

Running Head: UNAPPRECIATED HETEROGENEITY OF EFFECTS

The Unappreciated Heterogeneity of Effect Sizes:
Implications for Power, Precision, Planning of Research, and Replication

David A. Kenny

University of Connecticut

Charles M. Judd

University of Colorado Boulder

David A. Kenny is a member of Department of Psychological Sciences and can be reached at david.kenny@uconn.edu and Charles M. Judd is a member of Department of Psychology and Neuroscience and can be reached at charles.judd@colorado.edu. We especially thank Gary McClelland and Christopher Rhoads, as well as Deborah Kashy, Joshua Correll, Vincent Yzerbyt, Blair Johnson, Betsy McCoach, and Thomas Ledermann who provided us with helpful feedback.

Abstract

Repeated investigations of the same phenomenon typically yield effect sizes that vary more than one would expect from sampling error alone. Although such heterogeneity of effect sizes is well documented, its implications for power analyses, the precision of effects, and the planning of research are not fully appreciated. In the presence of heterogeneity, there exists a range of effects that research studies might uncover, including effects that go in the opposite direction from the average effect reported in the literature. Additionally, the usual power calculations and confidence intervals are misleading, and the preference for definitive large- n studies is misguided. Effect size heterogeneity arises not only from measureable factors that moderate effects, but also possibly from random variation between studies. Our usual notions about what constitutes a “failure to replicate” an effect need to be rethought, given the ubiquity of heterogeneous effect sizes in the literature.

When multiple studies of the same effect exist, one expects variation in the estimated effect sizes because of sampling error, even if they are all estimating the same true effect. Typically, however, in meta-analyses, there exists more variation in effect sizes than can be attributed to sampling error alone, leading to the conclusion of heterogeneous effect sizes across those studies.

To index heterogeneous effect sizes, one can estimate their true standard deviation, over and above what might be expected from sampling error. For instance, if the effect size is a Cohen's d , its true standard deviation might be denoted as σ_δ . We subscript with δ , rather than d , to make clear that we are referring to variation in effects sizes after removing sampling error. Meta-analysts typically test whether this differs from zero using a chi-square test of homogeneity, symbolized as Q . Additionally, they commonly report I^2 : the percentage of variance in study effects due to heterogeneity, the remaining proportion being sampling error. Less commonly reported is the estimated true variance, generically called τ^2 in the meta-analysis literature.

Meta-analysts consistently find evidence of heterogeneity. Higgins (2008) wrote:

Heterogeneity is to be expected in a meta-analysis: it would be surprising if multiple studies, performed by different teams in different places with different methods, all ended up estimating the same underlying parameter (p. 1158).

Richard, Bond, and Stokes Zoota (2003) conducted a meta-analysis of meta-analyses in 18 different domains of social psychology (a total of 322 meta-analyses summarizing 33,912 individual studies). They reported an average effect size (r) of .21 and an average true standard deviation of those effect sizes of .15. In fact, in two of the domains (expectancy effects and social influence), the standard deviation of the effects was larger

than the average effect size, suggesting that many studies in those domains were estimating effects that were actually in the opposite direction from most of the other effects.

Individual examples of meta-analyses with heterogeneity include the following:

1) Decoster and Claypool (2004) studied the effect of priming on assimilation in impression formation across 45 studies and found an average d of 0.35 and an I^2 of 43.

2) Haelermans and Borghans (2012) studied the effectiveness of job training programs on wages across 163 studies and found an average percent increase in wages of 3.9 percent and an I^2 of 83.

3) Gershoff and Grogan-Kaylor (2016) studied the effect of physical punishment on children's externalizing behavior across 14 studies and found an average d of 0.41 and an I^2 of 88.

4) Hagger and Chatzisarantis (2016) conducted a meta-analysis of the ego-depletion effect, using only studies that employed the same dependent variable: reaction times. They found an average d of 0.04 and an estimate of σ_δ of 0.13.

This last example illustrates an important point. There are probably many good reasons for finding variation in effect sizes. For instance, different studies may use rather different operationalizations of the independent and dependent variables. However, even when we control for that difference, as did Hagger and Chatzisarantis (2016), effect size heterogeneity still occurs. There are many factors, both subtle and obvious, that are likely responsible for heterogeneity. At a later point, we discuss in greater detail why it is

reasonable to anticipate heterogeneity of effects sizes. Importantly, it is not just in the behavioral sciences that meta-analyses typically find evidence of heterogeneity. Studies in medicine (Turner, Davey, Clarke, Thompson, & Higgins, 2012) and cell biology (Aird, Candella, & Mantis, 2017) also find considerable heterogeneity.

