

# Heterogeneity of Research Results: A New Perspective From Which to Assess and Promote Progress in Psychological Science

Perspectives on Psychological Science  
1–19

© The Author(s) 2021



Article reuse guidelines:

[sagepub.com/journals-permissions](http://sagepub.com/journals-permissions)

DOI: 10.1177/1745691620964193

[www.psychologicalscience.org/PPS](http://www.psychologicalscience.org/PPS)



**Audrey Helen Linden and Johannes Hönekopp** 

Department of Psychology, Northumbria University

## Abstract

Heterogeneity emerges when multiple close or conceptual replications on the same subject produce results that vary more than expected from the sampling error. Here we argue that unexplained heterogeneity reflects a lack of coherence between the concepts applied and data observed and therefore a lack of understanding of the subject matter. Typical levels of heterogeneity thus offer a useful but neglected perspective on the levels of understanding achieved in psychological science. Focusing on continuous outcome variables, we surveyed heterogeneity in 150 meta-analyses from cognitive, organizational, and social psychology and 57 multiple close replications. Heterogeneity proved to be very high in meta-analyses, with powerful moderators being conspicuously absent. Population effects in the average meta-analysis vary from small to very large for reasons that are typically not understood. In contrast, heterogeneity was moderate in close replications. A newly identified relationship between heterogeneity and effect size allowed us to make predictions about expected heterogeneity levels. We discuss important implications for the formulation and evaluation of theories in psychology. On the basis of insights from the history and philosophy of science, we argue that the reduction of heterogeneity is important for progress in psychology and its practical applications, and we suggest changes to our collective research practice toward this end.

## Keywords

meta-analysis, heterogeneity, replication, statistical power, philosophy of science, psychological research

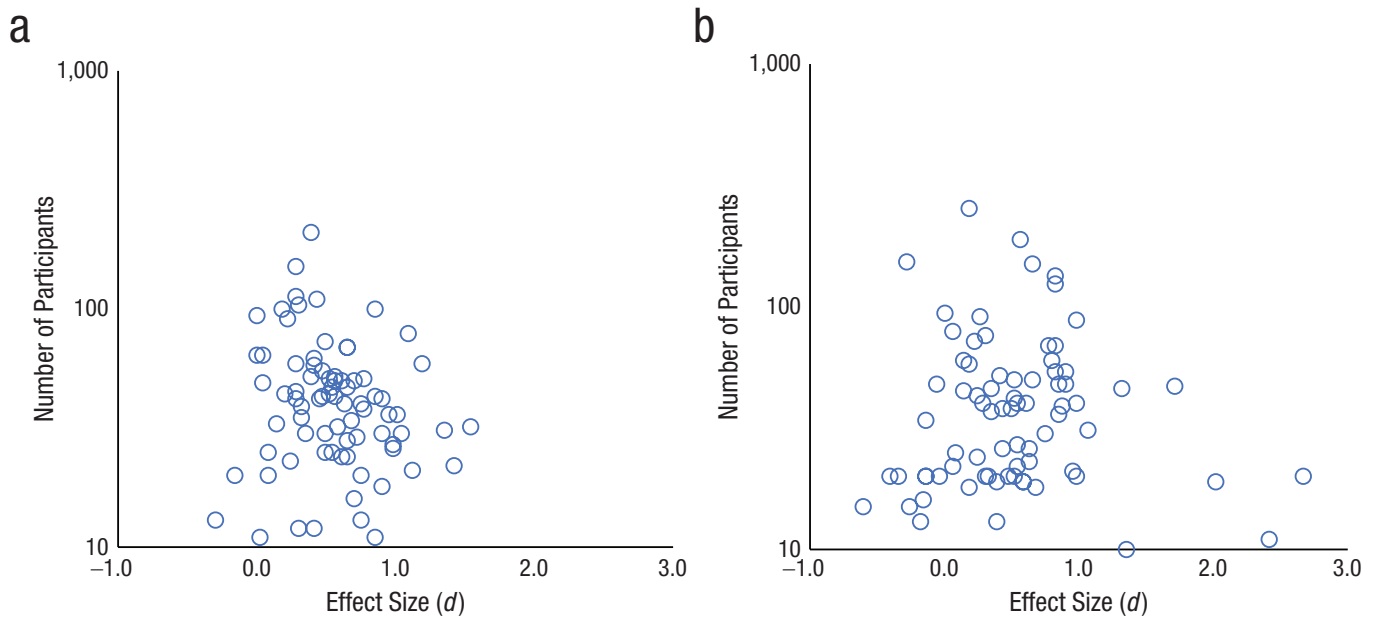
Meta-analysis, which seeks to summarize the results of multiple studies into the same phenomenon, has become an indispensable tool in contemporary research. In pioneering work, Smith and Glass (1977) showed that psychotherapy has a strong positive effect on the average patient studied, and F. L. Schmidt and Hunter (1977) demonstrated that the validity of employment tests generalize more readily across different job types than previously believed. Influential surveys of meta-analyses have demonstrated the effectiveness of psychological interventions (Lipsey & Wilson, 1993), provided effect-size benchmarks for social psychology (Richard et al., 2003), and summarized findings on psychological gender similarities (Hyde, 2014). Here we provide a survey of meta-analyses that shifts the perspective from the mean effect size in a population of studies (i.e., the size of the average effect in a particular domain) to the heterogeneity of results (i.e., the degree to which results

differ across studies into the same issue). In any meta-analysis, heterogeneity indicates the extent to which the summarized studies tap into the same population effect size. If the same population effect size is investigated, heterogeneity will be zero. Even in this case, the sampling error will create differences in observed effects across studies. Zero heterogeneity is inferred when these observed differences do not exceed the level expected as a result of the sampling error. Consider the effectiveness of psychotherapy as an example. If heterogeneity was zero, the effectiveness of psychotherapy would be the same across all studies, regardless of the issue patients present with (e.g., anorexia, depression, specific phobia), the type of therapy they receive (e.g.,

---

### Corresponding Author:

Johannes Hönekopp, Department of Psychology, Northumbria University  
E-mail: [johannes.honekopp@unn.ac.uk](mailto:johannes.honekopp@unn.ac.uk)



**Fig. 1.** Funnel plots for two meta-analyses. Linck et al. (2014) investigated the link between working memory and second-language comprehension (a). The estimated mean of the population of effect sizes, standard deviation of observed effect sizes, and estimated heterogeneity of true effect sizes were  $d = 0.51$ ,  $SD = 0.36$ , and  $T = 0.11$  ( $I^2 = 11$ ), respectively. Baker et al. (2014) investigated the link between intelligence and performance on the Reading the Mind in the Eyes test (b). The estimated mean of the population of effect sizes, standard deviation of observed effect sizes, and estimated heterogeneity of true effect sizes were  $d = 0.49$ ,  $SD = 0.59$ , and  $T = 0.35$  ( $I^2 = 53$ ), respectively.

cognitive-behavioral, psychoanalytic), and other differences. This is obviously unrealistic; for example, some conditions are treated more successfully than others (Huhn et al., 2014). Heterogeneity thus reflects how much the population effect sizes differ across studies. We provide a formal treatment of heterogeneity later, but Figure 1 provides examples with high and low heterogeneity.

Heterogeneity tends to receive little attention from researchers (Aytug et al., 2012; Dieckmann et al., 2009; Ioannidis, 2008), but we argue here that much is to be gained from its study because (a) heterogeneity reflects the degree of understanding of the subject matter being investigated and (b) its study offers useful suggestions regarding the improvement of our collective research practice.

### Why Heterogeneity Matters

Low (as opposed to high) heterogeneity reflects a more advanced understanding of the subject matter being studied. This is because high heterogeneity, at least as long as it remains unexplained, suggests the lack of a strong coherence between the concepts applied and the data observed. Take visuospatial skills in people with autism spectrum conditions (ASCs) as an example (Muth et al., 2014). In line with current theorizing, the average study found (moderately) better visuospatial

performance in people with ASCs than in IQ-matched control subjects on a number of standardized tasks. At the same time, the heterogeneity of the results proved to be high, even for the same task. Not accounted for by any theory, this random variation in study results (which might have resulted from unrecognized variability in ASCs, unreliability in diagnosis, or other factors) points to a shortcoming in our understanding. It also implies that the result of the next study into the same question is highly unpredictable (i.e., over and above the uncertainty arising from sampling error).

Moreover, low heterogeneity should facilitate future progress for two reasons. First, a clear structure in observable data can in itself guide understanding—a point stressed by 17th-century luminaries Francis Bacon and Isaac Newton as well as modern philosophers of science such as Hans Reichenbach, Norwood Russell Hanson, and Herbert Simon (Schickore, 2018; Simon, 1973). For example, the 19th-century astronomer William Huggins observed that the light of different stars, when seen through a prism, shows the same set of spectral lines; however, he also observed that these lines are collectively shifted to varying degrees. The observation of this systematic redshift pattern led to the discovery that stars move away from us and at different speeds (Schneider, 2014). Sixty years later, Edwin Hubble observed that the degree of stars' redshift is linearly related to their distance from us, which led to

the discovery that the universe is expanding (Schneider, 2014). Skinner (1956) and Stevens (1957) provide prominent examples for a guiding role of orderly observation data in psychology. Second, the systematic violation of expectations has often proved crucial for scientific discovery (Kuhn, 1970). Thus, the failure of an increasingly convoluted Ptolemaic system to further improve the predictions of astronomic events motivated Copernicus to devise a new, heliocentric model of the cosmos; and the failure to detect expected changes in the speed of light—derived from the idea that light propagates through a medium—led Einstein to abandon the idea of a luminiferous ether and to fundamentally rethink physics. As captured in Bacon's dictum that "truth emerges more readily from error than from confusion" (Kuhn, 1970, p. 57), such anomalies cannot emerge when theoretical concepts and observed data lack a clear connection in the first place.

We therefore propose heterogeneity as a useful perspective from which to judge the success of psychological science, alongside other yardsticks such as the generation of good theories (Wallis, 2015), the design of successful interventions (Lipsey & Wilson, 1993), and beneficial contributions to policy design (Fischhoff, 1990). Thus, heterogeneity is of considerable intrinsic value, which is why we seek to systematically measure it in the psychological-research results presented here. What are typical levels? Do they differ across domains, and if so, can we make sense of these differences? Apart from its intrinsic value, knowledge of actual levels of heterogeneity has immediate practical implications: Heterogeneity has been demonstrated to typically decrease the statistical power of studies<sup>1</sup>; that is, any real effect under investigation is less likely to produce a statistically significant result (Kenny & Judd, 2019; McShane & Böckenholt, 2014; Shrout & Rodgers, 2018). For sample-size planning to take this into account, reliable estimates of heterogeneity are needed, which we supply here. Finally, and perhaps most importantly, our findings have clear implications for improving our collective research practice, as we discuss at the end of this article. Before we can address the details of our study, it is necessary to deal with a number of critical points, which we address in the next sections.

## Moderators

Heterogeneity reflects a lack of understanding only when it remains unaccounted for. Let us reconsider our example of psychotherapy effectiveness. A meta-analysis that summarizes all sensible studies should find large heterogeneity because these studies will differ in key variables such as the issue being treated, the therapy

being used, and so on. If this heterogeneity can be explained by moderators (e.g., that effectiveness differs strongly across treated disorders or across types of psychotherapy), this obviously no longer indicates a lack of knowledge. (On the contrary, it might be argued that explained heterogeneity reflects an increase in understanding.) We are not aware of any study to date that has systematically investigated the extent to which the heterogeneity that is observed in a set of studies is accounted for by moderators. We therefore investigate it here.

