

Why Are Replication Rates So Low?

Patrick Vu[†]

Abstract

Many explanations have been offered for why replication rates are low in the social sciences, including selective publication, p -hacking, and treatment effect heterogeneity. This article emphasizes that issues with the most commonly used approach for setting sample sizes in replication studies may also play an important role. Theoretically, I show in a simple model of the publication process that we should expect the replication rate to fall below its nominal target, even when original studies are unbiased. The main mechanism is that the most commonly used approach for setting the replication sample size does not properly account for the fact that original effect sizes are estimated. Specifically, it sets the replication sample size to achieve a nominal power target under the assumption that estimated effect sizes correspond to fixed true effects. However, since there are non-linearities in the replication power function linking original effect sizes to power, ignoring the fact that effect sizes are estimated leads to systematically lower replication rates than intended. Empirically, I find that a parsimonious model accounting only for these issues can fully explain observed replication rates in experimental economics and social science, and two-thirds of the replication gap in psychology. I conclude with practical recommendations for replicators. (*JEL* C18, C53, C90)

Keywords: Replications, statistical power, experiments.

1. Introduction

In a 2016 survey conducted by *Nature*, 90% of researchers across various fields agreed that the scientific community faces a ‘reproducibility crisis’ (Baker, 2016). Growing consensus has been supported by high-profile replication projects which find that the replication rate – i.e. the fraction of replications that are significant with the same sign as the original finding – is just

[†]*This version:* October 16, 2024. University of New South Wales. Email: patrick.vu@unsw.edu.au.

36% in psychology, 61% in experimental economics, and 62% in experimental social science (Open Science Collaboration, 2015; Camerer et al., 2016, 2018).

Understanding the underlying cause of low replication rates is important for researchers and reformers aiming to improve the credibility of published research. There is a large and growing literature examining a wide range of explanations, including selective publication against null results (Franco et al., 2014; Open Science Collaboration, 2015; Camerer et al., 2016, 2018); p -hacking and other questionable research practices (Ioannidis, 2005, 2008; Simonsohn et al., 2014; Brodeur et al., 2016, 2020, 2022; Elliott et al., 2022); and heterogeneity across original studies and replications in research design and experimental subjects (Higgins and Thompson, 2002; Cesario, 2014; Simons, 2014; Stanley et al., 2018; Bryan et al., 2019).

This article focuses on an alternative issue that has received less attention. Its main theoretical result shows that issues with the most commonly used approach of setting sample sizes in replication studies implies that the replication rate should be expected to fall short of its intended target. Importantly, this is true even if original studies are unbiased and free from common concerns over selective publication, p -hacking, treatment effect heterogeneity, and low statistical power in original studies. To see why this is the case, note that the sample size in replication studies is commonly set as a (decreasing) function of the observed effect size in the original study (Maxwell et al., 2015; Open Science Collaboration, 2015; Camerer et al., 2016). The actual approach, which I refer to as the common power rule, is designed to deliver a prespecified replication probability target (e.g. 0.9) if the observed effect is exactly equal to the unobserved true effect. In practice, replication rates consistently fall below the intended power target, which is commonly interpreted as an indicator that original effects are biased due to factors such as selective publication, p -hacking, or treatment effect heterogeneity (Open Science Collaboration, 2015; Camerer et al., 2016). However, this article highlights that replication power is in fact a non-linear, locally concave function of the original estimate. Thus, even if original estimates were unbiased, Jensen’s inequality implies that the expected replication rate must fall below the replication probability evaluated at the expectation, which, with unbiased original estimates, equals the nominal target. To further develop the intuition behind this central idea, I present a simple illustrative example in Section 2, before showing that it applies to more general settings in Section 3.

Practically, the main result means that stated replication rate targets in large-scale replication studies using the common power rule do not in fact set an attainable benchmark against which to judge observed replication rates; even if original studies were unbiased, such targets are not reachable in expectation. I also show that the gap between the expected replication rate and its intended power target is larger when the original published studies have low power, a problem that we might expect to be severe in practice given evidence of low power in various em-

pirical literatures (Button et al., 2013; Ioannidis et al., 2017; Stanley et al., 2018; Arel-Bundock et al., 2023).

The main theoretical result applies to studies using what I refer to as the common power rule, which sets replication power to detect the original estimated effect size. More recently, some studies have begun to use a higher-power variant which I refer to as the fractional power rule, wherein replication power is set to detect some fraction of the estimated effect size (Camerer et al., 2018). Building on results in Andrews and Kasy (2019), I show that the expected replication rate using the fractional power rule can be either above or below the stated power target.

To what extent can these theoretical insights explain the low replication rates actually observed in large-scale replication studies? Although the theory predicts that the actual replication rate will always fall below the target when using the common power rule, the magnitude of this gap is an empirical question. Likewise, for replication studies using the fractional power rule, both the sign and the magnitude of the gap is an empirical question.

To evaluate the importance of replication power issues in practice, I empirically investigate the results of three replication studies, two of which use the common power rule (Open Science Collaboration, 2015; Camerer et al., 2016) and one of which uses the fractional power rule (Camerer et al., 2018). In each application, I estimate the empirical model in Andrews and Kasy (2019) using a ‘metastudy approach’ that corrects for publication bias to obtain the underlying distribution of latent studies prior to screening by the publication process. I then use the estimated latent distribution of studies to simulate what we should expect the replication rate to be based on the power calculations actually implemented in replications. Importantly, the model and its predictions are based only on data from original studies and assume away researcher manipulation and heterogeneous treatment effects. The empirical exercise asks, in effect, whether observed replication rates could have been predicted by issues with replication power alone, before the replication studies themselves were actually undertaken and in a parsimonious model without treatment effect heterogeneity or p -hacking.

I find that the predicted replication rate is almost identical to observed replication rates in experimental economics (60% vs. 61%) and experimental social science (54% vs. 57%). Replications in experimental economics implemented the common power rule, while those in experimental social science used a fractional power rule.¹ These empirical results are consistent with the null hypothesis that observed replication rates in these studies are driven entirely by issues with power calculations, rather than other issues such as p -hacking or treatment effect

¹In the experimental social science replications (Camerer et al., 2018), replicators used a fractional power rule in the first stage of replications predicted here, where replication power was set to detect 75% of the original effect size with 90% intended power.

heterogeneity. Of course, failure to reject a hypothesis does not mean that it is true, and thus we should not necessarily conclude that these other factors are not present in these settings. Nevertheless, other evidence has also suggested a relatively limited role for p -hacking in the context of lab experiments studied here (Brodeur et al., 2016, 2020; Imai et al., 2020).

In psychology, the predicted replication rate is 55%, whereas the observed replication rate is 35%. Since the intended power target was 92%, issues with power calculations explain only two-thirds of the gap in psychology. In the case of psychology, we can therefore reject the null that the replication gap is entirely explained by issues with power calculations. This provides strong evidence that some other factors are important in psychology. Some possibilities discussed in the literature include heterogeneous treatment effects, p -hacking, and differences in replicability across subfields.

In an extension, I examine the relative effect size (defined as the mean of the ratio of the replication effect size and the original effect size), a common complementary continuous measure of replication. I generate relative effect size predictions in each field using a similar method as for the replication rate. I once again find that the predictions are quite similar to observed outcomes in economics (0.70 vs. 0.66). The model is somewhat farther off for social sciences (0.53 vs. 0.44), perhaps suggesting some role for other factors, although the difference is not statistically distinguishable from zero. In psychology, predictions are quite far off (0.64 vs. 0.37), again providing strong evidence for alternative factors.

When analyzing relative effect sizes, it is important to note that the common practice of only choosing statistically significant results to replicate will mechanically induce upward bias in original estimates, irrespective of whether or not the published literature is biased. Since replication estimates must regress to the mean, relative effect sizes should be expected to be below the nominal target of one. This means that the nominal target is actually unattainable in expectation, much like for the replication rate. Note that selecting only significant results to replicate might also be expected to lead to lower replication rates. However, in Subsection 3.3, I show that the impact of selection on the replication rate is actually theoretically ambiguous, as it also tends to select studies with larger true effects that have higher replication probabilities. In Subsection 4.4, I find that the overall impact of selection on the replication rate is empirically small.

The results in this article highlight that the common power rule and selection on significance make both the binary measure of replication and the relative effect size measure difficult to interpret. In light of these limitations, I make several practical recommendations to replicators in the final section. First, I recommend focusing on relative effect sizes over the binary measure of replication, and caution against the practice of selecting only significant results to replicate. This removes the regression-to-the-mean issue and restores the nominal target of one as a

meaningful benchmark. I then suggest formally testing whether relative effect sizes deviate from one using the prediction interval approach in Patil et al. (2016). A key advantage of this approach is that it is more informative about the level of selective publication than the binary measure of replication.²

This article contributes to the large metascience literature and the growing literature on predicting research outcomes (Ioannidis, 2005; Franco et al., 2014; Gelman and Carlin, 2014; Dreber et al., 2015; Maxwell et al., 2015; Anderson and Maxwell, 2017; Stanley et al., 2018; Miguel and Christensen, 2018; Altmejd et al., 2019; Amrhein et al., 2019; DellaVigna et al., 2020; Gordon et al., 2020; Frankel and Kasy, 2022; DellaVigna and Linos, 2022; Nosek et al., 2022). Andrews and Kasy (2019) and Kasy (2021) provide stylized examples showing that the replication rate can vary widely depending on the latent distribution of studies (i.e. the joint distribution of true effects and standard errors for published and unpublished studies). Theoretically, this article builds on this observation by establishing that the expected replication rate is bounded above by its nominal target owing to issues with common power calculations in replication studies. This result holds for any distribution of latent studies. Empirically, I provide evidence that among the profusion of explanations for low replication rates, a parsimonious model accounting only for issues with replication power calculations and low power in original studies can adequately account for observed replication rates in experimental economics and social science.

The remainder of the article is organized as follows. Section 2 introduces the key intuition behind the main theoretical result in a minimalistic setting. Section 3 extends these ideas to more general settings. Section 4 presents the empirical applications. Finally, Section 5 offers some practical recommendations to replicators.

2. Illustrative Example

In this section, I present a simple example to illustrate the main ideas. Consider a study examining the impact of framing on consumer choice. Suppose that the type of framing being studied boosts purchases of a particular product of interest, and denote this unobserved true effect by $\theta > 0$. Researchers hoping to learn about this effect conduct an experiment and publish their findings in a leading journal. Their main results show that the estimated effect of framing on purchases is equal to $X > 0$ and is statistically significant at the 5% level.

Several years after publication, a different team of researchers becomes interested in testing

²The replication rate is not very informative about selective publication or the ‘file-drawer’ problem (Andrews and Kasy, 2019; Kasy, 2021). This is for the simple reason that the replication rate only includes significant results in its definition, and hence tells us very little, or nothing, about the degree to which null results are censored. Subsection 3.3 expands on this point.

the reliability of this finding. To do this, they conduct a replication using identical procedures to the original study, but using a new set of experimental participants. Of primary interest is whether the original result ‘replicates’, in the sense that the estimated effect in the replication has the same sign as the original study and is statistically significant. Assuming the replication estimate is approximately normally distributed, the probability of this outcome – or replication power – is equal to $1 - \Phi\left(1.96 - \frac{\theta}{\sigma_r}\right)$, where the replication standard error σ_r is determined by the choice of the replication sample size and $\Phi(\cdot)$ denotes the normal cdf.³

A central problem faced by replicators is how to choose the replication standard error σ_r – or equivalently the replication sample size – to deliver a prespecified intended power target of, say, 0.9. If θ were known, then setting $\sigma_r = \theta / (1.96 - \Phi^{-1}(1 - 0.9))$ would yield replication power equal to target power of 0.9. However, in reality, θ is unobserved. The most common solution has been to replace the unobserved true effect θ with the observed original estimate X , leading to what I refer to in this paper as the common power rule: $\sigma_r(X) = |X| / (1.96 - \Phi^{-1}(1 - 0.9))$. Under this rule for setting the replication standard error, replication power to detect a positive effect is given by

$$RP(X) = 1 - \Phi\left(1.96 - \frac{\theta}{|X|} (1.96 - \Phi^{-1}(1 - 0.9))\right) \quad (1)$$

Clearly, if the original estimate X coincides exactly with the true effect θ , then $RP(X) = 0.9$, as intended. However, this reasoning fails to account for the fact that the original estimate is a random variable and that $RP(\cdot)$ is a non-linear function. This article shows that the replication function $RP(\cdot)$ is in fact a non-linear, locally concave function. Thus, even if original estimates were unbiased, by Jensen’s inequality we have that $\mathbb{E}[RP(X)] < RP(\mathbb{E}[X]) = RP(\theta) = 0.9$. In other words, the expected replication rate under the common power rule is below its nominal target, even in the ‘ideal’ scenario where original studies are unbiased and there is no publication bias, p -hacking, or treatment effect heterogeneity.

How much might this matter in practice? Section 4 provide a fuller answer to this question by analyzing three large-scale replication studies. However, for illustrative purposes, suppose that original estimates are drawn from an $N(\frac{1}{2}, 1)$ distribution, and hence unbiased and free from concerns over distortions from selective publication and p -hacking. If this were the case, then the expected replication rate based on the equation above would be $\mathbb{E}[RP(X)] = 0.59$, which, despite the unbiasedness of original estimates, is far below its nominal target of 0.9.⁴

³Approximate normality holds under mild conditions and is widely assumed in practice.

⁴This calculation assumes that all original estimates are replicated, irrespective of whether or not they are statistically significant. In practice, replicators typically choose only significant results to replicate. Performing the same calculation but conditioning on statistically significant original results with a positive gives a much lower replication rate: $\mathbb{E}[RP(X)|1.96 \leq X] = 0.10$; note also that this number is practically unchanged if significant results are replicated irrespective of the sign of the original estimate. Subsection 3.3 and Appendix

3. Theory

3.1. Model of Large-Scale Replication Studies

This section shows that the idea illustrated in the previous section applies to much more general settings. Appendices in the supplementary materials contain proofs for all theoretical results. For the general framework, I consider the model in [Andrews and Kasy \(2019\)](#). Suppose a large-scale replication study is conducted in an empirical literature of interest and we observe the estimated effect sizes and standard errors for original studies and their replications. Let upper case letters denote random variables, lower case letters realizations. Latent studies (published or unpublished) have a superscript $*$ and published studies have no superscript. The model of the DGP has five steps:

1. **Draw a population parameter and standard error:** Draw a research question with population parameter (Θ^*) and standard error (Σ^*) :

$$(\Theta^*, \Sigma^*) \sim \mu_{\Theta, \Sigma}$$

where $\mu_{\Theta, \Sigma}$ is the joint distribution of latent true effects and latent standard errors.