Finding that effect sizes in a domain vary over and above what might be expected from sampling error implies that there is a range of true effect sizes that exist, rather than a single one. Within the meta-analysis literature, the recognition of heterogeneity has led to the development of procedures for “random effects” meta-analyses (Hedges & Vevea, 1998). However, the implications of heterogeneity, outside of the methodological literature on how to conduct meta-analyses, have not been fully explored. The goal of this paper is to examine the implications of effect size heterogeneity for power analysis, the precision of effect estimation, and the conduct of research.

Before we begin, we introduce notation and simplifying assumptions. Consider an effect that is measured using Cohen’s d . We assume that all studies utilize two equal-sized independent groups of participants and their standardized mean difference yields the effect size estimate. We assume each study has a total of n persons, with n likely varying across studies. The effect size estimate from the i^{th} individual study is d_i and it is an estimate of the true effect size for that study, δ_i . Across studies, there is a distribution of these true effect sizes, with a mean, denoted as δ , and a standard deviation, denoted as σ_δ . To be clear, σ_δ refers the standard deviation after sampling error has been removed and is denoted as τ in the meta-analysis literature. For a particular study i , the effect size d_i , has two parts: its true effect size, δ_i , and its sampling error, $d_i - \delta_i$. In sum, there is a mean effect size (δ), and there are multiple *true* effect sizes rather than a single

one. We assume that in any particular literature, we have a sample of all methodologically-sound studies, and thus this sample provides estimates of both δ and σ_δ . In Supplemental Material (davidakenny.net/doc/KJ_Sup17.pdf), we discuss the difficulties underlying these assumptions.

In the next two sections we discuss two related consequences of heterogeneity of effect sizes: the statistical power of detecting a significant effect size and the precision of any effect size estimate.

Power Given Heterogeneity

In a conditional power analysis in which the effect size is known and not an estimate (Maxwell, Lau, & Howard, 2015), one computes the power for a given effect from that effect size. This conventional analysis, where the only variability derives from sampling error, treats the effect size as fixed. However, if there is effect size variation over and above sampling error, then the effect size does not take on one value, but rather is a random variable. McShane and Böckenholdt (2014) have suggested a way to estimate power in this case, essentially using multilevel modeling. An alternative approach that we employ is to average the power estimates¹ across all possible values of effect sizes. We presume that the distribution of effect sizes is normal.

To estimate power in the presence of effect heterogeneity, we need to estimate the area under the conjunction of two distributions, within study sampling error and between study variance. To do so we compute power in the conventional manner for each of the different possible true effect sizes, and then weight these power estimates by the likelihood of each δ_j . In this way, we estimate power in the presence of effect

heterogeneity, given δ , σ_δ , α and n . A web-based program that implements this is available at <https://davidakenny.shinyapps.io/SVPower>.

Table 1 presents power estimates for a range of positive values of δ (0.2, 0.5, and 0.8, corresponding to small, medium and large average effects), a range of values of the standard deviation of effect sizes, σ_δ (0.0, 0.1, 0.2, and 0.3), and a range of values of sample sizes, n (100, 200, 500, 1000, and ∞). Alpha is set at .05. The standard deviations of the effect sizes, σ_δ , are in the range of plausible values based on past meta-analyses (Richard et al., 2003). Note that when σ_δ is equal to 0.0, there is no effect size heterogeneity and the results in the table are consistent with a conventional power analysis. For each combination of values two probabilities are given: The one denoted by “+” is the probability of rejecting the null hypothesis when δ_i is positive (i.e., in the same direction as δ), and the one denoted by “-” is the probability of rejecting the null hypothesis when δ_i is negative (i.e., in the opposite direction).

Examining first the results from conventional power analyses ($\sigma_\delta = 0.0$), power of finding a positive effect increases as the effect size and sample size increase. Also there is almost no chance of finding a significant negative effect. If we compare these results to what we find when there is heterogeneity, there are several dramatic differences.

First, whenever conventional power analyses yields a power value less than .50, the estimate that allows for heterogeneity is greater than the estimate based on the absence of heterogeneity. For instance, when $\delta = 0.2$ and $n = 100$, power is .166 with no heterogeneity and rises to .294 when σ_δ is 0.3. The increase in power is due to the fact that half the time the value of δ_i is larger than 0.2, its mean, which results in a boost in power.

Second when conventional power analyses yield values greater than .50, the opposite happens: Power declines once an allowance is made for heterogeneity. For instance, when $\delta = 0.5$ and $n = 200$, power is an impressive .940 with no heterogeneity but sinks to .740 when σ_δ is 0.3. In the presence of study variation, power declines because half of the δ_i values are smaller than the mean of .5, resulting in a net loss of power.

Third, Table 1 gives the probability of finding a significant effect in the negative direction, opposite in sign from the average effect. When there is no study variation, this probability is negligible. However, with increasing heterogeneity, not too surprisingly this probability increases. Perhaps surprisingly, this probability actually increases as the sample size increases. For instance, with an n of 1,000, $\delta = 0.2$, and σ_δ of 0.3, there is nearly a 15 percent chance of finding a significant result in the opposite direction from the average effect size.