## Conceptual Versus Close Replications

Heterogeneity as a concept makes sense only if the set of studies for which it is computed can, in some sense, be conceived as replications of each other. In this context, the differentiation between close and conceptual replications has become fruitful (S. Schmidt, 2009; Zwaan et al., 2018). The former seek to replicate an earlier study as faithfully as possible. The Open Science Collaboration (2015) project is a famous example. In a massive collaborative effort, the authors sought to replicate 100 studies published in high-profile psychology journals. The replications sought to copy study materials, data analyses, and other key aspects of the original studies as closely as possible and can therefore be considered close replications. In contrast, the studies summarized in a meta-analysis can typically be considered to be conceptual replications (F. L. Schmidt & Oh, 2016); that is, although they address the same topic or mechanism, they often differ markedly in their design, study materials, participants, data analysis, and other key aspects. Heterogeneity should thus tend to be larger in conceptual replications than in close replications.

A systematic comparison of heterogeneity in close and conceptual replications should be instructive. For example, Stanley et al. (2018) argued that the low replicability observed in Open Science Collaboration (2015) might reflect low power caused by high heterogeneity. However, the heterogeneity data that they presented in support of this argument stemmed almost exclusively from conceptual replications. Their assumption that heterogeneity in close replication attempts might be similar rested on only two examples for the latter.

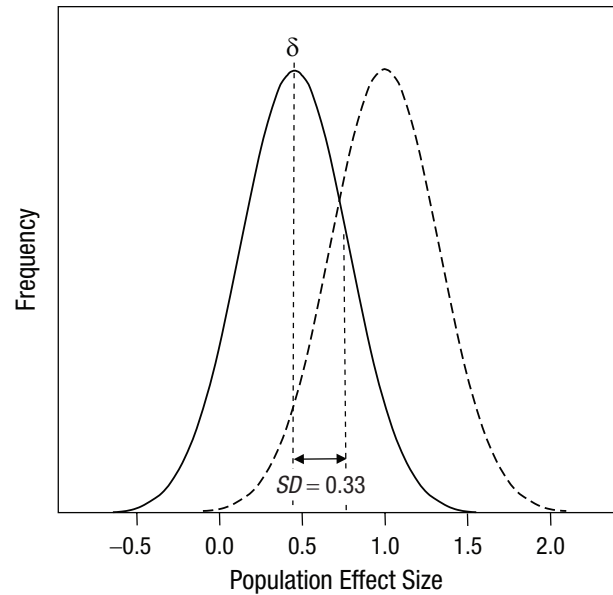
Note that the heterogeneity for each of the 100 twin studies in the Open Science Collaboration (original and replication) cannot be reliably estimated. Instead of a single replication, this would require multiple close replications of the same effect (e.g., Klein et al., 2014). We thus use Many Labs-type replications to study heterogeneity in close replications.

## Measuring Heterogeneity

So far, we have not addressed how heterogeneity can be quantified. In psychology the idea of heterogeneity is usually discussed in the context of standardized effect sizes (e.g., Cohen's  $d$  instead of a difference between group means in raw scores), and we stick to this perspective here. Two established approaches to quantify heterogeneity are  $I^2$  and  $\tau$ . Before we deal with these approaches, it is helpful to consider what we should expect to see in the absence of heterogeneity. Even if all primary studies tap into the same population effect size, we expect to see differences in the observed effect sizes as a result of sampling error. Thus, observed differences between effect sizes do not necessarily point to heterogeneity.

The first approach,  $I^2$ , estimates the percentage of observed effect-size variability that reflects real differences in effect sizes. When  $I^2$  is near zero, the observed variability is mostly down to sampling error; when  $I^2$  is near 100, most of the observed variability reflects differences in population effect sizes. However,  $I^2$  does not discriminate well when heterogeneity is large (Huedo-Medina et al., 2006). Moreover,  $I^2$  depends on the sample size in the primary studies (Borenstein et al., 2017; IntHout et al., 2015). Imagine that all studies used small samples. Individual effect sizes will scatter widely around their population effect size. A large percentage of observed variability thus reflects the sampling error, and  $I^2$  will be low. Now imagine that all studies used very large samples. Each study will provide a highly accurate estimate of its population effect size. Thus, only a small percentage of observed variability reflects the sampling error, and  $I^2$  will be high. Finally, for the current study, our approach to heterogeneity involves summarizing heterogeneity estimates across multiple meta-analyses. Using  $I^2$  in this way strikes us as questionable unless average sample sizes are similar.

The second approach directly estimates the variability in population effect sizes. It is generally assumed that population effect sizes relating to a given phenomenon follow a normal distribution;  $\tau$  refers to their standard deviation (Borenstein et al., 2009) and can be calculated when individual study effect sizes and standard errors are available. As an example, consider the meta-analysis in Figure 1a. The standard deviation of the observed effect sizes in the primary studies is 0.36. (For the sake of consistency, we use Cohen's  $d$  as a measure of effect size in this example and throughout this article.) Some of the observed effect-size variability must be due to sampling error. When the sampling error is removed, heterogeneity is estimated to be only 0.11. To better understand heterogeneity, consider Figure 2. Here, the mean for the population of true effect sizes is 0.45, and their standard deviation is 0.33; therefore,



**Fig. 2.** Two distributions of population effect sizes (standardized mean differences). The distribution on the left (solid line) shows a population effect size with a mean ( $\delta$ ) of 0.45 and a standard deviation ( $T$ ) of 0.33. The distribution on the right (dashed line) shows a population effect size with a mean ( $\delta$ ) of 1 and a standard deviation ( $T$ ) of 0.33.

$\tau = 0.33$ . The standard deviation roughly reflects how far data points are, on average, away from the mean. Any study's population effect size will thus typically deviate from 0.45 by approximately 0.33; moreover, the 95% credibility interval ranges from  $-0.20$  to  $1.10$ , estimating that 95% of all population effect sizes fall into this bracket (Hunter & Schmidt, 2004). Because  $\tau$  is independent from  $N$  in primary studies (which differ markedly between meta-analyses) and  $\tau$ , unlike  $I^2$ , is on an equal interval scale that allows meaningful computations of means, we use it here.<sup>2</sup>

Because  $\tau$  is an unknown population parameter, it must be estimated. Its estimator  $T$  often comes with considerable uncertainty, especially when a meta-analysis is based on few studies (e.g., Huedo-Medina et al., 2006). For individual meta-analyses this can be a serious issue, especially when heterogeneity is wrongly estimated to be zero (Chung et al., 2013). However, this is less of a concern for our current purpose because our focus is not on individual meta-analyses but on aggregates of 50 or more, and we do not believe there is reason to suggest that heterogeneity estimates will be consistently biased in one direction (see Hönekopp & Linden, 2019).

## A Sensible Sampling Frame

What is a sensible sampling frame for a survey of heterogeneity? One potential strategy would be to use a

representative sample of meta-analyses across all of psychology. However, our heterogeneity measure  $T$  is not suitable for odds ratios and similar effect sizes, which are frequently used in clinical psychology. A sample of meta-analyses amenable to our heterogeneity measure would thus fail to be representative. We therefore decided to focus on a number of subdisciplines instead in which effect-size measures  $d$  and  $r$ , for which our heterogeneity measure works, predominate. We chose cognitive, social, and organizational psychology because these subdisciplines differ in their relative emphasis on fundamental versus applied research and because they were the focus of previous metascientific inquiries (Mitchell, 2012; Open Science Collaboration, 2015). Mitchell (2012) compared effect sizes from laboratory-based and field-based studies in organizational and social psychology. The correlation between lab- and field-based effect sizes was higher in the former ( $r = .89$ ) than in the latter ( $r = .53$ ). The Open Science Collaboration (2015) found that findings from cognitive psychology were substantially more replicable than those from social psychology. These observations could point to greater heterogeneity within social psychology in general. Indeed, Stanley et al. (2018) argued that the low replication rate observed in the Open Science Collaboration might reflect low power caused by high heterogeneity. From this perspective, the difference in observed Open Science Collaboration replication rates would then suggest higher heterogeneity in social psychology than cognitive psychology. We test this idea here.

## Aims

Our aims are as follows: Given the intrinsic value of heterogeneity as an indicator of a lack of understanding, we seek to establish a typical level of heterogeneity in conceptual replications. We compare these levels across the subdisciplines of cognitive, organizational, and social psychology and against heterogeneity observed in close replications. We also investigate the extent to which heterogeneity in a set of studies can typically be accounted for by moderators. We further explore whether any characteristics explain differences in heterogeneity.

To foreshadow our key results, we find that heterogeneity tends to be very large in conceptual replications but moderate in close replications. Our investigations regarding the drivers of heterogeneity show that moderators do little to account for heterogeneity. We also find a previously unexplored strong relationship between heterogeneity and effect size, which allows us, for the first time, to make predictions about expected levels of heterogeneity for a given phenomenon. These findings have clear implications for the improvement

of our collective research practice, as we discuss at the end of this article.

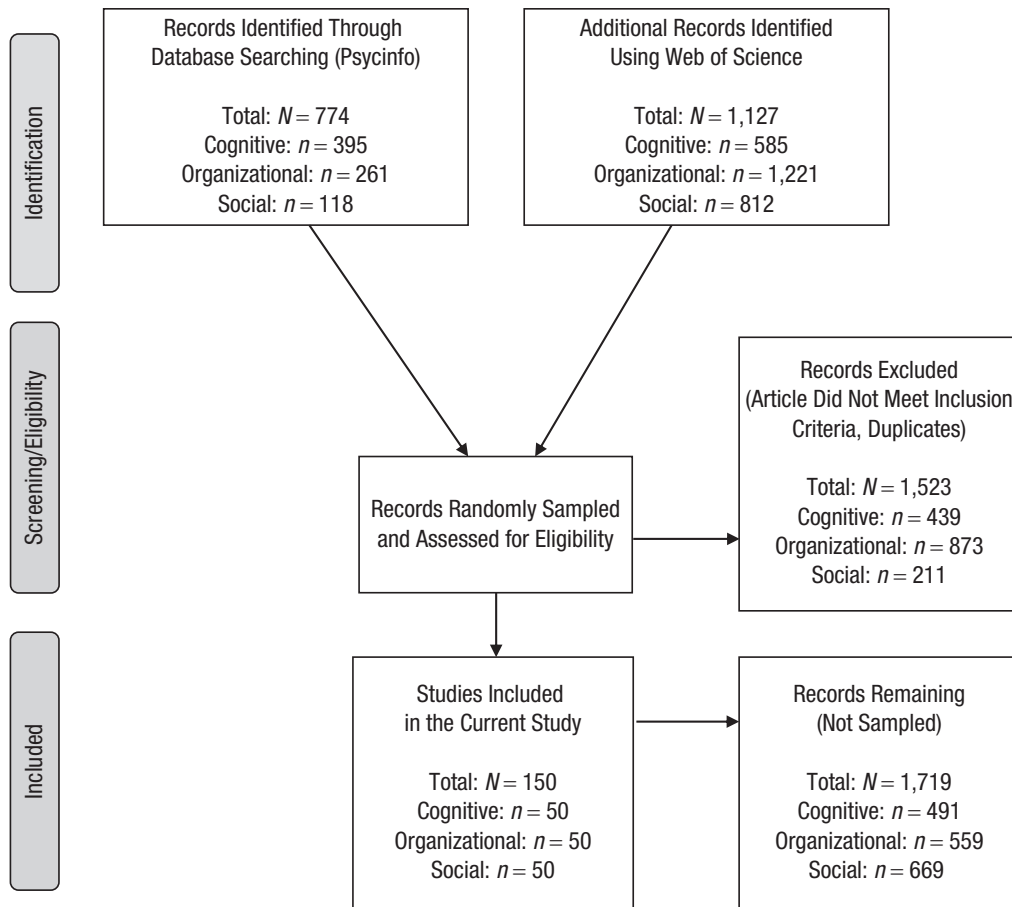
## Method

### *Study search and selection strategy*

We aimed to investigate all available Many Labs-type replications. We searched CurateScience.org for relevant reports in April 2017 and added studies from Many Labs 2 (Klein et al., 2018) at a later stage. Further, we investigated 50 meta-analyses each from cognitive, organizational, and social psychology. Feasibility, rather than power considerations, determined this choice. Our preregistered study protocol is available at <http://aspredicted.org/blind.php?x=bf46k8>.