2. **Estimate the effect:** Draw an estimated effect from a normal distribution with parameters from step 1:

$$X^* | \Theta^*, \Sigma^* \sim N(\Theta^*, \Sigma^{*2})$$

3. **Publication selection:** Selective publication is modeled by the function $p(\cdot)$, which returns the probability of publication for any given t -ratio. Let D be a Bernoulli random variable equal to 1 if the study is published and 0 otherwise:

$$\mathbb{P}(D = 1 | X^*/\Sigma^*) = p\left(\frac{X^*}{\Sigma^*}\right) \quad (2)$$

4. **Replication selection:** Replications are sampled from published studies (X, Σ, Θ) (i.e. latent studies $(X^*, \Sigma^*, \Theta^*)$ conditional on publication $(D = 1)$). Replication selection is modeled by the function $r(\cdot)$, which returns the probability of being chosen for replication for any given t -ratio. Let R be a Bernoulli random variable equal to 1 if the study is chosen for replication and 0 otherwise:

$$\mathbb{P}(R = 1 | X/\Sigma) = r\left(\frac{X}{\Sigma}\right) \quad (3)$$

C explore the interaction between selection and non-linearities in more detail.

5. **Replication:** A replication draw is made with:

$$X_r | \Theta, X, \Sigma, \Sigma_r, D = 1, R = 1 \sim N\left(\Theta, \Sigma_r^2\right)$$

We observe i.i.d draws of $(X, \Sigma, X_r, \Sigma_r)$ from the conditional distribution of $(X^*, \Sigma^*, X_r, \Sigma_r)$ given $D = 1$ and $R = 1$. I consider what happens in this model when the replication standard error, Σ_r , is set to detect the original estimate X with a prespecified power level $1 - \beta$, where β is the target probability of Type II error. This approach is implemented, for example, in [Open Science Collaboration \(2015\)](#) and [Camerer et al. \(2016\)](#), and a survey of replications the psychology literature by [Anderson and Maxwell \(2017\)](#) shows that it is the most commonly implemented approach. I refer to this as the common power rule, which is formalized as follows:

Definition 1 (Common Power Rule). *The common power rule to detect original effect size x with intended power $1 - \beta$ sets the replication standard error to*

$$\sigma_r(x, \beta) = \frac{|x|}{1.96 - \Phi^{-1}(\beta)} \quad (4)$$

This is equivalent to setting the replication sample size to $N \times \left[\frac{\sigma}{|x|}(1.96 - \Phi^{-1}(\beta))\right]^2$, where N and σ are the original study's sample size and standard deviation, respectively.

The justification for the common power rule is that the power in any given replication study will equal its intended power target of $1 - \beta$ when the original estimate coincides exactly with the true effect, $x = \theta$ (Lemma B1). In practice, replication rates consistently fall below this benchmark, which is typically taken as evidence that original estimates are biased because of selective publication or p -hacking. While this argument has intuitive appeal, it does not account for the fact that replication power is a non-linear function of the random original estimate X ; thus, even if $\mathbb{E}[X|\Theta = \theta] = \theta$, the replication probability evaluated at the expectation (which equals the intended target) will not, in general, be equal to the expected replication rate.

This argument is developed more formally in the following subsections. We make two assumptions in the following analysis. First, we impose that the publication probability $p(\cdot)$ is weakly increasing in the absolute t -ratio over the statistically significant region of the support, and symmetric around zero. Intuitively, this means that studies that are ‘more significant’ have a (weakly) higher probability of being selected for publication.

Assumption 1 (Publication Selection Function). *Let $p(t)$ be weakly increasing for all $t \geq 1.96$, and $p(t) = p(-t)$ for all $t \geq 1.96$. Allow $p(\cdot)$ to take any form when $t \in (-1.96, 1.96)$. Finally, let $p(t) \neq 0$ for all $|t| \geq 1.96$.*

This allows for very general forms of publication bias (or lack thereof). Assuming symmetry in the publication selection function may be appropriate in some settings but not others. For example, symmetric publication bias may not be appropriate in the literature examining the impact of the minimum wage on employment, since there are priors about the sign of the effect. However, for large-scale replication studies, which are the primary focus of this paper, the assumption is more plausible. This is because large-scale replication typically include a range of studies examining different outcomes and hence, the relative signs of effects across studies are arbitrary.

Next, we make the same assumptions about the replication selection mechanism, $r(\cdot)$:

Assumption 2 (Replication Selection Function). *Let $r(t)$ be weakly increasing for all $t \geq 1.96$, and $r(t) = r(-t)$ for all $t \geq 1.96$. Allow $r(\cdot)$ to take any form when $t \in (-1.96, 1.96)$. Finally, let $r(t) \neq 0$ for all $|t| \geq 1.96$.*

This includes, for example, random sampling from published, significant studies. This assumption is plausible in all three empirical applications examined in this paper, which I discuss in further detail in Subsection 4.1.

Finally, note that the article uses three distinct concepts of statistical power. First, *power in an original study* is defined as the probability of obtaining a statistically significant estimate in the same direction as the true effect: $1 - \Phi(1.96 - \frac{\theta}{\sigma})$ when $\theta > 0$; and $\Phi(-1.96 - \frac{\theta}{\sigma})$ when $\theta \leq 0$.⁵ Second, *power in a replication study* (or the ‘replication probability’; Definition 2 below) is defined as the probability of obtaining a significant effect with the same sign as the original study, and will depend on the rule for setting replication power. Finally, the *intended power target* of any given rule for setting replication power is denoted by $1 - \beta$, where β is the target probability of Type II error.

3.2. Common Power Calculations and Low Replication Rates

This subsection defines the replication rate and then discusses the main result. First, we define the replication probability of a single study and then use this to define the expected replication rate over multiple studies.

Definition 2 (Replication Probability of Individual Study). *The replication probability of a published study (X, Σ, Θ) chosen for replication ($R = 1$) is*

$$RP(X, \Theta, \sigma_r(X, \beta)) = \mathbb{P}\left(\frac{|X_r|}{\sigma_r(X, \beta)} \geq 1.96, \text{sign}(X_r) = \text{sign}(X) \mid X, \Theta, \beta, R = 1\right) \quad (5)$$

⁵The arguments made throughout are essentially unchanged if we consider the alternative definition of obtaining a statistically significant estimate irrespective of the sign.

This definition captures the dual requirement that the replication estimate is statistically significant and has the same sign as the original study.

Definition 3 (Expected Replication Rate). *The expected replication probability is defined over published studies (X, Σ, Θ) which are chosen for replication ($R = 1$) and statistically significant ($S_X = 1$). It is equal to*

$$\mathbb{E}[RP(X, \Theta, \sigma_r(X, \beta)) | R = 1, S_X = 1] \quad (6)$$

Substituting the common power rule in Definition 1 for the replication standard error gives the expected replication rate under the common power rule. Note that while insignificant results may be replicated, they are not included in the replication rate in Definition 3. This matches the main definition reported in most large-scale replication studies (Klein et al., 2014; Open Science Collaboration, 2015; Camerer et al., 2016, 2018; Klein et al., 2018).⁶ With this, we can state the main theoretical result:

Proposition 1 (Common Power Rule Implies the Expected Replication Rate is Always Below Target). *Consider the model in 3.1. Under Assumptions 1 and 2, if replication standard errors are set by the common power rule to detect original estimates with intended power $1 - \beta \geq 0.8314$, then*

$$\mathbb{E}[RP(X, \Theta, \sigma_r(X, \beta)) | R = 1, S_X = 1] < 1 - \beta \quad (7)$$

From a practical perspective, Proposition 1 means that replicators who set the replication sample size to detect original effect sizes should not expect the replication rate to reach its intended target, regardless of whether or not there is selective publication, and even under ‘ideal’ conditions with no researcher manipulation, replications with identical designs and comparable samples (i.e. no heterogeneity in true effects), no measurement error, random sampling in replication selection, and high-powered original studies. That the intended target is not in fact attainable in expectation underscores fundamental difficulties in interpreting replication rate gaps observed in large-scale replication studies.

Figure 1 provides the key intuition underlying this result. It plots the replication probability of a single study in Definition 2 as a function of the original effect X , for a fixed true effect θ and assuming that the common power rule is applied with an intended power target of $1 - \beta = 0.9$. Under a prespecified common power rule with a fixed true effect, the replication

⁶Replication power calculations themselves are typically designed with this definition in mind. Complementary replication measures include: the relative effect size; whether the 95% confidence interval of the replication covers the original estimate; replication based on meta-analytic estimates; the 95% prediction interval approach (Patil et al., 2016); the ‘small telescopes’ approach (Simonsohn, 2015); and the one-sided default Bayes factor (Wagenmakers et al., 2016).

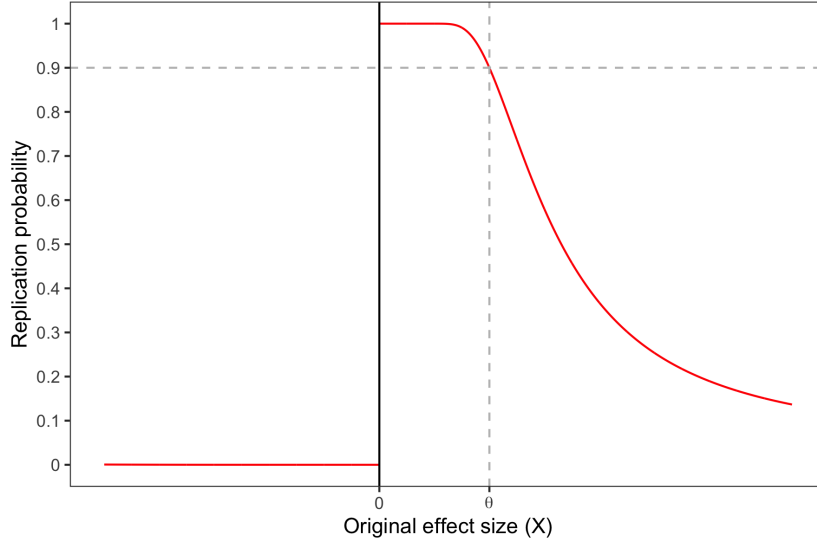


FIGURE 1. REPLICATION PROBABILITY FUNCTION CONDITIONAL ON Θ

Notes: Replication probability function in Definition 2 conditional on a fixed θ . The replication standard error is calculated using the common power rule in Definition 1 to detect original effect sizes with 90% power (i.e. $\sigma_r(X, 0.1) = |X|/3.242$).

probability is only a function of the original estimate X . Denote this conditional replication probability function as $RP(X|\theta)$. It is clear that $RP(X|\theta)$ is non-linear in X , which implies that $\mathbb{E}[RP(X|\theta)] \neq RP(\mathbb{E}[X|\theta])$, even if X is unbiased. If $RP(\cdot|\theta)$ were globally concave, Proposition 1 would immediately follow from Jensen's inequality. However, it is only locally concave around the true effect θ . The proof of Proposition 1 shows that when $1 - \beta \geq 0.8314$, local concavity is sufficient to arrive at the same result for any distribution of latent studies.

The difference between the expected replication rate and its intended target is larger when power in original studies is low. This is because the concavity of $RP(\cdot|\theta)$ is more pronounced when power in original studies is low. As an illustration, Figure 2 plots the relationship between the expected replication rate and power in original studies, again assuming the intended power target in replications is set to 90%, close to mean reported intended replication power in [Open Science Collaboration \(2015\)](#) and [Camerer et al. \(2016\)](#). To highlight the impact of power in original studies, the relationship is derived assuming no p -hacking, no selective publication, and no heterogeneity (i.e. assuming exact replications). The plot shows that the expected replication rate is bounded above by its intended target of 90%, in line with Proposition 1, and is especially low when power in original studies is low. For instance, the expected probability of replicating an original study with 33% power is around 50%. With relatively low estimates of power across various empirical literatures, this provides strong theoretical grounds for expecting low replication rates in practice, even in the absence of issues with p -hacking or treatment effect

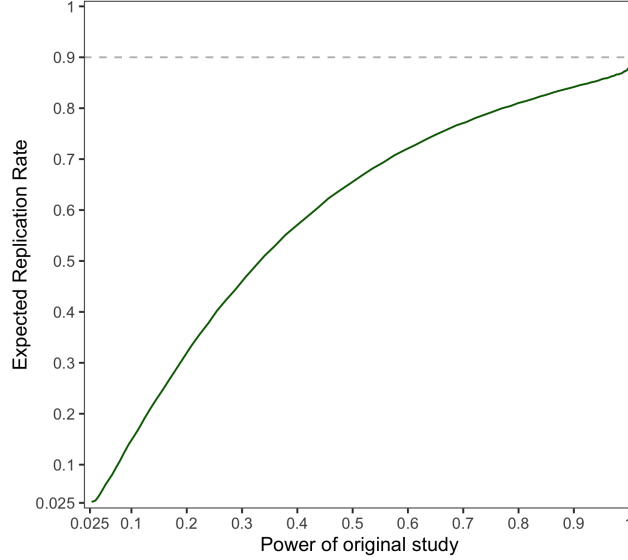


FIGURE 2. ORIGINAL POWER AND THE EXPECTED REPLICATION RATE UNDER THE COMMON POWER RULE

Notes: Power of original study and the expected replication rate under the common power rule are both functions of $\omega = \theta/\sigma$ (normalized to be positive). Power in the original study to obtain a significant effect with the same sign as the true effect is equal to $1 - \Phi(1.96 - \omega)$. The expected replication rate is calculated by taking 10^6 draws of Z from $N(\omega, 1)$; simulating replications under the common power rule to detect the original effect size with 90% power (Definition 1); and then calculating the fraction which ‘replicate’ by being significant and having the same sign as the original study. This figure assumes no p -hacking, no heterogeneity in true effects, no selective publication and random replication selection.

heterogeneity. For intuition, note that if the true effect is zero, the replication probability is 0.025 (regardless of the how the replication standard error is chosen). Continuity implies that when original studies have true effects close to zero (and therefore power in original studies is low), replication probabilities will also be very low.

