we refer to the replications as 'replication of [study reference]' though.

For Cheon and Hong⁵⁴, the result chosen for replication is reported as part of a 2×2 ANOVA in the original article; since the paper does not report the main effect, the original authors kindly provided us with the corresponding estimates. For Gheorghiu et al.⁶⁰, the result to be replicated is only reported with its P value in the paper; a precise estimate of the test statistic has been obtained from a re-analysis of the original data, which the original authors kindly provided. For the study by Kraus et al.⁷⁵, we could not reproduce the result reported in the original article using the original data. The original authors acknowledged that there had been a reporting error in the original article. For the replication, we use the analysis described in the paper; the effect size and the test statistic reported in the original paper were replaced by the re-estimated result. For the study by Williams et al.⁸⁶, the focal hypothesis test in the replication is based on a composite score of five suites of behaviour (which are tested separately in the original article) to have a single test. The original authors also report tests on composite measures in the Supplementary Information of their article, and they approved the choice to investigate the replicability of the focal hypothesis using a composite score. These changes are transparently reported in the replication reports for each study (see <https://osf.io/sejyp> for details).

Decision market and prediction survey

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

In the survey, we asked participants to assess, for each replication study, (1) the likelihood that the hypothesis will successfully replicate (on a scale from 0% to 100%); (2) their stated expertise for the study/

the hypothesis (on a scale from 1 to 7); and (3) whether they believe that the pandemic has affected the likelihood of replication. The question about the pandemic was measured on a scale from -3 ('the pandemic has definitely decreased the probability of replication') to 3 ('the pandemic has definitely increased the probability of replication'); the 0 midpoint implies that they do not think that the pandemic has affected the probability of replication. The survey was not incentivized.

The decision market opened on 8 November (2021) and closed after 2 weeks on 22 November (and before the decision market opened, participants had at least 1 week to complete the prediction survey). In the decision market, participants could trade (bet) on whether they expected the 41 studies to replicate. While participants had the opportunity to bet on the replication outcome of the 41 studies, we did not carry out replications for all 41 studies, but the final decision market prices determined which studies to replicate. We replicated the 12 studies that had the highest and the 12 studies that had the lowest market prices when the market closed. In addition, 2 out of the remaining 17 studies were randomly selected for replication to ensure a non-zero probability for each study to be replicated (that is, we replicated $12 + 12 + 2 = 26$ studies in total). Since payoffs are only determined based on forecasts of studies that were eventually replicated, payoffs were scaled up by the inverse of their probability of being selected for replication in the decision rule (see below for details). This incentive scheme encourages trading based on traders' true beliefs, even though some studies will not be replicated. Consequently, participants have the incentive to buy shares of a particular study whenever they believe that the likelihood of replication is higher than the current market price; likewise, participants have an incentive to (short) sell shares whenever they believe that the likelihood of replication is lower than the current market price. Thus, as long as the market price differs from the predicted likelihood of replication for a participant, the participant has an incentive to buy or (short) sell shares of a particular study and realizing a trade according to a trader's belief will move the market price in the direction of the trader's belief. The decision rule for which studies to replicate was based on final market prices and was common knowledge to the market participants; the instructions (provided to participants who completed the prediction survey by the due date) are available at <https://osf.io/a24zq/>.

We chose 12 studies with the lowest predicted probability and 12 studies with the highest predicted probability based on a power calculation using the pooled data from our previous prediction market studies⁴⁵. The power calculations were conducted by randomly sampling 41 studies from the dataset described in Gordon et al.⁴⁵ in a simulation with 10,000 iterations and then selecting the forecasts and outcomes from the 12 studies with the lowest predicted probability, the 12 studies with the highest predicted probability and 2 random studies. We failed to set a random seed for the simulation when the study was conducted, implying that the preregistered power estimates could not be numerically reproduced when we wrote up the study results. For full transparency, we report the power estimates included in the PAP in parentheses below. The median point-biserial correlation coefficient across the 10,000 runs is 0.671 (reported as 0.66 in the PAP), and we have 91.0% power (reported as >90% in the PAP) to detect a statistically significant correlation ($n = 26$) between the decision market prices and the replication outcomes at the 0.5% level and 99.4% power (reported as >95% in the PAP) to detect a statistically significant correlation at the 5% level, which is our first primary hypothesis test. As a secondary hypothesis test, we test if the fraction of studies that successfully replicate differs between the 12 studies with the highest and the 12 studies with the lowest predicted replication probabilities using Fisher's exact test. Applying the same sampling approach as for primary hypothesis 1, the median difference in replication rates between the 12 studies with the highest and the 12 studies with the lowest market prices is 0.663; the secondary test ($n = 24$) has 66.5% power at the 0.5% level (reported as 66% in the PAP) and 94.9% power

at the 5% level (reported as 95% in the PAP). The code for the power simulations is available at <https://osf.io/47drs>.

Implementation of the decision market. We used a web-based trading platform, similar to the ones used in Camerer et al.^{16,17} and identical to the one used in Botvinik-Nezer et al.⁶. The trading platform involves two main views: (1) the market overview and (2) the trading page. The market overview listed the 41 assets (that is, one corresponding to each study) in tabular format, including information on the current price for buying a share and the number of shares held (separated for long and short positions). Via the trading page, which was shown after clicking on any of the assets, participants could make investment decisions (that is, buy or sell shares) and view price developments in graphical format for the particular asset.

Trading and incentivization. Decision market participants received an endowment of 100 tokens corresponding to USD 50. Once the markets opened, market participants could use the tokens to trade shares of the assets available in the market. An automated market maker, implementing a logarithmic market scoring rule¹³³, determined the assets' prices. At the beginning of the markets, all assets were valued at 0.50 tokens per share. The market maker calculated the share price for each infinitesimal transaction and updated the price based on the scoring rule. With this mechanism, participants had incentives to invest according to their beliefs^{43,44}. With the logarithmic scoring rule, the price p for an infinitesimal trade is determined as $p = e^{s/b} \div (e^{s/b} + 1)$, where s denotes the net sales (shares held - shares borrowed) that the market maker has done so far in a market; the liquidity parameter b determines how strongly the market price is affected by trade and was set to $b = 100$, implying that by investing ten tokens, traders could move the price of a single asset from 0.50 to about 0.55. We opted for the same value of b as the one used in the prediction markets in the SSRP¹⁷, which appears intuitively sensible in terms of striking a good balance between price sensitivity and liquidity. Notwithstanding, it is worth noting that it is unclear whether or not our results are sensitive to the choice of this parameter. Decision market participants were paid only for studies chosen for replication (based on their final holdings). Participants received one token per correct share for the replications with the 12 lowest and 12 highest final market prices. For the two randomly selected replications, participants received $17 \div 2 = 8.5$ tokens for each share; for replications that were not chosen for replication, participants received no compensation for their holdings. We followed this procedure to keep information revelation in the decision market incentive compatible, with the increased payouts for the randomly selected studies compensating for the 'voided' shares in studies not chosen for replication. Participants were paid after all 26 replications had been completed.

Participation. A total of 193 participants completed the prediction survey (a prerequisite to participate in market trading) after providing consent to participate and were subsequently invited to trade on the decision market. Of these 193 participants, 162 (83.9%) traded in the market at least once. During the 2-week trading period, a total of 4,412 transactions were recorded. On average, each trader prompted 27.2 transactions (s.d. = 30.7; min = 1, max = 185). The average number of traders per hypothesis was 65.1 (s.d. = 15.3; min = 35, max = 98); the average number of transactions recorded per hypothesis was 107.6 (s.d. = 35.2; min = 56, max = 213). See Supplementary Table 2 for descriptive statistics on the trading activity for each market.

