

## RESEARCH ARTICLE SUMMARY

## PSYCHOLOGY

# Estimating the reproducibility of psychological science

Open Science Collaboration\*

**INTRODUCTION:** Reproducibility is a defining feature of science, but the extent to which it characterizes current research is unknown. Scientific claims should not gain credence because of the status or authority of their originator but by the replicability of their supporting evidence. Even research of exemplary quality may have irreproducible empirical findings because of random or systematic error.

**RATIONALE:** There is concern about the rate and predictors of reproducibility, but limited evidence. Potentially problematic practices include selective reporting, selective analysis, and insufficient specification of the conditions necessary or sufficient to obtain the results. Direct replication is the attempt to recreate the conditions believed sufficient for obtaining a pre-

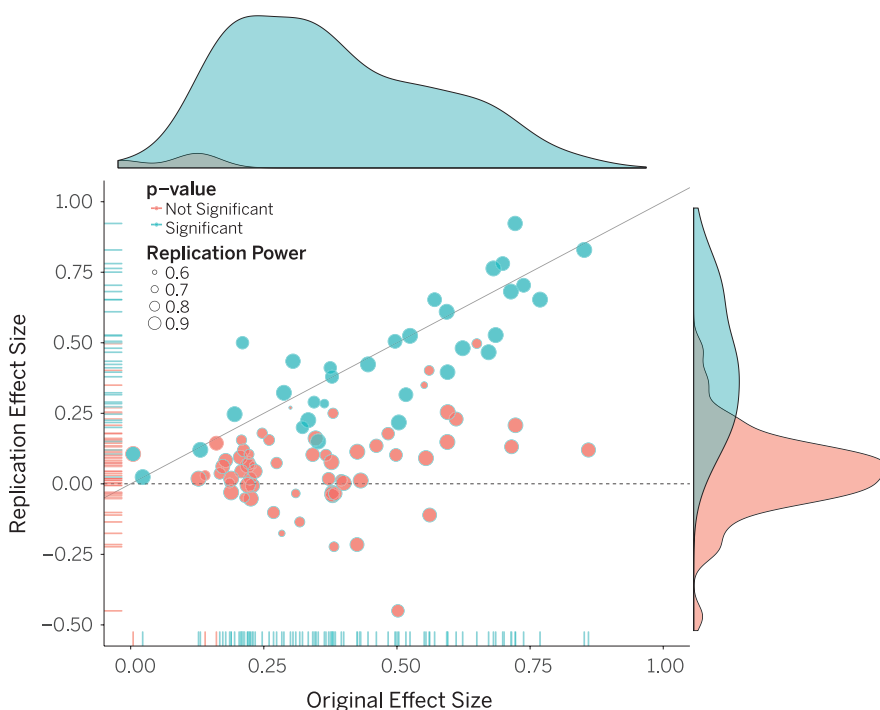
viously observed finding and is the means of establishing reproducibility of a finding with new data. We conducted a large-scale, collaborative effort to obtain an initial estimate of the reproducibility of psychological science.

**RESULTS:** We conducted replications of 100 experimental and correlational studies published in three psychology journals using high-powered designs and original materials when available. There is no single standard for evaluating replication success. Here, we evaluated reproducibility using significance and *P* values, effect sizes, subjective assessments of replication teams, and meta-analysis of effect sizes. The mean effect size ( $r$ ) of the replication effects ( $M_r = 0.197$ ,  $SD = 0.257$ ) was half the magnitude of the mean effect size of the original effects ( $M_o = 0.403$ ,  $SD = 0.188$ ), representing a

substantial decline. Ninety-seven percent of original studies had significant results ( $P < .05$ ). Thirty-six percent of replications had significant results; 47% of original effect sizes were in the 95% confidence interval of the replication effect size; 39% of effects were subjectively rated to have replicated the original result; and if no bias in original results is assumed, combining original and replication results left 68% with statistically significant effects. Correlational tests suggest that replication success was better predicted by the strength of original evidence than by characteristics of the original and replication teams.

**CONCLUSION:** No single indicator sufficiently describes replication success, and the five indicators examined here are not the only ways to evaluate reproducibility. Nonetheless, collectively these results offer a clear conclusion: A large portion of replications produced weaker evidence for the original findings despite using materials provided by the original authors, review in advance for methodological fidelity, and high statistical power to detect the original effect sizes. Moreover, correlational evidence is consistent with the conclusion that variation in the strength of initial evidence (such as original *P* value) was more predictive of replication success than variation in the characteristics of the teams conducting the research (such as experience and expertise). The latter factors certainly can influence replication success, but they did not appear to do so here.

Reproducibility is not well understood because the incentives for individual scientists prioritize novelty over replication. Innovation is the engine of discovery and is vital for a productive, effective scientific enterprise. However, innovative ideas become old news fast. Journal reviewers and editors may dismiss a new test of a published idea as unoriginal. The claim that “we already know this” belies the uncertainty of scientific evidence. Innovation points out paths that are possible; replication points out paths that are likely; progress relies on both. Replication can increase certainty when findings are reproduced and promote innovation when they are not. This project provides accumulating evidence for many findings in psychological research and suggests that there is still more work to do to verify whether we know what we think we know. ■



**Original study effect size versus replication effect size (correlation coefficients).** Diagonal line represents replication effect size equal to original effect size. Dotted line represents replication effect size of 0. Points below the dotted line were effects in the opposite direction of the original. Density plots are separated by significant (blue) and nonsignificant (red) effects.

The list of author affiliations is available in the full article online.  
\*Corresponding author. E-mail: nosek@virginia.edu  
Cite this article as Open Science Collaboration, *Science* 349, aac4716 (2015). DOI: 10.1126/science.aac4716

## RESEARCH ARTICLE

## PSYCHOLOGY

# Estimating the reproducibility of psychological science

Open Science Collaboration\*†

Reproducibility is a defining feature of science, but the extent to which it characterizes current research is unknown. We conducted replications of 100 experimental and correlational studies published in three psychology journals using high-powered designs and original materials when available. Replication effects were half the magnitude of original effects, representing a substantial decline. Ninety-seven percent of original studies had statistically significant results. Thirty-six percent of replications had statistically significant results; 47% of original effect sizes were in the 95% confidence interval of the replication effect size; 39% of effects were subjectively rated to have replicated the original result; and if no bias in original results is assumed, combining original and replication results left 68% with statistically significant effects. Correlational tests suggest that replication success was better predicted by the strength of original evidence than by characteristics of the original and replication teams.

Reproducibility is a core principle of scientific progress (1–6). Scientific claims should not gain credence because of the status or authority of their originator but by the replicability of their supporting evidence. Scientists attempt to transparently describe the methodology and resulting evidence used to support their claims. Other scientists agree or disagree whether the evidence supports the claims, citing theoretical or methodological reasons or by collecting new evidence. Such debates are meaningless, however, if the evidence being debated is not reproducible.

Even research of exemplary quality may have irreproducible empirical findings because of random or systematic error. Direct replication is the attempt to recreate the conditions believed sufficient for obtaining a previously observed finding (7, 8) and is the means of establishing reproducibility of a finding with new data. A direct replication may not obtain the original result for a variety of reasons: Known or unknown differences between the replication and original study may moderate the size of an observed effect, the original result could have been a false positive, or the replication could produce a false negative. False positives and false negatives provide misleading information about effects, and failure to identify the necessary and sufficient conditions to reproduce a finding indicates an incomplete theoretical understanding. Direct replication provides the opportunity to assess and improve reproducibility.

There is plenty of concern (9–13) about the rate and predictors of reproducibility but limited evidence. In a theoretical analysis, Ioannidis estimated that publishing and analytic practices make it likely that more than half of research

results are false and therefore irreproducible (9). Some empirical evidence supports this analysis. In cell biology, two industrial laboratories reported success replicating the original results of landmark findings in only 11 and 25% of the attempted cases, respectively (10, 11). These numbers are stunning but also difficult to interpret because no details are available about the studies, methodology, or results. With no transparency, the reasons for low reproducibility cannot be evaluated.