Proposition 1 applies to replications implementing the common power rule. Some more recent replication studies have used a higher-power variant which I refer to as the fractional power rule, wherein replication power is set to detect some fraction ψ of the estimated effect size (Camerer et al., 2018, 2022). In Proposition B3 in Appendix B, I show that the expected replication rate under the fractional power rule can be either above or below the stated power target $1 - \beta$. More specifically, the expected replication rate can range anywhere between 0.025 and $\Phi[1.96 - \frac{1}{\psi}(1.96 - \Phi^{-1}(\beta))] > 1 - \beta$ depending on the statistical power of original studies. For instance, if $\psi = \frac{3}{4}$ and $1 - \beta = 0.9$, as in the first-stage in Camerer et al. (2018), then the expected replication rate could range anywhere between 0.025 and 0.99. These results build on those in Andrews and Kasy (2019), who argue that replication rates may vary widely depending on the latent distribution of studies. Finally, note that as with Proposition 1, these conclusions hold whether or not there is selective publication, and even in the absence of p -hacking or

treatment effect heterogeneity.

3.3. Additional Factors

The primary focus thus far has been on the impact of the concavity of the replication probability function on expected replication rates. The reason for emphasizing this over other factors is because this specific issue alone turns out to largely explain observed replication rates, at least in economics and social sciences. In this subsection, I explore three additional factors that might also be expected to contribute to low replication rates, but whose impact turns out to be limited in practice.

Consider first the common perception that selective publication favoring significant results – either by authors or journals – produces more ‘false-positives’ in the published literature, which are in turn harder to replicate. This theory is important to address because it enjoys substantial support: over 90% of researchers cite ‘selective reporting’ as a contributing factor to irreproducibility, more than any other factor (Baker, 2016). However, Andrews and Kasy (2019) and Kasy (2021) point out that the replication rate in fact tells us very little about selective publication. Both provide examples showing that the replication rate can take on almost any value depending on the latent distribution of true effects, irrespective of how selective publication is. In fact, in Proposition B4 in Appendix B, I show that the the replication rate in the Andrews and Kasy (2019) model is completely insensitive to selective publication against null results.⁷ This follows from the simple fact that the replication rate definition does not include statistically insignificant results. Thus, even if it were the case that insignificant results were being widely published, they would not be included in the replication rate.^{8,9} Consequently, the replication rate is ill-suited to measuring the extent of the ‘file-drawer’ problem with respect to statistically insignificant results.

The second factor is related to the first: the replication rate induces upward bias in original estimates because it is, by definition, calculated on a selected sample of significant findings. Replication estimates will therefore mechanically regress to the mean (Galton, 1886).¹⁰ This

⁷The proposition proves this more generally for measures $g(\cdot)$ that condition on statistical significance. Setting $g(x, \sigma, x_r, \beta) = \mathbb{1}\left[\frac{|x_r|}{\sigma_r(x, \sigma, \beta)} \geq 1.96, \text{sign}(x_r) = \text{sign}(x)\right]$ gives the result for the replication rate measure.

⁸A caveat is that the model assumes a fixed distribution of latent studies, whereas in practice it may be endogenous, for example, if researchers engage in more specification searching when publication bias against null results is high (Simonsohn et al., 2014; Brodeur et al., 2016, 2020, 2022).

⁹Appendix D examines measures of replication which may be more sensitive to changes in selective publication than the replication rate. For evaluating efforts to reduce selective publication, simulation results show that the prediction interval approach (Patil et al., 2016), when calculated over both significant and insignificant results, provides a useful alternative to the replication rate, the confidence interval measure, and the meta-analysis approach.

¹⁰For a formal statement and proof, see Proposition B2 in Appendix B.

distortion immediately implies that relative effect sizes, which measure the ratio of replication estimates to original estimates, should be below one in expectation, even when original estimates in the published literature are unbiased. While it might be intuitive that this selection would lower the replication rate, the impact is actually theoretically ambiguous when true effects differ across studies, and hence an empirical question. This is because conditioning on significance also tends to select studies for replication with larger true effects, which, all else equal, have higher replication probabilities. In economics and psychology, using the empirical methodology described in the next section, I find that these two effects broadly offset one another, such that the overall impact of selection into significance is relatively small. This is discussed in greater detail in Appendix C.

Finally, a third factor potentially affecting replication rates is that when original estimates are significant but with the ‘wrong’ sign, the probability of replication is very low because it requires the highly unlikely event that the replication estimate also has the wrong sign and is statistically significant. This can be seen in Figure 1, where the replication probability when the original estimate has the wrong sign is bounded above by 0.025.¹¹ Empirically, this turns out to have a small impact on replication rates because the probability of observing a significant original effect with the wrong sign is relatively low.

4. Empirical Applications

In this section, I test the null hypothesis that the expected replication rate from the model in Section 3 matches the replication rate actually observed in three large-scale replication studies. Since the model does not include p -hacking or heterogeneity, the null tests whether observed replication rates can be entirely explained by issues with common power calculations emphasized in Proposition 1. To test this hypothesis, the theory requires that we estimate the latent distribution of studies. This can then be used to generate replication rate predictions which can be compared to observed replication rates. The procedure is as follows:

1. Estimate the latent distribution of studies, $\mu_{\Theta, \Sigma}$ using an augmented version of the [Andrews and Kasy \(2019\)](#) model applied to three large-scale replications.¹² Estimation does not use any data from replications.
2. Use the estimated model to simulate replications and predict what fraction of significant results would replicate, absent any other issues such as p -hacking or heterogeneity.

¹¹Lemma A1.4 in Appendix A shows that the replication probability under the common power rule approaches zero as x approaches zero from below, and approaches 0.025 as x approaches negative infinity.

¹²Note that estimating the latent distribution of studies requires modeling selective publication. However, with estimates of the latent distribution in hand, replication rate predictions in step 2 will not depend on the degree to which null results are suppressed, since the replication rate is defined only over significant results.

3. Compare these predictions (which do not use any data from the replications) to actual replication outcomes.

4.1. Replication Studies

I examine three replication studies. [Camerer et al. \(2016\)](#) replicate results from all 18 between subjects laboratory experiments published in *American Economic Review* and *Quarterly Journal of Economics* between 2011 and 2014. [Open Science Collaboration \(2015\)](#) replicate results from 100 psychology studies in 2008 from *Psychological Science*, *Journal of Personality and Social Psychology*, and *Journal of Experimental Psychology: Learning, Memory, and Cognition*. Following [Andrews and Kasy \(2019\)](#), I consider a subsample of 73 studies with test statistics that are well-approximated by z -statistics. [Camerer et al. \(2018\)](#) replicate 21 experimental studies in the social sciences published between 2010 and 2015 in *Science* and *Nature*.

In [Camerer et al. \(2016\)](#), replicators used the common power rule to detect original effects with at least 90% power. In [Open Science Collaboration \(2015\)](#), replication teams were instructed to achieve at least 80% power using the common power rule, and encouraged to obtain higher power if feasible. Reported mean intended power was 92% in both cases. [Camerer et al. \(2018\)](#) implemented a higher-powered fractional power rule consisting of two stages. In the first stage, replicators aimed to detect 75% of the original effect with 90% power. In the second stage, further data collection was undertaken for insignificant results from the first stage, such that the pooled sample from both stages was calibrated to detect half of the original effect size with 90% power. I predict replication outcomes in the first stage.¹³

Note that the theoretical result in [Proposition 1](#) showing that the expected replication rate is bounded above by its intended target applies to the common power rule and not to the fractional power rule. For the fractional power rule, the expected replication rate can either be above or below the stated power target ([Proposition B3](#)). In both cases, the magnitude of the gap is an empirical question.

Finally, note that [Proposition 1](#) assumes that the probability of selecting significant studies for replication is weakly decreasing in the p -value ([Assumption 2](#)). This assumption plausibly holds in all three empirical applications. In psychology, replicators simply chose the last experiment reported in each article for replication. In economics and social science experiments, replicators selected the ‘most important statistically significant result’ within a study, as emphasized by the authors. In these two applications, [Assumption 2](#) would be violated if, among the set of statistically significant results with p -values below 0.05, authors systematically emphasize a result as more important if it has a *higher* p -value, which seems implausible.

¹³Predicting second-stage outcomes is complicated by the fact that one study that was ‘successfully’ replicated in the first stage was erroneously included in the second stage.

4.2. Estimation

To calculate the expected replication rate, it is necessary to estimate the latent distribution of studies $\mu_{\Theta, \Sigma}$. To do this, I estimate an augmented version of the empirical model in [Andrews and Kasy \(2019\)](#). Specifically, [Andrews and Kasy \(2019\)](#) develop an empirical model to estimate the marginal distribution of true effects Θ^* , but not of standard errors Σ^* . Since predictions of the replication rate also require knowledge of the distribution of Σ^* , I augment the model to estimate the joint distribution of (Θ^*, Σ^*) . For details on the likelihood, see [Appendix E](#).

Estimation is based on the ‘metastudy approach’, which only uses data from original studies, and not from the replication studies whose outcomes are being predicted. I assume that Σ^* follows a gamma distribution with shape and scale parameters denoted by $(\kappa_\sigma, \lambda_\sigma)$. For all other aspects of the model, I implement identical model specifications as [Andrews and Kasy \(2019\)](#), whose focus is on estimating publication bias. Matching their specifications, I assume that $|\Theta^*|$ follows a gamma distribution with shape and scale parameters $(\kappa_\theta, \lambda_\theta)$; and that the joint probability of being published and chosen for replication, $p(X/\Sigma) \times r(X/\Sigma)$, is a step-function parameterized by $\beta_{\mathbf{p}}$. The inclusion of steps at common significance levels (1.64, 1.96, 2.58) varies slightly across applications owing to different approaches for choosing which studies to replicate.¹⁴ In estimation, I normalize the sign of the original estimates to be positive.

Identification for the empirical model requires that latent true effects are statistically independent of latent standard errors, a common assumption in meta-analyses. This assumption is not required for [Proposition 1](#), but only for estimating the empirical model since it is used to identify the publication selection function. [Andrews and Kasy \(2019\)](#) present an alternative ‘systematic replication studies approach’ to estimation, which they suggest can be used as a check the reliability of the meta-study estimates. This is because the ‘systematic replication studies approach’ does not rely on the independence assumption (although it does use data from replication studies which makes it undesirable for the purposes of prediction). Across applications, this alternative approach yields broadly similar estimates of the selection parameters as the ‘meta-study approach’, lending support to the reliability of these estimates.

[Table 1](#) presents the maximum likelihood estimates from the meta-study approach, together with reproduced estimates from [Andrews and Kasy \(2019\)](#) for comparison.¹⁵ For common pa-

¹⁴Details on mechanisms for replication selection are outlined in [Appendix F](#). With $Z = X/\Sigma$, the selection functions in each application are: $r(X/\Sigma) \times p(X/\Sigma) \propto \mathbb{1}(1.64 \leq |Z| < 1.96)\beta_{p2} + \mathbb{1}(|Z| \geq 1.96)$ in economics; $r(X/\Sigma) \times p(X/\Sigma) \propto \mathbb{1}(|Z| < 1.64)\beta_{p1} + \mathbb{1}(1.64 \leq |Z| < 1.96)\beta_{p2} + \mathbb{1}(|Z| \geq 1.96)$ in psychology; and $r(X/\Sigma) \times p(X/\Sigma) \propto \mathbb{1}(1.96 \leq |Z| < 2.58)\beta_{p3} + \mathbb{1}(|Z| \geq 2.58)$ for social science experiments.

¹⁵Estimates for psychology in this article are slightly different to the meta-study estimates reported in [Andrews and Kasy \(2019\)](#) (their [Table 2](#)). The difference is due to a misreported p -value in the raw psychology data for one study, which leads to an erroneous outlier in the distribution of original study standard errors. [Table 1](#) in this article reproduces estimates of their model with the corrected data. Excluding this study in the augmented model leads to very similar replication rate predictions.

TABLE 1 – MAXIMUM LIKELIHOOD ESTIMATES

	Latent true effects $ \Theta^* $		Latent standard errors Σ^*		Selection parameters		
	κ_θ	λ_θ	κ_σ	λ_σ	β_{p1}	β_{p2}	β_{p3}
<i>Economics experiments</i>							
Augmented model	1.426 (1.282)	0.148 (0.072)	2.735 (0.536)	0.103 (0.031)	0.000 (0.000)	0.039 (0.050)	–
Andrews and Kasy (2019)	1.343 (1.285)	0.157 (0.075)	– –	– –	0.000 (0.000)	0.038 (0.050)	–
<i>Psychology experiments</i>							
Augmented model	0.782 (0.423)	0.179 (0.055)	4.698 (0.605)	0.044 (0.008)	0.012 (0.007)	0.303 (0.134)	–
Andrews and Kasy (2019)	0.734 (0.405)	0.185 (0.056)	– –	– –	0.012 (0.007)	0.300 (0.134)	–
<i>Social science experiments</i>							
Augmented model	0.077 (0.106)	0.644 (0.333)	6.249 (1.828)	0.028 (0.009)	0.000 (0.000)	0.000 (0.000)	0.611 (0.425)
Andrews and Kasy (2019)	0.070 (0.091)	0.663 (0.327)	– –	– –	0.000 (0.000)	0.000 (0.000)	0.583 (0.418)

Notes: Maximum likelihood estimates for economics (Camerer et al., 2016), psychology (Open Science Collaboration, 2015) and social sciences (Camerer et al., 2018). Robust standard errors are in parentheses. Latent true effects and standard errors are assumed to follow a gamma distribution; parameters (κ, λ) are the shape and scale parameters, respectively. In economics and psychology, joint publication and replication probability coefficients are measured relative to the omitted category of studies significant at 5 percent level. Parameters β_{p1} , β_{p2} in this case are the relative publication probabilities of studies that are insignificant at the 10% level; and significant at the 10% level but not at the 5% level. For example, in experimental economics, an estimate of $\beta_{p2} = 0.039$ implies that results which are significant at the 5% level are about 26 times more likely to be published and chosen for replication than results that are significant at the 5% level. Note that in economics, results which were insignificant at the 10% level were not selected for replication and hence $\beta_{p1} = 0$. In social sciences, the omitted category is studies significant at the 1% level. Results below the 5% significance level were not chosen for replication so that $\beta_{p1} = \beta_{p2} = 0$, and β_{p3} measures the publication probability of a result that is significant at the 5% level but not at the 1% level, relative to that of a significant result at the 1% level. Andrews and Kasy (2019) estimates are reproduced from accessible data and code from their analysis.

parameters, estimates are very close. Appendix E examines the sensitivity of the main results to the parametric assumptions for the distribution of latent true effect and standard errors. Overall, different parametric assumptions (e.g. log-normal) give rise to replication rate predictions with very similar accuracy to those presented in the main results below.