Replications

We carry out close replications¹⁰⁷ as closely as possible following the experimental design, sample restrictions, exclusion criteria and analysis as used in the original studies and carried out in the same population (Amazon Mechanical Turk) as the original studies.

The replications started in January 2022 and were completed in October 2023. The replications were planned and preregistered by five replication teams: a team at CalTech, LMU and Wharton; a team at the Stockholm School of Economics; a team at the National University of Singapore; a team at the University of Amsterdam; and a team at the University of Innsbruck.

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

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

The individual replication studies sometimes also used additional exclusion criteria that are detailed in the preregistered replication report for each replication (we tried to use the same exclusion criteria for the replications as used in the original studies as much as possible). The replication sample sizes defined below are the sample sizes after any exclusions of participants.

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

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

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

Incentivization in the replication experiments. We standardized payments across all replications such that studies had a certain show-up fee depending on the expected length of the study. In particular, we paid an hourly fee of USD 8.00 for all studies, and we calculated the show-up fee for each study based on the expected length of the study. For all studies, we implemented a minimum payoff of USD 1.00. For studies with incentive payments, we used the same incentive payment as in the original study, paid in addition to the show-up fee. If we faced problems in recruiting participants, we increased the show-up fee, which happened for two studies^{61,65}.

Replication indicators

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

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

Small-telescopes approach (secondary indicator). We also used the small-telescopes approach¹¹². For this indicator, we estimated whether the replication effect size was significantly smaller (using a

one-sided test at the 5% level) than a ‘small effect’, defined as the effect size the original study would have had 33% power to detect. For studies using *t*-tests (or *F*-tests converted to a *t*-test statistic), we based ‘the small effect size’ on the effect size that a *t*-test had 33% power to detect (at the 5% level in a two-sided test); for studies using *z*-test statistics (or chi-square tests converted to a *z*-test statistic), we based ‘the small effect size’ on the effect size that a *z*-test had 33% power to detect (at the 5% level in a two-sided test). To test whether the replication effect size was significantly smaller than ‘the small effect size’ in a one-sided test at the 5% level, we estimated a 90% CI of the replication effect size. We tested if the 90% CI overlapped the small effect size with CIs constructed as described in Supplementary Notes. If the effect size in the replication was significantly smaller than this ‘small effect size’, the result was considered a failed replication; otherwise, it was considered successful. We report the fraction of studies that replicate according to this criterion and the 95% Clopper–Pearson CI of this fraction. The small-telescopes results are reported in Fig. 4 and Supplementary Table 4.

Bayes factors (secondary indicators). We also compute the one-sided default Bayes factors on the replication data, allowing us to obtain the strength of evidence in favour of the hypothesis that stipulates an effect in the direction of the original experiment (where a default prior in terms of a truncated Cauchy distribution with scale 0.707 was assigned to the size of the effect) versus the null hypothesis that stipulates the effect to be absent¹¹³. In addition, we also computed (one-sided) replication Bayes factors, which quantifies the additional evidence for the hypothesis given the evidence already provided by the original study¹¹⁴. (We are counting the one-sided default and replication as Bayes factors as two separate indicators, which they are.) These results are reported in Fig. 5 and Supplementary Table 4. We use the evidence categories proposed by Jeffreys¹¹⁵ to interpret the Bayes factors. A detailed report on the estimation of the Bayes factors is available at <https://osf.io/47drs/>.

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

Reporting summary

Further information on research design is available in the Nature Portfolio Reporting Summary linked to this article.

Data availability

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

Code availability

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

References

- Ioannidis, J. P. A. Why most published research findings are false. *PLoS Med.* **2**, e124 (2005).

- Leamer, E. E. Let’s take the con out of econometrics. *Am. Econ. Rev.* **73**, 31–43 (1983).
- Begley, C. G. & Ellis, L. M. Raise standards for preclinical cancer research. *Nature* **483**, 531–533 (2012).
- McNutt, M. Reproducibility. *Science* **343**, 229–229 (2014).
- Gertler, P., Galiani, S. & Romero, M. How to make replication the norm. *Nature* **554**, 417–419 (2018).
- Botvinik-Nezer, R. et al. Variability in the analysis of a single neuroimaging dataset by many teams. *Nature* **582**, 84–88 (2020).
- Breznau, N. et al. Observing many researchers using the same data and hypothesis reveals a hidden universe of uncertainty. *Proc. Natl Acad. Sci. USA* **119**, e2203150119 (2022).
- Delios, A. et al. Examining the generalizability of research findings from archival data. *Proc. Natl Acad. Sci. USA* **119**, e2120377119 (2022).
- Huber, C. et al. Competition and moral behavior: a meta-analysis of forty-five crowd-sourced experimental designs. *Proc. Natl Acad. Sci. USA* **120**, e2215572120 (2023).
- Dreber, A. et al. Using prediction markets to estimate the reproducibility of scientific research. *Proc. Natl Acad. Sci. USA* **112**, 15343–15347 (2015).
- Maniadi, Z., Tufano, F. & List, J. A. To replicate or not to replicate? Exploring reproducibility in economics through the lens of a model and a pilot study. *Econ. J.* **127**, F209–F235 (2017).
- Klein, R. A. et al. Investigating variation in replicability: a ‘many labs’ replication project. *Soc. Psychol.* **45**, 142–152 (2014).
- 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).
- Klein, R. A. et al. Many Labs 2: investigating variation in replicability across samples and settings. *Adv. Methods Pract. Psychol. Sci.* **1**, 443–490 (2018).
- Open Science Collaboration Estimating the reproducibility of psychological science. *Science* **349**, aac4716 (2015).
- Camerer, C. F. et al. Evaluating replicability of laboratory experiments in economics. *Science* **351**, 1433–1436 (2016).
- Camerer, C. F. et al. Evaluating the replicability of social science experiments in *Nature* and *Science* between 2010 and 2015. *Nat. Hum. Behav.* **2**, 637–644 (2018).
- Ioannidis, J. P. A. Why most discovered true associations are inflated. *Epidemiology* **19**, 640 (2008).
- Button, K. S. et al. Power failure: why small sample size undermines the reliability of neuroscience. *Nat. Rev. Neurosci.* **14**, 365–376 (2013).
- Wegener, D. T., Fabrigar, L. R., Pek, J. & Hoisington-Shaw, K. Evaluating research in personality and social psychology: considerations of statistical power and concerns about false findings. *Pers. Soc. Psychol. Bull.* **48**, 1105–1117 (2022).
- 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).
- John, L. K., Loewenstein, G. & Prelec, D. Measuring the prevalence of questionable research practices with incentives for truth telling. *Psychol. Sci.* **23**, 524–532 (2012).
- Finkel, E. J., Eastwick, P. W. & Reis, H. T. Replicability and other features of a high-quality science: toward a balanced and empirical approach. *J. Pers. Soc. Psychol.* **113**, 244–253 (2017).
- Flake, J. K., Davidson, I. J., Wong, O. & Pek, J. Construct validity and the validity of replication studies: a systematic review. *Am. Psychol.* **77**, 576–588 (2022).
- Pittelkow, M.-M. et al. The process of replication target selection in psychology: what to consider? *R. Soc. Open Sci.* **10**, 210586 (2023).