Fourth, related to the previous point, there are non-obvious results as n gets large. As expected, when there is no variation in effect size, power goes to a value of one with increasing sample sizes. However, with study variation, we see that power to detect an effect in the same direction as the average effect size goes to a value less than one; how much less than one depends on the σ_δ to δ ratio. As it gets larger, there is greater chance of obtaining a significant negative effect and this leads to a decrease in power for detecting a significant positive effect. For instance, given $\delta = 0.2$ and $\sigma_\delta = 0.3$, power in the same direction as the average effect size never reaches the traditionally desired level of .80.

To summarize, given effect heterogeneity, the power in testing an effect in any particular study is different from what conventional power analyses suggest. Whenever a conventional power analyses yields a power value less than .50, an estimate that allows for heterogeneity is greater; and when a conventional analysis yields a power value greater than .50, the estimate given heterogeneity is less.

Second, given some heterogeneity and a small to moderate average effect size, there is a non-trivial chance of finding a significant effect in the opposite direction from the average effect size reported in the literature. Perhaps, even more surprisingly, the power to detect an effect in the wrong direction (e.g., δ is positive, but the test shows a significant negative effect) can be quite large. This probability increases as n increases.

Precision in Estimating Effect Sizes Given Heterogeneity

The effect size in a study provides an estimate of the true effect size. Its standard error permits an estimation of the confidence interval for that true effect size. Assuming a fixed effect size, the standard error derives solely from the sampling error within a study. For δ , this can be closely approximated² by $2/\sqrt{n}$. If there was study variation, this is the standard error for δ_i , the true effect size for the particular study, and not δ , the mean of all possible effect sizes. In the presence of effect size variation, the proper standard error for

$$\delta \text{ is } \sqrt{4/n + \sigma_\delta^2}.$$

Table 2 presents the 95% confidence interval for δ , given an estimated d from a study with the indicated n , assuming varying degrees of known heterogeneity. The values in this table indicate unsurprisingly that the confidence interval becomes narrower as the study n increases. They also show, again perhaps unsurprisingly, that as effect heterogeneity increases, the confidence interval for the true effect size becomes wider.

This difference can be quite dramatic. Looking at the last row of Table 2, the width of confidence interval with no heterogeneity, a large effect size of 0.8, and a sample size of 1,000 is 0.248, a value much narrower than that with smaller sample sizes, but still relatively wide (Simonsohn, 2014). However, if σ_δ is 0.3, the width of the interval widens by over a factor of four, to 1.202. With large effect sizes and sample sizes, we might have high power with heterogeneity, but we still have quite a bit of uncertainty about the size of the effect.

The confidence intervals in Table 2 assume a single study. Both Maxwell et al. (2015) and ShROUT and Rodgers (2018) have argued that when conducting replication studies it may make sense to conduct multiple such studies to narrow the confidence interval. If multiple studies were run, all estimating the same average effect, then the confidence interval for the average effect size decreases as a function of essentially pooling the observations from all studies into a single standard error, the approximate formula being $\sqrt{(4/n + \sigma_\delta^2)/k}$ where k is the number of studies. In Table 3, we present the confidence intervals for δ if five studies were run, all examining an effect in the same domain but with heterogeneity in effect sizes as indicated by the value of σ_δ .

To see the precision benefits of running five studies, as opposed to one, let us first compare the confidence intervals for the first columns in Tables 2 and 3, where there is no heterogeneity of effect sizes, i.e., $\sigma_\delta = 0.0$. If one runs a single study, with an n of 100 there is considerably less precision than if one runs five such studies, each with an n of 100. In fact, the confidence interval is exactly the same with one study having an n of 500 as for five studies each with an n of 100.

Importantly, however, if there is effect size heterogeneity, then there are substantial precision benefits that accrue from multiple smaller studies compared to a single large study. Compare again the rows in Table 2 where n equals 500 with the rows in Table 3 where the n equals 100 in each study, for a combined n across the studies of 500. If there is effect size heterogeneity, the confidence interval for δ is substantially narrower with five studies, each with an n of 100, than for a single study with an n of 500. Parallel conclusions are found when comparing $n = 1,000$ in Table 2 to $n = 200$ in Table 3.

Many analysts recommend what might be called a *one-basket strategy*. They put all their eggs in the one basket of a very large n study. It is also now common for psychologists to dismiss a study as having too small a sample size and pay attention to only large n studies. As Tables 2 and 3 make clear, given that effect sizes vary across studies, such a strategy is misguided.