4.3. The Predicted Replication Rate

Model parameter estimates in Table 1 can be used to generate replication rate predictions by simulating replications using the following procedure:

1. Draw 10^6 latent (published or unpublished) research questions and standard errors $(\theta^{*sim}, \sigma^{*sim})$ from the estimated joint distribution $\hat{\mu}_{\Theta, \Sigma}(\hat{\kappa}_\theta, \hat{\lambda}_\theta, \hat{\kappa}_\sigma, \hat{\lambda}_\sigma)$.
2. Draw original estimates $x^{*sim} | \theta^{*sim}, \sigma^{*sim} \sim N(\theta^{*sim}, \sigma^{*sim2})$ for each latent study.

3. Use the estimated selection parameters $\hat{\beta}_{\mathbf{p}}$ to determine the subset of studies that are published and chosen for replication.
4. For studies chosen for replication, calculate the replication standard error σ_r^{sim} according to the following rule

$$\sigma_r^{sim}(x^{sim}, \beta, \psi) = \frac{\psi \cdot |x^{sim}|}{1.96 - \Phi^{-1}(\beta)} \quad (8)$$

where $\psi = 1$ and $1 - \beta = 0.92$ in economics and psychology, which corresponds to the common power rule; and $\psi = \frac{3}{4}$ and $1 - \beta = 0.9$ in social science experiments, which corresponds to a fractional power rule.¹⁶

5. Simulate replications by drawing replication estimates $x_r^{sim} | \theta^{sim}, \sigma_r^{sim} \sim N(\theta^{sim}, \sigma_r^{sim2})$

The replication rate prediction is then formed by calculating the fraction of simulated replications in which the original study is replicated with the same sign and statistical significance. More formally, let $\{x_i, \sigma_i, x_{r,i}, \sigma_{r,i}\}_{i=1}^{M_{sig}}$ denote the (simulated) set of published, replicated original studies that are significant at the 5% level, and their corresponding replication results. M_{sig} is the number of replicated originally-significant studies. Further restricting attention to the set of statistically significant results is done to match the procedure in large-scale replications where only significant results were included in the replication rate calculation (or chosen for replication in the first place).^{17,18} The predicted replication rate is the simulated analogue of the expected replication rate in Definition 3, and thus equal to the share of replication estimates $x_{r,i}$ that are statistically significant and have the same sign as the corresponding original estimate x_i

$$\frac{1}{M_{sig}} \sum_{i=1}^{M_{sig}} \mathbb{1}\left(|x_{r,i}| \geq 1.96\sigma_{r,i}, \text{sign}(x_{r,i}) = \text{sign}(x_i)\right) \quad (9)$$

¹⁶This assumes all simulated replications set intended power equal to the mean of reported intended power. In practice, there was some variation in the application of the power rule around the mean. Appendix G reports predicted replication rates allowing for variation in intended power across studies that matches the empirical variation in each application. Results are very similar and in fact slightly more accurate in all three applications (61.5% in economics; 52.2% in psychology; and 55.5% in social science).

¹⁷In both experimental economics and psychology, a small number of original results whose p -values were slightly above 0.05 were treated as ‘positive’ results and included in the replication rate calculation. To match this, I set the cutoff for significant findings for the purposes of replication equal to the smallest z -statistic that was treated as a ‘positive’ result for replication. Predictions are almost identical with a strict 0.05 significance threshold.

¹⁸Note that original estimates are normalized to be positive in estimation but not in the calculation of the predicted replication rate.

4.4. Results

Table 2 presents the results. In experimental economics, the predicted replication rate is 60%, which is very close to the observed rate of 61.1%. This is an “out-of-sample” prediction in the sense that the model is estimated only using information from the original studies, and does not incorporate any information from the replications. The accuracy of this prediction is consistent with the null hypothesis that the observed replication rate in economics can be explained entirely by a parsimonious model accounting only for issues with power calculations, and not other issues such as p -hacking or treatment effect heterogeneity. Failure to reject the null hypothesis does not, of course, imply that it is true, and thus we should not necessarily conclude that these other factors are not present. Nonetheless, other evidence points to a relatively limited role for p -hacking in the context of lab experiments studied here, perhaps due to fewer researcher degrees of freedom as compared with observational settings (Brodeur et al., 2016, 2020; Imai et al., 2020). Note that despite the very accurate point estimate, the standard error is relatively large, which implies limited power to reject the model’s prediction (perhaps owing to the fact that there are only 18 replicated studies).

In psychology, the model predicts a replication rate of 54.5%. This is well below mean intended power of 92%, but higher than the observed replication rate of 34.8%. In this case, the model accounts for around two-thirds of the replication rate gap, and we can reject the null hypothesis that the replication gap is entirely explained by issues with common power calculations. The unexplained portion of the gap in psychology provides evidence that other factors discussed in the literature and not incorporated in the model may be important, including heterogeneity in true effects, p -hacking, and measurement error. Another possibility is that the model should account for differences in replicating main effects and interaction effects, and differences across subfields (Open Science Collaboration, 2015; Altmejd et al., 2019).

Low predicted replication rates can largely be attributed to the issue that replication power calculations do not account for the concavity of the power function. Appendix C shows that the two other factors affecting replication rates in the model, and discussed in Subsection 3.3, turn out to have a relatively small impact empirically. The first factor is that original estimates with the wrong sign have very low replication probabilities, as the chance of obtaining a significant replication estimate with the wrong sign is very low. This has a small impact on expected replication rates because the probability that significant original estimates have the wrong sign is relatively low in both economics (3.0%) and psychology (5.4%). The second factor is that the replication rate is calculated on a selected sample of significant findings, which mechanically induces upward bias in original estimates that can make successful replications less probable. Perhaps surprisingly, this has only a small negative impact on the expected replication rate

	Economics experiments	Psychology	Social sciences
Nominal target (intended power)	0.92	0.92	–
Observed replication rate	0.611	0.348	0.571
Predicted replication rate	0.600	0.545	0.543
	(0.122)	(0.054)	(0.134)

TABLE 2 – REPLICATION RATE PREDICTIONS

Notes: Economics experiments refers to [Camerer et al. \(2016\)](#), psychology experiments to [Open Science Collaboration \(2015\)](#) and social sciences to [Camerer et al. \(2018\)](#). The replication rate is defined as the share of original estimate whose replications have statistically significant findings of the same sign. Figures in the first row report the mean intended power reported in both applications. The second row shows observed replication rates. The third row reports the predicted replication rate in equation (9) calculated using parameter estimates Table 1. The fourth row shows standard errors for the predicted replication rate which are calculated using the delta method. In social sciences, power is set to detect three-quarters of the original effect size with 90% power. This approach does not have a fixed nominal target for the replication rate. See Appendix E for robustness to alternative parametric assumptions in the estimated model.

in economics while actually slightly *increasing* it in psychology. As discussed in Subsection 3.3, this is because selection on significance also leads to the replication of more studies with larger true effects (since they are more likely to be significant), which have higher replication probabilities (Figure 2). These two effect broadly offset one another such that the overall impact of selection is relatively small.

A popular variant for the common power rule is the fractional power rule, where replication power is set to detect some fraction of the original effect size with a given level of statistical power (e.g. [Camerer et al. \(2018\)](#) and [Camerer et al. \(2022\)](#)). Theoretically, under the specific rule applied in [Camerer et al. \(2018\)](#), the expected replication rate can range anywhere between 0.025 and 0.99 depending on the power in original studies.¹⁹ Empirically, the predicted replication rate for the experimental social sciences is 54.3%, which is very close to the observed rate of 57.1%. The difference is statistically indistinguishable from zero, although the standard error of the prediction is quite large. Similarly to experimental economics, the accuracy of the point estimate of the prediction implies that we cannot reject the null hypothesis that the observed replication rate can be explained by a parsimonious model accounting only for issues with power calculations.

Finally, for the social sciences, note that I do not estimate the relative importance of alternative factors in explaining low replication rates as I do for the other applications e.g. selection on significance and original estimates with the wrong sign. This is because the approach described in Appendix C is based on decomposing the gap between the replication rate and its

¹⁹Proposition B3 shows that the expected replication rate can range between 0.025 and $1 - \Phi[1.96 - \frac{1}{\psi}(1.96 - \Phi^{-1}(\beta))]$. With the fraction of original effect size to detect equal to $\psi = 3/4$, and intended power set to $1 - \beta = 0.9$, the upper range equals 0.99.

nominal target, and the fractional power rule does not provide a clear nominal target.

4.5. Extensions

I examine three extensions. In Appendix H, I use the empirical models estimated in Table 1 to generate predicted average relative effect sizes, using a similar procedure to the replication rate predictions. I find that the predicted relative effect size is quite similar to the observed value in economics (0.70 vs. 0.66). In the social sciences, the model is somewhat farther off (0.53 vs. 0.44), which may suggest a role for other factors such as p -hacking or heterogeneity, although the difference is not statistically distinguishable from zero. Finally, in psychology, the prediction is quite far off (0.64 vs. 0.37), again providing strong evidence for alternative factors. Note that relative effect sizes are affected both by selection of significant results for replication and the level of statistical power in original studies.²⁰

A second extension considers the proposed rule of setting replication power equal to original power in Appendix G. In a review of 108 psychology replications by Anderson and Maxwell (2017), 19 (17.6%) implemented this approach. In all three applications, this approach leads to lower predicted replication rates than under the common power rule.

Given the issues that stem from conditioning on statistical significance, the third extension in Appendix I examines the suggestion of extending the replication rate definition to include null results that are ‘replicated’ if their replications are also insignificant. For empirical models in economics and psychology, this ‘extended’ replication rate remains below intended power under the common power rule.

5. Practical Recommendations

The results thus far highlight some important limitations in the common approaches used for conducting and interpreting replication outcomes. In particular, (i) expected replication rates cannot reach stated nominal power targets under the common power rule used for setting replication power (Proposition 1); and (ii) relative effect sizes are below their nominal target of one in expectation when replicators only select significant results to replicate, because selection implies that replication estimates must regress to the mean (Proposition B2). This makes observed replication rates and relative effect sizes difficult to interpret because, in both cases, there is no clear benchmark against which to judge them. Importantly, issues (i) and (ii) are present even in the absence of p -hacking, publication bias, and treatment effect heterogeneity.

²⁰Figure H1 in Appendix H shows that the expected relative effect size is an increasing function of power in original studies and approaches one as original power approach 100%.

In light of these limitations, what should researchers running replications do? In this section, I make four practical recommendations relevant to different stages of the replication process, discussing each in greater detail below.

1. **Replication selection:** for large-scale studies replicating a number of findings, do not impose a selection rule that only chooses significant results for replication.
2. **Replication sample sizes:** set replication sample sizes independently of the magnitude of original estimates.
3. **Replication measures:** focus on the relative effect size and formally test whether it deviates from one using the prediction interval approach (Patil et al., 2016).
4. **Interpreting replication outcomes:** report the relative effect size in conjunction with prediction intervals and include an explicit discussion addressing the power of the test.

The motivation for the first recommendation is clear from previous discussions, namely, that samples selected on extreme characteristics, such as statistical significance, will regress to the mean in repeated samples. Selection of this kind implies that relative effect sizes will be below one in expectation, even when original findings are unbiased, and should therefore be avoided where possible. For example, replicating a randomly chosen set of results would restore the nominal target of one as a meaningful benchmark.

Second, as an alternative to the common power rule or the fractional power rule, replicators should consider setting replication sample sizes independently of the magnitude of original estimates. This approach was taken, for instance, in Protzko et al. (2024), where all replications were conducted using a target sample of at least 1,500 participants. Another approach would be to determine the replication sample size by scaling the original sample size by some constant factor (e.g. double the original sample size). There are several reasons for this. First, the common power rule and the fractional power rule are designed to detect original estimates with a prespecified level of power, and hence are not properly suited to replicating null results (as is suggested in the first recommendation). Second, setting the replication sample size independently of the original effect size is required for validity of the recommended prediction interval approach.²¹

The third recommendation is to focus on relative effect sizes over the binary measure of replication. The relative effect size can also be connected to a formal statistical test based on the prediction interval approach in Patil et al. (2016), which tests the null hypothesis that

²¹More formally, the prediction interval approach requires that original and replication estimates are statistically independent: $X \perp\!\!\!\perp X_r$.

$X \sim N(\theta, \Sigma^2)$ and $X_r \sim N(\theta, \Sigma_r^2)$ — i.e. that both the original and replicated estimates are (uncensored) normal draws centered around the same true effect. Under this null, the probability that the relative effect size X/X_r lies within the prediction interval $(1 - 1.96\frac{\tilde{\Sigma}}{X}, 1 + 1.96\frac{\tilde{\Sigma}}{X})$ equals 0.95, where $\tilde{\Sigma} \equiv \sqrt{\Sigma^2 + \Sigma_r^2}$.²² In other words, when the relative effect size deviates sufficiently far from one, we can reject the null that both the original and replication estimates come from normal distributions with the same true effect. See Appendix J for details on the derivation, and some guidance on implementation.²³

There are at least three advantages to using prediction interval approach. First, it connects the relative effect size to a well-specified null hypothesis test which is tightly linked to common concerns of replicators over publication bias, p -hacking, and treatment effect heterogeneity. Note also that choosing to replicate only significant results immediately violates the null hypothesis. The null in the prediction interval approach can have advantages over the null in the binary replication measure of no true effect. For instance, under the null in the binary measure, the replication probability can be high if the replication sample size is very large, even if the original study is severely biased due to p -hacking. By contrast, with large sample sizes, the prediction interval approach is likely to reject the null. This is because the prediction interval approach captures not only the sign and significance of the replication estimate, but also its magnitude.