26. Makel, M. C., Plucker, J. A. & Hegarty, B. Replications in psychology research: how often do they really occur? *Perspect. Psychol. Sci.* **7**, 537–542 (2012).
27. Lindsay, D. S. Replication in psychological science. *Psychol. Sci.* **26**, 1827–1832 (2015).
28. Block, J. & Kuckertz, A. Seven principles of effective replication studies: strengthening the evidence base of management research. *Manag. Rev. Q.* **68**, 355–359 (2018).
29. Coles, N. A., Tiokhin, L., Scheel, A. M., Isager, P. M. & Lakens, D. The costs and benefits of replication studies. *Behav. Brain Sci.* **41**, e124 (2018).
30. Alipourfard, N. et al. Systematizing confidence in open research and evidence (SCORE). Preprint at SocArXiv <https://doi.org/10/hn4g> (2021).
31. Hardwicke, T. E., Tessler, M. H., Peloquin, B. N. & Frank, M. C. A Bayesian decision-making framework for replication. *Behav. Brain Sci.* **41**, e132 (2018).
32. Field, S. M., Hoekstra, R., Bringmann, L. & van Ravenzwaaij, D. When and why to replicate: as easy as 1, 2, 3? *Collabra Psychol.* **5**, 46 (2019).
33. Isager, P. M. et al. Deciding what to replicate: a decision model for replication study selection under resource and knowledge constraints. *Psychol. Methods* **28**, 438–451 (2023).
34. O'Donnell, M. et al. Empirical audit and review and an assessment of evidentiary value in research on the psychological consequences of scarcity. *Proc. Natl Acad. Sci. USA* **118**, e2103313118 (2021).
35. Kuehberger, A. & Schulte-Mecklenbeck, M. Selecting target papers for replication. *Behav. Brain Sci.* **41**, e139 (2018).
36. Hoogeveen, S., Sarafoglou, A. & Wagenmakers, E.-J. Laypeople can predict which social-science studies will be replicated successfully. *Adv. Methods Pract. Psychol. Sci.* **3**, 267–285 (2020).
37. Marcoci, A. et al. Predicting the replicability of social and behavioural science claims from the COVID-19 Preprint Replication Project with structured expert and novice groups. Preprint at MetaArXiv <https://doi.org/10.31222/osf.io/xdsjf> (2023).
38. Wolfers, J. & Zitzewitz, E. Prediction markets. *J. Econ. Perspect.* **18**, 107–126 (2004).
39. Arrow, K. J. et al. The promise of prediction markets. *Science* **320**, 877–878 (2008).
40. Tziralis, G. & Tatsiopoulos, I. Prediction markets: an extended literature review. *J. Predict. Mark.* **1**, 75–91 (2012).
41. Hanson, R. Decision markets. *IEEE Intell. Syst.* **14**, 16–19 (1999).
42. Hanson, R. Combinatorial information market design. *Inf. Syst. Front.* **5**, 107–119 (2003).
43. Chen, Y., Kash, I., Ruberry, M. & Shnayder, V. Decision markets with good incentives. In *Proc. Internet and Network Economics* (eds Chen, N. et al.) 72–83 (Springer, 2011).
44. Wang, W. & Pfeiffer, T. Securities based decision markets. In *Proc. Distributed Artificial Intelligence* Vol. 13170 (eds Chen, J. et al.) 79–92 (Springer, 2022).
45. Gordon, M., Viganola, D., Dreber, A., Johannesson, M. & Pfeiffer, T. Predicting replicability—analysis of survey and prediction market data from large-scale forecasting projects. *PLoS ONE* **16**, e0248780 (2021).
46. Ames, D. L. & Fiske, S. T. Perceived intent motivates people to magnify observed harms. *Proc. Natl Acad. Sci. USA* **112**, 3599–3605 (2015).
47. Atir, S. & Ferguson, M. J. How gender determines the way we speak about professionals. *Proc. Natl Acad. Sci. USA* **115**, 7278–7283 (2018).
48. Baldwin, M. & Lammers, J. Past-focused environmental comparisons promote proenvironmental outcomes for conservatives. *Proc. Natl Acad. Sci. USA* **113**, 14953–14957 (2016).
49. Bear, A., Fortgang, R. G., Bronstein, M. V. & Cannon, T. D. Mistiming of thought and perception predicts delusionality. *Proc. Natl Acad. Sci. USA* **114**, 10791–10796 (2017).
50. Boswell, R. G., Sun, W., Suzuki, S. & Kober, H. Training in cognitive strategies reduces eating and improves food choice. *Proc. Natl Acad. Sci. USA* **115**, E11238–E11247 (2018).
51. Caruso, E. M., Burns, Z. C. & Converse, B. A. Slow motion increases perceived intent. *Proc. Natl Acad. Sci. USA* **113**, 9250–9255 (2016).
52. Casella, A., Kartik, N., Sanchez, L. & Turban, S. Communication in context: interpreting promises in an experiment on competition and trust. *Proc. Natl Acad. Sci. USA* **115**, 933–938 (2018).
53. Chao, M. Demotivating incentives and motivation crowding out in charitable giving. *Proc. Natl Acad. Sci. USA* **114**, 7301–7306 (2017).
54. Cheon, B. K. & Hong, Y.-Y. Mere experience of low subjective socioeconomic status stimulates appetite and food intake. *Proc. Natl Acad. Sci. USA* **114**, 72–77 (2017).
55. Clarkson, J. J. et al. The self-control consequences of political ideology. *Proc. Natl Acad. Sci. USA* **112**, 8250–8253 (2015).
56. Cooney, G., Gilbert, D. T. & Wilson, T. D. When fairness matters less than we expect. *Proc. Natl Acad. Sci. USA* **113**, 11168–11171 (2016).
57. Côté, S., House, J. & Willer, R. High economic inequality leads higher-income individuals to be less generous. *Proc. Natl Acad. Sci. USA* **112**, 15838–15843 (2015).
58. Flesch, T., Balaguer, J., Dekker, R., Nili, H. & Summerfield, C. Comparing continual task learning in minds and machines. *Proc. Natl Acad. Sci. USA* **115**, E10313–E10322 (2018).
59. Genschow, O., Rigoni, D. & Brass, M. Belief in free will affects causal attributions when judging others' behavior. *Proc. Natl Acad. Sci. USA* **114**, 10071–10076 (2017).
60. Gheorghiu, A. I., Callan, M. J. & Skylark, W. J. Facial appearance affects science communication. *Proc. Natl Acad. Sci. USA* **114**, 5970–5975 (2017).
61. Guilbeault, D., Becker, J. & Centola, D. Social learning and partisan bias in the interpretation of climate trends. *Proc. Natl Acad. Sci. USA* **115**, 9714–9719 (2018).
62. Halevy, N. & Halali, E. Selfish third parties act as peacemakers by transforming conflicts and promoting cooperation. *Proc. Natl Acad. Sci. USA* **112**, 6937–6942 (2015).
63. Handley, I. M., Brown, E. R., Moss-Racusin, C. A. & Smith, J. L. Quality of evidence revealing subtle gender biases in science is in the eye of the beholder. *Proc. Natl Acad. Sci. USA* **112**, 13201–13206 (2015).
64. Hoffman, K. M., Trawalter, S., Axt, J. R. & Oliver, M. N. Racial bias in pain assessment and treatment recommendations, and false beliefs about biological differences between blacks and whites. *Proc. Natl Acad. Sci. USA* **113**, 4296–4301 (2016).
65. Hofstetter, R., Rüppell, R. & John, L. K. Temporary sharing prompts unrestrained disclosures that leave lasting negative impressions. *Proc. Natl Acad. Sci. USA* **114**, 11902–11907 (2017).
66. Horne, Z., Powell, D., Hummel, J. E. & Holyoak, K. J. Countering antivaccination attitudes. *Proc. Natl Acad. Sci. USA* **112**, 10321–10324 (2015).
67. Isley, S. C., Stern, P. C., Carmichael, S. P., Joseph, K. M. & Arent, D. J. Online purchasing creates opportunities to lower the life cycle carbon footprints of consumer products. *Proc. Natl Acad. Sci. USA* **113**, 9780–9785 (2016).
68. Jachimowicz, J. M., Chafik, S., Munrat, S., Prabhu, J. C. & Weber, E. U. Community trust reduces myopic decisions of low-income individuals. *Proc. Natl Acad. Sci. USA* **114**, 5401–5406 (2017).
69. John, L. K., Barasz, K. & Norton, M. I. Hiding personal information reveals the worst. *Proc. Natl Acad. Sci. USA* **113**, 954–959 (2016).
70. Jordan, J. J., Hoffman, M., Nowak, M. A. & Rand, D. G. Uncalculating cooperation is used to signal trustworthiness. *Proc. Natl Acad. Sci. USA* **113**, 8658–8663 (2016).