Knowing the Magnitude and Distribution of Heterogeneity

We have just seen that the degree of heterogeneity in effect sizes has substantial consequences for statistical power and precision. We have so far assumed that the magnitude of heterogeneity is known. However, knowing exactly the magnitude of σ_δ is difficult. We might use estimates from prior research, but simulation studies (e.g., Chung, Rabe-Hesketh, & Choi, 2013) have shown that estimates of heterogeneity are not very accurate, especially when the number of studies is small.

Power analyses always rest on a series of informed guesses. To conduct a conditional power analysis, we start with an informed guess of the effect size. Similarly,

in the presence of heterogeneity of effect sizes, an informed guess for that heterogeneity is also needed.

How might a researcher make an informed guess? One might surmise that research domains with larger average effect sizes have larger effect size variances. To examine this, we correlated the two for the 18 domains of research in Richard et al. (2003) and found little or no correlation ($r = .06$).

Another idea would be to use the average effect size variance values reported by Richard et al. (2003): about .15 for r and 0.30 for d . Perhaps, we might use 0.1 for σ_δ as a value for small heterogeneity, 0.3 for medium, and 0.5 for large. Obviously these are just initial proposals that need further evaluation. (See McShane and Böckenholdt (2014) for alternative suggestions.)

For precision, knowing the value of σ_δ is more problematic as it needs to be integrated with other statistical information (i.e., the amount of sampling error within studies). Even if we have multiple studies and so have a statistical estimate of heterogeneity, that estimate has a great deal of sampling error. One could just guess at the value and treat it as a population value. Alternatively, a Bayesian analysis, as outlined by Higgins, Thompson, and Spiegelhalter (2009) and Maxwell et al. (2015), might be attempted.

There are certainly difficulties of knowing the extent to which there is effect size variance in a given domain. That said, we strongly feel those difficulties are no excuse for just assuming that it is zero. Effect size variation is both widespread and consequential.

We have also assumed that true effect sizes are normally distributed. One consequence of this is that finding a significant effect in the opposite direction from the average effect size is possible. There has been some work on specifying alternatives to the normal distribution for random effects (Lee & Thompson, 2007). One candidate that we have explored is a log-normal one. The advantage of such a distribution is that it has only positive values, and so finding negative effect sizes can only happen because of sampling error. A table of power results using the log normal distribution, paralleling Table 1, is available in Supplemental Material.

The Origins of Effect Size Heterogeneity

Consistent with the considerable evidence for the heterogeneity of effect sizes, there are good reasons why such heterogeneity is to be expected. First and most obviously study characteristics affect the magnitude of effects that are likely to be found even in a well-defined and limited research domain. Many factors vary between studies, such as sample characteristics, particular measurement techniques, investigators, locations, and historical events; all of these may be moderators of the effect studied.

However even when we exhaust all the factors that we can think of that may moderate the effect, there likely remains additional purely random heterogeneity. Consider an analogy with random variance associated with participants in how they respond to some treatment. Participants typically vary for a variety of potentially knowable reasons that might be measured and controlled. But over and above these, there is also simply random variance in people's responses that we are unlikely to ever explain completely. The same holds true, we believe, for effects shown in different

studies searching for a common effect. There is random variation due to study in the effects produced and this is not entirely reducible to a finite set of effect moderators.

One (see, for instance, Simonsohn, 2015) might question the idea of studies as random, as they are surely not randomly sampled from some known population. We would suggest, however, that just as participants in most experiments are not randomly sampled, yet always treated as random, so too should studies be considered as random.

Planning and Replicating Research Given Heterogeneity

We have shown that effect size heterogeneity has important consequences for statistical power and for the precision of effect size estimates. These consequences deserve attention in planning research. We first explore these consequences in planning new research to demonstrate an effect. We then turn to implications for replication research, a topic that is particularly important in the context of recent concerns about replicability (e.g., Open Science Collaboration, 2015).

Research to Demonstrate an Effect

Conventional wisdom suggests that one is generally better off doing a single very large study to demonstrate an effect rather than doing a series of smaller and more modest studies. The results we have shown in the discussion surrounding Tables 2 and 3 lead us to take issue with this conventional wisdom.

We present formulas in online Supplemental Material for the estimation of how many studies and how many participants per study are needed to minimize the standard error of the effect, and so the width of the confidence interval, given available limited resources. What is likely to result is several studies with fewer participants than we normally might think would be enough.

To illustrate, in one example discussed in Supplemental Material, we start with a single study with a modest number of participants. Anticipating an effect size, δ , of 0.4 and 138 participants, the conventional conditional power estimate is .645, which is not very good. Even worse, if we allow for heterogeneity in effect sizes of 0.25, the power is only .584. We are faced with a dilemma. Conventional advice is that one should conduct only high powered studies. But given heterogeneity, any given study, no matter how large its sample size, might be far away from the mean of the effect sizes. Moreover, given heterogeneity, the power of any given study is not as great as might be thought. What then is the alternative? We see it as conducting a series of studies, each of which might be only moderately powered, but the combination of those studies would have decent power.