Other investigations point to practices and incentives that may inflate the likelihood of obtaining false-positive results in particular or irreproducible results more generally. Potentially problematic practices include selective reporting, selective analysis, and insufficient specification of the conditions necessary or sufficient to obtain the results (12–23). We were inspired to address the gap in direct empirical evidence about reproducibility. In this Research Article, we report a large-scale, collaborative effort to obtain an initial estimate of the reproducibility of psychological science.

## Method

Starting in November 2011, we constructed a protocol for selecting and conducting high-quality replications (24). Collaborators joined the project, selected a study for replication from the available studies in the sampling frame, and were guided through the replication protocol. The replication protocol articulated the process of selecting the study and key effect from the available articles, contacting the original authors for study materials, preparing a study protocol and analysis plan, obtaining review of the protocol by the original authors and other members within the present project, registering the protocol publicly, conducting the replication, writing the final report, and auditing the process and analysis for quality control. Project coordinators

facilitated each step of the process and maintained the protocol and project resources. Replication materials and data were required to be archived publicly in order to maximize transparency, accountability, and reproducibility of the project (<https://osf.io/ezcuq>).

In total, 100 replications were completed by 270 contributing authors. There were many different research designs and analysis strategies in the original research. Through consultation with original authors, obtaining original materials, and internal review, replications maintained high fidelity to the original designs. Analyses converted results to a common effect size metric [correlation coefficient ( $r$ )] with confidence intervals (CIs). The units of analysis for inferences about reproducibility were the original and replication study effect sizes. The resulting open data set provides an initial estimate of the reproducibility of psychology and correlational data to support development of hypotheses about the causes of reproducibility.

## Sampling frame and study selection

We constructed a sampling frame and selection process to minimize selection biases and maximize generalizability of the accumulated evidence. Simultaneously, to maintain high quality, within this sampling frame we matched individual replication projects with teams that had relevant interests and expertise. We pursued a quasi-random sample by defining the sampling frame as 2008 articles of three important psychology journals: *Psychological Science* (PSCD), *Journal of Personality and Social Psychology* (JPSP), and *Journal of Experimental Psychology: Learning, Memory, and Cognition* (JEP: LMC). The first is a premier outlet for all psychological research; the second and third are leading disciplinary-specific journals for social psychology and cognitive psychology, respectively [more information is available in (24)]. These were selected a priori in order to (i) provide a tractable sampling frame that would not plausibly bias reproducibility estimates, (ii) enable comparisons across journal types and sub-disciplines, (iii) fit with the range of expertise available in the initial collaborative team, (iv) be recent enough to obtain original materials, (v) be old enough to obtain meaningful indicators of citation impact, and (vi) represent psychology sub-disciplines that have a high frequency of studies that are feasible to conduct at relatively low cost.

The first replication teams could select from a pool of the first 20 articles from each journal, starting with the first article published in the first 2008 issue. Project coordinators facilitated matching articles with replication teams by interests and expertise until the remaining articles were difficult to match. If there were still interested teams, then another 10 articles from one or more of the three journals were made available from the sampling frame. Further, project coordinators actively recruited teams from the community with relevant experience for particular articles. This approach balanced competing goals: minimizing selection bias by

\*All authors with their affiliations appear at the end of this paper.

†Corresponding author. E-mail: [nosek@virginia.edu](mailto:nosek@virginia.edu)

having only a small set of articles available at a time and matching studies with replication teams' interests, resources, and expertise.

By default, the last experiment reported in each article was the subject of replication. This decision established an objective standard for study selection within an article and was based on the intuition that the first study in a multiple-study article (the obvious alternative selection strategy) was more frequently a preliminary demonstration. Deviations from selecting the last experiment were made occasionally on the basis of feasibility or recommendations of the original authors. Justifications for deviations were reported in the replication reports, which were made available on the Open Science Framework (OSF) (<http://osf.io/ezcu>). In total, 84 of the 100 completed replications (84%) were of the last reported study in the article. On average, the to-be-replicated articles contained 2.99 studies ( $SD = 1.78$ ) with the following distribution: 24 single study, 24 two studies, 18 three studies, 13 four studies, 12 five studies, 9 six or more studies. All following summary statistics refer to the 100 completed replications.

For the purposes of aggregating results across studies to estimate reproducibility, a key result from the selected experiment was identified as the focus of replication. The key result had to be represented as a single statistical inference test or an effect size. In most cases, that test was a  $t$  test,  $F$  test, or correlation coefficient. This effect was identified before data collection or analysis and was presented to the original authors as part of the design protocol for critique. Original authors occasionally suggested that a different effect be used, and by default, replication teams deferred to original authors' judgments. Nonetheless, because the single effect came from a single study, it is not necessarily the case that the identified effect was central to the overall aims of the article. In the individual replication reports and subjective assessments of replication outcomes, more than a single result could be examined, but only the result of the single effect was considered in the aggregate analyses [additional details of the general protocol and individual study methods are provided in the supplementary materials and (25)].

In total, there were 488 articles in the 2008 issues of the three journals. One hundred fifty-eight of these (32%) became eligible for selection for replication during the project period, between November 2011 and December 2014. From those, 111 articles (70%) were selected by a replication team, producing 113 replications. Two articles had two replications each (supplementary materials). And 100 of those (88%) replications were completed by the project deadline for inclusion in this aggregate report. After being claimed, some studies were not completed because the replication teams ran out of time or could not devote sufficient resources to completing the study. By journal, replications were completed for 39 of 64 (61%) articles from *PSCI*, 31 of 55 (56%) articles from *JPSP*, and 28 of 39 (72%) articles from *JEP:LMC*.

The most common reasons for failure to match an article with a team were feasibility constraints for conducting the research. Of the 47 articles from the eligible pool that were not claimed, six (13%) had been deemed infeasible to replicate because of time, resources, instrumentation, dependence on historical events, or hard-to-access samples. The remaining 41 (87%) were eligible but not claimed. These often required specialized samples (such as macaques or people with autism), resources (such as eye tracking machines or functional magnetic resonance imaging), or knowledge making them difficult to match with teams.

### Aggregate data preparation

Each replication team conducted the study, analyzed their data, wrote their summary report, and completed a checklist of requirements for sharing the materials and data. Then, independent reviewers and analysts conducted a project-wide audit of all individual projects, materials, data, and reports. A description of this review is available on the OSF (<https://osf.io/xtine>). Moreover, to maximize reproducibility and accuracy, the analyses for every replication study were reproduced by another analyst independent of the replication team using the R statistical programming language and a standardized analytic format. A controller R script was created to regenerate the entire analysis of every study and recreate the master data file. This R script, available at <https://osf.io/fkmwg>, can be executed to reproduce the results of the individual studies. A comprehensive description of this reanalysis process is available publicly (<https://osf.io/a2eyg>).

### Measures and moderators

We assessed features of the original study and replication as possible correlates of reproducibility and conducted exploratory analyses to inspire further investigation. These included characteristics of the original study such as the publishing journal; original effect size,  $P$  value, and sample size; experience and expertise of the original research team; importance of the effect, with indicators such as the citation impact of the article; and rated surprisingness of the effect. We also assessed characteristics of the replication such as statistical power and sample size, experience and expertise of the replication team, independently assessed challenge of conducting an effective replication, and self-assessed quality of the replication effort. Variables such as the  $P$  value indicate the statistical strength of evidence given the null hypothesis, and variables such as "effect surprisingness" and "expertise of the team" indicate qualities of the topic of study and the teams studying it, respectively. The master data file, containing these and other variables, is available for exploratory analysis (<https://osf.io/5wup8>).