71. Jun, Y., Meng, R. & Johar, G. V. Perceived social presence reduces fact-checking. *Proc. Natl Acad. Sci. USA* **114**, 5976–5981 (2017).
72. KC, R. P., Kunter, M. & Mak, V. The influence of a competition on noncompetitors. *Proc. Natl Acad. Sci. USA* **115**, 2716–2721 (2018).
73. Klein, N. & O'Brien, E. People use less information than they think to make up their minds. *Proc. Natl Acad. Sci. USA* **115**, 13222–13227 (2018).
74. Kouchaki, M. & Gino, F. Memories of unethical actions become obfuscated over time. *Proc. Natl Acad. Sci. USA* **113**, 6166–6171 (2016).
75. Kraus, M. W., Rucker, J. M. & Richeson, J. A. Americans misperceive racial economic equality. *Proc. Natl Acad. Sci. USA* **114**, 10324–10331 (2017).
76. McCall, L., Burk, D., Laperrière, M. & Richeson, J. A. Exposure to rising inequality shapes Americans' opportunity beliefs and policy support. *Proc. Natl Acad. Sci. USA* **114**, 9593–9598 (2017).
77. Morris, A., MacGlashan, J., Littman, M. L. & Cushman, F. Evolution of flexibility and rigidity in retaliatory punishment. *Proc. Natl Acad. Sci. USA* **114**, 10396–10401 (2017).
78. Mummolo, J. Militarization fails to enhance police safety or reduce crime but may harm police reputation. *Proc. Natl Acad. Sci. USA* **115**, 9181–9186 (2018).
79. Payne, B. K., Brown-Iannuzzi, J. L. & Hannay, J. W. Economic inequality increases risk taking. *Proc. Natl Acad. Sci. USA* **114**, 4643–4648 (2017).
80. Phillips, J. & Cushman, F. Morality constrains the default representation of what is possible. *Proc. Natl Acad. Sci. USA* **114**, 4649–4654 (2017).
81. Rai, T. S., Valdesolo, P. & Graham, J. Dehumanization increases instrumental violence, but not moral violence. *Proc. Natl Acad. Sci. USA* **114**, 8511–8516 (2017).
82. Reeck, C., Wall, D. & Johnson, E. J. Search predicts and changes patience in intertemporal choice. *Proc. Natl Acad. Sci. USA* **114**, 11890–11895 (2017).
83. Schilke, O., Reimann, M. & Cook, K. S. Power decreases trust in social exchange. *Proc. Natl Acad. Sci. USA* **112**, 12950–12955 (2015).
84. Stern, C., West, T. V. & Rule, N. O. Conservatives negatively evaluate counterstereotypical people to maintain a sense of certainty. *Proc. Natl Acad. Sci. USA* **112**, 15337–15342 (2015).
85. Vacharkulksemsuk, T. et al. Dominant, open nonverbal displays are attractive at zero-acquaintance. *Proc. Natl Acad. Sci. USA* **113**, 4009–4014 (2016).
86. Williams, K. E. G., Sng, O. & Neuberg, S. L. Ecology-driven stereotypes override race stereotypes. *Proc. Natl Acad. Sci. USA* **113**, 310–315 (2016).
87. Schimmelpfennig, R. et al. The moderating role of culture in the generalizability of psychological phenomena. *Adv. Methods Pract. Psychol. Sci.* **7**, 25152459231225163 (2024).
88. Zhou, H. & Fishbach, A. The pitfall of experimenting on the web: how unattended selective attrition leads to surprising (yet false) research conclusions. *J. Pers. Soc. Psychol.* **111**, 493–504 (2016).
89. Chmielewski, M. & Kucker, S. C. An MTurk crisis? Shifts in data quality and the impact on study results. *Soc. Psychol. Pers. Sci.* **11**, 464–473 (2020).
90. Aguinis, H., Villamor, I. & Ramani, R. S. MTurk research: review and recommendations. *J. Manag.* **47**, 823–837 (2021).
91. Brodeur, A., Cook, N. & Heyes, A. *We Need to Talk About Mechanical Turk: What 22,989 Hypothesis Tests Tell Us About Publication Bias and P-Hacking in Online Experiments* Discussion Paper No. 15478 (IZA Institute of Labor Economics, 2022).
92. Peer, E., Rothschild, D., Gordon, A., Evernden, Z. & Damer, E. Data quality of platforms and panels for online behavioral research. *Behav. Res. Methods* **54**, 1643–1662 (2022).
93. Webb, M. A. & Tangney, J. P. Too good to be true: bots and bad data from Mechanical Turk. *Perspect. Psychol. Sci.* <https://doi.org/10.1177/17456916221120027> (2022).
94. Douglas, B. D., Ewell, P. J. & Brauer, M. Data quality in online human-subjects research: comparisons between MTurk, Prolific, CloudResearch, Qualtrics, and SONA. *PLoS ONE* **18**, e0279720 (2023).
95. Abelson, R. P. *Statistics as Principled Argument* (Psychology Press, 1995).
96. Macdonald, R. R. Statistical inference and Aristotle's Rhetoric. *Br. J. Math. Stat. Psychol.* **57**, 193–203 (2004).
97. Scheel, A. M., Schijen, M. R. M. J. & Lakens, D. An excess of positive results: comparing the standard psychology literature with registered reports. *Adv. Methods Pract. Psychol. Sci.* **4**, 25152459211007467 (2021).
98. Soderberg, C. K. et al. Initial evidence of research quality of registered reports compared with the standard publishing model. *Nat. Hum. Behav.* **5**, 990–997 (2021).
99. Brodeur, A., Cook, N., Hartley, J. & Heyes, A. Do pre-registration and pre-analysis plans reduce p-hacking and publication bias? Evidence from 15,992 test statistics and suggestions for improvement. *JPE Micro.* **2**, 527–561 (2024).
100. Yamada, Y. How to crack pre-registration: toward transparent and open science. *Front. Psychol.* **9**, 1831 (2018).
101. Flis, I. The function of literature in psychological science. *Rev. Gen. Psychol.* **26**, 146–156 (2022).
102. Rubin, M. Questionable metascience practices. *JOTE* <https://doi.org/10.36850/mr4> (2023).
103. Wagenmakers, E.-J., Wetzels, R., Borsboom, D., van der Maas, H. L. J. & Kievit, R. A. An agenda for purely confirmatory research. *Perspect. Psychol. Sci.* **7**, 632–638 (2012).
104. Nosek, B. A., Ebersole, C. R., DeHaven, A. C. & Mellor, D. T. The preregistration revolution. *Proc. Natl Acad. Sci. USA* **115**, 2600–2606 (2018).
105. Maxwell, S. E., Lau, M. Y. & Howard, G. S. Is psychology suffering from a replication crisis? What does 'failure to replicate' really mean? *Am. Psychol.* **70**, 487–498 (2015).
106. Shrout, P. E. & Rodgers, J. L. Psychology, science, and knowledge construction: broadening perspectives from the replication crisis. *Annu. Rev. Psychol.* **69**, 487–510 (2018).
107. Dreber, A. & Johannesson, M. A framework for evaluating reproducibility and replicability in economics. *Econ. Inq.* <https://doi.org/10.1111/ecin.13244> (2024).
108. Benjamin, D. J. et al. Redefine statistical significance. *Nat. Hum. Behav.* **2**, 6–10 (2018).
109. Barnard, G. A. Significance tests for 2 × 2 tables. *Biometrika* **34**, 123–138 (1947).
110. Mehrotra, D. V., Chan, I. S. F. & Berger, R. L. A cautionary note on exact unconditional inference for a difference between two independent binomial proportions. *Biometrics* **59**, 441–450 (2003).
111. Boschloo, R. D. Raised conditional level of significance for the 2x2-table when testing the equality of two probabilities. *Stat. Neerl.* **24**, 1–9 (1970).
112. Simonsohn, U. Small telescopes: detectability and the evaluation of replication results. *Psychol. Sci.* **26**, 559–569 (2015).
113. Ly, A., Verhagen, J. & Wagenmakers, E.-J. Harold Jeffreys's default Bayes factor hypothesis tests: explanation, extension, and application in psychology. *J. Math. Psychol.* **72**, 19–32 (2016).
114. Ly, A., Etz, A., Marsman, M. & Wagenmakers, E.-J. Replication Bayes factors from evidence updating. *Behav. Res. Methods* **51**, 2498–2508 (2019).
115. Jeffreys, H. *The Theory of Probability* (Oxford Univ. Press, 1961).
116. 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).