A second advantage is that it explicitly incorporates sampling variation from both the original study and the replication study. In particular, the range of estimates consistent with the null decreases as the sample sizes in the original and replication studies increase i.e. as $\tilde{\Sigma}$ decreases. By contrast, the binary measure of replication under the common power rule only considers the original estimate, not its standard error. This is an undesirable property, since, intuitively, with very noisy original estimates, it should be more difficult to make precise conclusions about whether a replication study reproduces the original result. That power to reject the null in the prediction interval approach depends on the power of original studies also highlights the usefulness of well-powered original studies for making precise conclusions about replicability.

A third advantage of the prediction interval approach is that it allows us to test for the presence of publication bias (conditional on following the first recommendation). This is not

²²Note that I recast the prediction interval from the original Patil et al. (2016) study in terms of relative effect sizes. The motivation behind this reformulation is simply that relative effect sizes are commonly reported in the replication literature, making them more convenient to interpret and compare across studies. However, for the replication of null results, it may be more appropriate to use the unnormalized prediction interval approach, namely, to reported whether $X_r \in (X - 1.96\tilde{\Sigma}, X + 1.96\tilde{\Sigma})$. This is because dividing by an original estimate that is close to zero might lead to extreme values which are difficult to interpret.

²³An illustrative example provides steps on how to covert the prediction interval from Fisher transformation units – which are used for inference (Fisher, 1915) – to correlation coefficient units.

the case for the replication rate. As discussed in Subsection 3.3, the replication rate is defined only over significant original studies and is therefore insensitive to whether or not null results are censored by publication bias. Moreover, a definition of replication “success” which is based on significance does not extend naturally to replicating original null results. By contrast, the prediction interval measure is low when selective publication is high, and approaches 95% as the probability of publishing null results approach one. For more details, see Appendix D. This is arguably a very important property for a replication metric, since publication bias remains a central concern for replicators and the scientific community more broadly.

The fourth and final recommendation concerns how the relative effect size should be reported and interpreted. When interpreting replication outcomes, researchers should discuss both: (i) the distance of relative effect sizes from one; and (ii) whether prediction intervals contain scientifically significant deviations from one. Jointly considering these metrics is important for making an overall assessment of the reproducibility of the original finding. For example, it is possible to observe a relative effect size well below one yet still fail to reject the null due to large prediction intervals. In this case, the test has limited power to reject the null and suggests a need to focus on better powered original and/or replication studies in order to reach more robust conclusions.

6. Conclusion

The prominence of the replication rate stems in part from its apparent transparency and ease of interpretation. However, caution should be applied when interpreting the replication rate from large-scale replication studies using the common power rule for setting replication power. In general, intended replication targets are not attainable in expectation. Moreover, the replication rate gap will be particularly large when original power is low. Empirical evidence supports the importance of these theoretical insights. In a parsimonious model with neither heterogeneity nor p -hacking, predicted replication rates in experimental economics and social science are very close to observed values. This is consistent with the null hypothesis that problems with power calculations alone are sufficient to explain observed replication rates in these fields.

As an alternative to focusing on the binary measure of replication, replicators might consider focusing primarily on relative effect sizes and formally testing whether they deviate from one. Moreover, replicators should be cautious about selecting only significant results for replication, as this induces distortions which complicate the interpretation of replication outcomes.

7. Acknowledgements

I am especially grateful for the feedback, advice, and encouragement of Jonathan Roth. For helpful comments, suggestions and conversations, I thank Johannes Abeler, Abel Brodeur, Daniel Björkegren, Pedro Dal Bó, Anna Dreber, Peter Hull, Toru Kitagawa, Soonwoo Kwon, and Jesse Shapiro, as well as seminar participants at Brown University, University College London, 2024 BITTS Annual Meeting, and the AIMOS 2022 conference.

References

- Altmejd, A., A. Dreber, E. Forsell, et al. (2019). Predicting the Replicability of Social Science Lab Experiments. *PLoS ONE* 14(12).
- Amrhein, V., D. Trafimow, and S. Greenland (2019). Inferential Statistics as Descriptive Statistics: There Is No Replication Crisis if We Don't Expect Replication. *The American Statistician* 73(1), 262–270.
- Anderson, S. F. and S. E. Maxwell (2017). Addressing the “Replication Crisis”: Using Original Studies to Design Replication Studies with Appropriate Statistical Power. *Multivariate Behavioral Research* 52(3), 305–324.
- Andrews, I. and M. Kasy (2019). Identification of and Correction for Publication Bias. *American Economic Review* 109(8), 2766–2794.
- Arel-Bundock, V., R. C. Briggs, H. Doucouliagos, et al. (2023). Quantitative Political Science Research is Greatly Underpowered. *OSF Preprint*.
- Baker, M. (2016). 1,500 Scientists Lift the Lid on Reproducibility. *Nature* 533, 452–454.
- Brodeur, A., N. Cook, and A. Heyes (2020). Methods Matter: p-Hacking and Publication Bias in Causal Analysis in Economics. *American Economic Review* 110(11), 3634–3660.
- Brodeur, A., N. Cook, and A. Heyes (2022). We Need to Talk About Mechanical Turk: What 22,989 Hypothesis Tests Tell Us About Publication Bias and p-Hacking in Online Experiments. *IZA Discussion Paper* 15478.
- Brodeur, A., M. Lé, M. Sangnier, and Y. Zylberberg (2016). Star Wars: The Empirics Strike Back. *American Economic Journal: Applied Economics* 8(1), 1–32.
- Bryan, C. J., D. S. Yeager, and J. M. O'Brien (2019). Replicator Degrees of Freedom Allow Publication of Misleading Failures to Replicate. *Proceedings of the National Academy of Sciences of the United States of America* 116(51), 25535–25545.

- Button, K. S., J. P. Ioannidis, C. Mokrysz, et al. (2013). Power Failure: Why Small Sample Size Undermines the Reliability of Neuroscience. *Nature Reviews Neuroscience* 14(5), 365–376.
- Camerer, C., Y. Chen, A. Dreber, et al. (2022). Mechanical Turk Replication Project.
- Camerer, C. F., A. Dreber, E. Forsell, et al. (2016). Evaluating Replicability of Laboratory Experiments in Economics. *Science* 351(6280), 1433–1436.
- Camerer, C. F., A. Dreber, F. Holzmeister, et al. (2018). Evaluating the replicability of social science experiments in Nature and Science between 2010 and 2015. *Nature Human Behaviour* 2(9), 637–644.
- Cesario, J. (2014). Priming, Replication, and the Hardest Science. *Perspectives on Psychological Science* 9(1), 40–48.
- DellaVigna, S. and E. Linos (2022). RCTs to Scale: Comprehensive Evidence From Two Nudge Units. *Econometrica* 90(1), 81–116.
- DellaVigna, S., N. Otis, and E. Vivaldi (2020). Forecasting the Results of Experiments: Piloting an Elicitation Strategy. *AEA Papers and Proceedings* 110, 75–79.
- Dreber, A., T. Pfeiffer, J. Almenberg, et al. (2015). Using Prediction Markets to Estimate the Reproducibility of Scientific Research. *Proceedings of the National Academy of Sciences of the United States of America* 112(50), 15343–15347.
- Elliott, G., N. Kudrin, and K. Wüthrich (2022). Detecting p-Hacking. *Econometrica* 90(2), 887–906.
- Fisher, R. A. (1915). Frequency Distribution of the Values of the Correlation Coefficient in Samples from an Indefinitely Large Population. *Biometrika* 10(4), 507–521.
- Franco, A., N. Malhotra, and G. Simonovits (2014). Publication Bias in the Social Sciences: Unlocking the File Drawer. *Science* 345(6203), 1502–1505.
- Frankel, A. and M. Kasy (2022). Which Findings Should Be Published? *American Economic Journal: Microeconomics* 14(1), 1–38.
- Galton, F. (1886). Regression Towards Mediocrity in Hereditary Stature. *The Journal of the Anthropological Institute of Great Britain and Ireland* 15, 246–263.
- Gelman, A. and J. Carlin (2014). Beyond Power Calculations: Assessing Type S (Sign) and Type M (Magnitude) Errors. *Perspectives on Psychological Science* 9(6), 641–651.
- Gordon, M., D. Viganola, M. Bishop, et al. (2020). Are Replication Rates the Same Across Academic Fields? Community Forecasts from the DARPA SCORE Programme. *Royal Society Open Science* 7.

- Higgins, J. P. and S. G. Thompson (2002). Quantifying heterogeneity in a meta-analysis. *Statistics in Medicine* 21(11), 1539–1558.
- Imai, T., K. Zemlianova, N. Kotecha, et al. (2020). How Common are False Positives in Laboratory Economics Experiments? Evidence from the P-Curve Method. *Working Paper*.
- Ioannidis, J. P. (2005). Why Most Published Research Findings Are False. *PLoS Medicine* 2(8).
- Ioannidis, J. P. (2008). Why Most Discovered True Associations Are Inflated. *Epidemiology* 19(5), 640–648.
- Ioannidis, J. P., T. D. Stanley, and H. Doucouliagos (2017). The Power of Bias in Economics Research. *The Economic Journal* 127(605), 236–265.
- Kasy, M. (2021). Of Forking Paths and Tied Hands: Selective Publication of Findings, and What Economists Should Do about It. *Journal of Economic Perspectives* 35(3), 175–192.
- Klein, R. A., K. A. Ratliff, M. Vianello, et al. (2014). Investigating Variation in Replicability: A “Many Labs” Replication Project. *Social Psychology* 45(3), 142–152.
- Klein, R. A., M. Vianello, F. Hasselman, et al. (2018). Many Labs 2: Investigating Variation in Replicability Across Samples and Settings. *Advances in Methods and Practices in Psychological Science* 1(4), 443–490.
- Laird, N. M. and F. Mosteller (1990). Some Statistical Methods for Combining Experimental Results. *International Journal of Technology Assessment in Health Care* 6(1), 5–30.
- Maxwell, S. E., M. Y. Lau, and G. S. Howard (2015). Is Psychology Suffering from a Replication Crisis? What Does “Failure to Replicate” Really Mean? . *American Psychologist* 70(6), 487–498.
- Miguel, E. and G. Christensen (2018). Transparency, Reproducibility, and the Credibility of Economics Research. *Journal of Economic Literature* 56(3), 920–980.
- Nosek, B. A., T. E. Hardwicke, H. Moshontz, et al. (2022). Replicability, Robustness, and Reproducibility in Psychological Science. *Annual Review of Psychology* 73, 719–748.
- Open Science Collaboration (2015). Estimating the Reproducibility of Psychological Science. *Science* 349(6251).
- Patil, P., R. D. Peng, and J. T. Leek (2016). What Should Researchers Expect When They Replicate Studies? A Statistical View of Replicability in Psychological Science. *Perspectives on Psychological Science* 11(4), 539–544.
- Protzko, J., J. Krosnick, L. Nelson, et al. (2024). High Replicability of Newly Discovered Social-Behavioural findings is Achievable. *Nature Human Behaviour* 9, 311–319.

- Simons, D. J. (2014). The Value of Direct Replication. *Perspectives on Psychological Science* 9(1), 76–80.
- Simonsohn, U. (2015). Small Telescopes: Detectability and the Evaluation of Replication Results. *Psychological Science* 26(5), 559–69.
- Simonsohn, U., L. D. Nelson, and J. P. Simmons (2014). P-Curve: A Key to the File-Drawer. *Journal of Experimental Psychology: General* 143(2), 534–547.
- Stanley, T. D., E. C. Carter, and H. Doucouliagos (2018). What Meta-Analyses Reveal About the Replicability of Psychological Research. *Psychological Bulletin* 144(12), 1325–1346.
- Vu, P. (2022). Replication data for: Can the Replication Rate Tell Us About Selective Publication? *American Economic Association [publisher], Inter-university Consortium for Political and Social Research [distributor]*.
- Wagenmakers, E.-J., J. Verhagen, and A. Ly (2016). How to Quantify the Evidence for the Absence of a Correlation. *Behavior Research Methods* 48(2), 413–26.

Online Appendix

This appendix contain proofs and supplementary materials for “Why Are Replication Rates So Low?”

A. Properties of the Replication Probability Function

This Appendix derives a number of properties of the replication probability function (Definition 1). The simply provides a convenient, compact notation. The remaining properties consider the replication probability function under the common power rule to detect original effect sizes with $1 - \beta$ intended power (Definition 1). Recall that the replication probability for original study (x, σ, θ) is equal to

$$RP(x, \theta, \sigma_r(x, \sigma, \beta)) = \mathbb{P}\left(\frac{|X_r|}{\sigma_r(x, \beta)} \geq 1.96, \text{sign}(X_r) = \text{sign}(x)\right) \quad (10)$$

To provide intuition of the properties, Figure A1 provides an illustration of the replication probability function for different values of x under the common power rule for $1 - \beta = 0.9$ and a fixed value of $\theta > 0$. Note that Lemma A1 assumes $\theta > 0$, although flipping the sign simply reflects the replication function about the y-axis.

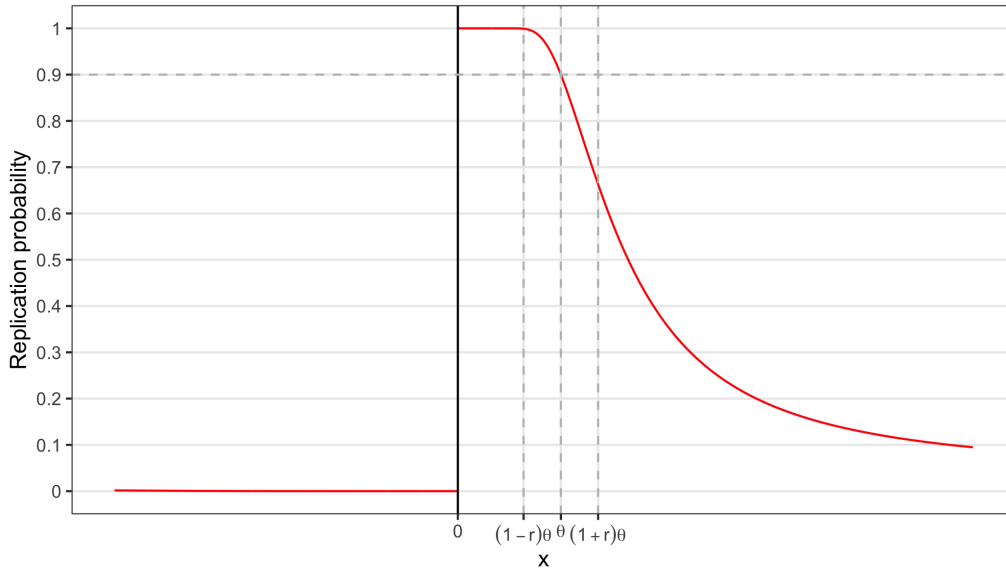


FIGURE A1. Replication Probability Function

Notes: Example of the replication probability function under the common power rule with intended power $(1 - \beta) = 0.9$. The two vertical lines around θ marks the open interval over which the replication probability function is strictly concave, where r^* is given by equation (12).