For instance, let us return to the case in which δ is 0.4 and σ_δ is 0.25, and 7 studies have been conducted, each with a sample size of 138. The power of finding a significant effect in any one study is only .584, making the power of finding all 7 tests significant only .023. To test the null hypothesis that δ is zero, one conducts a random effects meta-analysis of the 7 studies (Maxwell et al., 2015). Denoting n_p is the number of persons per study and k the number of studies and setting $k = 7$, $n_p = 138$, and $\sigma_\delta = 0.25$, the power of a one-sample t -test of mean d is .829. Thus, although the power of any one study is not very impressive, the power of the test of the mean is quite acceptable. Moreover, that power is greater than if one had done fewer studies with the same total number of participants. Finally, across studies one can critically examine heterogeneity and begin to test factors responsible for variation in effect sizes.

However, if there are few studies, say just 3 or 4, a random effects meta-analysis is impractical as there are too few studies to have a reliable estimate of the variance of effect sizes. Our earlier discussion of how to determine the level of heterogeneity applies. However, just pretending that there is no heterogeneity should not be seen as a defensible option.

We have suggested that multiple smaller studies are preferable to a single large one, given effect size heterogeneity. But what exactly does it mean to conduct multiple smaller studies? Clearly, it would not do to conduct one large study, say with an n of 1000 and break it up, acting as if one had done five studies each with an n of 200. Conducting multiple studies must allow for the existing effect size heterogeneity which, as we have already discussed, accrues randomly from a multitude of sources, including experimenters, samples, and so forth. The point is simply that we are better served by a number of studies that permit one to examine the existing variability of effect sizes in a domain.

Research to Replicate an Effect

Recently concerns have been raised about the replicability of effects in psychology (Ioannides, 2005). In one well-publicized examination of replicability (Open Science Collaboration, 2015), 100 published psychology studies were each replicated one time. The results were interpreted to be relatively disturbing, as less than half of the studies were successfully replicated.

What can be learned from a single replication study? Table 1 can help provide an answer. Imagine that the initial study to be replicated yields an estimated effect size of 0.5. In an effort to conduct the replication with sufficient power, we assume δ , the true

mean effect size, is 0.5, and we plan on a sample size of 200. This gives rise to an estimate of .94 power based on a conventional conditional power analysis. If the study fails to replicate, it seems reasonable to question the initial study result.

Let us, however, assume heterogeneity of effect sizes in the effect to be replicated, with σ_{δ} equal to 0.3. In this case, then the actual power is much less than .94, roughly about .73. Thus, over 25 percent of the time the study will fail to replicate. There is even a chance, albeit a very small one, of finding a significant effect in the opposite direction from the original effect.

In fact, the power in the case we have just explored is certainly worse than we have portrayed it, for two reasons. First, we assumed that δ is the same as effect size we estimated in the original study. But that initial effect size has sampling error in it that has not been factored in (Maxwell et al., 2015). Second, over and above the sampling error in the original effect size estimate, due to publication biases, the actual true effect size is likely smaller than the typical reported estimated effect size.

In the presence of heterogeneity, our results show that power is not nearly as high as it would seem and that even large n studies may have a non-trivial chance of finding a result in the opposite direction from the original study. This makes us question the wisdom of placing a great deal of faith in a single replication study. The presence of heterogeneity implies that there are a variety of true effects that could be produced.

Additionally, the presence of heterogeneity makes us question the common practice of seeing whether zero is in the confidence interval of the difference between the effect in the original study and the effect in the replication study.³ Doing so presumes that the only source of variance between the two studies is sampling error. However,

given heterogeneity, the width of the confidence interval would be greater than that based solely on sampling error. For instance, consider two studies each with an n of 200 and estimated effect sizes of 0.60 and 0.15. The 95 percent confidence interval for the difference between these two effect sizes, assuming no heterogeneity, is from 0.053 to 0.847. Because this interval does not include zero, it appears that the two studies are statistically different. However, if we allow for study variance of $\sigma_{\delta} = 0.15$, the confidence interval actually goes from -0.044 to 0.944, which now includes zero. Ignoring study variation leads to too narrow a confidence interval and sometimes the mistaken conclusion that the original and replication study results have produced inconsistent results.

Part of the recent focus on replication is based on the implicit belief that if procedures could be fully standardized, the only difference between study effects would then be sampling error. Such a view is likely mistaken (Maxwell et al., 2015). Even in a well-conducted replication there are still many factors that may lead to effect heterogeneity. For instance, studies are conducted in different locations, with different experimenters, in different historical moments, and with different non-randomly selected participants. All of these, and a variety of other randomly varying factors, likely lead to heterogeneity. And this heterogeneity leads to concerns about the utility of any single replication study.