117. Gelman, A. & Stern, H. The difference between 'significant' and 'not significant' is not itself statistically significant. *Am. Stat.* **60**, 328–331 (2006).
118. Cumming, G. Replication and P intervals: P values predict the future only vaguely, but confidence intervals do much better. *Perspect. Psychol. Sci.* **3**, 286–300 (2008).
119. Muradchianian, J., Hoekstra, R., Kiers, H. & van Ravenzwaaij, D. How best to quantify replication success? A simulation study on the comparison of replication success metrics. *R. Soc. Open Sci.* **8**, 201697 (2021).
120. Altmeld, A. et al. Predicting the replicability of social science lab experiments. *PLoS ONE* **14**, e0225826 (2019).
121. Yang, Y., Youyou, W. & Uzzi, B. Estimating the deep replicability of scientific findings using human and artificial intelligence. *Proc. Natl Acad. Sci. USA* **117**, 10762–10768 (2020).
122. Rajtmajer, S. et al. A synthetic prediction market for estimating confidence in published work. *Proc. AAAI Conf. Artif. Intell.* **36**, 13218–13220 (2022).
123. Youyou, W., Yang, Y. & Uzzi, B. A discipline-wide investigation of the replicability of psychology papers over the past two decades. *Proc. Natl Acad. Sci. USA* **120**, e2208863120 (2023).
124. Agle, J., Xiao, Y., Nolan, R. & Goltzari-Arroyo, L. Quality control questions on Amazon's Mechanical Turk (MTurk): a randomized trial of impact on the USAUDIT, PHQ-9, and GAD-7. *Behav. Res. Methods* **54**, 885–897 (2022).
125. Veselovsky, V., Ribeiro, M. H. & West, R. Artificial artificial intelligence: crowd workers widely use large language models for text production tasks. Preprint at <https://arxiv.org/abs/2306.07899> (2023).
126. Olsson-Collentine, A., Wicherts, J. M. & van Assen, M. A. L. M. Heterogeneity in direct replications in psychology and its association with effect size. *Psychol. Bull.* **146**, 922–940 (2020).
127. Linden, A. H. & Hönemann, J. Heterogeneity of research results: a new perspective from which to assess and promote progress in psychological science. *Perspect. Psychol. Sci.* **16**, 358–376 (2021).
128. Holzmeister, F. et al. Heterogeneity in effect size estimates. *Proc. Natl Acad. Sci. USA* **121**, e2403490121 (2024).
129. Moshontz, H. et al. The Psychological Science Accelerator: advancing psychology through a distributed collaborative network. *Adv. Methods Pract. Psychol. Sci.* **1**, 501–515 (2018).
130. Epstein, R. & Robertson, R. E. The search engine manipulation effect (SEME) and its possible impact on the outcomes of elections. *Proc. Natl Acad. Sci. USA* **112**, E4512–E4521 (2015).
131. Gallo, E. & Yan, C. The effects of reputational and social knowledge on cooperation. *Proc. Natl Acad. Sci. USA* **112**, 3647–3652 (2015).
132. Li, V., Michael, E., Balaguer, J., Hecce Castañón, S. & Summerfield, C. Gain control explains the effect of distraction in human perceptual, cognitive, and economic decision making. *Proc. Natl Acad. Sci. USA* **115**, E8825–E8834 (2018).
133. Hanson, R. Logarithmic market scoring rules for modular combinatorial information aggregation. *J. Predict. Mark.* **1**, 3–15 (2007).

Acknowledgements

We thank A. Andevall for helping with the data collection and programming of experiments and R. Willer for helpful advice on defining the IP address check and exclusion criteria used to exclude individuals from participating to minimize low-quality participant data. For financial support, we thank the Austrian Science FWF

(grant SFB F63 to J.H. and M.K.), Jan Wallander and Tom Hedelius Foundation (grants P21-0091 and P23-0098 to A.D.), Knut and Alice Wallenberg Foundation and Marianne and Marcus Wallenberg Foundation (Wallenberg Scholar grant to A.D.) and Riksbankens Jubileumsfond (grant P21-0168 to M.J.). The funders had no role in study design, data collection and analysis, decision to publish or preparation of the paper. One author (V.W.) is currently employed by Sveriges Riksbank but did this work before being employed by Sveriges Riksbank; the opinions expressed in this article are the sole responsibility of the authors and should not be interpreted as reflecting the views of Sveriges Riksbank.

Author contributions

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 the study. Y.C., A.D., F.H., M.J. and T.P. designed and implemented the decision market. A.D., F.H., M.J., B.M. and V.W. selected articles and critical findings for (potential) replication. A.D., 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. and V.W. designed the replications and collected replication data. F.H., S.H., T.I., L.J., A.S., R.S. and V.W. conducted the preregistered statistical tests on the individual replications. A.L. computed the Bayes factors. F.H. conducted all analyses reported in the paper. A.D., F.H. and M.J. wrote the paper. All authors reviewed and approved the final paper.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41562-024-02062-9>.

Correspondence and requests for materials should be addressed to Anna Dreber.

Peer review information *Nature Human Behaviour* thanks Michael Varnum and the other, anonymous, reviewer(s) for their contribution to the peer review of this work.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2024

Reporting Summary

Nature Portfolio wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Portfolio policies, see our [Editorial Policies](#) and the [Editorial Policy Checklist](#).

Statistics

For all statistical analyses, confirm that the following items are present in the figure legend, table legend, main text, or Methods section.

n/a | Confirmed

- The exact sample size (n) for each experimental group/condition, given as a discrete number and unit of measurement
- A statement on whether measurements were taken from distinct samples or whether the same sample was measured repeatedly
- The statistical test(s) used AND whether they are one- or two-sided
Only common tests should be described solely by name; describe more complex techniques in the Methods section.
- A description of all covariates tested
- A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons
- A full description of the statistical parameters including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals)
- For null hypothesis testing, the test statistic (e.g. F , t , r) with confidence intervals, effect sizes, degrees of freedom and P value noted
Give P values as exact values whenever suitable.
- For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings
- For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes
- Estimates of effect sizes (e.g. Cohen's d , Pearson's r), indicating how they were calculated

Our web collection on [statistics for biologists](#) contains articles on many of the points above.