Lemma A1 (Replication Probability Function Properties). *Let $\theta > 0$. The replication probability function satisfies the following properties:*

1. *For any replication standard error $\sigma_r(x, \sigma, \beta)$, the replication probability for an original study (x, σ, θ) can be written compactly as*

$$RP(x, \theta, \sigma_r(x, \sigma, \beta)) = 1 - \Phi\left(1.96 - \text{sign}(x) \frac{\theta}{\sigma_r(x, \sigma, \beta)}\right) \quad (11)$$

The remaining properties assume the replication standard error $\sigma_r(x, \beta)$ is set using the common power rule in Definition 1 with intended power $1 - \beta$:

2. *If $1 - \beta > 0.025$, then $RP(x, \theta, \sigma_r(x, \beta))$ is strictly decreasing in x over $(-\infty, 0)$ and $(0, \infty)$.*
3. *If $(1 - \beta) > 0.6628$, then $RP(x, \theta, \sigma_r(x, \beta))$ is strictly concave with respect to x over the open interval $(\max\{0, [1 - r^*(\beta)]\theta\}, [1 + r^*(\beta)]\theta)$, where*

$$r^*(\beta) = -(2 + 1.96.h(\beta)) + \sqrt{\frac{(2 + 1.96.h(\beta))^2 - 4 \times (1 + 1.96.h(\beta) - h(\beta)^2)}{2}} > 0 \quad (12)$$

with $h(\beta) = (1.96 - \Phi^{-1}(\beta))$.

4. *The limits of the replication probability function with respect to x are*

$$\lim_{x \rightarrow \infty} RP(x, \theta, \sigma_r(x, \beta)) = 0.025 \text{ and } \lim_{x \rightarrow -\infty} RP(x, \theta, \sigma_r(x, \beta)) = 0.025 \quad (13)$$

$$\lim_{x \uparrow 0} RP(x, \theta, \sigma_r(x, \beta)) = 0 \text{ and } \lim_{x \downarrow 0} RP(x, \theta, \sigma_r(x, \beta)) = 1 \quad (14)$$

5. *Suppose $X^* \sim N(\theta, \sigma^2)$. Then $\mathbb{E}[RP(X, \theta, \sigma_r(X, \beta))] \rightarrow 1 - \beta$ as $\theta \rightarrow \infty$ for fixed σ .*

Proof of 1.

The probability in equation (10) equals $[\mathbb{1}(x/\sigma \geq 1.96) \times (1 - \Phi(1.96 - \frac{\theta}{\sigma_r}))] + [\mathbb{1}(x/\sigma \leq -1.96) \times \Phi(-1.96 - \frac{\theta}{\sigma_r})]$. This captures the two requirements for ‘successful’ replication: the replication estimate must attain statistical significance and have the same sign as the original estimate. Equation (11) is obtained using the symmetry of the normal distribution, which implies that $\Phi(t) = 1 - \Phi(-t)$ for any t . \square

Proof of 2.

The first derivative of the replication probability function with the common power rule is

$$\frac{\partial RP(x, \theta, \sigma_r(x, \beta))}{\partial x} = \begin{cases} -\frac{\theta}{x^2}(1.96 - \Phi^{-1}(\beta)) \times \phi\left(1.96 - \frac{\theta}{x}(1.96 - \Phi^{-1}(\beta))\right), & x > 0 \\ -\frac{\theta}{x^2}(1.96 - \Phi^{-1}(\beta)) \times \phi\left(-1.96 - \frac{\theta}{|x|}(1.96 - \Phi^{-1}(\beta))\right), & x < 0 \end{cases} \quad (15)$$

These are strictly negative whenever $(1.96 - \Phi^{-1}(\beta)) > 0 \iff (1 - \beta) > 0.025$. \square

Proof of 3.

First, note that for $x > 0$, the second derivative of the replication probability function with the common power rule is

$$\frac{\partial^2 RP(x, \theta, \sigma_r(x, \beta))}{\partial x^2} = \left(\frac{h(\beta)\theta}{x^3}\right)\phi\left(1.96 - \frac{h(\beta)\theta}{x}\right)\left[1 + \left(\frac{h(\beta)\theta}{x}\right)\left(1.96 - \frac{h(\beta)\theta}{x}\right)\right] \quad (16)$$

Let $x = (1 + r)\theta$. Substituting this into the previous equation and simplifying shows that equation (16) is strictly negative when the following inequality is satisfied

$$r^2 + (2 + 1.96h(\beta)).r + (1 + 1.96h(\beta) - h(\beta)^2) < 0 \quad (17)$$

The solution to the quadratic equation has a unique positive solution $r^*(\beta)$ whenever $(1 - \beta) > 0.6628$. To see this, note that there exists a unique positive solution when $(1 + 1.96h(\beta) - h(\beta)^2) < 0$. This quadratic equation in $h(\beta)$ must have a unique positive and negative solution in turn, since the parabola opens downwards and equals 1 when $h(\beta) = 0$. The positive root can be obtained from the quadratic formula, which gives 2.38014. Since the quadratic function opens downward, this implies that for any $h(\beta) > 2.38014$, we have $(1 + 1.96h(\beta) - h(\beta)^2) < 0$. Thus, a unique positive solution to equation (17) exists whenever this condition is satisfied. In particular, a unique positive solution exists whenever

$$\begin{aligned} h(\beta) &= 1.96 - \Phi^{-1}(\beta) > 2.38014 \\ \iff \Phi(1.96 - 2.38014) &> \beta \\ \iff (1 - \beta) &> 0.6628 \end{aligned} \quad (18)$$

The unique positive solution for equation (17) can again be obtained by the quadratic

formula, which gives equation (12). Note that for any $r > 0$ where the inequality for concavity in equation (17) is satisfied, the same must also be true of $-r$, since it makes the left-hand-side strictly smaller. This implies that the replication probability function is strictly concave (since its second derivative is strict negative) over $(\max\{0, [1 - r^*(\beta)]\theta\}, [1 + r^*(\beta)]\theta)$, where the maximum is taken because the replication probability function is discontinuous at 0. This follows because of the properties of the quadratic function. Specifically, suppose $f(x)$ is a parabola that opens upward and intersects the y-axis at a negative value. Then for any two points (a, b) with $a < b$ and $f(a), f(b) < 0$, it must be that $f(c) < 0$ for any $c \in (a, b)$. \square

Proof of 4.

Substituting the common power rule into the replication probability function gives

$$RP(x, \theta, \sigma_r(x, \beta)) = 1 - \Phi\left(1.96 - \frac{\theta}{x}(1.96 - \Phi^{-1}(\beta))\right) \quad (19)$$

The values of the limits can be seen immediately from this expression. \square

Proof of 5.

This proof consists of two steps. In the first step, I show that the replication probability function approaches linearity in x in an even interval around θ , as $\theta \rightarrow \infty$ for fixed σ . To see this, fix $r \in (0, 1)$. Then the second derivative evaluated at any point $c\theta \in (r\theta, (1+r)\theta)$ equals

$$\left. \frac{\partial^2 RP(x, \theta, \sigma_r(x, \beta))}{\partial x^2} \right|_{x=c\theta} = \left(\frac{h(\beta)}{c^3\theta^2} \right) \phi\left(1.96 - \frac{h(\beta)}{c}\right) \left[1 + \left(\frac{h(\beta)}{c} \right) \left(1.96 - \frac{h(\beta)}{c}\right) \right] \quad (20)$$

This approaches zero as $\theta \rightarrow \infty$, which implies that $RP(x, \theta, \sigma_r(x, \beta))$ approaches linearity in x over the interval $(r\theta, (1+r)\theta)$ in the limit.

For the second step, see that as $\theta \rightarrow \infty$ with fixed σ , we have that

$$\mathbb{P}[X^* \in (r\theta, (1+r)\theta) | \theta, \sigma] = \Phi\left(\frac{(1+r)\theta - \theta}{\sigma}\right) - \Phi\left(\frac{r\theta - \theta}{\sigma}\right) \rightarrow 1 \quad (21)$$

That is, the probability of drawing X^* inside of the range $(r\theta, (1+r)\theta)$ approaches one in the limit. But from the first step we know that the replication probability function is linear over this range as $\theta \rightarrow \infty$ with fixed σ . This implies in the limit that $\mathbb{E}[RP(X, \theta, \sigma_r(X, \beta))] = RP(\mathbb{E}[X], \theta, \sigma_r(X, \beta)) = RP(\theta, \theta, \sigma_r(X, \beta)) = 1 - \beta$, as shown in Lemma B1.

B. Proofs of Propositions

For convenience, the proofs sometimes use notation distinguishing selection functions over significant and insignificant regions. For example, the publication probability function $p(\cdot)$ is equal to

$$p(X^*/\Sigma^*) = \begin{cases} p_{sig}(X^*/\Sigma^*) & \text{if } S_X^* = 1 \\ p_{insig}(X^*/\Sigma^*) & \text{if } S_X^* = 0 \end{cases}$$

where S_X^* is an indicator variable that equals one if $|X^*/\Sigma^*| \geq 1.96$ and zero otherwise. Similar notation is applied to the replication selection function $r(\cdot)$ and the joint publication and replication function $g(\cdot) \equiv p(\cdot) \times r(\cdot)$.

Lemma B1 (Justification of the common power rule). *Consider a published study (x, σ, θ) . If $x = \theta$ and a replication uses the common power rule to detect the original effect with intended power $1 - \beta$, then*

$$RP(\theta, \theta, \sigma_r(\theta, \beta)) = 1 - \beta \quad (22)$$

Proof. Substitute the common power rule in the replication probability function derived in Lemma A1.1 in Appendix A. If $x = \theta$, then

$$RP(\theta, \theta, \sigma_r(\theta, \beta)) = 1 - \Phi\left(1.96 - \text{sign}(\theta) \frac{\theta}{\sigma_r(\theta, \beta)}\right) = 1 - \Phi\left(1.96 - \frac{\theta}{\theta}(1.96 - \Phi^{-1}(\beta))\right) = 1 - \beta \quad (23)$$

□

Proof of Proposition 1: For notational convenience, let $(X_{sig}, \Sigma_{sig}, \Theta_{sig})$ denote the distribution of latent studies $(X^*, \Sigma^*, \Theta^*)$ conditional on being statistically significant at the 5% level ($|X^*/\Sigma^*| \geq 1.96$), published ($D = 1$), and selected for replication ($R = 1$). The expected replication probability (Definition 2) under the common power rule (Definition 1) is equal to

$$\begin{aligned} & \mathbb{E}_{X^*, \Sigma^*, \Theta^* | D, R, S_X^*} \left[RP(X^*, \Theta^*, \sigma_r(X^*, \beta)) \Big| D = 1, R = 1, |X^*/\Sigma^*| \geq 1.96 \right] \\ &= \mathbb{E}_{X_{sig}, \Sigma_{sig}, \Theta_{sig}} \left[RP(X_{sig}, \Theta_{sig}, \sigma_r(X_{sig}, \Sigma_{sig}, \beta)) \right] \\ &= \mathbb{E}_{\Sigma_{sig}, \Theta_{sig}} \left[\mathbb{E}_{X_{sig} | \Sigma_{sig}, \Theta_{sig}} \left[RP(X_{sig}, \Theta_{sig}, \sigma_r(X_{sig}, \beta)) \Big| \Theta_{sig} = \theta, \Sigma_{sig} = \sigma \right] \right] \end{aligned} \quad (24)$$

where the last equality uses the Law of Iterated Expectations. The proof shows that the conditional expected replication probability satisfies $\mathbb{E}_{X_{sig}|\Sigma_{sig},\Theta_{sig}}[RP(X_{sig}, \Theta_{sig}, \sigma_r(X_{sig}, \beta))|\Theta_{sig} = \theta, \Sigma_{sig} = \sigma] < 1 - \beta$, which implies that the expected replication probability is also less than intended power $1 - \beta$. For greater clarity in what follows, let $\mathbb{E}[RP(X_{sig}|\theta, \sigma, \beta)]$ be shorthand for $\mathbb{E}_{X_{sig}|\Sigma_{sig},\Theta_{sig}}[RP(X_{sig}, \Theta_{sig}, \sigma_r(X_{sig}, \beta))|\Theta_{sig} = \theta, \Sigma_{sig} = \sigma]$.

Note that the conditional expected replication probability can be written explicitly as

$$\mathbb{E}[RP(X_{sig}|\theta, \sigma, \beta)] = \int RP(X_{sig}|\theta, \sigma, \beta) \cdot \frac{g\left(\frac{x}{\sigma}\right)\frac{1}{\sigma}\phi\left(\frac{x-\theta}{\sigma}\right)\mathbb{1}\left(\left|\frac{x}{\sigma}\right| \geq 1.96\right)dx}{\int_{x'} g\left(\frac{x'}{\sigma}\right)\frac{1}{\sigma}\phi\left(\frac{x'-\theta}{\sigma}\right)\mathbb{1}\left(\left|\frac{x'}{\sigma}\right| \geq 1.96\right)dx'} \quad (25)$$

where $g(\cdot) \equiv p(\cdot) \times r(\cdot)$ is the joint selection function for both publication and replication. Note that the density differs from a normal density in two respects: (1) the selection function $g(\cdot)$ reweights the distribution; and (2) conditioning on statistical significance truncates original effects falling in the insignificant region $(-1.96\sigma, 1.96\sigma)$.