In November 2016, we searched PsycINFO (journals only) for “meta-analy\*” in the abstract field. We restricted searches to PsycINFO classifications “3000 Social Psychology,” “3600 Industrial and Organizational Psychology,” “2340 Cognitive Processes,” “2343 Learning and Memory,” and “2346 Attention.” Because this search did not yield a sufficient number of eligible meta-analyses (see below for inclusion criteria), we also searched the Web of Science (articles only) for “meta-analy\*” in the categories “Psychology Social,” “Psychology Applied,” and “Psychology” (excluding meta-analyses that fell outside our target subdisciplines; see Fig. 3). All eligible meta-analyses were inspected in random order until we reached the desired number of 50 meta-analyses.

**Inclusion criteria.** Meta-analyses for the three subdisciplines were included if they met all of the following criteria. First, they had to address a substantive psychological effect (rather than, e.g., the psychometric properties of a questionnaire). Second, the analyzed effects had to be described as standardized mean differences (Cohen’s  $d$  or Hedges’s  $g$ ) or correlations (Pearson’s  $r$  or Fisher’s  $z$ ). Standardized differences and correlations are closely related concepts, and one can easily be converted into the other. Similar conversions are less sensible if categorical dependent variables are used (Ferguson, 2009), and our heterogeneity measure  $T$  is also not suitable for these types of effect sizes. For this reason, we excluded meta-analyses that used odds ratios, risk ratios, and similar measures. Third, the effect-size and sample-size information had to have been provided for the original studies. This was necessary to calculate heterogeneity. When only the sample sizes or effect sizes were available, an attempt was made to obtain missing data from the corresponding author. Finally, for practical reasons, the full article had to be available in English. All close replication reports that met the same criteria (Many Labs 1–3 and Registered Replication Reports 3–6) were included (Cheung et al.,



**Fig. 3.** Sampling of meta-analyses.

2016; Ebersole et al., 2016; Eerland et al., 2016; Hagger et al., 2016; Klein et al., 2014, 2018; Wagenmakers et al., 2016).<sup>3</sup>

In this way, we identified 50 meta-analyses for cognitive psychology, 50 for organizational psychology, 50 for social psychology, and 57 for close replications (see Table S1 in the Supplemental Material available online).

**Data extraction and analysis.** If the results of more than one meta-analysis were reported, the one including the largest number of studies was extracted. If multiple meta-analyses included the same number of studies, the first was used.

Heterogeneity for each meta-analysis was computed using the DerSimonian-Laird estimator in the metafor package (Version 2.1-0; Viechtbauer, 2010) for the R software environment (Version 3.4.1; R Core Team, 2017).<sup>4</sup> To keep effect sizes and levels of heterogeneity consistent across studies, all effect sizes were input as Cohen's  $d$ . All other effect sizes were converted accordingly.

It turned out that the frequency distributions for some of our observed outcome variables were skewed

to the right. For example, among the 150 meta-analyses,  $T$  had a skewness 0.99 (the largest  $Z$  score being 3.57). We therefore report Winsorized means ( $M_{\text{win}}$ ) and respective standard deviations ( $SD_{\text{win}}$ ). Winsorizing replaces the smallest and largest values in a distribution (in this case the smallest 10% and largest 10%) with the observations that are closest to them. If frequency distributions are normal in nature, this will not alter the results.  $M_{\text{win}}$  therefore removes the undue effect of outliers but retains much more information than the median, which trims all scores but the one in the middle of the distribution (Erceg-Hurn & Miroseovich, 2008). Specifically for  $T$ ,  $M_{\text{win}}$  should also counteract the likely overestimation resulting from setting negative heterogeneity estimates ( $T$ ) to 0. We used the Yuen-Welch method (Wilcox, 2005)—which is similar to the  $t$  test but based on  $M_{\text{win}}$ —for group comparisons of  $T$ . Likewise, we used Winsorized correlations ( $r_{\text{win}}$ ; Wilcox, 2005). Winsorized correlations limit the effect of outliers but retain more information than Spearman's rank-based correlation ( $r_s$ ). All data and further materials can be found at <https://osf.io/yr3xd>.

**Table 1.** Descriptive Statistics for Study 1

Study type	Number of meta-analyses	Mean $k$ per meta-analysis	Mean $T$ ( $SD$ )
Close replications	57	35.2	0.09 (0.07)
Subdiscipline			
Social	50	35.7	0.31 (0.11)
Cognitive	50	36.5	0.32 (0.13)
Organizational	50	38.3	0.35 (0.10)
Total	150	36.9	0.33 (0.11)

Note: All means are Winsorized.

## Results

### *How meta-analyses address heterogeneity*

Of 150 meta-analyses, 123 tested moderators, but only 83 (55%) reported a measure of heterogeneity. In 2009, the influential Preferred Reporting Items for Systematic Reviews and Meta-Analyses (PRISMA) guidelines (Moher et al., 2009) recommended that meta-analyses should address heterogeneity. Even for meta-analyses published after 2009, heterogeneity was reported in only 60% of cases. Note that the statistical significance of heterogeneity, for example,  $Q$ , was widely reported (77 times); of those meta-analyses, 45% did not report quantifying information. This focus on statistical significance and neglect of quantifying information runs counter to the meta-analysis estimation perspective (Hunter, 1997).

Overall, heterogeneity was quantified in less than a third of cases (43 times out of 150):  $I^2$  was reported in 33 cases,  $T^2$  in 9, and another measure was reported once. In addition to the observed neglect of quantification, it is interesting that authors unanimously reported  $T^2$  (the heterogeneity variance) instead of  $T$  (the standard deviation). Whereas standard deviation has a meaning that is comparatively easy to grasp (it approximates the average difference from the mean), variance does not have a similarly accessible interpretation. (This is why researchers most commonly report standard deviations, and not variances, in their descriptive statistics. Likewise, as shown by Fig. 2, humans have an intuitive understanding of the concept of standard deviation, so it makes sense to graph it. But humans lack an intuitive understanding of variance, so graphing it would be pointless.)

### *Heterogeneity observed in close replications and meta-analyses*

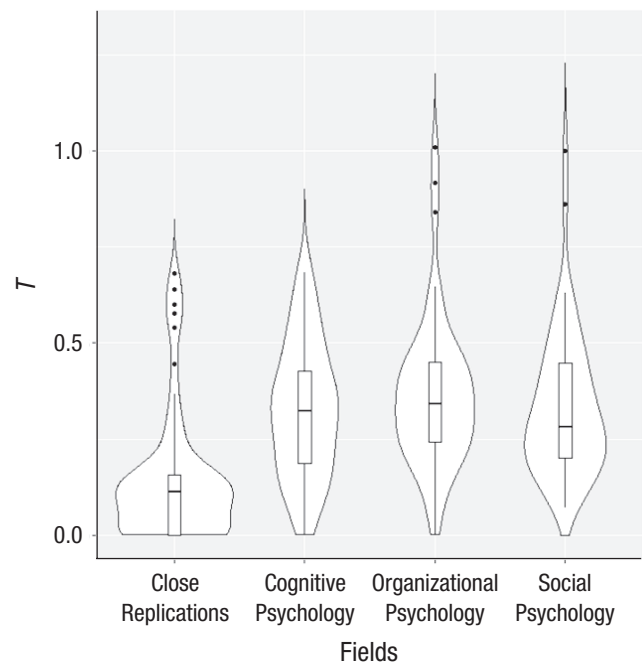
Table 1 shows descriptive statistics for close replications and meta-analyses. As expected, heterogeneity was much lower ( $M_{\text{win}} = 0.09$ ) in close replications than in

the meta-analyses ( $M_{\text{win}} = 0.33$ ),  $t(94.9) = 10.43$ ,  $p < .001$  (see Fig. 4).

Levels of heterogeneity were unexpectedly similar across all three subdisciplines—cognitive versus social psychology:  $t(57.3) = 0.33$ ,  $p = .370$  (one-tailed); social versus organizational:  $t(57.0) = 1.17$ ,  $p = .125$  (one-tailed); and organizational versus cognitive psychology:  $t(54.9) = 0.74$ ,  $p = .463$  (two-tailed).

### *Moderators*

To investigate the extent to which moderators account for heterogeneity, we looked at all 36 meta-analyses with a  $k \geq 60$  because moderators will be most reliably



**Fig. 4.** Observed levels of heterogeneity for 57 close replications and 50 meta-analyses in each subdiscipline (cognitive, organizational, and social psychology). In each plot, the horizontal line indicates the Winsorized mean, the top and bottom of the box indicates the top and bottom of the interquartile range (IQR), the whiskers represent values above and below the IQR, and dots represent outliers.

identified in large meta-analyses. We looked at meta-analyses only in which moderators were reported by the original authors. Where possible, we used the strongest moderator for which sufficient information was reported for further analyses. All moderators thus identified were grouping variables (i.e., none was continuous). We used these moderators to partition studies into appropriate subsets, which left us with 22 meta-analyses. We then excluded broad subsets. For example, Baker et al. (2014) investigated whether a relationship exists between intelligence and performance on the Reading the Mind in the Eyes test. The strongest moderator they examined was the type of intelligence test used. Studies were split into two subsets on the basis of this examination: IQ measured using the Wechsler IQ test and IQ measured using any other test. We excluded the broad “other” subset and compared this  $T$  against the  $T$  in the initial overall meta-analysis. The average  $T$  in the 22 subsets ( $M_{\text{win}} = 0.33$ ) was very similar to the average  $T$  in the corresponding 22 overall meta-analyses ( $M_{\text{win}} = 0.37$ ),  $t(13) = 1.39$ ,  $p = .187$ . Powerful moderators might only emerge when they are based on theoretical considerations (Tipton et al., 2019a, 2019b). We therefore looked at 10 of the 22 meta-analyses that presented a theoretical rationale for the moderator. Again, we used these moderators to split the studies from each meta-analysis into appropriate subsets and repeated the previous analysis. We found that the average  $T$  in the 10 moderator-based subsets ( $M_{\text{win}} = 0.37$ ) was again very similar to the average  $T$  in the 10 corresponding overall meta-analyses ( $M_{\text{win}} = 0.37$ ),  $t(5) = 0.73$ ,  $p = .499$ .

This moderator analysis does not suggest that the large heterogeneity in our meta-analysis sample is readily explained by mixing apples and oranges. Still, the possibility remains that authors (potentially unwisely) combine highly diverse studies and then fail to address relevant moderators. To address this point, we rated the broadness or narrowness of the inclusion criteria for each meta-analysis, using a single, global five-point scale ranging from *low* to *high*. Ratings considered the extent to which the addressed question was broad (e.g., “How effective is psychotherapy?”) versus narrow (e.g., “How effective is cognitive behavioral therapy in treating simple phobias?”), the extent to which the manipulation of the independent variable and the measurement of the dependent variable followed a standard protocol, and the similarity of the samples included. Ratings were conducted by the second author without knowledge of the actual levels of heterogeneity; to establish reliability, a random sample of 30 meta-analyses were independently rerated by the first author. We computed interrater agreement as Cohen’s  $\kappa$  using quadratic weights and observed a  $\kappa_w$  of .73, which is typically interpreted as good (Jakobsson & Westergren, 2005). For the 58

meta-analyses whose broadness of inclusion criteria was rated low or low to medium, average heterogeneity was still very high ( $M_{\text{win}} = 0.29$ ). In other words, if authors generally avoided meta-analyses that integrate fairly diverse studies, the levels of observed heterogeneity would probably not be much lower. In sum, our analyses do not support the view that the unwise mixing of apples and oranges is a strong driver of observed heterogeneity in meta-analyses.