Software and code

Policy information about [availability of computer code](#)

Data collection Qualtrics (qsf file from 2021) based prediction survey (files available on OSF). We replicated 26 original studies using the software of the original study whenever possible; in the replications where we used other software this is stated in the SI and the Replication Report for each replication (the Replication Reports are available at OSF at <https://osf.io/sk82q>).

Data analysis We have posted code for all data analyses carried out in the Replication Reports for each replication and for all the analyses in the manuscript and SI at OSF (<https://osf.io/sk82q>). In the paper there is a separate code availability entitled section with the OSF web-link.

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors and reviewers. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Portfolio [guidelines for submitting code & software](#) for further information.

Data

Policy information about [availability of data](#)

All manuscripts must include a [data availability statement](#). This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A description of any restrictions on data availability
- For clinical datasets or third party data, please ensure that the statement adheres to our [policy](#)

The data availability statement is provided in the manuscript, and contains the following: "The data reported in this paper is tabulated in Supplementary Tables 1–

8. The replication reports (both the pre-replication and the post-replication versions), the pre-analysis plan, the data from the survey and the decision market, and the data for each of the 26 replications are available at the project's OSF repository (<https://osf.io/sk82q>). The separate code availability section contains the following: "The analysis scripts generating all results, figures, and tables reported in the main text and the Supplementary Information are available at the project's OSF repository (<https://osf.io/sk82q>)."

Research involving human participants, their data, or biological material

Policy information about studies with [human participants or human data](#). See also policy information about [sex, gender \(identity/presentation\), and sexual orientation](#) and [race, ethnicity and racism](#).

Reporting on sex and gender

Some of the replications involved collecting socio-economic data because the original studies included these variables, including sex and gender, but these are not used in any analysis and all shared data has been anonymized.

Reporting on race, ethnicity, or other socially relevant groupings

Some of the replications involved collecting socio-economic data because the original studies included these variables, including race or ethnicity, but these are not used in any analysis and all shared data has been anonymized.

Population characteristics

Some of the replications involved collecting socio-economic data like age (if the original studies had included such data) but this is not used in any analysis and all shared data has been anonymized. All participants were adults.

Recruitment

Decision market participants were recruited through public mailing lists (ESA and JDM lists) and social media (e.g., Twitter/X); we also emailed colleagues asking them to distribute the call for participants within their professional networks. This sample is thus self-selected and it is not clear to what extent the results would generalize to other samples of researchers making predictions. Participants in the replication studies were recruited through Amazon Mechanical Turk, similarly to the original studies. The self-selection on Amazon Mechanical Turk in the replications is thus similar to the self-selection in the original studies.

Ethics oversight

We sought ethical approval from the Swedish Ethical Review Authority who had no ethical objections to the decision market part of the project and judged the replication part of the project to not be covered by the Swedish ethical review law (Dnr 2019-06501).

Note that full information on the approval of the study protocol must also be provided in the manuscript.

Field-specific reporting

Please select the one below that is the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

Life sciences

Behavioural & social sciences

Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see nature.com/documents/nr-reporting-summary-flat.pdf

Life sciences study design

All studies must disclose on these points even when the disclosure is negative.

Sample size

Describe how sample size was determined, detailing any statistical methods used to predetermine sample size OR if no sample-size calculation was performed, describe how sample sizes were chosen and provide a rationale for why these sample sizes are sufficient.

Data exclusions

Describe any data exclusions. If no data were excluded from the analyses, state so OR if data were excluded, describe the exclusions and the rationale behind them, indicating whether exclusion criteria were pre-established.

Replication

Describe the measures taken to verify the reproducibility of the experimental findings. If all attempts at replication were successful, confirm this OR if there are any findings that were not replicated or cannot be reproduced, note this and describe why.

Randomization

Describe how samples/organisms/participants were allocated into experimental groups. If allocation was not random, describe how covariates were controlled OR if this is not relevant to your study, explain why.

Blinding

Describe whether the investigators were blinded to group allocation during data collection and/or analysis. If blinding was not possible, describe why OR explain why blinding was not relevant to your study.

Behavioural & social sciences study design

All studies must disclose on these points even when the disclosure is negative.

Study description

This is a quantitative study: we have a decision market, a survey with researchers and we perform 26 replication studies of experimental studies in the social sciences using quantitative methods (with pre-registration of all the replications).

Research sample

For the decision markets and survey we recruited participants who were required to hold a Ph.D. degree or to be a Ph.D. student. The sample is unlikely to be representative of researchers. We did not collect demographic information such as age, sex etc. For the

	replication studies we recruited participants through Amazon Mechanical Turk. The replication samples are similar to the ones used in the original studies and are described in the Replication Reports for each replication posted at OSF (https://osf.io/sk82q). Most of the replication studies do not contain information on demographic information, and the samples are unlikely to be representative of the general population.
Sampling strategy	For the decision markets and survey this is not applicable. We rely on self-selection of researchers who choose to participate in the decision markets. We did not predetermine a sample size; no sample size calculation was performed for the sample size of the participants (but we do have a power calculation in our pre-analysis plan for the number of studies included in the decision markets) and we aimed to include as many participants as possible. For the replication studies, the sampling strategy is similar to the one used in the original studies and is described in the Replication Reports for each replication posted at OSF (https://osf.io/sk82q). The sampling procedure for the replications is convenience sampling. For the replications, we determined sample sizes ex ante and pre-registered the sample sizes we needed to have 90% power to detect 2/3 of the original effect sizes at the 5% level (two-sided test).
Data collection	The data collection for all replications was done as similarly as possible to the data collection in the original studies and are described in detail in the Replication Report for each replication posted at OSF (https://osf.io/sk82q). All data collection was done online and participants in the replications were blinded to the experimental conditions and study hypotheses.
Timing	The data collection for the decision markets and survey with researchers was done in October-November 2021. The data collection for the replications were done between January 2022 and October 2023.
Data exclusions	For the survey with researchers we exclude those that did not participate in the decision markets (this was pre-registered). 193 participants completed the survey by the due date and were invited to participate in the markets. 162 of these 193 participants traded in the markets at least once and these 162 participants are thus our analysis sample also for the survey. As pre-registered, we used the same criteria for data exclusions as in the original studies and any deviations from this are stated in the Replication Reports for each replication posted at OSF (https://osf.io/sk82q). In addition, as pre-registered we excluded participants who failed an IP check (this resulted in 29% of participants who accepted to participate in the replications to be excluded).
Non-participation	289 participants signed up to participate in the survey and were forwarded the link to the survey. 193 participants completed the survey, and 162 participants participated in the decision markets. Those that dropped out at various stages did not provide a reason for it. In the 26 replications, there is variation across the replications in how many dropped out, and they did not provide reasons for dropping out. Any non-participation is detailed in the Replication Reports for each replication posted at OSF (https://osf.io/sk82q).
Randomization	For the decision market and survey there is no randomization. Participants in the replication studies were randomized to conditions.

Ecological, evolutionary & environmental sciences study design

All studies must disclose on these points even when the disclosure is negative.