In their classic paper on “The Law of Small Numbers,” Tversky and Kahneman (1971) described an experimenter who does the same study twice and in the first study, he or she obtains a significant effect whereas in the second study the effect is no longer significant. When asked what they would do if faced with this situation, a plurality of

psychologists said they would “try to find an explanation for the difference between the two groups” (p. 27). Perhaps even more perplexing, we have shown that the second study may even come up with a significant effect in the opposite direction from the first. The second study, does not necessarily “disconfirm” the first; rather it may well lead to the conclusion of considerable random variance in the effect in question.

Conclusion

Effect size heterogeneity is found nearly everywhere in what we study. However, in power analyses, computing confidence intervals, and the planning of research, the field acts as if the results of studies are homogeneous. We have shown that heterogeneity leads to both lower and higher power than expected, a surprising high probability of finding the “wrong” results, and the conclusion that multiple smaller studies are preferable to a single large one. All of this leads to very different ideas about the conduct of research and the quest to establish the true effect in the presence of random variation. Replication research, it seems to us, should not be about confirming or disconfirming earlier results in the literature. Replication researchers should not strive to conduct the definitive large n study in an effort to establish once and for all whether a given effect exists. The goal of replication research should instead be to establish both typical effects in a domain and the range of possible effects, given all of what Campbell called the “heterogeneity of irrelevancies” (Cook, 1990) that affect studies and their results. Many smaller studies that vary those irrelevancies likely better serve us than one single large study.

Most researchers tend to believe that in any given domain, when evaluating any given effect, there really is only one effect and one should strive to uncover it in studies

that are undertaken. It can be disconcerting, at best, to believe that there really are a variety of effects that exist and that might be found. However, that is what it means to have study variation in effect sizes and, as we emphasized early on, that is what we typically find in meta-analyses. As a field, we need to begin to understand what it means for effects to vary.

We have raised the issue of heterogeneity and explored some of its implications, while nevertheless highlighting some difficult issues that require further attention. These include the nature of the underlying distribution of effect sizes and how to estimate their variability. Another important one, that we only mention here and in Supplemental Material, is the presence of methodological problems in studies (e.g., biased effect size estimates, and p -hacking) and their consequences for the estimation of the mean and variance of effect size estimates. Finally, we have limited our discussion of effects sizes to d ; a full treatment of the topic would require extending the discussion to other effect size measures, e.g., correlations and odd ratios.

These issues notwithstanding, we firmly believe that we need to accept and, in fact, embrace heterogeneity. If there truly exist multiple effect sizes in a given domain, then power analyses and confidence intervals need to allow for that. Moreover, research should also examine that variability, and the factors that can partly explain it, rather than focusing solely on whether an effect exists or does not.

References

- Aird, F., Candela, I., & Mantis, C. (2017) Replication study: BET bromodomain inhibition as a therapeutic Strategy to target c-Myc. *eLife*, *6*, e21253.
- Biesanz, J. C., & Schrage, S. M. (2010). *Sample size planning with effect size estimates*. Unpublished paper, University of British Columbia.
- 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.
- Cook, T. D. (1990). The generalization of causal connections: Multiple theories in search of clear practice. In L. Sechrest, E. Perrin, & J. Bunker (Eds.), *Research methodology: Strengthening causal interpretations of nonexperimental data* (DHHS Publication No. PHS 90-3454, pp. 9–31). Rockville MD: Department of Health and Human Services.
- Dallow, N., & Fina, P. (2011). The perils with the misuse of predictive power. *Pharmaceutical Statistics*, *10*, 311–317.
- Decoster, J., & Claypool, H. M. (2004). A meta-analysis of priming effects on impression formation supporting a general model of informational biases. *Personality and Social Psychology Review*, *8*, 2-27.
- Gershoff, E. T., & Grogan-Kaylor, A. (2016) Spanking and child outcomes: Old controversies and new meta-analyses. *Journal of Family Psychology*, *30*, 453-469.
- Gillett, R. (1994). An average power criterion for sample size estimation. *The Statistician*, *43*, 389-394.

Gillett, R. (2002). The unseen power loss: Stemming the flow. *Educational and Psychological Measurement*, 62, 960-968.

Haelermans, C., & Borghans, L. (2012). Wage effects of on-the-job training: A meta-analysis. *British Journal of Industrial Relations, London School of Economics*, 50, 502-528.

Hagger, M. S., & Chatzisarantis, N. L. D. (2016). A multilab preregistered replication of the ego-depletion effect. *Perspectives on Psychological Sciences*, 11, 546-573.

Hedges, L. V., & Olkin, I. (1985) *Statistical methods for meta-analysis*. London: Academic Press.