Without loss of generality, assume $\theta > 0$. We first introduce some notation. Define $(l^*, u^*) = ((1 - r^*)\theta, (1 + r^*)\theta)$ when $r^* \in (0, 1)$ and $(l^*, u^*) = (0, 2\theta)$ when $r^* \geq 1$, and where r^* depends on the value of β as specified in equation (12). In both cases, the replication probability function is strictly concave over the interval with mid-point θ . Concavity is guaranteed by Lemma A1.3, which states that if $(1 - \beta) > 0.6628$, then $RP(x, |\theta, \sigma, \beta)$ is strictly concave over the open interval $(\max\{0, [1 - r^*]\theta\}, [1 + r^*]\theta)$. This Proposition assumes $(1 - \beta) > 0.8314$, so the condition is satisfied.

Consider first the case where $r^* \geq 1$ so that $(l^*, u^*) = (0, 2\theta)$. The conditional replication probability can be written as

$$\begin{aligned} \mathbb{E}[RP(X_{sig}|\theta, \sigma, \beta)] &= \mathbb{P}(X_{sig} < l^*)\mathbb{E}[RP(X_{sig}|\theta, \sigma, \beta)|X_{sig} < l^*] \\ &+ \mathbb{P}(l^* \leq X_{sig} \leq u^*)\mathbb{E}[RP(X_{sig}|\theta, \sigma, \beta)|l^* \leq X_{sig} \leq u^*] + \mathbb{P}(X_{sig} > u^*)\mathbb{E}[RP(X_{sig}|\theta, \sigma, \beta)|X_{sig} > u^*] \\ &< \mathbb{P}(X_{sig} < l^*)0.025 + \mathbb{P}(l^* \leq X_{sig} \leq u^*)\mathbb{E}[RP(X_{sig}|\theta, \sigma, \beta)|l^* \leq X_{sig} \leq u^*] + \mathbb{P}(X_{sig} > u^*)(1 - \beta) \end{aligned} \quad (26)$$

In the last line, the first term in the sum uses the fact that the maximum value of the replication probability when $x < l^* = 0$ is 0.025 (Lemma A1.2 and Lemma A1.4 in Appendix A). The third term follows because $RP(2\theta|\theta, \sigma, \beta)$ is the maximum value the function takes over $x > u^* = 2\theta$, since the function is strictly decreasing over $x > 0$ (Lemma A1.2); and therefore that $RP(2\theta|\theta, \sigma, \beta) < RP(\theta|\theta, \sigma, \beta) = 1 - \beta$, where the equality is shown in Lemma B1. From equation (26), we can see that $\mathbb{E}[RP(X_{sig}|\theta, \sigma, \beta)|l^* \leq X_{sig} \leq u^*] < 1 - \beta$ is a sufficient condition for $\mathbb{E}[RP(X_{sig}|\theta, \sigma, \beta)] < 1 - \beta$.

Before showing that this sufficient condition is satisfied, we show that the identical sufficient

condition is also applicable in the second case, where $r^* \in (0, 1)$ so that $(l^*, u^*) = ((1-r^*)\theta, (1+r^*)\theta)$. First, we can again express the conditional replication probability as a weighted sum

$$\begin{aligned} \mathbb{E}\left[RP(X_{sig}|\theta, \sigma, \beta)\right] &= \mathbb{P}\left(X_{sig} \leq l^*\right)\mathbb{E}\left[RP(X_{sig}|\theta, \sigma, \beta)\middle|X_{sig} \leq l^*\right] \\ &+ \mathbb{P}\left(l^* \leq X_{sig} \leq u^*\right)\mathbb{E}\left[RP(X_{sig}|\theta, \sigma, \beta)\middle|l^* \leq X_{sig} \leq u^*\right] + \mathbb{P}\left(X_{sig} \geq u^*\right)\mathbb{E}\left[RP(X_{sig}|\theta, \sigma, \beta)\middle|X_{sig} \geq u^*\right] \\ &< \mathbb{P}\left(X_{sig} \leq l^*\right) + \mathbb{P}\left(l^* \leq X_{sig} \leq u^*\right)\mathbb{E}\left[RP(X_{sig}|\theta, \sigma, \beta)\middle|l^* \leq X_{sig} \leq u^*\right] + \mathbb{P}\left(X_{sig} \geq u^*\right)RP(u^*|\theta, \sigma, \beta) \end{aligned} \quad (27)$$

The strict inequality follows for two reasons. For the first term in the sum, one is the maximum value the function can take for any x . For the third term, $RP(u^*|\theta, \sigma, \beta)$ is the function's maximum value over $x \geq u^*$, since the integrand is strictly decreasing over positive values (Lemma A1.2). With an additional step, we can write this inequality as

$$\begin{aligned} \mathbb{E}\left[RP(X_{sig}|\theta, \sigma, \beta)\right] &< \frac{1}{2}\left(1 - \mathbb{P}\left(l^* \leq X_{sig} \leq u^*\right)\right)\left(1 + RP(u^*|\theta, \sigma, \beta)\right) \\ &+ \mathbb{P}\left(l^* \leq X_{sig} \leq u^*\right)\mathbb{E}\left[RP(X_{sig}|\theta, \sigma, \beta)\middle|l^* \leq X_{sig} \leq u^*\right] \end{aligned} \quad (28)$$

This follows because $\mathbb{P}(X_{sig} \leq l^*) \leq \mathbb{P}(X_{sig} \geq u^*)$ and $RP(u^*|\theta, \sigma, \beta) < 1$. That is, increasing the relative weight on the maximum value of one, such that both tails are equally weighted, must lead to a (weakly) larger value. The weak inequality $\mathbb{P}(X_{sig} \leq l^*) \leq \mathbb{P}(X_{sig} \geq u^*)$ required for this simplification is shown below:

Lemma B2. *Suppose $X|\theta, \sigma$ follows the truncated normal pdf in equation (25). Then for any $r^* \in (0, 1)$, the following inequality holds: $\mathbb{P}(X_{sig} \leq (1-r^*)\theta) < \mathbb{P}(X_{sig} \geq (1+r^*)\theta)$.*

Proof. First, note that $((1-r^*)\theta, (1+r^*)\theta)$ is an interval over the positive real line centered at θ . Consider two cases:

Case 1: Let $(1-r^*)\theta \leq 1.96\sigma$. Define the normalization constant $C = \int_{x'} g\left(\frac{x'}{\sigma}\right)\frac{1}{\sigma}\phi\left(\frac{x'-\theta}{\sigma}\right)\mathbb{1}\left(\left|\frac{x'}{\sigma}\right| \geq 1.96\right)dx'$. Then

$$\begin{aligned} \mathbb{P}\left(X_{sig} \leq (1-r^*)\theta\right) &= \frac{1}{C} \int_{-\infty}^{-1.96\sigma} g_{sig}\left(\frac{x}{\sigma}\right)\frac{1}{\sigma}\phi\left(\frac{x-\theta}{\sigma}\right)dx' \leq \frac{1}{C} \int_{2\theta+1.96\sigma}^{\infty} g_{sig}\left(\frac{x}{\sigma}\right)\frac{1}{\sigma}\phi\left(\frac{x-\theta}{\sigma}\right)dx' \\ &< \frac{1}{C} \int_{2\theta+1.96\sigma}^{\infty} g_{sig}\left(\frac{x}{\sigma}\right)\frac{1}{\sigma}\phi\left(\frac{x-\theta}{\sigma}\right)dx' + \frac{1}{C} \int_{\max\{1.96\sigma, (1+r^*)\theta\}}^{2\theta+1.96\sigma} g_{sig}\left(\frac{x}{\sigma}\right)\frac{1}{\sigma}\phi\left(\frac{x-\theta}{\sigma}\right)dx' = \mathbb{P}\left(X_{sig} \geq (1+r^*)\theta\right) \end{aligned} \quad (29)$$

Consider the weak inequality. Note that the mid-point between -1.96σ and $2\theta + 1.96\sigma$ is θ . Thus, with no selective publication (i.e. $p(t) = 1$ for all t), we would have equality owing to the symmetry of the normal distribution. However, recall that $g_{sig}(\cdot)$ is symmetric about zero and weakly increasing in absolute value (Assumptions 1 and 2). It follows that $|2\theta + 1.96\sigma| > |-1.96\sigma|$ implies $g_{sig}(|2\theta + 1.96\sigma|) \geq g_{sig}(|-1.96\sigma|)$; using this fact and symmetry of the normal distribution about θ gives the weak inequality. The strict inequality follows because the additional term is strictly positive, since $g_{sig}(\cdot)$ is assumed to be non-zero.

Case 2: Let $(1 - r^*)\theta > 1.96\sigma$. The argument is similar to the first case:

$$\begin{aligned} \mathbb{P}\left(X_{sig} \leq (1 - r^*)\theta\right) &= \frac{1}{C} \int_{-\infty}^{-1.96\sigma} g_{sig}\left(\frac{x}{\sigma}\right) \frac{1}{\sigma} \phi\left(\frac{x - \theta}{\sigma}\right) dx' + \frac{1}{C} \int_{1.96\sigma}^{(1-r^*)\theta} g_{sig}\left(\frac{x}{\sigma}\right) \frac{1}{\sigma} \phi\left(\frac{x - \theta}{\sigma}\right) dx' \\ &< \frac{1}{C} \int_{2\theta+1.96\sigma}^{\infty} g_{sig}\left(\frac{x}{\sigma}\right) \frac{1}{\sigma} \phi\left(\frac{x - \theta}{\sigma}\right) dx' + \frac{1}{C} \int_{(1+r^*)\theta}^{2\theta-1.96\sigma} g_{sig}\left(\frac{x}{\sigma}\right) \frac{1}{\sigma} \phi\left(\frac{x - \theta}{\sigma}\right) dx' \\ &\quad + \frac{1}{C} \int_{2\theta-1.96\sigma}^{2\theta+1.96\sigma} g_{sig}\left(\frac{x}{\sigma}\right) \frac{1}{\sigma} \phi\left(\frac{x - \theta}{\sigma}\right) dx' = \mathbb{P}\left(X_{sig} \geq (1 + r^*)\theta\right) \end{aligned} \quad (30)$$

□

The inequality in equation (28) can be further simplified by placing restrictions on intended power. In particular, if intended power satisfies $1 - \beta \geq 0.8314$, then

$$\begin{aligned} \mathbb{E}\left[RP(X_{sig}|\theta, \sigma, \beta)\right] &< \left(1 - \mathbb{P}(l^* \leq X_{sig} \leq u^*)\right)(1 - \beta) \\ &\quad + \mathbb{P}(l^* \leq X_{sig} \leq u^*) \mathbb{E}\left[RP(X_{sig}|\theta, \sigma, \beta) \mid l^* \leq X_{sig} \leq u^*\right] \end{aligned} \quad (31)$$

This follows because with $u^* = (1 + r^*)\theta$, we have

$$\begin{aligned} \frac{1}{2}\left(1 + RP(u^*|\theta, \sigma, \beta)\right) &= \frac{1}{2}\left(1 + \left(1 - \Phi\left(1.96 - \frac{1.96 - \Phi^{-1}(\beta)}{1 + r^*(\beta)}\right)\right)\right) \\ &\leq 1 - \beta \iff 1 - \beta \geq 0.8314 \end{aligned} \quad (32)$$

From equation (31), we can see that $\mathbb{E}[RP(X_{sig}|\theta, \sigma, \beta) \mid l^* \leq X_{sig} \leq u^*] < 1 - \beta$ is a sufficient condition for $\mathbb{E}[RP(X_{sig}|\theta, \sigma, \beta)] < 1 - \beta$. Thus, in both cases $-r^* \geq 1$ and $r^* \in (0, 1)$ – the sufficient condition for the desired result is the same.

This sufficient condition is shown in two steps. In the first, I show that this inequality holds even in the case where there is no selection such that all published results are replicated, and thus $X \sim N(\Theta, \Sigma^2)$. In the second, I show that this inequality remains true once we allow for selection and truncation of the distribution due to conditioning on statistical significance.

Lemma B3 states the first step and is of independent interest. It shows that even in the optimistic scenario where original estimates are unbiased, there is no selective publication, and all results are published and replicated, that the expected replication probability still falls below intended power.

Lemma B3. *Let published effects be distributed according to $X|\theta, \sigma \sim N(\theta, \sigma^2)$. Suppose $p(t) = 1$ and $r(t) = 1$ for all $t \in \mathbb{R}$. Assume all results are included in the replication rate calculation. Let power in replications is set according to the common power rule with intended power $1 - \beta \geq 0.8314$. Then $\mathbb{E}[RP(X|\theta, \sigma, \beta)] < 1 - \beta$.*

Proof. Recall that $RP(x|\theta, \sigma, \beta)$ is strictly concave with respect to x over the interval (l^*, u^*) , where $(l^*, u^*) = ((1 - r^*)\theta, (1 + r^*)\theta)$ when $r^* \in (0, 1)$ and $(l^*, u^*) = (0, 2\theta)$; in both cases, the mid-point of the interval is θ . We have that

$$\mathbb{E}[RP(X|\theta, \sigma, \beta)|l^* \leq X \leq u^*] = RP(\mathbb{E}[X|l^* \leq X \leq u^*]|\theta, \sigma, \beta) < RP(\theta|\theta, \sigma, \beta) = 1 - \beta \quad (33)$$

where the strict inequality follows from Jensen's inequality and the fact that $\mathbb{E}[X|l^* \leq X \leq u^*] = \theta$. The final equality is a property of the replication probability function shown in Lemma B1. This is the sufficient condition required for the desired result.