Study description	<i>Briefly describe the study. For quantitative data include treatment factors and interactions, design structure (e.g. factorial, nested, hierarchical), nature and number of experimental units and replicates.</i>
Research sample	<i>Describe the research sample (e.g. a group of tagged <i>Passer domesticus</i>, all <i>Stenocereus thurberi</i> within Organ Pipe Cactus National Monument), and provide a rationale for the sample choice. When relevant, describe the organism taxa, source, sex, age range and any manipulations. State what population the sample is meant to represent when applicable. For studies involving existing datasets, describe the data and its source.</i>
Sampling strategy	<i>Note the sampling procedure. Describe the statistical methods that were used to predetermine sample size OR if no sample-size calculation was performed, describe how sample sizes were chosen and provide a rationale for why these sample sizes are sufficient.</i>
Data collection	<i>Describe the data collection procedure, including who recorded the data and how.</i>
Timing and spatial scale	<i>Indicate the start and stop dates of data collection, noting the frequency and periodicity of sampling and providing a rationale for these choices. If there is a gap between collection periods, state the dates for each sample cohort. Specify the spatial scale from which the data are taken</i>
Data exclusions	<i>If no data were excluded from the analyses, state so OR if data were excluded, describe the exclusions and the rationale behind them, indicating whether exclusion criteria were pre-established.</i>
Reproducibility	<i>Describe the measures taken to verify the reproducibility of experimental findings. For each experiment, note whether any attempts to repeat the experiment failed OR state that all attempts to repeat the experiment were successful.</i>
Randomization	<i>Describe how samples/organisms/participants were allocated into groups. If allocation was not random, describe how covariates were controlled. If this is not relevant to your study, explain why.</i>
Blinding	<i>Describe the extent of blinding used during data acquisition and analysis. If blinding was not possible, describe why OR explain why blinding was not relevant to your study.</i>

Did the study involve field work? Yes No

Field work, collection and transport

Field conditions	<i>Describe the study conditions for field work, providing relevant parameters (e.g. temperature, rainfall).</i>
Location	<i>State the location of the sampling or experiment, providing relevant parameters (e.g. latitude and longitude, elevation, water depth).</i>
Access & import/export	<i>Describe the efforts you have made to access habitats and to collect and import/export your samples in a responsible manner and in compliance with local, national and international laws, noting any permits that were obtained (give the name of the issuing authority, the date of issue, and any identifying information).</i>
Disturbance	<i>Describe any disturbance caused by the study and how it was minimized.</i>

Reporting for specific materials, systems and methods

We require information from authors about some types of materials, experimental systems and methods used in many studies. Here, indicate whether each material, system or method listed is relevant to your study. If you are not sure if a list item applies to your research, read the appropriate section before selecting a response.

Materials & experimental systems

n/a	Involvement in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> Antibodies
<input checked="" type="checkbox"/>	<input type="checkbox"/> Eukaryotic cell lines
<input checked="" type="checkbox"/>	<input type="checkbox"/> Palaeontology and archaeology
<input checked="" type="checkbox"/>	<input type="checkbox"/> Animals and other organisms
<input checked="" type="checkbox"/>	<input type="checkbox"/> Clinical data
<input checked="" type="checkbox"/>	<input type="checkbox"/> Dual use research of concern
<input checked="" type="checkbox"/>	<input type="checkbox"/> Plants

Methods

n/a	Involvement in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> ChIP-seq
<input checked="" type="checkbox"/>	<input type="checkbox"/> Flow cytometry
<input checked="" type="checkbox"/>	<input type="checkbox"/> MRI-based neuroimaging

Antibodies

Antibodies used	<i>Describe all antibodies used in the study; as applicable, provide supplier name, catalog number, clone name, and lot number.</i>
Validation	<i>Describe the validation of each primary antibody for the species and application, noting any validation statements on the manufacturer's website, relevant citations, antibody profiles in online databases, or data provided in the manuscript.</i>

Eukaryotic cell lines

Policy information about [cell lines and Sex and Gender in Research](#)

Cell line source(s)	<i>State the source of each cell line used and the sex of all primary cell lines and cells derived from human participants or vertebrate models.</i>
Authentication	<i>Describe the authentication procedures for each cell line used OR declare that none of the cell lines used were authenticated.</i>
Mycoplasma contamination	<i>Confirm that all cell lines tested negative for mycoplasma contamination OR describe the results of the testing for mycoplasma contamination OR declare that the cell lines were not tested for mycoplasma contamination.</i>
Commonly misidentified lines (See ICLAC register)	<i>Name any commonly misidentified cell lines used in the study and provide a rationale for their use.</i>

Palaeontology and Archaeology

Specimen provenance	<i>Provide provenance information for specimens and describe permits that were obtained for the work (including the name of the issuing authority, the date of issue, and any identifying information). Permits should encompass collection and, where applicable, export.</i>
Specimen deposition	<i>Indicate where the specimens have been deposited to permit free access by other researchers.</i>

Dating methods

If new dates are provided, describe how they were obtained (e.g. collection, storage, sample pretreatment and measurement), where they were obtained (i.e. lab name), the calibration program and the protocol for quality assurance OR state that no new dates are provided.

Tick this box to confirm that the raw and calibrated dates are available in the paper or in Supplementary Information.

Ethics oversight

Identify the organization(s) that approved or provided guidance on the study protocol, OR state that no ethical approval or guidance was required and explain why not.

Note that full information on the approval of the study protocol must also be provided in the manuscript.

Animals and other research organisms

Policy information about [studies involving animals](#); [ARRIVE guidelines](#) recommended for reporting animal research, and [Sex and Gender in Research](#)

Laboratory animals

For laboratory animals, report species, strain and age OR state that the study did not involve laboratory animals.

Wild animals

Provide details on animals observed in or captured in the field; report species and age where possible. Describe how animals were caught and transported and what happened to captive animals after the study (if killed, explain why and describe method; if released, say where and when) OR state that the study did not involve wild animals.

Reporting on sex

Indicate if findings apply to only one sex; describe whether sex was considered in study design, methods used for assigning sex. Provide data disaggregated for sex where this information has been collected in the source data as appropriate; provide overall numbers in this Reporting Summary. Please state if this information has not been collected. Report sex-based analyses where performed, justify reasons for lack of sex-based analysis.

Field-collected samples

For laboratory work with field-collected samples, describe all relevant parameters such as housing, maintenance, temperature, photoperiod and end-of-experiment protocol OR state that the study did not involve samples collected from the field.

Ethics oversight

Identify the organization(s) that approved or provided guidance on the study protocol, OR state that no ethical approval or guidance was required and explain why not.

Note that full information on the approval of the study protocol must also be provided in the manuscript.

Clinical data

Policy information about [clinical studies](#)

All manuscripts should comply with the ICMJE [guidelines for publication of clinical research](#) and a completed [CONSORT checklist](#) must be included with all submissions.

Clinical trial registration

Provide the trial registration number from ClinicalTrials.gov or an equivalent agency.

Study protocol

Note where the full trial protocol can be accessed OR if not available, explain why.

Data collection

Describe the settings and locales of data collection, noting the time periods of recruitment and data collection.

Outcomes

Describe how you pre-defined primary and secondary outcome measures and how you assessed these measures.