Hedges, L. V., & Vevea, J. L. (1998). Fixed-and random-effects models in meta-analysis. *Psychological Methods*, 3(4), 486-504.

Higgins, J. P. T. (2008) Commentary: Heterogeneity in meta-analysis should be expected and appropriately quantified. *International Journal of Epidemiology*, 37, 1158-1160.

Higgins, J. P. T., Thompson, S. G., & Spiegelhalter, D. J. (2009). A re-evaluation of random-effects meta-analysis. *Journal of the Royal Statistical Society, Series A*, 172, Part 1,137–159.

Ioannidis, J. P. (2005). Why most published research findings are false. *PLoS med*, 2(8), e124.

Lee, K. J., & Thompson, S.G. (2007). Flexible parametric models for random-effects distributions. *Statistics in Medicine*, 27, 418–434.

Maxwell, S. E., Lau, M. Y., & Howard, G. S. (2015). Is psychology suffering from a replication crisis? What does “failure to replicate” really mean? *American Psychologist, 79*, 487-498.

McShane, B. B., & Böckenholdt, U. (2014). You cannot step into the same river twice: When power analyses are optimistic. *Perspectives on Psychological Sciences, 9*, 612-625.

Open Science Collaboration (2015) Estimating the reproducibility of psychological science. *Science, 349* (6251), aac4716.

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.

Shrout, P., & Rodgers, J. (2018). Research and statistical practices that promote accumulation of scientific findings. *Annual Review of Psychology*, in press.

Simonsohn, U. (May 1, 2014). Data Colada: [20] We cannot afford to study effect size in the lab, datacolada.org/20.

Simonsohn, U. (February 9, 2015). Data Colada: [33] “The” effect size does not exist, datacolada.org/33.

Turner, R. M., Davey, J., Clarke, M. J., Thompson, S. G., & Higgins, J. P. T. (2012). Predicting the extent of heterogeneity in meta-analysis, using empirical data from the Cochrane Database of Systematic Reviews. *International Journal of Epidemiology, 41*, 818-27.

Tversky, A., & Kahneman, D. (1971). Belief in the law of small numbers. *Psychological Bulletin, 76*, 105-110.

Table 1

Power for a positive “+” and a negative effect “-” given an effect size (δ), study variation (σ_δ), total sample size (n) and an alpha of .05

| δ | n | σ_δ | | | | | | | |
|----------|----------|-----------------|------|-------|------|-------|------|------|------|
| | | 0.0 | | 0.1 | | 0.2 | | 0.3 | |
| | | + | - | + | - | + | - | + | - |
| 0.2 | 100 | .166 | .002 | .192 | .004 | .245 | .018 | .294 | .050 |
| | 200 | .290 | .000 | .326 | .003 | .374 | .026 | .406 | .075 |
| | 500 | .607 | .000 | .572 | .003 | .544 | .043 | .531 | .115 |
| | 1,000 | .885 | .000 | .739 | .003 | .641 | .061 | .598 | .145 |
| | ∞ | 1.000 | .000 | .977 | .023 | .841 | .159 | .748 | .252 |
| 0.5 | 100 | .697 | .000 | .678 | .000 | .643 | .001 | .613 | .007 |
| | 200 | .940 | .000 | .899 | .000 | .817 | .001 | .748 | .010 |
| | 500 | 1.000 | .000 | .992 | .000 | .931 | .001 | .850 | .015 |
| | 1,000 | 1.000 | .000 | .999 | .000 | .963 | .001 | .890 | .021 |
| | ∞ | 1.000 | .000 | 1.000 | .000 | .994 | .006 | .952 | .048 |
| 0.8 | 100 | .977 | .000 | .964 | .000 | .922 | .000 | .868 | .000 |
| | 200 | 1.000 | .000 | .999 | .000 | .983 | .000 | .942 | .001 |
| | 500 | 1.000 | .000 | 1.000 | .000 | .998 | .000 | .977 | .001 |
| | 1,000 | 1.000 | .000 | 1.000 | .000 | .999 | .000 | .986 | .001 |
| | ∞ | 1.000 | .000 | 1.000 | .000 | 1.000 | .000 | .996 | .004 |

^aProbability of detecting a positive effect, i.e., one consistent with the average effect size.

^bProbability of detecting a negative effect, i.e., one inconsistent with the average effect size.