It is possible to derive a variety of hypotheses about predictors of reproducibility. To reduce the likelihood of false positives due to many tests, we aggregated some variables into summary indicators: experience and expertise of original team, experience and expertise of replication team, challenge of replication, self-assessed quality of repli-

cation, and importance of the effect. We had no a priori justification to give some indicators stronger weighting over others, so aggregates were created by standardizing [mean ( $M$ ) = 0,  $SD = 1$ ] the individual variables and then averaging to create a single index. In addition to the publishing journal and subdiscipline, potential moderators included six characteristics of the original study and five characteristics of the replication (supplementary materials).

### Publishing journal and subdiscipline

Journals' different publishing practices may result in a selection bias that covaries with reproducibility. Articles from three journals were made available for selection: *JPSP* ( $n = 59$  articles), *JEP:LMC* ( $n = 40$  articles), and *PSCI* ( $n = 68$  articles). From this pool of available studies, replications were selected and completed from *JPSP* ( $n = 32$  studies), *JEP:LMC* ( $n = 28$  studies), and *PSCI* ( $n = 40$  studies) and were coded as representing cognitive ( $n = 43$  studies) or social-personality ( $n = 57$  studies) subdisciplines. Four studies that would ordinarily be understood as "developmental psychology" because of studying children or infants were coded as having a cognitive or social emphasis. Reproducibility may vary by subdiscipline in psychology because of differing practices. For example, within-subjects designs are more common in cognitive than social psychology, and these designs often have greater power to detect effects with the same number of participants.

### Statistical analyses

There is no single standard for evaluating replication success (25). We evaluated reproducibility using significance and  $P$  values, effect sizes, subjective assessments of replication teams, and meta-analyses of effect sizes. All five of these indicators contribute information about the relations between the replication and original finding and the cumulative evidence about the effect and were positively correlated with one another ( $r$  ranged from 0.22 to 0.96, median  $r = 0.57$ ). Results are summarized in Table 1, and full details of analyses are in the supplementary materials.

### Significance and $P$ values

Assuming a two-tailed test and significance or  $\alpha$  level of 0.05, all test results of original and replication studies were classified as statistically significant ( $P \leq 0.05$ ) and nonsignificant ( $P > 0.05$ ). However, original studies that interpreted nonsignificant  $P$  values as significant were coded as significant (four cases, all with  $P$  values  $< 0.06$ ). Using only the nonsignificant  $P$  values of the replication studies and applying Fisher's method (26), we tested the hypothesis that these studies had "no evidential value" (the null hypothesis of zero-effect holds for all these studies). We tested the hypothesis that the proportions of statistically significant results in the original and replication studies are equal using the McNemar test for paired nominal data and calculated a CI of the reproducibility parameter. Second, we compared the central tendency of the distribution of  $P$  values of original and

**Table 1. Summary of reproducibility rates and effect sizes for original and replication studies overall and by journal/discipline.** *df/N* refers to the information on which the test of the effect was based (for example, *df* of *t* test, denominator *df* of *F* test, sample size  $-3$  of correlation, and sample size for *z* and  $\chi^2$ ). Four original results had *P* values slightly higher than 0.05 but were considered positive results in the original article and are treated that way here. Exclusions (explanation provided in supplementary materials, A3) are “replications  $P < 0.05$ ” (3 original nulls excluded;  $n = 97$  studies); “mean original and replication effect sizes” (3 excluded;  $n = 97$  studies); “meta-analytic mean estimates” (27 excluded;  $n = 73$  studies); “percent meta-analytic ( $P < 0.05$ )” (25 excluded;  $n = 75$  studies); and, “percent original effect size within replication 95% CI” (5 excluded,  $n = 95$  studies).

	Effect size comparison					Original and replication combined						
	Replications $P < 0.05$ in original direction	Percent	Mean (SD) original effect size	Median original <i>df/N</i>	Mean (SD) replication effect size	Median replication <i>df/N</i>	Average replication power	Meta- analytic mean (SD) estimate	Percent meta- analytic ( $P < 0.05$ )	Percent original effect size within replication 95% CI	Percent subjective “yes” to “Did it replicate?”	
Overall	35/97	36	0.403 (0.188)	54	0.197 (0.257)	68	0.92	0.309 (0.223)	68	47	39	
<i>JPSP</i> , social	7/31	23	0.29 (0.10)	73	0.07 (0.11)	120	0.91	0.138 (0.087)	43	34	25	
<i>JEP:LMC</i> , cognitive	13/27	48	0.47 (0.18)	36.5	0.27 (0.24)	43	0.93	0.393 (0.209)	86	62	54	
<i>PSCI</i> , social	7/24	29	0.39 (0.20)	76	0.21 (0.30)	122	0.92	0.286 (0.228)	58	40	32	
<i>PSCI</i> , cognitive	8/15	53	0.53 (0.2)	23	0.29 (0.35)	21	0.94	0.464 (0.221)	92	60	53	

**Table 2. Spearman's rank-order correlations of reproducibility indicators with summary original and replication study characteristics.** Effect size difference computed after converting *r* to Fisher's *z*. *df/N* refers to the information on which the test of the effect was based (for example, *df* of *t* test, denominator *df* of *F* test, sample size  $-3$  of correlation, and sample size for *z* and  $\chi^2$ ). Four original results had *P* values slightly higher than 0.05 but were considered positive results in the original article and are treated that way here. Exclusions (explanation provided in supplementary materials, A3) are “replications  $P < .05$ ” (3 original nulls excluded;  $n = 97$  studies), “effect size difference” (3 excluded;  $n = 97$  studies); “meta-analytic mean estimates” (27 excluded;  $n = 73$  studies); and, “percent original effect size within replication 95% CI” (5 excluded,  $n = 95$  studies).

	Replications $P < 0.05$ in original direction	Effect size difference	Meta-analytic estimate	Original effect size within replication 95% CI	Subjective “yes” to “Did it replicate?”
Original study characteristics					
Original <i>P</i> value	-0.327	-0.057	-0.468	0.032	-0.260
Original effect size	0.304	0.279	0.793	0.121	0.277
Original <i>df/N</i>	-0.150	-0.194	-0.502	-0.221	-0.185
Importance of original result	-0.105	0.038	-0.205	-0.133	-0.074
Surprising original result	-0.244	0.102	-0.181	-0.113	-0.241
Experience and expertise of original team	-0.072	-0.033	-0.059	-0.103	-0.044
Replication characteristics					
Replication <i>P</i> value	-0.828	0.621	-0.614	-0.562	-0.738
Replication effect size	0.731	-0.586	0.850	0.611	0.710
Replication power	0.368	-0.053	0.142	-0.056	0.285
Replication <i>df/N</i>	-0.085	-0.224	-0.692	-0.257	-0.164
Challenge of conducting replication	-0.219	0.085	-0.301	-0.109	-0.151
Experience and expertise of replication team	-0.096	0.133	0.017	-0.053	-0.068
Self-assessed quality of replication	-0.069	0.017	0.054	-0.088	-0.055

replication studies using the Wilcoxon signed-rank test and the *t* test for dependent samples. For both tests, we only used study-pairs for which both *P* values were available.

#### Effect sizes

We transformed effect sizes into correlation coefficients whenever possible. Correlation coefficients have several advantages over other effect size measures, such as Cohen's *d*. Correlation coefficients are bounded, well known, and therefore more readily interpretable. Most important for our purposes, analysis of correlation coefficients is straightforward because, after ap-

plying the Fisher transformation, their standard error is only a function of sample size. Formulas and code for converting test statistics *z*, *F*, *t*, and  $\chi^2$  into correlation coefficients are provided in the appendices at <http://osf.io/ezum7>. To be able to compare and analyze correlations across study-pairs, the original study's effect size was coded as positive; the replication study's effect size was coded as negative if the replication study's effect was opposite to that of the original study.