Dual use research of concern

Policy information about [dual use research of concern](#)

Hazards

Could the accidental, deliberate or reckless misuse of agents or technologies generated in the work, or the application of information presented in the manuscript, pose a threat to:

No	Yes	
<input checked="" type="checkbox"/>	<input type="checkbox"/>	Public health
<input checked="" type="checkbox"/>	<input type="checkbox"/>	National security
<input checked="" type="checkbox"/>	<input type="checkbox"/>	Crops and/or livestock
<input checked="" type="checkbox"/>	<input type="checkbox"/>	Ecosystems
<input checked="" type="checkbox"/>	<input type="checkbox"/>	Any other significant area

Experiments of concern

Does the work involve any of these experiments of concern:

No	Yes
<input checked="" type="checkbox"/>	<input type="checkbox"/> Demonstrate how to render a vaccine ineffective
<input checked="" type="checkbox"/>	<input type="checkbox"/> Confer resistance to therapeutically useful antibiotics or antiviral agents
<input checked="" type="checkbox"/>	<input type="checkbox"/> Enhance the virulence of a pathogen or render a nonpathogen virulent
<input checked="" type="checkbox"/>	<input type="checkbox"/> Increase transmissibility of a pathogen
<input checked="" type="checkbox"/>	<input type="checkbox"/> Alter the host range of a pathogen
<input checked="" type="checkbox"/>	<input type="checkbox"/> Enable evasion of diagnostic/detection modalities
<input checked="" type="checkbox"/>	<input type="checkbox"/> Enable the weaponization of a biological agent or toxin
<input checked="" type="checkbox"/>	<input type="checkbox"/> Any other potentially harmful combination of experiments and agents

Plants

Seed stocks	<i>Report on the source of all seed stocks or other plant material used. If applicable, state the seed stock centre and catalogue number. If plant specimens were collected from the field, describe the collection location, date and sampling procedures.</i>
Novel plant genotypes	<i>Describe the methods by which all novel plant genotypes were produced. This includes those generated by transgenic approaches, gene editing, chemical/radiation-based mutagenesis and hybridization. For transgenic lines, describe the transformation method, the number of independent lines analyzed and the generation upon which experiments were performed. For gene-edited lines, describe the editor used, the endogenous sequence targeted for editing, the targeting guide RNA sequence (if applicable) and how the editor was applied.</i>
Authentication	<i>Describe any authentication procedures for each seed stock used or novel genotype generated. Describe any experiments used to assess the effect of a mutation and, where applicable, how potential secondary effects (e.g. second site T-DNA insertions, mosaicism, off-target gene editing) were examined.</i>

ChIP-seq

Data deposition

- Confirm that both raw and final processed data have been deposited in a public database such as [GEO](#).
- Confirm that you have deposited or provided access to graph files (e.g. BED files) for the called peaks.

Data access links <i>May remain private before publication.</i>	<i>For "Initial submission" or "Revised version" documents, provide reviewer access links. For your "Final submission" document, provide a link to the deposited data.</i>
Files in database submission	<i>Provide a list of all files available in the database submission.</i>
Genome browser session (e.g. UCSC)	<i>Provide a link to an anonymized genome browser session for "Initial submission" and "Revised version" documents only, to enable peer review. Write "no longer applicable" for "Final submission" documents.</i>

Methodology

Replicates	<i>Describe the experimental replicates, specifying number, type and replicate agreement.</i>
Sequencing depth	<i>Describe the sequencing depth for each experiment, providing the total number of reads, uniquely mapped reads, length of reads and whether they were paired- or single-end.</i>
Antibodies	<i>Describe the antibodies used for the ChIP-seq experiments; as applicable, provide supplier name, catalog number, clone name, and lot number.</i>
Peak calling parameters	<i>Specify the command line program and parameters used for read mapping and peak calling, including the ChIP, control and index files used.</i>
Data quality	<i>Describe the methods used to ensure data quality in full detail, including how many peaks are at FDR 5% and above 5-fold enrichment.</i>
Software	<i>Describe the software used to collect and analyze the ChIP-seq data. For custom code that has been deposited into a community repository, provide accession details.</i>

Flow Cytometry

Plots

Confirm that:

- The axis labels state the marker and fluorochrome used (e.g. CD4-FITC).
- The axis scales are clearly visible. Include numbers along axes only for bottom left plot of group (a 'group' is an analysis of identical markers).
- All plots are contour plots with outliers or pseudocolor plots.
- A numerical value for number of cells or percentage (with statistics) is provided.

Methodology

Sample preparation

Describe the sample preparation, detailing the biological source of the cells and any tissue processing steps used.

Instrument

Identify the instrument used for data collection, specifying make and model number.

Software

Describe the software used to collect and analyze the flow cytometry data. For custom code that has been deposited into a community repository, provide accession details.

Cell population abundance

Describe the abundance of the relevant cell populations within post-sort fractions, providing details on the purity of the samples and how it was determined.

Gating strategy

Describe the gating strategy used for all relevant experiments, specifying the preliminary FSC/SSC gates of the starting cell population, indicating where boundaries between "positive" and "negative" staining cell populations are defined.

- Tick this box to confirm that a figure exemplifying the gating strategy is provided in the Supplementary Information.

Magnetic resonance imaging

Experimental design

Design type

Indicate task or resting state; event-related or block design.

Design specifications

Specify the number of blocks, trials or experimental units per session and/or subject, and specify the length of each trial or block (if trials are blocked) and interval between trials.

Behavioral performance measures

State number and/or type of variables recorded (e.g. correct button press, response time) and what statistics were used to establish that the subjects were performing the task as expected (e.g. mean, range, and/or standard deviation across subjects).

Acquisition

Imaging type(s)

Specify: functional, structural, diffusion, perfusion.

Field strength

Specify in Tesla

Sequence & imaging parameters

Specify the pulse sequence type (gradient echo, spin echo, etc.), imaging type (EPI, spiral, etc.), field of view, matrix size, slice thickness, orientation and TE/TR/flip angle.

Area of acquisition

State whether a whole brain scan was used OR define the area of acquisition, describing how the region was determined.

Diffusion MRI

Used

Not used

Preprocessing

Preprocessing software

Provide detail on software version and revision number and on specific parameters (model/functions, brain extraction, segmentation, smoothing kernel size, etc.).

Normalization

If data were normalized/standardized, describe the approach(es): specify linear or non-linear and define image types used for transformation OR indicate that data were not normalized and explain rationale for lack of normalization.

Normalization template

Describe the template used for normalization/transformation, specifying subject space or group standardized space (e.g. original Talairach, MNI305, ICBM152) OR indicate that the data were not normalized.

Noise and artifact removal

Describe your procedure(s) for artifact and structured noise removal, specifying motion parameters, tissue signals and physiological signals (heart rate, respiration).

Volume censoring

Define your software and/or method and criteria for volume censoring, and state the extent of such censoring.

Statistical modeling & inference

Model type and settings

Specify type (mass univariate, multivariate, RSA, predictive, etc.) and describe essential details of the model at the first and second levels (e.g. fixed, random or mixed effects; drift or auto-correlation).

Effect(s) tested

Define precise effect in terms of the task or stimulus conditions instead of psychological concepts and indicate whether ANOVA or factorial designs were used.

Specify type of analysis: Whole brain ROI-based Both

Statistic type for inference

Specify voxel-wise or cluster-wise and report all relevant parameters for cluster-wise methods.

(See [Eklund et al. 2016](#))

Correction

Describe the type of correction and how it is obtained for multiple comparisons (e.g. FWE, FDR, permutation or Monte Carlo).

Models & analysis

n/a | Involved in the study

 Functional and/or effective connectivity Graph analysis Multivariate modeling or predictive analysis

Functional and/or effective connectivity

Report the measures of dependence used and the model details (e.g. Pearson correlation, partial correlation, mutual information).

Graph analysis

Report the dependent variable and connectivity measure, specifying weighted graph or binarized graph, subject- or group-level, and the global and/or node summaries used (e.g. clustering coefficient, efficiency, etc.).

Multivariate modeling and predictive analysis

Specify independent variables, features extraction and dimension reduction, model, training and evaluation metrics.