Table 2

95 Percent Confidence Interval for the Effect with Lower (L) and Upper (U) Limits from a Single Study for Different Sample Sizes (n), Effect Sizes (d), and Level of Heterogeneity (σ_δ)

| d | n | σ_δ | | | | | | | |
|-----|-------|-----------------|-------|--------|-------|--------|-------|--------|-------|
| | | 0.0 | | 0.1 | | 0.2 | | 0.3 | |
| | | L | U | L | U | L | U | L | U |
| 0.2 | 100 | -0.192 | 0.592 | -0.238 | 0.638 | -0.354 | 0.754 | -0.507 | 0.907 |
| | 200 | -0.077 | 0.477 | -0.139 | 0.539 | -0.280 | 0.680 | -0.450 | 0.850 |
| | 500 | 0.025 | 0.375 | -0.063 | 0.463 | -0.229 | 0.629 | -0.414 | 0.814 |
| | 1,000 | 0.076 | 0.324 | -0.032 | 0.432 | -0.211 | 0.611 | -0.401 | 0.801 |
| 0.5 | 100 | 0.108 | 0.892 | 0.062 | 0.938 | -0.054 | 1.054 | -0.207 | 1.207 |
| | 200 | 0.223 | 0.777 | 0.161 | 0.839 | 0.020 | 0.980 | -0.150 | 1.150 |
| | 500 | 0.325 | 0.675 | 0.237 | 0.763 | 0.071 | 0.929 | -0.114 | 1.114 |
| | 1,000 | 0.376 | 0.624 | 0.268 | 0.732 | 0.089 | 0.911 | -0.101 | 1.101 |
| 0.8 | 100 | 0.408 | 1.192 | 0.362 | 1.238 | 0.246 | 1.354 | 0.093 | 1.507 |
| | 200 | 0.523 | 1.077 | 0.461 | 1.139 | 0.320 | 1.280 | 0.150 | 1.450 |
| | 500 | 0.625 | 0.975 | 0.537 | 1.063 | 0.371 | 1.229 | 0.186 | 1.414 |
| | 1,000 | 0.676 | 0.924 | 0.568 | 1.032 | 0.389 | 1.211 | 0.199 | 1.401 |

Table 3

95 Percent Confidence Interval for the Mean Effect with Lower (L) and Upper (U) Limits from Five Studies for Different Sample Sizes (n), Effect Sizes (d), and Level of Heterogeneity (σ_δ)

| d | n | σ_δ | | | | | | | |
|-----|-------|-----------------|-------|-------|-------|--------|-------|--------|-------|
| | | 0.0 | | 0.1 | | 0.2 | | 0.3 | |
| | | L | U | L | U | L | U | L | U |
| 0.2 | 100 | 0.025 | 0.375 | 0.004 | 0.396 | -0.048 | 0.448 | -0.116 | 0.516 |
| | 200 | 0.076 | 0.324 | 0.048 | 0.352 | -0.015 | 0.415 | -0.091 | 0.491 |
| | 500 | 0.122 | 0.278 | 0.082 | 0.318 | 0.008 | 0.392 | -0.074 | 0.474 |
| | 1,000 | 0.145 | 0.255 | 0.096 | 0.304 | 0.016 | 0.384 | -0.069 | 0.469 |
| 0.5 | 100 | 0.325 | 0.675 | 0.304 | 0.696 | 0.252 | 0.748 | 0.184 | 0.816 |
| | 200 | 0.376 | 0.624 | 0.348 | 0.652 | 0.285 | 0.715 | 0.209 | 0.791 |
| | 500 | 0.422 | 0.578 | 0.382 | 0.618 | 0.308 | 0.692 | 0.226 | 0.774 |
| | 1,000 | 0.445 | 0.555 | 0.396 | 0.604 | 0.316 | 0.684 | 0.231 | 0.769 |
| 0.8 | 100 | 0.625 | 0.975 | 0.604 | 0.996 | 0.552 | 1.048 | 0.484 | 1.116 |
| | 200 | 0.676 | 0.924 | 0.648 | 0.952 | 0.585 | 1.015 | 0.509 | 1.091 |
| | 500 | 0.722 | 0.878 | 0.682 | 0.918 | 0.608 | 0.992 | 0.526 | 1.074 |
| | 1,000 | 0.745 | 0.855 | 0.696 | 0.904 | 0.616 | 0.984 | 0.531 | 1.069 |

Footnotes

¹Gillett (1994; 2002) has also suggested computing power across a distribution of effect sizes. However, the distribution in this work is a prior distribution of effect sizes and not the degree of heterogeneity of the given effect. Additionally, Biesanz and Schragger (2010) and Dallow and Fina (2011) have suggested in replication studies using a distribution of effect sizes based on its confidence interval.

²Hedges and Olkin (p. 86, 1985) give the standard deviation for d as $\sqrt{\frac{4}{n} + \frac{\delta^2}{2n}}$, where δ is the population value of d for the study and n is the study sample size with the assumption that $n_1 = n_2$. This is the formula that we use throughout this paper, but in cases where δ is unknown, as it is here, or varies, we drop the second term.

³Sometimes researchers mistakenly check to see if the original effect size is in the confidence interval of the replication effect size. Such a practice is flawed because it ignores that the original effect has sampling error (Maxwell et al., 2015).