We compared effect sizes using four tests. We compared the central tendency of the effect size distributions of original and replication studies using both a paired two-sample *t* test and the

Wilcoxon signed-rank test. Third, we computed the proportion of study-pairs in which the effect of the original study was stronger than in the replication study and tested the hypothesis that this proportion is 0.5. For this test, we included findings for which effect size measures were available but no correlation coefficient could be computed (for example, if a regression coefficient was reported but not its test statistic). Fourth, we calculated “coverage,” or the proportion of study-pairs in which the effect of the original study was in the CI of the effect of the replication study, and compared this with the expected proportion using a goodness-of-fit  $\chi^2$  test. We carried

out this test on the subset of study pairs in which both the correlation coefficient and its standard error could be computed [we refer to this data set as the meta-analytic (MA) subset]. Standard errors could only be computed if test statistics were  $r$ ,  $t$ , or  $F(1,df_2)$ . The expected proportion is the sum over expected probabilities across study-pairs. The test assumes the same population effect size for original and replication study in the same study-pair. For those studies that tested the effect with  $F(df_1 > 1, df_2)$  or  $\chi^2$ , we verified coverage using other statistical procedures (computational details are provided in the supplementary materials).

#### Meta-analysis combining original and replication effects

We conducted fixed-effect meta-analyses using the R package metafor (27) on Fisher-transformed correlations for all study-pairs in subset MA and on study-pairs with the odds ratio as the dependent variable. The number of times the CI of all these meta-analyses contained 0 was calculated. For studies in the MA subset, estimated effect sizes were averaged and analyzed by discipline.

#### Subjective assessment of “Did it replicate?”

In addition to the quantitative assessments of replication and effect estimation, we collected subjective assessments of whether the replication provided evidence of replicating the original result. In some cases, the quantitative data anticipate a straightforward subjective assessment of replication. For more complex designs, such as multivariate interaction effects, the quantitative analysis may not provide a simple interpretation. For subjective assessment, replication teams answered “yes” or “no” to the question, “Did your results replicate the original effect?” Additional subjective variables are available for analysis in the full data set.

#### Analysis of moderators

We correlated the five indicators evaluating reproducibility with six indicators of the origi-

nal study (original  $P$  value, original effect size, original sample size, importance of the effect, surprising effect, and experience and expertise of original team) and seven indicators of the replication study (replication  $P$  value, replication effect size, replication power based on original effect size, replication sample size, challenge of conducting replication, experience and expertise of replication team, and self-assessed quality of replication) (Table 2). As follow-up, we did the same with the individual indicators comprising the moderator variables (tables S3 and S4).

## Results

### Evaluating replication effect against null hypothesis of no effect

A straightforward method for evaluating replication is to test whether the replication shows a statistically significant effect ( $P < 0.05$ ) with the same direction as the original study. This dichotomous vote-counting method is intuitively appealing and consistent with common heuristics used to decide whether original studies “worked.” Ninety-seven of 100 (97%) effects from original studies were positive results (four had  $P$  values falling a bit short of the 0.05 criterion— $P = 0.0508, 0.0514, 0.0516, \text{ and } 0.0567$ —but all of these were interpreted as positive effects). On the basis of only the average replication power of the 97 original, significant effects [ $M = 0.92$ , median ( $Mdn$ ) = 0.95], we would expect approximately 89 positive results in the replications if all original effects were true and accurately estimated; however, there were just 35 [36.1%; 95% CI = (26.6%, 46.2%)], a significant reduction [McNemar test,  $\chi^2(1) = 59.1, P < 0.001$ ].

A key weakness of this method is that it treats the 0.05 threshold as a bright-line criterion between replication success and failure (28). It could be that many of the replications fell just short of the 0.05 criterion. The density plots of  $P$  values for original studies (mean  $P$  value = 0.028) and replications (mean  $P$  value = 0.302) are shown in Fig. 1, left. The 64 nonsignificant

$P$  values for replications were distributed widely. When there is no effect to detect, the null distribution of  $P$  values is uniform. This distribution deviated slightly from uniform with positive skew, however, suggesting that at least one replication could be a false negative,  $\chi^2(128) = 155.83, P = 0.048$ . Nonetheless, the wide distribution of  $P$  values suggests against insufficient power as the only explanation for failures to replicate. A scatterplot of original compared with replication study  $P$  values is shown in Fig. 2.

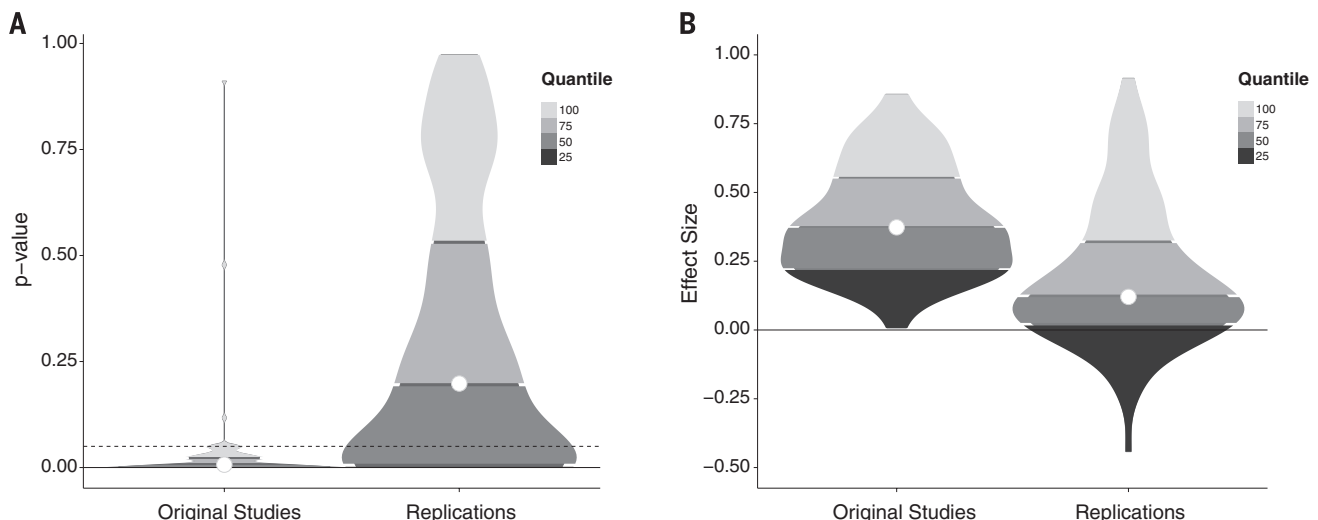
### Evaluating replication effect against original effect size

A complementary method for evaluating replication is to test whether the original effect size is within the 95% CI of the effect size estimate from the replication. For the subset of 73 studies in which the standard error of the correlation could be computed, 30 (41.1%) of the replication CIs contained the original effect size (significantly lower than the expected value of 78.5%,  $P < 0.001$ ) (supplementary materials). For 22 studies using other test statistics [ $F(df_1 > 1, df_2)$  and  $\chi^2$ ], 68.2% of CIs contained the effect size of the original study. Overall, this analysis suggests a 47.4% replication success rate.