### ***Exploratory analyses on what drives heterogeneity***

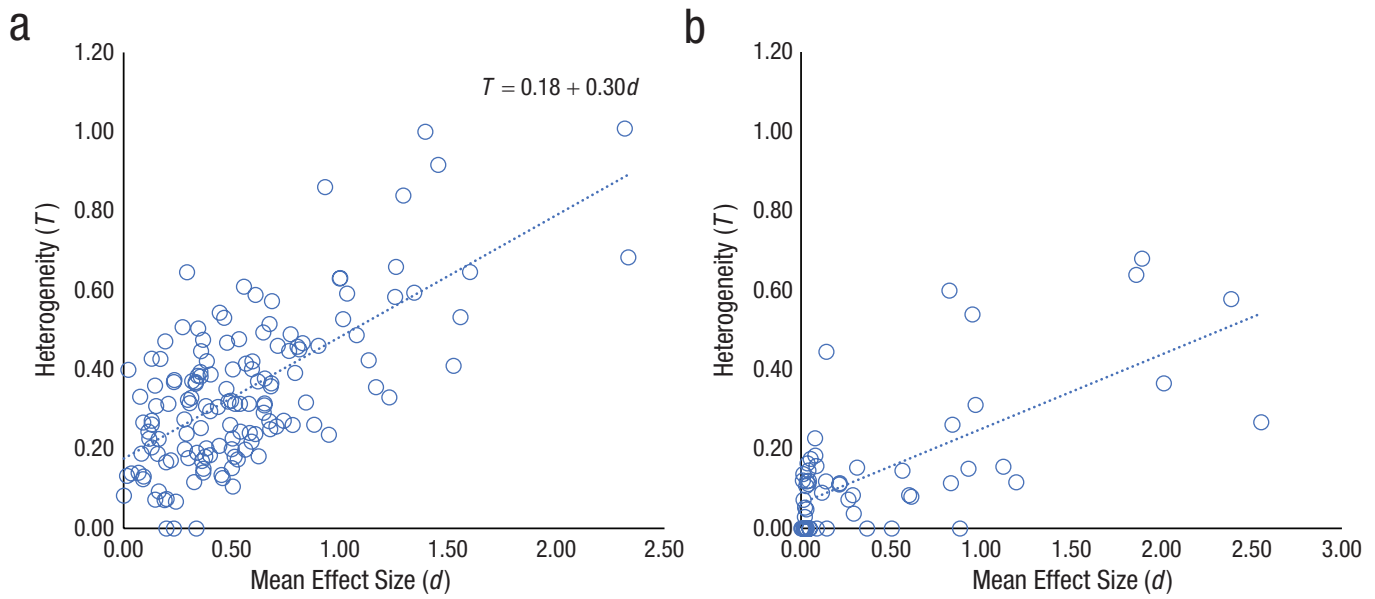
Heterogeneity differed substantially between meta-analyses ( $SD_{\text{win}} = 0.11$ ). A number of ideas have been proposed to explain this difference. Kenny and Judd (2019) suggested that research areas with larger average effect sizes should have greater levels of heterogeneity. We therefore correlated mean  $d$  with  $T$ . In line with Kenny and Judd (2019), we found a strong correlation ( $r_{\text{win}} = .49$ ,  $p < .001$ ) for the set of 150 meta-analyses (see Fig. 5). This correlation was replicated across all three subdisciplines (cognitive:  $r_{\text{win}} = .34$ ,  $p = .019$ ; social:  $r_{\text{win}} = .52$ ,  $p < .001$ ; and organizational:  $r_{\text{win}} = .62$ ,  $p < .001$ ). As shown by Figure 5, the relationship between mean  $d$  with  $T$  also held for the set of 57 close replications ( $r_{\text{win}} = .48$ ,  $p < .001$ ). The observed relationship between  $d$  and  $T$  might, at least partly, arise from differences in participants’ motivation across studies. When the treatment has a strong effect, it should typically make a big difference as to whether participants engage well with the study or not; when the treatment has only a small effect, such differences should be less consequential (Weiss et al., 2017).

In light of the link between mean  $d$  and  $T$ , the direct comparison of heterogeneity in close replications versus meta-analyses might appear doubtful. This is because mean  $d$  proved considerably lower in close replications (0.31) than in meta-analyses (0.47). Correction via the relevant regression equation (Fig. 5) suggests an average  $T$  of 0.27 for meta-analyses at  $d$  of 0.31 (the average for close replications). This is still much larger than that observed for close replications ( $M_{\text{win}} = 0.09$ ).

Looking at systematic reviews in health care, IntHout et al. (2015) found that smaller studies are more heterogeneous (measured as  $T^2$ ) than larger ones. We therefore worked out the median sample size for each meta-analysis and correlated this with  $T$ . This resulted in an  $r_{\text{win}}$  of  $-.10$  ( $p = .108$ , one-tailed), which does not suggest an important role for average sample size.

Richard et al. (2003) proposed that as a research field matures, the focus shifts from establishing an effect to exploring its boundaries, and this should increase heterogeneity in findings. If we accept the number of studies ( $k$ ) as a proxy for the maturity of a research





**Fig. 5.** Heterogeneity as a function of meta-analyses' mean effect size. The funnel plot in (a) shows 150 meta-analyses from cognitive, social, and organizational psychology; the funnel plot in (b) shows 57 meta-analyses for close replications.

field, we should expect a positive correlation between  $k$  and  $T$ . Indeed, Richard et al. (2003) found a correlation ( $r$ ) of .11 in a large survey of meta-analyses in social psychology. For our 150 meta-analyses, we found an  $r_{\text{win}}$  of .23 ( $p = .005$ ), consistent with this idea. An obvious alternative interpretation is that meta-analyses with broader inclusion criteria cast a wider net and will therefore include more studies than those that use narrow inclusion criteria (Murphy, 2017).

We thought to test these two competing explanations (exploring boundaries vs. broader inclusion criteria). If later research into an effect tends to explore its boundaries, we would expect to see higher heterogeneity in studies conducted late than in those conducted early. We therefore looked at all meta-analyses that seemed to capture a sufficiently mature research area and included those 82 with a  $k$  of 30 or greater and a time range for included studies of at least 10 years. For each of these meta-analyses, we then determined  $T$  separately for the earlier and the latter half of the included studies. Although the difference was in the expected direction (early:  $M_{\text{win}} = 0.336$ ,  $SD_{\text{win}} = 0.097$ ; late:  $M_{\text{win}} = 0.342$ ,  $SD_{\text{win}} = 0.110$ ), it was small and not statistically significant,  $t(49) = 1.01$ ,  $p = .319$ , and therefore did not support the idea of boundary exploration.

If the observed correlation between  $k$  and  $T$  resulted from a broader inclusion criteria, we would expect to see that meta-analyses with broader inclusion criteria show a higher  $T$ . We therefore correlated our ratings of the broadness of inclusion criteria with  $T$  and found an  $r_{\text{win}}$  of .12 ( $p = .142$ ). This does not offer strong evidence

that the observed correlation between  $k$  and  $T$  reflects the broadness of inclusion criteria. In sum, it is not clear from our data why  $T$  tends to increase with  $k$ .

## Discussion

We found that the quantification of heterogeneity in meta-analyses is uncommon. When it is undertaken, authors rarely rely on the measure that we argue is most informative. Average heterogeneity proved to be  $T$  of 0.33 for meta-analyses, with powerful (or even decent) moderators being conspicuously absent, and  $T$  of 0.09 for close replications. Heterogeneity showed a strong positive association with average effect size. Although based on exploratory analyses, this finding is credible because of the strong consistency shown across all three subdisciplines and close replications.

The effect of heterogeneity on statistical power and its implications for the interpretation of low replicability rates in the Open Science Collaboration (2015) project has received considerable attention (Shrout & Rodgers, 2018; Stanley et al., 2018). We address these more specific issues here. We discuss more general implications for the progress of psychological science in the General Discussion.

### *The meaning of observed heterogeneity levels*

It is helpful to first consider the meaning of average heterogeneity levels ( $T = 0.09$  for close replications and

$T = 0.33$  for meta-analyses). As we discussed earlier, this partly depends on the effect size ( $d$ ), which averaged 0.31 and 0.47 (Winsorized means for close replications and meta-analyses, respectively).<sup>5</sup> The average result for meta-analyses is depicted in Figure 2 (solid line). Remember that Cohen's  $d$  values of 0.2, 0.5, and 0.8 are often considered benchmarks for small, medium, and large effects, respectively. All of these values occur frequently in the distribution of population effect sizes for the average meta-analysis. Therefore, the observed heterogeneity in conceptual replications appears to be large. In contrast, the average close replication showed moderate variability in population effect sizes.

We further illustrate the meaning of heterogeneity with two examples from cognitive psychology. To further understand the importance of working memory for second-language proficiency development and processing, Linck et al. (2014) investigated the strength of this link in a meta-analysis. Included studies used a range of working memory tasks and second-language comprehension measures in diverse samples. The strength of the relationship proved medium in size ( $d = 0.51$ ), and heterogeneity was estimated to be low ( $T = 0.11$ ; see Fig. 1a). The latter implies high consistency of the relationship and ready generalizability across paradigms. In line with this idea, most population effect sizes (mean  $\pm 1$   $SD$ ) should fall in a narrow range of medium-sized effects ( $d = 0.40$ – $0.62$ ).

Baker et al. (2014) used meta-analysis to investigate the degree of independence between general intelligence and mental-state understanding. Included studies used a range of established intelligence tests in diverse samples; however, all studies used the same widely used test of mental-state understanding (Reading the Mind in the Eyes test). As in the previous example, the strength of the relationship proved medium in size ( $d = 0.49$ ); however, heterogeneity was estimated to be much higher ( $T = 0.35$ ; Fig. 1b), although the same test of mental-state understanding was used throughout. Although the observed level of heterogeneity was medium and not large, it already implies low consistency of the studied relationship and a lack of generalizability across paradigms. Most population effect sizes (mean  $\pm 1$   $SD$ ) should fall in a wide range of very small to large effects ( $d = 0.14$ – $0.84$ ). Baker et al. (2014) reported a statistically significant moderator, but given its atheoretical nature and the number of moderators tested, it remains debatable whether this reflects progress in understanding or successful capitalization on chance (Ioannidis, 2008).

In sum, it appears that the relationship between working memory and second-language proficiency is better understood than that between intelligence and performance on the Reading the Mind in the Eyes test. More

generally, everything else being equal, meta-analyses with lower heterogeneity will be more informative.

### **Potential biases in heterogeneity estimates**

Before we address the implications of these findings in greater detail, it is necessary to highlight a number of points regarding the trustworthiness of our estimates.

**Representativeness of our samples.** Our sampling of meta-analyses in cognitive, organizational, and social psychology was rigorous, and perusal of the topics (see Table S1) confirms a broad coverage of topics typical for these subdisciplines. We did not find evidence for heterogeneity differences across these subdisciplines, which might indicate that our results generalize more broadly across psychology. This is supported by the fact that Stanley et al. (2018), in a broader sample of meta-analyses from *Psychological Bulletin*, found heterogeneity levels (median  $T = 0.35$ ) very similar to ours ( $T = 0.33$ ). In contrast, representativeness cannot be claimed for our sample of close replications. Because of the novelty of the concept and the enormous effort involved, such Many Labs-type studies are relatively rare and focus on effects that are relatively easy to study; for this reason, the set of original studies that motivated these replications cannot be considered to be representative of the three subdisciplines we studied or psychological research in general. The observed difference in average effect size (close replications:  $d = 0.31$ ; meta-analyses:  $d = 0.47$ ) reinforces this point. It therefore remains unknown how readily the low heterogeneity observed in close replications would generalize to findings in psychological research in general.

**Publication bias and questionable research practices.** Publication bias (Sterling, 1959) and questionable research practices (QRPs; Simmons et al., 2011) are problems in psychological research (John et al., 2012; McShane et al., 2016; Simmons et al., 2011; Sterling, 1959). As a result, only a biased sample of all conducted studies appears in the published literature; “unsuccessful” studies typically remain invisible. Given that larger compared with smaller observed effects are more likely to be statistically significant (and thus “successful”), publication bias leads to upwardly biased effect sizes in published studies and meta-analyses (e.g., McShane et al., 2016). To achieve statistically significant, and therefore publishable, results, researchers might resort to QRPs (e.g., collect a number of similar dependent variables but report findings from only the most successful one). QRPs can dramatically increase the rate of false-positive results (Simmons et al., 2011) and thus lead again to inflated effect sizes in published studies and meta-analyses.

Mathematical modeling and computer simulations suggest that publication bias can lead to an underestimation or overestimation of heterogeneity; however, the former tends to be more prevalent than the latter (Augusteijn et al., 2019; Jackson, 2006). In addition, the overestimation of heterogeneity resulting from publication bias and QRPs tended to be much smaller than the levels of heterogeneity observed in conceptual replications here (Hönemann & Linden, 2019). From this viewpoint, the very large  $T$  in conceptual replications cannot be attributed to bias but instead seems to represent real heterogeneity that is not well understood. Our heterogeneity estimates for close replications are not affected in this way because the protocol for Many Labs-type replications precludes publication bias and QRPs (e.g., Shrout & Rodgers, 2018).

**Overreliance on WEIRD samples.** Studies in psychology, even if they seek to address human nature in general, rely almost exclusively on samples from Western, educated, industrialized, rich, and democratic (WEIRD) societies. Henrich et al. (2010) argued that WEIRD samples are among the least suitable to make general inferences about human nature and that many phenomena that are well established in WEIRD populations fail to generalize to other populations. Obviously, this is a concern only for those studies that seek to address human nature; however, this is frequently the case. For these cases, the findings from Henrich et al. imply that observed heterogeneity would often be higher if researchers did not rely almost entirely on WEIRD samples. This should hold equally for conceptual and close replications.

**Accuracy of meta-analyses.** For feasibility reasons, we had to rely on reported effect sizes for the underlying primary studies. One systematic investigation of meta-analyses in medicine found that about one in five effect-size computations for primary studies was erroneous (Gøtzsche et al., 2007). This should add (error) variance to the meta-analysis and consequently inflate observed heterogeneity. Given the strict protocols and high degree of transparency for Many Labs-type studies (e.g., Klein et al., 2014, 2018), erroneous effect sizes should be less of a concern for close replications.

**Summary.** In sum, our data for cognitive, organizational, and social psychology should be fairly representative for these disciplines, and results might generalize fairly well beyond. Publication bias, QRPs, and overreliance on WEIRD samples should artificially lower heterogeneity estimates; meta-analytic errors regarding the extraction of effect sizes from primary studies should have the opposite effect. On balance then, there is no strong evidence to suggest that our very-high heterogeneity estimates grossly

overestimate actual levels of heterogeneity. If anything, heterogeneity-deflating biases appear more numerous than heterogeneity-inflating biases. Thus, our results suggest that actual heterogeneity is typically very high in sets of conceptual replications. Although the representativeness of our close-replications sample is unclear, resulting heterogeneity estimates should, overall, be less prone to error than those for conceptual replications.

### ***Implications for the replicability of close replications***

As discussed previously, the Open Science Collaboration (2015) project famously attempted close replications of 100 studies. Although larger samples were used than in the original studies, statistical significance was achieved in only 36% of replications (25% in social psychology and 50% in cognitive psychology). This finding has become a catalyst of the controversial debate about the health of psychology research, which is still ongoing (e.g., Earp & Trafimow, 2015; Pashler & Harris, 2012; F. L. Schmidt & Oh, 2016; Simons, 2014; Stroebe & Strack, 2014). This is not the place to review this debate (for a comprehensive summary, see Zwaan et al., 2018), but one of its strands is of particular interest here. Stanley et al. (2018) suggested that heterogeneity accounts for the Open Science Collaboration's low replication rates. The authors estimated heterogeneity to be high ( $T = 0.25$ ) on the basis of only two Many Labs-type close replications. Subsequent power calculations demonstrated that heterogeneity should therefore decrease power in the typical psychological study to levels that are in line with the low replication rate observed in Open Science Collaboration (2015). However, heterogeneity of the order we observed in a much larger sample for close replications (mean  $T = 0.09$ ) reduces statistical power only marginally.<sup>6</sup> If we stick to sample sizes that generate 80% power at zero heterogeneity, power does not drop at all for large effects, drops to 78% for medium effects, and drops to 70% for small effects (McShane and Böckenholt, 2014). The mean effect size for the original studies included in Open Science Collaboration (2015) was large ( $d = 0.87$ ).<sup>7</sup> Therefore, replication power should not be greatly affected provided that the differences between the Open Science Collaboration replication studies and their original counterparts are comparable with the differences in multiple close replications.

Moreover, if replication failure reflects heterogeneity-driven low power, as Stanley et al. (2018) claimed, the large difference in replication rates between cognitive and social psychology (Open Science Collaboration, 2015) should be reflected in larger heterogeneity in the latter. Our finding of virtually identical heterogeneity

levels across cognitive and social psychology does not support this view. In conjunction with the low heterogeneity observed in close replications, it strengthens the interpretation that the low replication rate demonstrated in the Open Science Collaboration might be attributable to publication bias and QRPs. This is good news from our perspective because promising strategies to combat these biases have been developed (Munafò et al., 2017). On a more general level, one may note that the central issue with the results from the Open Science Collaboration is less about the percentage of original results that are true and more about the suggestion that a key plank in our common standards to accept evidence as valid ( $p < .05$ ) has little utility (Sedlmeier & Renkewitz, 2018, p. 621).

### ***How should heterogeneity be estimated for power calculations?***

Average levels of heterogeneity ( $T = 0.33$ ) have dramatic effects on power, which drops from 80% to 69% for large effects, 80% to 63% for medium effects, and 80% to 56% for small effects (McShane & Böckenholt, 2014). What level of heterogeneity should we expect for a new study that is not a close replication? This is an important question for proper sample-size planning. The researcher's informed judgment will always be necessary; however, the following suggestions might appear sensible. If a relevant meta-analysis reports  $T$ , use this. If such a meta-analysis reports only the effect size, use a  $T$  value of  $0.18 + 0.30d$  (see Fig. 5) to estimate heterogeneity. If there is no meta-analysis, the heterogeneity can still be estimated (although with lower precision) from the effect size of a single study. When we used all effect sizes from all 150 meta-analyses in our sample to predict heterogeneity,  $T = 0.28 + 0.11d$  was the resulting regression ( $R = .38$ ). Finally, when an effect-size estimate is not available, use the mean ( $T = 0.33$ ).

### **Conclusions**

We suggested that heterogeneity is a useful perspective for reflecting the degree of understanding psychology achieves. Science can be described as a quest to explain the apparent complexity of the natural world through simpler, fundamental principles. Empirical cumulativeness reflects the extent to which empirical findings fit such a simple or explicable pattern. All else being equal, high levels of (unexplained) heterogeneity indicate lower empirical cumulativeness (Asendorpf et al., 2013; Hedges, 1987; Murphy, 2017; Richard et al., 2003; Sells, 1963). For conceptual replications in three of psychology's core disciplines (and plausibly beyond; see Stanley et al., 2018; van Erp et al., 2017), we found

that heterogeneity is typically large (see Fig. 2) and unexplained, with little reason to believe that our estimates are inflated. To add some perspective, we can compare typical levels of heterogeneity (variability within a specific topic) with the variability in mean effect sizes across meta-analyses (variability between topics). Whereas we found a  $T$  value of 0.33 for the former, for the latter we observed an  $SD$  of 0.42 across all 150 meta-analyses. In other words, variability within phenomena measured in this way is only about 20% less than variability between phenomena. This large heterogeneity is sobering, as it reflects low empirical cumulativeness and therefore low coherence between the concepts researchers use and the data observed. On a brighter note, our findings also showed that large heterogeneity is not inevitable—in close replications, it was typically of moderate magnitude ( $T = 0.09$ )—and even hard sciences face some heterogeneity in their measurements (Hedges, 1987).

Before we explore important implications of this twin finding and possible improvements for our collective research practice, we address a likely objection to our argument that heterogeneity meaningfully reflects the degree of understanding psychology achieves.

### ***Reply to an objection***

A likely objection is that progress is driven by theories and that effect sizes tend to be irrelevant for most psychological theories (e.g., Baumeister, 2016; Strack, 2017); if effect sizes are largely irrelevant, their variability (i.e., heterogeneity) is likewise of little consequence. We think that such a perspective is mistaken for a number of reasons. First, even if effect sizes were largely irrelevant, the direction of effects remains important: In the face of large heterogeneity, the direction of an effect might be difficult to predict. Second, effect sizes are by no means irrelevant for increasing understanding; therefore, their degree of variability is also important. Although some psychological theories are not rooted in quantitative concepts (e.g., Piaget's stages in cognitive development), most psychological research is rooted in measurement. Given that measurement is regarded as a practically indispensable tool for investigation, it seems inconsistent to be disinterested in its result. In general, strong theories tend to be specific in the sense that they declare a large range of potential observations to be contrary to theory, thereby creating ample scope for the theory to be empirically challenged (Kuhn, 1970). Likewise, the ability to make precise predictions is often a hallmark of more mature science (Schickore, 2018). Effect sizes are obviously not the only route to achieve such specificity, but they may often provide a viable way forward. If heterogeneity is high, such specificity is difficult to achieve.

Finally, effect sizes are highly relevant for both explanations and practical applications. Psychological explanations typically rely on probabilistic relationships (e.g., in mate choice, men tend to put more emphasis on a partner's physical attractiveness than women; Feingold, 1990), and, all else being equal, stronger effects convey better explanations (Woodward, 2014). For example, the sex difference in height (approximately  $d = 2$ ; see Gustafsson & Lindenfors, 2009) is much stronger than the sex difference in relevance of attractiveness for mate choice (approximately  $d = 0.5$ ). Thus, "because she is a woman" more suitably explains why Aminah is shorter than Muhammad than why physical attractiveness matters less in her mate search than in his. If heterogeneity is large, it becomes unclear how powerful particular explanations are, which is obviously undesirable. Likewise, effect sizes are also highly relevant for practical applications. For example, sleep quality is a particularly strong predictor of adolescents' well-being, and for this reason it is a particularly promising lever for improving young people's well-being (Gireesh et al., 2018). Again, if heterogeneity is large, the effect of any intervention, which might be thought of as another conceptual replication, becomes more difficult to predict; and unless the average effect size is large, even the direction of the effect of the intervention could be uncertain (Fig. 2). In line with this idea, successful interventions can rarely be delineated from research findings but need to be tested (Cowen et al., 2017).

### ***Implications for testing theories***

Our twin finding of large heterogeneity in conceptual replications and moderate heterogeneity in close replications has important implications for the testing of theories.

***Knowledge as a tool.*** One relates to the use of knowledge as a tool. Imagine a situation in which the test of a psychological theory X requires inducing a particular mood. If this mood induction is based on a general principle that shows large heterogeneity, a negative finding of the test can be blamed on (unreliable) methods, and theory X is protected from failure. If heterogeneity thus precludes the meaningful empirical scrutiny of theories, theoretical progress will be limited (Ferguson & Heene, 2012; Greenwald, 2012; Kerr, 1998; LeBel & Peters, 2011; Meehl, 1978). In this context, the moderate heterogeneity observed in close replications ( $T = 0.09$ ) is encouraging, and it has a clear implication: A test of theory X should not rely on a general principle of mood induction; it should stick closely to a particular, successful protocol instead. This should typically bring about the expected change in mood.<sup>8</sup> In this context, we note that any

systematic investigation about which psychological effects are particularly reliable (i.e., strong and with low heterogeneity) is curiously absent.

***Theories' boundaries.*** Another implication of our findings is that the evaluation of theories also requires a broad exploration of the "research space" (Asendorpf et al., 2013), that is, the space defined by the combination of different manipulations of the independent variable, different dependent variables, different study populations, and so on. As an example, consider the set of stimuli used. If only a single standard set is used in a research domain to evoke the expected effect, some theory-irrelevant feature of that set might drive the observed effect (Fiedler, 2011). This problem can be detected only by using diverse (but theory-conforming) sets of stimuli. Consider also the case in which a theory offers a narrow explanation to account for an observation (e.g., memory for a word list is improved when the survival value of its items is to be judged). If a more general and thus more parsimonious explanation holds (e.g., memory for a word list is improved by any judgments that trigger self-referent encoding), this can be discovered only by testing instances of the research space that violate the overly narrow theory while still holding for the more general account (Fiedler et al., 2012; Shrouf & Rodgers, 2018).