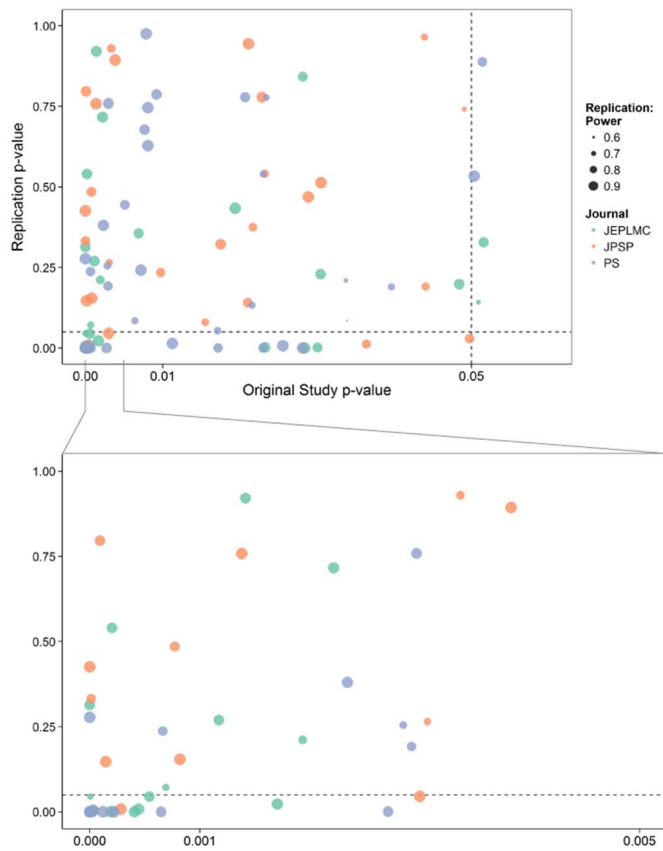
This method addresses the weakness of the first test that a replication in the same direction and a  $P$  value of 0.06 may not be significantly different from the original result. However, the method will also indicate that a replication “fails” when the direction of the effect is the same but the replication effect size is significantly smaller than the original effect size (29). Also, the replication “succeeds” when the result is near zero but not estimated with sufficiently high precision to be distinguished from the original effect size.

### Comparing original and replication effect sizes

Comparing the magnitude of the original and replication effect sizes avoids special emphasis on  $P$  values. Overall, original study effect sizes



**Fig. 1. Density plots of original and replication  $P$  values and effect sizes.** (A)  $P$  values. (B) Effect sizes (correlation coefficients). Lowest quantiles of  $P$  values are not visible because they are clustered near zero.



**Fig. 2. Scatterplots of original study and replication  $P$  values for three psychology journals.**

Data points scaled by power of the replication based on original study effect size. Dotted red lines indicate  $P = 0.05$  criterion. Subplot below shows  $P$  values from the range between the gray lines ( $P = 0$  to  $0.005$ ) in the main plot above.

( $M = 0.403$ ,  $SD = 0.188$ ) were reliably larger than replication effect sizes ( $M = 0.197$ ,  $SD = 0.257$ ), Wilcoxon's  $W = 7137$ ,  $P < 0.001$ . Of the 99 studies for which an effect size in both the original and replication study could be calculated (30), 82 showed a stronger effect size in the original study (82.8%;  $P < 0.001$ , binomial test) (Fig. 1, right). Original and replication effect sizes were positively correlated (Spearman's  $r = 0.51$ ,  $P < 0.001$ ). A scatterplot of the original and replication effect sizes is presented in Fig. 3.

### Combining original and replication effect sizes for cumulative evidence

The disadvantage of the descriptive comparison of effect sizes is that it does not provide information about the precision of either estimate or resolution of the cumulative evidence for the effect. This is often addressed by computing a meta-analytic estimate of the effect sizes by combining the original and replication studies (28). This approach weights each study by the inverse of its variance and uses these weighted estimates of effect size to estimate cumulative evidence and precision of the effect. Using a fixed-effect model, 51 of the 75 (68%) effects for which a meta-analytic estimate could be computed had 95% CIs that did not include 0.

One qualification about this result is the possibility that the original studies have inflated effect sizes due to publication, selection, reporting, or other biases (9, 12–23). In a discipline with low-powered research designs and an emphasis on positive results for publication, effect sizes will be systematically overestimated in the published literature. There is no publication bias in the replication studies because all results are reported. Also, there are no selection or reporting biases because all were confirmatory tests based on pre-analysis plans. This maximizes the interpretability of the replication  $P$  values and effect estimates. If publication, selection, and reporting biases completely explain the effect differences, then the replication estimates would be a better estimate of the effect size than would the meta-analytic and original results. However, to the extent that there are other influences, such as moderation by sample, setting, or quality of replication, the relative bias influencing original and replication effect size estimation is unknown.

### Subjective assessment of “Did it replicate?”

In addition to the quantitative assessments of replication and effect estimation, replication teams

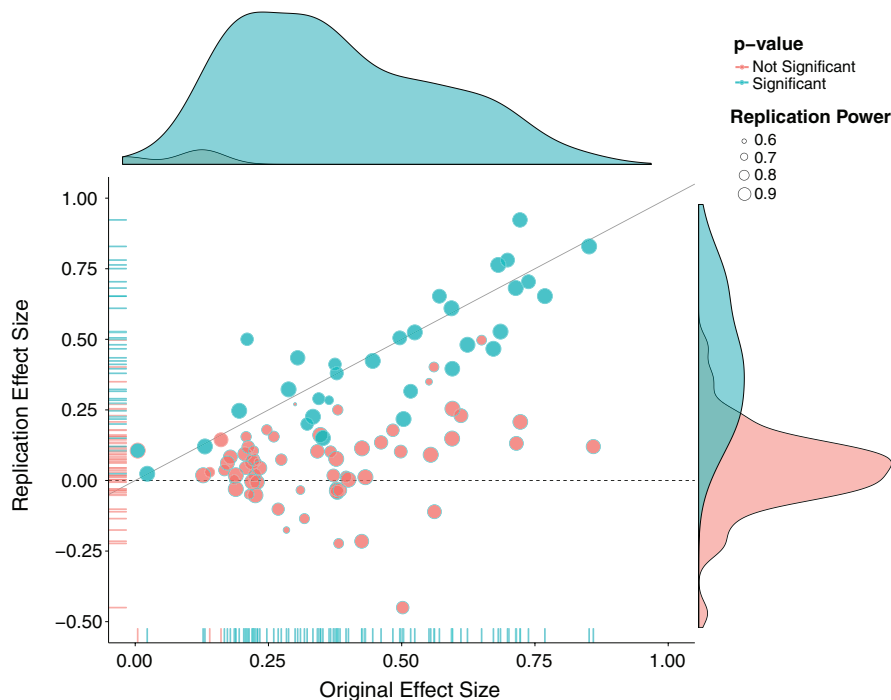
provided a subjective assessment of replication success of the study they conducted. Subjective assessments of replication success were very similar to significance testing results (39 of 100 successful replications), including evaluating “success” for two null replications when the original study reported a null result and “failure” for a  $P < 0.05$  replication when the original result was a null.

### Correlates of reproducibility

The overall replication evidence is summarized in Table 1 across the criteria described above and then separately by journal/discipline. Considering significance testing, reproducibility was stronger in studies and journals representing cognitive psychology than social psychology topics. For example, combining across journals, 14 of 55 (25%) of social psychology effects replicated by the  $P < 0.05$  criterion, whereas 21 of 42 (50%) of cognitive psychology effects did so. Simultaneously, all journals and disciplines showed substantial and similar [ $\chi^2(3) = 2.45$ ,  $P = 0.48$ ] declines in effect size in the replications compared with the original studies. The difference in significance testing results between fields appears to be partly a function of weaker original effects in social psychology studies, particularly in *JPSP*, and perhaps of the greater frequency of high-powered within-subjects manipulations and repeated measurement designs in cognitive psychology as suggested by high power despite relatively small participant samples. Further, the type of test was associated with replication success. Among original, significant effects, 23 of the 49 (47%) that tested main or simple effects replicated at  $P < 0.05$ , but just 8 of the 37 (22%) that tested interaction effects did.