***Meta-analysis and the testing of theories.*** A good theory should specify its scope. To evaluate the theory, meta-analysts must move beyond a narrow focus on the mean effect size and its statistical significance and take heterogeneity into account. This is obviously not a new insight (e.g., Higgins & Thompson, 2002; Hunter & Schmidt, 1990). However, our results regarding the reporting of heterogeneity in meta-analyses suggest this is rarely implemented in practice. One reason might be that frequently used approaches to heterogeneity fail to appeal to researchers' imagination: As shown earlier, quantification of heterogeneity is often missing or expressed in ways that might elude intuitive understanding ( $I^2$ ,  $T^2$ ). An increased focus on  $T$  might facilitate thinking about heterogeneity.

### ***Reducing unexplained heterogeneity as a sensible heuristic to advance understanding***

Given that unexplained heterogeneity tends to be both large and undesirable, its reduction should become an important goal. Among other advantages, this will increase coherence between the concepts we use and our observational data, facilitate empirical scrutiny of our theories, provide greater clarity regarding the power of the explanations we can offer, and facilitate the design of practical applications. Weiss et al. (2014)

offered a conceptual framework for heterogeneity in experiments, which is useful for discussing measures to either explain or reduce it.

**A conceptual framework for heterogeneity.** According to Weiss et al. (2014), heterogeneity in a set of experiments arises from three sources. First, studies can differ in their treatment contrasts, that is, the experimentally induced difference between the experimental and control group. The second source of heterogeneity are moderators that reside in the participants. Thus, if an effect is age-dependent, differences in participants' age across studies will induce heterogeneity. Finally, studies might differ on relevant context moderators; for example, an effect might vary across cultures or situations. Fruitful applications of this framework can be found in Weiss et al. (2017).

**Treatment contrasts.** Differences in studies' treatment contrasts will typically be driven by the strength of experimental manipulations. Stronger manipulations will often bring about stronger effects than weaker manipulations. Variability in the strength of manipulations across studies will thus induce heterogeneity in the results. If the strength of manipulations cannot be (or is not) properly expressed, it will be difficult to explain this heterogeneity. The unspecified or underspecified strength of experimental manipulations strikes us as a frequent issue across psychology that could often be avoided. We take the effect of bilateral symmetry on facial attractiveness as an arbitrary example. Correlational studies and experiments alike suggest that symmetry increases facial attractiveness (Rhodes, 2006). If the strength of experimental symmetry manipulations was described in relation to the natural variation in facial symmetry on which correlational studies rely, the variability of symmetry could be described on a common scale across all studies. These between-study differences in variability of symmetry (whether naturally occurring or experimentally induced) should be able to explain differences in results across studies and thus reduce heterogeneity. We are not aware of such attempts.

Our suggestion that systematically specifying the strength of manipulations of the independent variable will prove helpful is underpinned by the observation that many seminal insights in behavioral science relied on descriptions of the independent variable on a ratio scale. This is true for probabilities in classical conditioning (Rescorla & Wagner, 1972), operant conditioning (Herrnstein, 1961), perception under uncertainty (Tanner & Swets, 1954), and judgments and decision-making under uncertainty (Gigerenzer et al., 1991; Kahneman & Tversky, 1979); for the temporal relationship of stimuli or events and their effects on visual perception (Marcel, 1983), memory (Peterson & Peterson, 1959),

and the discounting of future outcomes (e.g., Frederick, 2002); the physical stimulus intensity and its relationship with perceived stimulus intensity (Stevens, 1957); and for degrees of genetic similarity, which underpin all estimates of the heritability of psychological traits (Plomin, 1990).

Differences in studies' treatment contrasts can also be affected by differences in the control groups, particularly in the case of real-world interventions. For these, "business as usual" (i.e., what it means *not* to be assigned to the intervention) will often differ in important ways between studies (Weiss et al., 2014, 2017). For example, an intervention promoting healthy behavior by providing information about health risks might create only a small treatment contrast in an environment in which ample information on health risks is at hand but a large treatment contrast in an environment in which such information is scarce.

**Person and context moderators.** The experimental test of a motivational intervention conducted by Yeager et al. (2019) provides an excellent illustration for both person and context moderators. Their short online intervention taught a nationally representative sample of U.S. students in secondary education that they can train their intellectual abilities like a muscle, which proved to have a small positive effect on students' grades. The authors hypothesized and confirmed that low-achieving students would benefit more from the intervention than high-achieving students (person moderator) and that the intervention would be most effective in schools with supportive peer norms (context moderator).

A meta-analytic search for moderators is most promising when it is driven by theory (Tipton et al., 2019a, 2019b). In this context it is noteworthy that psychologists have devoted great energy to describing individual differences in systematic ways (e.g., McCrae & Costa, 1997) but that comparable approaches to classify situations are, to the best of our knowledge, missing.

**Multisite experiments.** Meta-analyses are often limited in their ability to explain heterogeneity because relevant information on moderators or other sources of heterogeneity is unavailable for some or all of their primary studies. Multisite experiments, which directly address potential moderators in their design, are a promising alternative (e.g., Yeager et al., 2019). Such experiments are naturally arduous, but collaboration between many researchers through crowdsourcing holds great potential for such projects (Uhlmann et al., 2019).

**Standardized versus original-units effect sizes.** Finally, we want to draw attention to points outside of the heterogeneity framework proposed by Weiss et al.

**Table 2.** Irrelevant Differences in Standard Deviations Across Studies Negatively Affect the Suitability of Standardized Effect Sizes

Study	Control group ( <i>N</i> = 200)	Experimental group ( <i>N</i> = 200)	Difference	<i>d</i>
1	50.0 (10.0)	60.0 (10.0)	10.0	1.00
2	50.0 (15.0)	60.0 (15.0)	10.0	0.67
3	50.0 (10.0)	56.7 (10.0)	6.7	0.67

Note: Values are means with standard deviations in parentheses. Three similar, fictitious studies into the same phenomenon use the same dependent variable. Because Study 2 used a more diverse sample, the standardized effect size *d* misleadingly suggests that Studies 2 and 3 obtained the same results, whereas the difference in means shows that Studies 1 and 2 obtained the same results.

(2014). Our treatment of heterogeneity was based on descriptions of individual study results using standardized effect sizes. This is the norm for meta-analyses and conveys the obvious advantage that studies can be sensibly integrated even when they use different dependent variables. Nonetheless, standardized effect sizes might not be the best way to capture study results (Baguley, 2009; Bond et al., 2003; Tukey, 1969). Table 2 provides an example in which differences in sample means might be said to provide a more accurate description of individual results and of their differences across studies. This increase in accuracy might lead to reduced heterogeneity estimates and to the clearer emergence of informative moderators. The wealth of available data from Many Labs-type close replication studies (in which sets of close replications share the same dependent variable) provides rich opportunities for developing heterogeneity analyses on the basis of mean differences instead of standardized effect sizes and establishes whether this reduces heterogeneity estimates. If that is the case, we should also investigate the extent to which this can be fruitfully used for the analysis of conceptual replications.

Critics might argue that the portrayed shortcoming in standardized effect sizes (see Table 2) undermines our survey of heterogeneity. However, heterogeneity on the scale that we observed in conceptual replications cannot result from moderate inaccuracies in standardized effect sizes. Large heterogeneity is real, and its reduction should therefore become an important aim. To judge whether we make progress on this issue and to learn which strategies are best suited to reduce unexplained heterogeneity, its measurement is necessary. The approach we presented here strikes us as the most appropriate currently available.

## Outlook

Chemists in the 18th century, who did not yet understand the difference between compounds and mixtures, realized that substances often combine in fixed proportions

(e.g., you need 61.5 g of magnesia to neutralize 100 g of sulfuric acid; Leicester, 1965). Although useful for their daily practice, they did not attach much importance to this regularity because it appeared to lack universality (after all, you can mix one or three spoons of sugar in a cup of tea). Early in the 19th century, John Dalton parsed the seemingly incongruous observational data in a new way and realized the significance of fixed proportions, thus paving the way for the measurement of relative atomic weights and atomic theory, a major breakthrough in the history of chemistry (Kuhn, 1970). The linear relationship between stars' distance from Earth and the speed at which they move away from us was probably more obvious to perceive: Within a short time span, Georges Lemaître and Edwin Hubble independently discovered this law and, consequently, the expansion of the universe (Schneider, 2014). These examples illustrate that (a) regularity in observational data often acts as a lodestar for discovery (Simon, 1973) and (b) even the identification of pockets of regularity might be greatly beneficial. Reducing heterogeneity should make it easier for psychologists to perceive such regularities, and the prospect of new discoveries might be the strongest incentive to do so. We suggested some means to this end. We are sure that, once heterogeneity and its reduction receives more of the attention it deserves, the ingenuity of our colleagues will greatly add to our own ideas.

## Transparency

*Action Editor:* Laura A. King

*Editor:* Laura A. King

*Declaration of Conflicting Interests*

The author(s) declared that there were no conflicts of interest with respect to the authorship or the publication of this article.

## ORCID iD

Johannes Hönekopp  <https://orcid.org/0000-0002-0613-6801>

## Acknowledgments

We thank Frank Renkewitz and Peter Sedlmeier for critical comments on an earlier version of this article.

## Notes

1. This is because the gain in power for larger-than-expected effects is less than the loss in power for smaller-than-expected effects.
2. Indeed, across 141 meta-analyses reporting Cohen's  $d$ ,  $I^2$  and  $T$  correlated only moderately at  $r = .39$  (van Erp et al., 2017).
3. For Many Labs 2 (Klein et al., 2018), effect sizes for individual studies were not reported, but we could compute them from the published raw data.
4. We also computed analyses using the widely used Hunter-Schmidt, Paule-Mandel, and restricted maximum-likelihood estimators. These led to similar results (see Table S2 in the Supplemental Material) and the same conclusions.
5. We note that our value for meta-analyses corresponds closely to the average effect size ( $d$ ) of 0.39 observed in a sample of meta-analyses from *Psychological Bulletin* (Stanley et al., 2018).
6. Stanley et al. (2018) based their heterogeneity estimate for close replications on Eerland et al. (2016) and Hagger et al. (2016), both of which were part of our much larger sample.
7. The mean effect size ( $d$ ) of the original studies underlying the 57 Many Labs-type close replications was of 0.75 ( $SD = 0.37$ ). Regarding the effect size of the underlying original studies, Many Labs-type close replications and Open Science Collaboration replications are therefore comparable. This matters because of the observed link between  $T$  and mean effect size (see Fig. 5).
8. Note, however, that our heterogeneity estimate for close replications stems from preregistered studies published irrespective of their results. This precludes distortions of their effect sizes through publication bias and questionable research practices, which might often affect other published results (Ferguson & Brannick, 2012; Kühberger et al., 2014; LeBel & Peters, 2011; Levine et al., 2009; Pashler & Harris, 2012; Sterling, 1959). Thus, close replications based on a single published result might be less reliable (Open Science Collaboration, 2015).

## References

- Asendorpf, J. B., Conner, M., De Fruyt, F., De Houwer, J., Denissen, J. A., Fiedler, K., Fiedler, S., Funder, D. C., Kliegl, R., Nosek, B., Perugini, M., Roberts, B. W., Schmitt, M., van Aken, M. A. G., Weber, H., & Wicherts, J. M. (2013). Recommendations for increasing replicability in psychology. *European Journal of Personality*, *27*, 108–119.
- Augusteijn, H. E., van Aert, R., & van Assen, M. A. (2019). The effect of publication bias on the Q test and assessment of heterogeneity. *Psychological Methods*, *24*, 116–134.
- Aytug, Z. G., Rothstein, H. R., Zhou, W., & Kern, M. C. (2012). Revealed or concealed? Transparency of procedures, decisions, and judgment calls in meta-analyses. *Organizational Research Methods*, *15*, 103–133.
- Baguley, T. (2009). Standardized or simple effect size: What should be reported? *British Journal of Psychology*, *100*, 603–617.
- Baker, C. A., Peterson, E., Pulos, S., & Kirkland, R. A. (2014). Eyes and IQ: A meta-analysis of the relationship between intelligence and “Reading the Mind in the Eyes.” *Intelligence*, *44*, 78–92.
- Baumeister, R. F. (2016). Charting the future of social psychology on stormy seas: Winners, losers, and recommendations. *Journal of Experimental Social Psychology*, *66*, 153–158.
- Bond, C. F., Jr., Wiitala, W. L., & Richard, F. D. (2003). Meta-analysis of raw mean differences. *Psychological Methods*, *8*, 406–418.
- Borenstein, M., Hedges, L. V., Higgins, J., & Rothstein, H. R. (2009). *Introduction to meta-analysis*. Wiley Online Library.
- Borenstein, M., Higgins, J. P., Hedges, L. V., & Rothstein, H. R. (2017). Basics of meta-analysis:  $I^2$  is not an absolute measure of heterogeneity. *Research Synthesis Methods*, *8*(1), 5–18.
- Cheung, I., Campbell, L., LeBel, E. P., Ackerman, R. A., Aykutoğlu, B., Bahník, Š., Bowen, J. D., Bredow, C. A., Bromberg, C., Caprariello, P. A., Carcedo, R. J., Carson, K. J., Cobb, R. J., Collins, N. L., Corretti, C. A., Didonato, T. E., Ellithorpe, C., Fernández-Rouco, N., Fuglestad, P. T., . . . Yong, J. C. (2016). Registered replication report: Study 1 from Finkel, Rusult, Kumashiro, & Hannon (2002). *Perspectives on Psychological Science*, *11*, 750–764. <https://doi.org/10.1177/1745691616664694>
- Chung, Y., Rabe-Hesketh, S., & Choi, I. H. (2013). Avoiding zero between-study variance estimates in random-effects meta-analysis. *Statistics in Medicine*, *32*, 4071–4089.
- Cowen, N., Virk, B., Mascarenhas-Keyes, S., & Cartwright, N. (2017). Randomized controlled trials: How can we know “what works”? *Critical Review*, *29*, 265–292.
- Dieckmann, N. F., Malle, B. F., & Bodner, T. E. (2009). An empirical assessment of meta-analytic practice. *Review of General Psychology*, *13*, 101–115.
- Earp, B. D., & Trafimow, D. (2015). Replication, falsification, and the crisis of confidence in social psychology. *Frontiers in Psychology*, *6*, Article 621. <https://doi.org/10.3389/fpsyg.2015.00621>
- Ebersole, C. R., Atherton, O. E., Belanger, A. L., Skulborstad, H. M., Allen, J. M., Banks, J. B., Baranski, E., Bernstein, M. J., Bonfiglio, D. B. V., Boucher, L., Brown, E. R., Budiman, N. I., Cairo, A. H., Capaldi, C. A., Chartier, C. R., Chung, J. M., Cicero, D. C., Coleman, J. A., Conway, J. G., . . . Nosek, B. A. (2016). Many Labs 3: Evaluating participant pool quality across the academic semester via replication. *Journal of Experimental Social Psychology*, *67*, 68–82.
- Eerland, A., Sherrill, A. M., Magliano, J. P., Zwaan, R. A., Arnal, J. D., Aucoin, P., Berger, S. A., Birt, A. R., Capezza, N., Carlucci, M., Crocker, C., Ferretti, T. R., Kibbe, M. R., Knepp, M. M., Kurby, C. A., Melcher, J. M., Michael, S. W., Poirier, C., & Proulx, J. M. (2016). Registered replication report: Hart & Albarraçin (2011). *Perspectives on Psychological Science*, *11*, 158–171. <https://doi.org/10.1177/1745691616605826>
- Erceg-Hurn, D. M., & Mirosevich, V. M. (2008). Modern robust statistical methods: An easy way to maximize the accuracy



- and power of your research. *American Psychologist*, *63*, 591–601.
- Feingold, A. (1990). Gender differences in effects of physical attractiveness on romantic attraction: A comparison across five research paradigms. *Journal of Personality and Social Psychology*, *59*, 981–993.
- Ferguson, C. J. (2009). An effect size primer: A guide for clinicians and researchers. *Professional Psychology: Research and Practice*, *40*, 532–538.
- Ferguson, C. J., & Brannick, M. T. (2012). Publication bias in psychological science: Prevalence, methods for identifying and controlling, and implications for the use of meta-analyses. *Psychological Methods*, *17*, 120–128.
- Ferguson, C. J., & Heene, M. (2012). A vast graveyard of undead theories: Publication bias and psychological science's aversion to the null. *Perspectives on Psychological Science*, *7*, 555–561. <https://doi.org/10.1177/1745691612459059>
- Fiedler, K. (2011). Voodoo correlations are everywhere—not only in neuroscience. *Perspectives on Psychological Science*, *6*, 163–171. <https://doi.org/10.1177/1745691611400237>
- Fiedler, K., Kutzner, F., & Krueger, J. I. (2012). The long way from  $\alpha$ -error control to validity proper: Problems with a short-sighted false-positive debate. *Perspectives on Psychological Science*, *7*, 661–669. <https://doi.org/10.1177/1745691612462587>
- Fischhoff, B. (1990). Psychology and public policy: Tool or toolmaker? *American Psychologist*, *45*, 647–653.
- Frederick, S., Loewenstein, G., & O'Donoghue, T. (2002). Time discounting and time preference: A critical review. *Journal of Economic Literature*, *40*, 351–401.
- Gøtzsche, P. C., Hróbjartsson, A., Marić, K., & Tendal, B. (2007). Data extraction errors in meta-analyses that use standardized mean differences. *JAMA*, *298*, 430–437.
- Gigerenzer, G., Hoffrage, U., & Kleinbölting, H. (1991). Probabilistic mental models: A Brunswikian theory of confidence. *Psychological Review*, *98*, 506–528.
- Gireesh, A., Das, S., & Viner, R. M. (2018). Impact of health behaviours and deprivation on well-being in a national sample of English young people. *BMJ Paediatrics Open*, *2*(1), Article e000335. <https://doi.org/10.1136/bmjpo-2018-000335>
- Greenwald, A. G. (2012). There is nothing so theoretical as a good method. *Perspectives on Psychological Science*, *7*, 99–108. <https://doi.org/10.1177/1745691611434210>
- Gustafsson, A., & Lindenfors, P. (2009). Latitudinal patterns in human stature and sexual stature dimorphism. *Annals of Human Biology*, *36*(1), 74–87.
- Hagger, M. S., Chatzisarantis, N. L. D., Alberts, H., Anggono, C. O., Batailler, C., Birt, A. R., Brand, R., Brandt, M. J., Brewer, G., Bruyneel, S., Calvillo, D. P., Campbell, W. K., Cannon, P. R., Carlucci, M., Carruth, N. P., Cheung, T., Crowell, A., De Ridder, D. T. D., Dewitte, S., . . . Zwienerberg, M. (2016). A multilab preregistered replication of the ego-depletion effect. *Perspectives on Psychological Science*, *11*, 546–573. <https://doi.org/10.1177/1745691616652873>
- Hedges, L. V. (1987). How hard is hard science, how soft is soft science? The empirical cumulativeness of research. *American Psychologist*, *42*, 443–455.
- Henrich, J., Heine, S. J., & Norenzayan, A. (2010). The weirdest people in the world? *Behavioral and Brain Sciences*, *33*, 61–83.
- Herrnstein, R. J. (1961). Relative and absolute strength of response as a function of frequency of reinforcement. *Journal of the Experimental Analysis of Behavior*, *4*, 267–272.
- Higgins, J. P., & Thompson, S. G. (2002). Quantifying heterogeneity in a meta-analysis. *Statistics in Medicine*, *21*, 1539–1558.
- Hönekopp, J., & Linden, A. H. (2019). *Heterogeneity estimates in a biased world*. OSF. <https://osf.io/zx96p/>
- Huedo-Medina, T. B., Sánchez-Meca, J., Marín-Martínez, F., & Botella, J. (2006). Assessing heterogeneity in meta-analysis: Q statistic or  $I^2$  index? *Psychological Methods*, *11*, 193–206.
- Huhn, M., Tardy, M., Spineli, L. M., Kissling, W., Förstl, H., Pitschel-Walz, G., Leucht, C., Samara, M., Dold, M., & Davis, J. M. (2014). Efficacy of pharmacotherapy and psychotherapy for adult psychiatric disorders: A systematic overview of meta-analyses. *JAMA Psychiatry*, *71*, 706–715.
- Hunter, J. E. (1997). Needed: A ban on the significance test. *Psychological Science*, *8*, 3–7. <https://doi.org/10.1111/j.1467-9280.1997.tb00534.x>
- Hunter, J. E., & Schmidt, F. L. (1990). *Methods of meta-analysis: Correcting error and bias in research findings*. SAGE.
- Hunter, J. E., & Schmidt, F. (2004). *Methods of meta-analysis: Correcting error and bias in research findings* (2nd ed.). SAGE.
- Hyde, J. S. (2014). Gender similarities and differences. *Annual Review of Psychology*, *65*, 373–398.
- Int'Hout, J., Ioannidis, J. P., Borm, G. F., & Goeman, J. J. (2015). Small studies are more heterogeneous than large ones: A meta-meta-analysis. *Journal of Clinical Epidemiology*, *68*, 860–869.
- Ioannidis, J. (2008). Interpretation of tests of heterogeneity and bias in meta-analysis. *Journal of Evaluation in Clinical Practice*, *14*, 951–957.
- Jackson, D. (2006). The implications of publication bias for meta-analysis' other parameter. *Statistics in Medicine*, *25*, 2911–2921.
- Jakobsson, U., & Westergren, A. (2005). Statistical methods for assessing agreement for ordinal data. *Scandinavian Journal of Caring Sciences*, *19*, 427–431.
- John, L. K., Loewenstein, G., & Prelec, D. (2012). Measuring the prevalence of questionable research practices with incentives for truth telling. *Psychological Science*, *23*, 524–532. <https://doi.org/10.1177/0956797611430953>
- Kahneman, D., & Tversky, A. (1979). Prospect theory: An analysis of decision under risk. *Econometrica*, *47*, 263–291.
- Kenny, D. A., & Judd, C. M. (2019). The unappreciated heterogeneity of effect sizes: Implications for power, precision, planning of research, and replication. *Psychological Methods*, *24*, 578–589. <https://doi.org/10.1037/met000209>
- Kerr, N. L. (1998). HARKing: Hypothesizing after the results are known. *Personality and Social Psychology Review*, *2*, 196–217.