Correlations between reproducibility indicators and characteristics of replication and original studies are provided in Table 2. A negative correlation of replication success with the original study  $P$  value indicates that the initial strength of evidence is predictive of reproducibility. For example, 26 of 63 (41%) original studies with  $P < 0.02$  achieved  $P < 0.05$  in the replication, whereas 6 of 23 (26%) that had a  $P$  value between  $0.02 < P < 0.04$  and 2 of 11 (18%) that had a  $P$  value  $> 0.04$  did so (Fig. 2). Almost two thirds (20 of 32, 63%) of original studies with  $P < 0.001$  had a significant  $P$  value in the replication.

Larger original effect sizes were associated with greater likelihood of achieving  $P < 0.05$  ( $r = 0.304$ ) and a greater effect size difference between original and replication ( $r = 0.279$ ). Moreover, replication power was related to replication success via significance testing ( $r = 0.368$ ) but not with the effect size difference between original and replication ( $r = -0.053$ ). Comparing effect sizes across indicators, surprisingness of the original effect, and the challenge of conducting the replication were related to replication success for some indicators. Surprising effects were less reproducible, as were effects for which it was more challenging to conduct the replication. Last, there was little evidence that perceived importance of the effect, expertise of the original or replication teams, or self-assessed quality of the replication accounted for meaningful variation



**Fig. 3. Original study effect size versus replication effect size (correlation coefficients).** Diagonal line represents replication effect size equal to original effect size. Dotted line represents replication effect size of 0. Points below the dotted line were effects in the opposite direction of the original. Density plots are separated by significant (blue) and nonsignificant (red) effects.

in reproducibility across indicators. Replication success was more consistently related to the original strength of evidence (such as original  $P$  value, effect size, and effect tested) than to characteristics of the teams and implementation of the replication (such as expertise, quality, or challenge of conducting study) (tables S3 and S4).

## Discussion

No single indicator sufficiently describes replication success, and the five indicators examined here are not the only ways to evaluate reproducibility. Nonetheless, collectively, these results offer a clear conclusion: A large portion of replications produced weaker evidence for the original findings (31) despite using materials provided by the original authors, review in advance for methodological fidelity, and high statistical power to detect the original effect sizes. Moreover, correlational evidence is consistent with the conclusion that variation in the strength of initial evidence (such as original  $P$  value) was more predictive of replication success than was variation in the characteristics of the teams conducting the research (such as experience and expertise). The latter factors certainly can influence replication success, but the evidence is that they did not systematically do so here. Other investigators may develop alternative indicators to explore further the role of expertise and quality in reproducibility on this open data set.

### Insights on reproducibility

It is too easy to conclude that successful replication means that the theoretical understanding of

the original finding is correct. Direct replication mainly provides evidence for the reliability of a result. If there are alternative explanations for the original finding, those alternatives could likewise account for the replication. Understanding is achieved through multiple, diverse investigations that provide converging support for a theoretical interpretation and rule out alternative explanations.

It is also too easy to conclude that a failure to replicate a result means that the original evidence was a false positive. Replications can fail if the replication methodology differs from the original in ways that interfere with observing the effect. We conducted replications designed to minimize a priori reasons to expect a different result by using original materials, engaging original authors for review of the designs, and conducting internal reviews. Nonetheless, unanticipated factors in the sample, setting, or procedure could still have altered the observed effect magnitudes (32).

More generally, there are indications of cultural practices in scientific communication that may be responsible for the observed results. Lower-power research designs combined with publication bias favoring positive results together produce a literature with upwardly biased effect sizes (14, 16, 33, 34). This anticipates that replication effect sizes would be smaller than original studies on a routine basis—not because of differences in implementation but because the original study effect sizes are affected by publication and reporting bias, and the replications are not. Consistent with this expectation, most replication effects were

smaller than original results, and reproducibility success was correlated with indicators of the strength of initial evidence, such as lower original  $P$  values and larger effect sizes. This suggests publication, selection, and reporting biases as plausible explanations for the difference between original and replication effects. The replication studies significantly reduced these biases because replication preregistration and pre-analysis plans ensured confirmatory tests and reporting of all results.

The observed variation in replication and original results may reduce certainty about the statistical inferences from the original studies but also provides an opportunity for theoretical innovation to explain differing outcomes, and then new research to test those hypothesized explanations. The correlational evidence, for example, suggests that procedures that are more challenging to execute may result in less reproducible results, and that more surprising original effects may be less reproducible than less surprising original effects. Further, systematic, repeated replication efforts that fail to identify conditions under which the original finding can be observed reliably may reduce confidence in the original finding.

### Implications and limitations

The present study provides the first open, systematic evidence of reproducibility from a sample of studies in psychology. We sought to maximize generalizability of the results with a structured process for selecting studies for replication. However, it is unknown the extent to which these findings extend to the rest of psychology or other disciplines. In the sampling frame itself, not all articles were replicated; in each article, only one study was replicated; and in each study, only one statistical result was subject to replication. More resource-intensive studies were less likely to be included than were less resource-intensive studies. Although study selection bias was reduced by the sampling frame and selection strategy, the impact of selection bias is unknown.

We investigated the reproducibility rate of psychology not because there is something special about psychology, but because it is our discipline. Concerns about reproducibility are widespread across disciplines (9–21). Reproducibility is not well understood because the incentives for individual scientists prioritize novelty over replication (20). If nothing else, this project demonstrates that it is possible to conduct a large-scale examination of reproducibility despite the incentive barriers. Here, we conducted single-replication attempts of many effects obtaining broad-and-shallow evidence. These data provide information about reproducibility in general but little precision about individual effects in particular. A complementary narrow-and-deep approach is characterized by the Many Labs replication projects (32). In those, many replications of single effects allow precise estimates of effect size but result in generalizability that is circumscribed to those individual effects. Pursuing both strategies across disciplines, such as the ongoing effort in cancer biology (35), would yield insight about common and distinct

challenges and may cross-fertilize strategies so as to improve reproducibility.

Because reproducibility is a hallmark of credible scientific evidence, it is tempting to think that maximum reproducibility of original results is important from the onset of a line of inquiry through its maturation. This is a mistake. If initial ideas were always correct, then there would hardly be a reason to conduct research in the first place. A healthy discipline will have many false starts as it confronts the limits of present understanding.

Innovation is the engine of discovery and is vital for a productive, effective scientific enterprise. However, innovative ideas become old news fast. Journal reviewers and editors may dismiss a new test of a published idea as unoriginal. The claim that “we already know this” belies the uncertainty of scientific evidence. Deciding the ideal balance of resourcing innovation versus verification is a question of research efficiency. How can we maximize the rate of research progress? Innovation points out paths that are possible; replication points out paths that are likely; progress relies on both. The ideal balance is a topic for investigation itself. Scientific incentives—funding, publication, or awards—can be tuned to encourage an optimal balance in the collective effort of discovery (36, 37).

Progress occurs when existing expectations are violated and a surprising result spurs a new investigation. Replication can increase certainty when findings are reproduced and promote innovation when they are not. This project provides accumulating evidence for many findings in psychological research and suggests that there is still more work to do to verify whether we know what we think we know.

## Conclusion

After this intensive effort to reproduce a sample of published psychological findings, how many of the effects have we established are true? Zero. And how many of the effects have we established are false? Zero. Is this a limitation of the project design? No. It is the reality of doing science, even if it is not appreciated in daily practice. Humans desire certainty, and science infrequently provides it. As much as we might wish it to be otherwise, a single study almost never provides definitive resolution for or against an effect and its explanation. The original studies examined here offered tentative evidence; the replications we conducted offered additional, confirmatory evidence. In some cases, the replications increase confidence in the reliability of the original results; in other cases, the replications suggest that more investigation is needed to establish the validity of the original findings. Scientific progress is a cumulative process of uncertainty reduction that can only succeed if science itself remains the greatest skeptic of its explanatory claims.