- Klein, R. A., Ratliff, K. A., Vianello, M., Adams, R. B., Bahník, Š., Bernstein, M. J., Bocian, K., Brandt, M. J., Brooks, B., Brumbaugh, C. C., Cemalcilar, Z., Chandler, J., Cheong, W., Davis, W. E., Devos, T., Eisner, M., Frankowska, N., Furrow, D., Galliani, E. M., . . . Nosek, B. A. (2014). Investigating variation in replicability. *Social Psychology*, *45*, 142–152.
- Klein, R. A., Vianello, M., Hasselman, F., Adams, B. G., Adams, R. B., Alper, S., Aveyard, M., Axt, J. R., Babalola, M. T., Bahník, Š., Batra, R., Berkics, M., Bernstein, M. J., Berry, D. R., Bialobrzeska, O., Binan, E. D., Bocian, K., Brandt, M. J., Busching, R., Rédei, A. C., . . . Nosek, B. A. (2018). Many Labs 2: Investigating variation in replicability across samples and settings. *Advances in Methods and Practices in Psychological Science*, *1*, 443–490. <https://doi.org/10.1177/2515245918810225>
- Kühberger, A., Fritz, A., & Scherndl, T. (2014). Publication bias in psychology: A diagnosis based on the correlation between effect size and sample size. *PLOS ONE*, *9*(9), Article e105825. <https://doi.org/10.1371/journal.pone.0105825>
- Kuhn, T. S. (1970). *The structure of scientific revolutions* (2nd ed.). University of Chicago Press.
- LeBel, E. P., & Peters, K. R. (2011). Fearing the future of empirical psychology: Bem's (2011) evidence of psi as a case study of deficiencies in modal research practice. *Review of General Psychology*, *15*, 371–379.
- Leicester, H. M. (1965). *The historical background of chemistry*. Wiley.
- Levine, T. R., Asada, K. J., & Carpenter, C. (2009). Sample sizes and effect sizes are negatively correlated in meta-analyses: Evidence and implications of a publication bias against nonsignificant findings. *Communication Monographs*, *76*, 286–302.
- Linck, J. A., Osthus, P., Koeth, J. T., & Bunting, M. F. (2014). Working memory and second language comprehension and production: A meta-analysis. *Psychonomic Bulletin & Review*, *21*, 861–883.
- Lipsey, M. W., & Wilson, D. B. (1993). The efficacy of psychological, educational, and behavioral treatment: Confirmation from meta-analysis. *American Psychologist*, *48*, 1181–1209.
- Marcel, A. J. (1983). Conscious and unconscious perception: Experiments on visual masking and word recognition. *Cognitive Psychology*, *15*, 197–237.
- McCrae, R. R., & Costa, P. T., Jr. (1997). Personality trait structure as a human universal. *American Psychologist*, *52*, 509–516.
- McShane, B. B., & Böckenholt, U. (2014). You cannot step into the same river twice: When power analyses are optimistic. *Perspectives on Psychological Science*, *9*, 612–625. <https://doi.org/10.1177/1745691614548513>
- McShane, B. B., Böckenholt, U., & Hansen, K. T. (2016). Adjusting for publication bias in meta-analysis: An evaluation of selection methods and some cautionary notes. *Perspectives on Psychological Science*, *11*, 730–749. <https://doi.org/10.1177/1745691616662243>
- Meehl, P. E. (1978). Theoretical risks and tabular asterisks: Sir Karl, Sir Ronald, and the slow progress of soft psychology. *Journal of Consulting and Clinical Psychology*, *46*, 806–834.
- Mitchell, G. (2012). Revisiting truth or triviality: The external validity of research in the psychological laboratory. *Perspectives on Psychological Science*, *7*, 109–117. <https://doi.org/10.1177/1745691611432343>
- Moher, D., Liberati, A., Tetzlaff, J., & Altman, D. G. (2009). Preferred reporting items for systematic reviews and meta-analyses: The PRISMA statement. *Annals of Internal Medicine*, *151*, 264–269.
- Munafò, M. R., Nosek, B. A., Bishop, D. V. M., Button, K. S., Chambers, C. D., Percie Du Sert, N., Simonsohn, U., Wagenmakers, E., Ware, J. J., & Ioannidis, J. P. A. (2017). A manifesto for reproducible science. *Nature Human Behaviour*, *1*, Article 0021. <https://doi.org/10.1038/s41562-016-0021>
- Murphy, K. R. (2017). What inferences can and cannot be made on the basis of meta-analysis? *Human Resource Management Review*, *27*(1), 193–200.
- Muth, A., Hönekopp, J., & Falter, C. M. (2014). Visuo-spatial performance in autism: A meta-analysis. *Journal of Autism and Developmental Disorders*, *44*, 3245–3263.
- Open Science Collaboration. (2015). Estimating the reproducibility of psychological science. *Science*, *349*, 943–951.
- Pashler, H., & Harris, C. R. (2012). Is the replicability crisis overblown? Three arguments examined. *Perspectives on Psychological Science*, *7*, 531–536. <https://doi.org/10.1177/1745691612463401>
- Peterson, L., & Peterson, M. J. (1959). Short-term retention of individual verbal items. *Journal of Experimental Psychology*, *58*, 193–198.
- Plomin, R. (1990). The role of inheritance in behavior. *Science*, *248*(4952), 183–188.
- R Core Team. (2017). R: A language and environment for statistical computing (Version 3.4.1) [Computer software]. The R Project for Statistical Computing. <http://www.R-project.org>
- Rescorla, R. A., & Wagner, A. R. (1972). A theory of Pavlovian conditioning: Variations in the effectiveness of reinforcement and nonreinforcement. In A. H. Black & W. E. Prokasy (Eds.), *Classical conditioning II: Current research and theory* (pp. 64–99). Appleton-Century-Crofts.
- Rhodes, G. (2006). The evolutionary psychology of facial beauty. *Annual Review of Psychology*, *57*, 199–226.
- Richard, F. D., Bond, C. F., Jr., & Stokes-Zoota, J. J. (2003). One hundred years of social psychology quantitatively described. *Review of General Psychology*, *7*, 331–363.
- Schickore, J. (2018). Scientific discovery. In E. N. Zalta (Ed.), *The Stanford encyclopedia of philosophy* (Summer 2018). <https://plato.stanford.edu/archives/sum2018/entries/scientific-discovery>
- Schmidt, F. L., & Hunter, J. E. (1977). Development of a general solution to the problem of validity generalization. *Journal of Applied Psychology*, *62*, 529–540.
- Schmidt, F. L., & Oh, I.-S. (2016). The crisis of confidence in research findings in psychology: Is lack of replication the real problem? Or is it something else? *Archives of Scientific Psychology*, *4*(1), 32–37.

- Schmidt, S. (2009). Shall we really do it again? The powerful concept of replication is neglected in the social sciences. *Review of General Psychology, 13*, 90–100.
- Schneider, P. (2014). *Extragalactic astronomy and cosmology: An introduction*. Springer.
- Sedlmeier, P., & Renkewitz, F. (2018). *Forschungsmethoden und Statistik für Psychologen und Sozialwissenschaftler* [Research methods and statistics for psychologists and social scientists]. Pearson Studium München.
- Sells, S. B. (1963). An interactionist looks at the environment. *American Psychologist, 18*, 696–702.
- Shrout, P. E., & Rodgers, J. L. (2018). Psychology, science, and knowledge construction: Broadening perspectives from the replication crisis. *Annual Review of Psychology, 69*, 487–510.
- Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2011). False-positive psychology: Undisclosed flexibility in data collection and analysis allows presenting anything as significant. *Psychological Science, 22*, 1359–1366. <https://doi.org/10.1177/0956797611417632>
- Simon, H. A. (1973). Does scientific discovery have a logic? *Philosophy of Science, 40*, 471–480.
- Simons, D. J. (2014). The value of direct replication. *Perspectives on Psychological Science, 9*, 76–80. <https://doi.org/10.1177/1745691613514755>
- Skinner, B. F. (1956). A case history in scientific method. *American Psychologist, 11*, 221–233.
- Smith, M. L., & Glass, G. V. (1977). Meta-analysis of psychotherapy outcome studies. *American Psychologist, 32*, 752–760.
- Stanley, T., Carter, E. C., & Doucouliagos, H. (2018). What meta-analyses reveal about the replicability of psychological research. *Psychological Bulletin, 144*, 1325–1346.
- Sterling, T. D. (1959). Publication decisions and their possible effects on inferences drawn from tests of significance—or vice versa. *Journal of the American Statistical Association, 54*, 30–34.
- Stevens, S. S. (1957). On the psychophysical law. *Psychological Review, 64*, 153–181.
- Strack, F. (2017). From data to truth in psychological science. A personal perspective. *Frontiers in Psychology, 8*, Article 702. <https://doi.org/10.3389/fpsyg.2017.00702>
- Stroebe, W., & Strack, F. (2014). The alleged crisis and the illusion of exact replication. *Perspectives on Psychological Science, 9*, 59–71. <https://doi.org/10.1177/1745691613514450>
- Tanner, W. P., & Swets, J. A. (1954). A decision-making theory of visual detection. *Psychological Review, 61*, 401–409.
- Tipton, E., Pustejovsky, J. E., & Ahmadi, H. (2019a). Current practices in meta-regression in psychology, education, and medicine. *Research Synthesis Methods, 10*, 180–194.
- Tipton, E., Pustejovsky, J. E., & Ahmadi, H. (2019b). A history of meta-regression: Technical, conceptual, and practical developments between 1974 and 2018. *Research Synthesis Methods, 10*, 161–179.
- Tukey, J. W. (1969). Analyzing data: Sanctification or detective work? *American Psychologist, 24*, 83–91.
- Uhlmann, E. L., Ebersole, C. R., Chartier, C. R., Errington, T. M., Kidwell, M. C., Lai, C. K., McCarthy, R. J., Riegelman, A., Silberzahn, R., & Nosek, B. A. (2019). Scientific utopia III: Crowdsourcing science. *Perspectives on Psychological Science, 14*, 711–733. <https://doi.org/10.1177/1745691619850561>
- van Erp, S., Verhagen, J., Grasman, R. P., & Wagenmakers, E.-J. (2017). Estimates of between-study heterogeneity for 705 meta-analyses reported in *Psychological Bulletin* from 1990–2013. *Journal of Open Psychology Data, 5*(1), Article 4. <https://doi.org/10.5334/jopd.33>
- Viechtbauer, W. (2010). Conducting meta-analyses in R with the metafor package. *Journal of Statistical Software, 36*(3), 1–48.
- Wagenmakers, E., Beek, T., Dijkhoff, L., Gronau, Q. F., Acosta, A., Adams, R. B., Albohn, D. N., Allard, E. S., Benning, S. D., Blouin-Hudon, E., Bulnes, L. C., Caldwell, T. L., Calin-Jageman, R. J., Capaldi, C. A., Carfagno, N. S., Chasten, K. T., Cleeremans, A., Connell, L., Decicco, J. M., & Zwaan, R. A. (2016). Registered replication report: Strack, Martin, & Stepper (1988). *Perspectives on Psychological Science, 11*, 917–928. <https://doi.org/10.1177/1745691616674458>
- Wallis, S. E. (2015). Integrative propositional analysis: A new quantitative method for evaluating theories in psychology. *Review of General Psychology, 19*, 365–380.
- Weiss, M. J., Bloom, H. S., & Brock, T. (2014). A conceptual framework for studying the sources of variation in program effects. *Journal of Policy Analysis and Management, 33*, 778–808.
- Weiss, M. J., Bloom, H. S., Verbitsky-Savitz, N., Gupta, H., Vigil, A. E., & Cullinan, D. N. (2017). How much do the effects of education and training programs vary across sites? Evidence from past multisite randomized trials. *Journal of Research on Educational Effectiveness, 10*, 843–876.
- Wilcox, R. R. (2005). *Introduction to robust estimation and hypothesis testing* (2nd ed.). Elsevier Academic Press.
- Woodward, J. (2014). Scientific explanation. In E. N. Zalta (Ed.), *The Stanford encyclopedia of philosophy* (Fall 2017). <https://plato.stanford.edu/archives/fall2017/entries/scientific-explanation>
- Yeager, D. S., Hanselman, P., Walton, G. M., Murray, J. S., Crosnoe, R., Muller, C., Tipton, E., Schneider, B., Hulleman, C. S., Hinojosa, C. P., Paunesku, D., Romero, C., Flint, K., Roberts, A., Trott, J., Iachan, R., Buontempo, J., Yang, S. M., Carvalho, C. M., . . . Dweck, C. S. (2019). A national experiment reveals where a growth mindset improves achievement. *Nature, 573*, 364–369.
- Zwaan, R. A., Etz, A., Lucas, R. E., & Donnellan, M. B. (2018). Making replication mainstream. *Behavioral and Brain Sciences, 41*, Article e120. <https://doi.org/10.1017/S0140525X17001972>