The present results suggest that there is room to improve reproducibility in psychology. Any temptation to interpret these results as a defeat for psychology, or science more generally, must contend with the fact that this project demon-

strates science behaving as it should. Hypotheses abound that the present culture in science may be negatively affecting the reproducibility of findings. An ideological response would discount the arguments, discredit the sources, and proceed merrily along. The scientific process is not ideological. Science does not always provide comfort for what we wish to be; it confronts us with what is. Moreover, as illustrated by the Transparency and Openness Promotion (TOP) Guidelines (<http://cos.io/top>) (37), the research community is taking action already to improve the quality and credibility of the scientific literature.

We conducted this project because we care deeply about the health of our discipline and believe in its promise for accumulating knowledge about human behavior that can advance the quality of the human condition. Reproducibility is central to that aim. Accumulating evidence is the scientific community's method of self-correction and is the best available option for achieving that ultimate goal: truth.

## REFERENCES AND NOTES

- C. Hempel, Maximal specificity and lawlikeness in probabilistic explanation. *Philos. Sci.* **35**, 116–133 (1968). doi: [10.1086/288197](https://doi.org/10.1086/288197)
- C. Hempel, P. Oppenheim, Studies in the logic of explanation. *Philos. Sci.* **15**, 135–175 (1948). doi: [10.1086/286983](https://doi.org/10.1086/286983)
- I. Lakatos, in *Criticism and the Growth of Knowledge*, I. Lakatos, A. Musgrave, Eds. (Cambridge Univ. Press, London, 1970), pp. 170–196.
- P. E. Meehl, Appraising and amending theories: The strategy of Lakatosian defense and two principles that warrant it. *Psychol. Inq.* **1**, 108–141 (1990). doi: [10.1207/s15327965pli0102\\_1](https://doi.org/10.1207/s15327965pli0102_1)
- J. R. Platt, Strong Inference: Certain systematic methods of scientific thinking may produce much more rapid progress than others. *Science* **146**, 347–353 (1964). doi: [10.1126/science.146.3642.347](https://doi.org/10.1126/science.146.3642.347); PMID: [17739513](https://pubmed.ncbi.nlm.nih.gov/17739513/)
- W. C. Salmon, in *Introduction to the Philosophy of Science*, M. H. Salmon Ed. (Hackett Publishing Company, Indianapolis, 1999), pp. 7–41.
- B. A. Nosek, D. Lakens, Registered reports: A method to increase the credibility of published results. *Soc. Psychol.* **45**, 137–141 (2014). doi: [10.1027/1864-9335/a000192](https://doi.org/10.1027/1864-9335/a000192)
- S. Schmidt, Shall we really do it again? The powerful concept of replication is neglected in the social sciences. *Rev. Gen. Psychol.* **13**, 90–100 (2009). doi: [10.1037/a0015108](https://doi.org/10.1037/a0015108)
- J. P. A. Ioannidis, Why most published research findings are false. *PLoS Med.* **2**, e124 (2005). doi: [10.1371/journal.pmed.0020124](https://doi.org/10.1371/journal.pmed.0020124); PMID: [16060722](https://pubmed.ncbi.nlm.nih.gov/16060722/)
- C. G. Begley, L. M. Ellis, Drug development: Raise standards for preclinical cancer research. *Nature* **483**, 531–533 (2012). doi: [10.1038/483531a](https://doi.org/10.1038/483531a); PMID: [22460880](https://pubmed.ncbi.nlm.nih.gov/22460880/)
- F. Prinz, T. Schlange, K. Asadullah, Believe it or not: How much can we rely on published data on potential drug targets? *Nat. Rev. Drug Discov.* **10**, 712–713 (2011). doi: [10.1038/nrd3439-c1](https://doi.org/10.1038/nrd3439-c1); PMID: [21892149](https://pubmed.ncbi.nlm.nih.gov/21892149/)
- M. McNutt, Reproducibility. *Science* **343**, 229 (2014). doi: [10.1126/science.1250475](https://doi.org/10.1126/science.1250475); PMID: [24436391](https://pubmed.ncbi.nlm.nih.gov/24436391/)
- H. Pashler, E.-J. Wagenmakers, Editors' introduction to the special section on replicability in psychological science: A crisis of confidence? *Perspect. Psychol. Sci.* **7**, 528–530 (2012). doi: [10.1177/1745691612465253](https://doi.org/10.1177/1745691612465253); PMID: [26168108](https://pubmed.ncbi.nlm.nih.gov/26168108/)
- K. S. Button *et al.*, Power failure: Why small sample size undermines the reliability of neuroscience. *Nat. Rev. Neurosci.* **14**, 365–376 (2013). doi: [10.1038/nrn3475](https://doi.org/10.1038/nrn3475); PMID: [23571845](https://pubmed.ncbi.nlm.nih.gov/23571845/)
- D. Fanelli, “Positive” results increase down the hierarchy of the sciences. *PLOS ONE* **5**, e10068 (2010). doi: [10.1371/journal.pone.0010068](https://doi.org/10.1371/journal.pone.0010068); PMID: [20383332](https://pubmed.ncbi.nlm.nih.gov/20383332/)
- A. G. Greenwald, Consequences of prejudice against the null hypothesis. *Psychol. Bull.* **82**, 1–20 (1975). doi: [10.1037/h0076157](https://doi.org/10.1037/h0076157)
- G. S. Howard *et al.*, Do research literatures give correct answers? *Rev. Gen. Psychol.* **13**, 116–121 (2009). doi: [10.1037/a0015468](https://doi.org/10.1037/a0015468)
- J. P. A. Ioannidis, M. R. Munafò, P. Fusar-Poli, B. A. Nosek, S. P. David, Publication and other reporting biases in cognitive sciences: Detection, prevalence, and prevention. *Trends Cogn. Sci.* **18**, 235–241 (2014). doi: [10.1016/j.tics.2014.02.010](https://doi.org/10.1016/j.tics.2014.02.010); PMID: [24656991](https://pubmed.ncbi.nlm.nih.gov/24656991/)

- L. K. John, G. Loewenstein, D. Prelec, Measuring the prevalence of questionable research practices with incentives for truth telling. *Psychol. Sci.* **23**, 524–532 (2012). doi: [10.1177/0956797611430953](https://doi.org/10.1177/0956797611430953); PMID: [22508865](https://pubmed.ncbi.nlm.nih.gov/22508865/)
- B. A. Nosek, J. R. Spies, M. Motyl, Scientific utopia: II. Restructuring incentives and practices to promote truth over publishability. *Perspect. Psychol. Sci.* **7**, 615–631 (2012). doi: [10.1177/1745691612459058](https://doi.org/10.1177/1745691612459058); PMID: [26168121](https://pubmed.ncbi.nlm.nih.gov/26168121/)
- R. Rosenthal, The file drawer problem and tolerance for null results. *Psychol. Bull.* **86**, 638–641 (1979). doi: [10.1037/0033-2909.86.3.638](https://doi.org/10.1037/0033-2909.86.3.638)
- P. Rozin, What kind of empirical research should we publish, fund, and reward?: A different perspective. *Perspect. Psychol. Sci.* **4**, 435–439 (2009). doi: [10.1111/j.1745-6924.2009.01151.x](https://doi.org/10.1111/j.1745-6924.2009.01151.x); PMID: [26158991](https://pubmed.ncbi.nlm.nih.gov/26158991/)
- J. P. Simmons, L. D. Nelson, U. Simonsohn, False-positive psychology: Undisclosed flexibility in data collection and analysis allows presenting anything as significant. *Psychol. Sci.* **22**, 1359–1366 (2011). doi: [10.1177/0956797611417632](https://doi.org/10.1177/0956797611417632); PMID: [22006061](https://pubmed.ncbi.nlm.nih.gov/22006061/)
- Open Science Collaboration, An open, large-scale, collaborative effort to estimate the reproducibility of psychological science. *Perspect. Psychol. Sci.* **7**, 657–660 (2012). doi: [10.1177/1745691612462588](https://doi.org/10.1177/1745691612462588); PMID: [26168127](https://pubmed.ncbi.nlm.nih.gov/26168127/)
- Open Science Collaboration, in *Implementing Reproducible Computational Research (A Volume in The R Series)*, V. Stodden, F. Leisch, R. Peng, Eds. (Taylor & Francis, New York, 2014), pp. 299–323.
- R. A. Fisher, Theory of statistical estimation. *Math. Proc. Camb. Philos. Soc.* **22**, 700–725 (1925). doi: [10.1017/S0305004100009580](https://doi.org/10.1017/S0305004100009580)
- W. Viechtbauer, Conducting meta-analyses in R with the metafor package. *J. Stat. Softw.* **36**, 1–48 (2010).
- S. L. Braver, F. J. Thoemmes, R. Rosenthal, Continuously cumulating meta-analysis and replicability. *Perspect. Psychol. Sci.* **9**, 333–342 (2014). doi: [10.1177/1745691614529796](https://doi.org/10.1177/1745691614529796); PMID: [26173268](https://pubmed.ncbi.nlm.nih.gov/26173268/)
- U. Simonsohn, Small telescopes: Detectability and the evaluation of replication results. *Psychol. Sci.* **26**, 559–569 (2015). doi: [10.1177/0956797614567341](https://doi.org/10.1177/0956797614567341); PMID: [25800521](https://pubmed.ncbi.nlm.nih.gov/25800521/)
- D. Lakens, Calculating and reporting effect sizes to facilitate cumulative science: A practical primer for t-tests and ANOVAs. *Front. Psychol.* **4**, 863 (2013). doi: [10.3389/fpsyg.2013.00863](https://doi.org/10.3389/fpsyg.2013.00863); PMID: [24324449](https://pubmed.ncbi.nlm.nih.gov/24324449/)
- J. Lehrer, The truth wears off: Is there something wrong with the scientific method? *The New Yorker*, 52–57 (2010).
- R. Klein *et al.*, Investigating variation in replicability: A “many labs” replication project. *Soc. Psychol.* **45**, 142–152 (2014). doi: [10.1027/1864-9335/a000178](https://doi.org/10.1027/1864-9335/a000178)
- J. Cohen, The statistical power of abnormal-social psychological research: A review. *J. Abnorm. Soc. Psychol.* **65**, 145–153 (1962). doi: [10.1037/h0045186](https://doi.org/10.1037/h0045186); PMID: [13880271](https://pubmed.ncbi.nlm.nih.gov/13880271/)
- T. D. Sterling, Publication decisions and their possible effects on inferences drawn from tests of significance—or vice versa. *J. Am. Stat. Assoc.* **54**, 30–34 (1959).
- T. M. Errington *et al.*, An open investigation of the reproducibility of cancer biology research. *eLife* **3**, e04333 (2014). doi: [10.7554/eLife.04333](https://doi.org/10.7554/eLife.04333); PMID: [25490932](https://pubmed.ncbi.nlm.nih.gov/25490932/)
- J. K. Hartshorne, A. Schachner, Tracking replicability as a method of post-publication open evaluation. *Front. Comput. Neurosci.* **6**, 8 (2012). doi: [10.3389/fncom.2012.00008](https://doi.org/10.3389/fncom.2012.00008); PMID: [22403538](https://pubmed.ncbi.nlm.nih.gov/22403538/)
- B. A. Nosek *et al.*, Promoting an open research culture. *Science* **348**, 1422–1425 (2015). doi: [10.1126/science.aab2374](https://doi.org/10.1126/science.aab2374); PMID: [26113702](https://pubmed.ncbi.nlm.nih.gov/26113702/)

## ACKNOWLEDGMENTS

In addition to the coauthors of this manuscript, there were many volunteers who contributed to project success. We thank D. Acup, J. Anderson, S. Anzellotti, R. Araujo, J. D. Arnal, T. Bates, R. Battleday, R. Bauchwitz, M. Bernstein, B. Blohowiak, M. Boffo, E. Bruneau, B. Chabot-Hanowell, J. Chan, P. Chu, A. Dalla Rosa, B. Deen, P. DiGiacomo, C. Dogulu, N. Dufour, C. Fitzgerald, A. Foote, A. Garcia, E. Garcia, C. Gautreau, L. Germaine, T. Gill, L. Goldberg, S. D. Goldinger, H. Gweon, D. Haile, K. Hart, F. Hjorth, J. Hoening, A. Innes-Ker, B. Jansen, R. Jersakova, Y. Jie, Z. Kaldy, W. K. Kong, A. Kenney, J. Kingston, J. Koster-Hale, A. Lam, R. LeDonne, D. Lumian, E. Luong, S. Man-pui, J. Martin, A. Mauk, T. McElroy, K. McRae, T. Miller, K. Moser, M. Mullarkey, A. R. Munoz, J. Ong, C. Parks, D. S. Pate, D. Patron, H. J. M. Pennings, M. Penuliar, A. Pfammatter, J. P. Shanoltz, E. Stevenson, E. Pichler, H. Raudszus, H. Richardson, N. Rothstein, T. Scherndl, S. Schrager, S. Shah, Y. S. Tai, A. Skerry, M. Steinberg, J. Stoetrau, H. Tibboel, A. Tooley, A. Tullett, C. Vaccaro, E. Vergauwe, A. Watanabe, I. Weiss, M. H. White II,





## Estimating the reproducibility of psychological science

Open Science Collaboration

*Science* **349** (6251), aac4716.  
DOI: 10.1126/science.aac4716

### Empirically analyzing empirical evidence

One of the central goals in any scientific endeavor is to understand causality. Experiments that seek to demonstrate a cause/effect relation most often manipulate the postulated causal factor. Aarts *et al.* describe the replication of 100 experiments reported in papers published in 2008 in three high-ranking psychology journals. Assessing whether the replication and the original experiment yielded the same result according to several criteria, they find that about one-third to one-half of the original findings were also observed in the replication study.

*Science*, this issue 10.1126/science.aac4716

ARTICLE TOOLS	<a href="http://science.sciencemag.org/content/349/6251/aac4716">http://science.sciencemag.org/content/349/6251/aac4716</a>
SUPPLEMENTARY MATERIALS	<a href="http://science.sciencemag.org/content/suppl/2015/08/26/349.6251.aac4716.DC1">http://science.sciencemag.org/content/suppl/2015/08/26/349.6251.aac4716.DC1</a>
RELATED CONTENT	<a href="http://science.sciencemag.org/content/sci/349/6251/999.2.full">http://science.sciencemag.org/content/sci/349/6251/999.2.full</a> <a href="http://science.sciencemag.org/content/sci/349/6251/910.full">http://science.sciencemag.org/content/sci/349/6251/910.full</a> <a href="http://science.sciencemag.org/content/sci/348/6242/1422.full">http://science.sciencemag.org/content/sci/348/6242/1422.full</a> <a href="http://science.sciencemag.org/content/sci/343/6168/229.full">http://science.sciencemag.org/content/sci/343/6168/229.full</a> <a href="http://science.sciencemag.org/content/sci/348/6242/1403.full">http://science.sciencemag.org/content/sci/348/6242/1403.full</a> <a href="http://science.sciencemag.org/content/sci/351/6277/1037.2.full">http://science.sciencemag.org/content/sci/351/6277/1037.2.full</a> <a href="http://science.sciencemag.org/content/sci/351/6277/1037.3.full">http://science.sciencemag.org/content/sci/351/6277/1037.3.full</a>
REFERENCES	This article cites 37 articles, 3 of which you can access for free <a href="http://science.sciencemag.org/content/349/6251/aac4716#BIBL">http://science.sciencemag.org/content/349/6251/aac4716#BIBL</a>
PERMISSIONS	<a href="http://www.sciencemag.org/help/reprints-and-permissions">http://www.sciencemag.org/help/reprints-and-permissions</a>

Use of this article is subject to the [Terms of Service](#)