It follows that $\mathbb{E}[RP(X|\theta, \sigma, \beta)] < 1 - \beta$ due to the inequalities in equation (28) (for when $r^* \geq 1$) and equation (31) (for when $r^* \in (0, 1)$), which were derived under more general conditions. Specifically, these inequalities were derived assuming that normal distribution may be reweighted by $g(\cdot)$ and truncated based on significance. This setting is a special case with no selective publication (i.e. $g(t) = 1$ for all t) and no truncation. \square

The same conclusions hold when we introduce selective publication and replication (which reweights the normal distribution) and condition on statistical significance (which truncates the 'insignificant' regions of the density). Consider three cases. First, suppose that $u^* \leq 1.96\sigma$. Then $\mathbb{E}(RP(X_{sig}|\theta, \sigma, \beta)|l^* \leq X_{sig} \leq u^*) = 0 < 1 - \beta$ because of truncation. Second, suppose that $l^* \geq 1.96\sigma$. Then

$$\begin{aligned} \mathbb{E}\left[RP(X_{sig}|\theta, \sigma, \beta)|l^* \leq X_{sig} \leq u^*\right] &= \int_{l^*}^{u^*} RP(x|\theta, \sigma, \beta) \frac{g_{sig}\left(\frac{x}{\sigma}\right) \frac{1}{\sigma} \phi\left(\frac{x-\theta}{\sigma}\right) dx}{\int_{l^*}^{u^*} g_{sig}\left(\frac{x}{\sigma}\right) \frac{1}{\sigma} \phi\left(\frac{x'-\theta}{\sigma}\right) dx'} \\ &\leq \int_{l^*}^{u^*} RP(x|\theta, \sigma, \beta) \frac{\frac{1}{\sigma} \phi\left(\frac{x-\theta}{\sigma}\right) dx}{\int_{l^*}^{u^*} \frac{1}{\sigma} \phi\left(\frac{x'-\theta}{\sigma}\right) dx'} < RP(\theta|\theta, \sigma, \beta) = 1 - \beta \end{aligned} \quad (34)$$

Note that the distribution is invariant to the scale of $g_{sig}(\cdot)$. Consider first the weak inequality. This follows because $g_{sig}(\cdot)$ is assumed to be weakly increasing over (l^*, u^*) . When it is a

constant function over the interval, the equality holds. If $g_{sig}(x/\sigma) > 0$ for some $x \in (l^*, u^*)$ then the function redistributes weight to larger values of x . Since $RP(x|\theta, \sigma, \beta)$ is strictly decreasing over positive values of x (Lemma A1.2), placing higher relative weight on lower values implies that the weak inequality becomes strict. The strict inequality follows from Jensen's inequality, as shown in Lemma B3.

Finally, consider the case where $l^* < 1.96\sigma < u^*$. Then

$$\begin{aligned}
\mathbb{E}\left[RP(X_{sig}|\theta, \sigma, \beta) \mid l^* \leq X_{sig} \leq u^*\right] &= \int_{1.96\sigma}^{u^*} RP(x|\theta, \sigma, \beta) \frac{g_{sig}\left(\frac{x}{\sigma}\right) \frac{1}{\sigma} \phi\left(\frac{x-\theta}{\sigma}\right) dx}{\int_{1.96\sigma}^{u^*} g_{sig}\left(\frac{x'}{\sigma}\right) \frac{1}{\sigma} \phi\left(\frac{x'-\theta}{\sigma}\right) dx'} \\
&= \int_{1.96\sigma}^{2\theta-1.96\sigma} RP(x|\theta, \sigma, \beta) \frac{g_{sig}\left(\frac{x}{\sigma}\right) \frac{1}{\sigma} \phi\left(\frac{x-\theta}{\sigma}\right) dx}{\int_{1.96\sigma}^{u^*} g_{sig}\left(\frac{x'}{\sigma}\right) \frac{1}{\sigma} \phi\left(\frac{x'-\theta}{\sigma}\right) dx'} + \int_{2\theta-1.96\sigma}^{u^*} RP(x|\theta, \sigma, \beta) \frac{g_{sig}\left(\frac{x}{\sigma}\right) \frac{1}{\sigma} \phi\left(\frac{x-\theta}{\sigma}\right) dx}{\int_{1.96\sigma}^{u^*} g_{sig}\left(\frac{x'}{\sigma}\right) \frac{1}{\sigma} \phi\left(\frac{x'-\theta}{\sigma}\right) dx'} \\
&= \omega \int_{1.96\sigma}^{2\theta-1.96\sigma} RP(x|\theta, \sigma, \beta) \frac{g_{sig}\left(\frac{x}{\sigma}\right) \frac{1}{\sigma} \phi\left(\frac{x-\theta}{\sigma}\right) dx}{\int_{1.96\sigma}^{2\theta-1.96\sigma} g_{sig}\left(\frac{x'}{\sigma}\right) \frac{1}{\sigma} \phi\left(\frac{x'-\theta}{\sigma}\right) dx'} + (1-\omega) \int_{2\theta-1.96\sigma}^{u^*} RP(x|\theta, \sigma, \beta) \frac{g_{sig}\left(\frac{x}{\sigma}\right) \frac{1}{\sigma} \phi\left(\frac{x-\theta}{\sigma}\right) dx}{\int_{2\theta-1.96\sigma}^{u^*} g_{sig}\left(\frac{x'}{\sigma}\right) \frac{1}{\sigma} \phi\left(\frac{x'-\theta}{\sigma}\right) dx'} \\
&= \omega \int_{1.96\sigma}^{2\theta-1.96\sigma} RP(x|\theta, \sigma, \beta) \frac{\frac{1}{\sigma} \phi\left(\frac{x-\theta}{\sigma}\right) dx}{\int_{1.96\sigma}^{2\theta-1.96\sigma} \frac{1}{\sigma} \phi\left(\frac{x'-\theta}{\sigma}\right) dx'} + (1-\omega) \int_{2\theta-1.96\sigma}^{u^*} RP(x|\theta, \sigma, \beta) \frac{\frac{1}{\sigma} \phi\left(\frac{x-\theta}{\sigma}\right) dx}{\int_{2\theta-1.96\sigma}^{u^*} \frac{1}{\sigma} \phi\left(\frac{x'-\theta}{\sigma}\right) dx'} \\
&< \omega RP(\theta|\theta, \sigma, \beta) + (1-\omega) \cdot RP(2\theta-1.96\sigma|\theta, \sigma, \beta) < 1-\beta \tag{35}
\end{aligned}$$

with

$$\omega = \frac{\int_{1.96\sigma}^{2\theta-1.96\sigma} g_{sig}\left(\frac{x'}{\sigma}\right) \frac{1}{\sigma} \phi\left(\frac{x'-\theta}{\sigma}\right) dx'}{\int_{1.96\sigma}^{u^*} g_{sig}\left(\frac{x'}{\sigma}\right) \frac{1}{\sigma} \phi\left(\frac{x'-\theta}{\sigma}\right) dx'} \tag{36}$$

The second row simply breaks up the integral. The third row rearranges the sum so that the conditional expectation of the replication probability appears in both terms. The third line follows because, as in the previous case, the g_{sig} function redistributes weight to large values of x and hence lower values of $RP(x|\theta, \sigma, \beta)$. In the last line, the first term uses the concavity of $RP(x|\theta, \sigma, \beta)$ over $(1.96\sigma, 2\theta-1.96\sigma) \subset (l^*, u^*)$, Jensen's inequality, and the fact that the expected value of X over this interval is equal to θ . The second term follows because $2\theta-1.96\sigma$ is the maximum value the function can take because $RP(x|\theta, \sigma, \beta)$ is strictly decreasing in x over positive values. The final inequality follows because $RP(\theta|\theta, \sigma, \beta) = 1-\beta$ (Lemma B1) and $RP(2\theta-1.96\sigma|\theta, \sigma, \beta) < 1-\beta$ because $2\theta-1.96\sigma > \theta$ and the function is strictly decreasing over positive values.

This covers all cases, proving the proposition.

Proposition B2 (Regression to the Mean in Replications). *Under Assumption 1, when $\theta > 0$, we have that $\mathbb{E}[X^*|\Theta^* = \theta, S_X = 1, D = 1] > \theta = \mathbb{E}[X_r^*|\Theta^* = \theta, D = 1]$. When $\theta < 0$, the opposite inequality holds.*

Proof. Without loss of generality, assume $\theta > 0$. We have $\mathbb{E}[X_r^*|\Theta = \theta, D = 1] = \theta$ by assumption. Next, note that

$$\begin{aligned} \mathbb{E}_{X^*|\Theta^*, S_X^*, D} \left(X^*|\Theta^* = \theta, |X^*/\Sigma^*| \geq 1.96, D = 1 \right) &= \mathbb{E}_{X|\Theta, S_X} \left(X|\Theta = \theta, |X/\Sigma| \geq 1.96 \right) \\ &= \mathbb{E}_{\Sigma|\Theta, S_X} \left(\mathbb{E}_{X|\Theta, \Sigma, S_X} \left(X|\Theta = \theta, \Sigma = \sigma, |X/\sigma| \geq 1.96 \right) \right) \end{aligned} \quad (37)$$

where the second line uses the fact that X is defined as $X^*|D = 1$ and the last line uses the Law of Iterated Expectations. We will prove $\mathbb{E}_{X|\Theta, \Sigma, S_X^*} (X|\Theta = \theta, \Sigma = \sigma, |X/\sigma| \geq 1.96) > \theta$, which implies that the expression in equation (37) is also greater than θ . Recall that $X|\theta, \sigma$ is the effect size of published studies and follows a truncated normal distribution:

$$\frac{p\left(\frac{x}{\sigma}\right)\frac{1}{\sigma}\phi\left(\frac{x-\theta}{\sigma}\right)\mathbb{1}\left(\left|\frac{x}{\sigma}\right| \geq 1.96\right)}{\int p\left(\frac{x'}{\sigma}\right)\frac{1}{\sigma}\phi\left(\frac{x'-\theta}{\sigma}\right)\mathbb{1}\left(\left|\frac{x'}{\sigma}\right| \geq 1.96\right)dx'} \quad (38)$$

Define $X = \theta + \sigma Z$. Then the density for the transformed random variable Z is

$$\frac{p\left(z + \frac{\theta}{\sigma}\right)\phi(z)\mathbb{1}\left(\left|z + \frac{\theta}{\sigma}\right| \geq 1.96\right)}{\int p\left(z' + \frac{\theta}{\sigma}\right)\phi(z')\mathbb{1}\left(\left|z' + \frac{\theta}{\sigma}\right| \geq 1.96\right)dz'} \quad (39)$$

For notational convenience, define the following normalization constants:

$$\bar{\eta} = \mathbb{P}(X \leq -1.96\sigma) + \mathbb{P}(X \geq 1.96\sigma) = \mathbb{P}\left(Z \leq -1.96 - \frac{\theta}{\sigma}\right) + \mathbb{P}\left(Z \geq 1.96 - \frac{\theta}{\sigma}\right) \quad (40)$$

$$\eta_1 = \mathbb{P}(X \leq -1.96\sigma) = \mathbb{P}\left(Z \leq -1.96 - \frac{\theta}{\sigma}\right) \quad (41)$$

$$\eta_2 = \mathbb{P}(X \geq 2\theta + 1.96\sigma) = \mathbb{P}\left(Z \geq \frac{\theta}{\sigma} + 1.96\right) \quad (42)$$

$$\eta_3 = \mathbb{P}(1.96\sigma \leq X \leq 2\theta - 1.96\sigma) = \mathbb{P}\left(1.96 - \frac{\theta}{\sigma} \leq Z \leq \frac{\theta}{\sigma} - 1.96\right) \quad (43)$$

Case 1.

Consider two cases. First, suppose $\theta \in (0, 1.96\sigma)$. Conditional on (θ, σ) (where we suppress the conditional notation on (θ, σ) for clarity), the expected value of a published estimate conditional of statistical significance is

$$\begin{aligned} \mathbb{E}(X|1.96\sigma \leq |X|) &= \frac{1}{\bar{\eta}} \left(\eta_1 \mathbb{E}(X|X \leq -1.96\sigma) + \eta_2 \mathbb{E}(X|X \geq 2\theta + 1.96\sigma) \right. \\ &\quad \left. + (\bar{\eta} - \eta_1 - \eta_2) \mathbb{E}(X|1.96\sigma \leq X \leq 2\theta + 1.96\sigma) \right) \end{aligned} \quad (44)$$

First note that $\mathbb{E}(X|1.96\sigma \leq X \leq 2\theta + 1.96\sigma) > \theta$ since we assume that $\theta \in (0, 1.96\sigma)$ and $p_{sig}(\cdot) > 0$. If $\eta_1 \mathbb{E}(X|X \leq -1.96\sigma) + \eta_2 \mathbb{E}(X|X \geq 2\theta + 1.96\sigma) \geq (\eta_1 + \eta_2)\theta$, it follows that $\mathbb{E}(X|1.96\sigma \leq |X|) > \theta$, which is what we want to show. Consider the first expectation in this expression:

$$\mathbb{E}(X|X \leq -1.96\sigma) = \mathbb{E}\left(\theta + \sigma Z|Z \leq -1.96 - \frac{\theta}{\sigma}\right) = \theta + \sigma \mathbb{E}\left(Z|Z \leq -1.96 - \frac{\theta}{\sigma}\right) \quad (45)$$

Evaluating the expectation in the right-hand-side of equation (45) gives

$$\begin{aligned} \mathbb{E}\left(Z|Z \leq -1.96 - \frac{\theta}{\sigma}\right) &= \frac{1}{\eta_1} \int_{-\infty}^{-1.96 - \frac{\theta}{\sigma}} z p_{sig}\left(z + \frac{\theta}{\sigma}\right) \phi(z) dz = -\frac{1}{\eta_1} \int_{-\infty}^{-1.96 - \frac{\theta}{\sigma}} p_{sig}\left(z + \frac{\theta}{\sigma}\right) \phi'(z) dz \\ &= -\frac{1}{\eta_1} \left[p_{sig}(-1.96) \phi\left(-1.96 - \frac{\theta}{\sigma}\right) - p_{sig}(-\infty) \phi(-\infty) - \int_{-\infty}^{-1.96 - \frac{\theta}{\sigma}} p'_{sig}\left(z + \frac{\theta}{\sigma}\right) \phi(z) dz \right] \\ &= -\frac{1}{\eta_1} p_{sig}(-1.96) \phi\left(-1.96 - \frac{\theta}{\sigma}\right) + \frac{1}{\eta_1} \int_{-\infty}^{-1.96 - \frac{\theta}{\sigma}} p'_{sig}\left(z + \frac{\theta}{\sigma}\right) \phi(z) dz \end{aligned} \quad (46)$$

where the second equality uses $\phi'(z) = -z\phi(z)$; the third equality uses integration by parts; and the final equality follows because $p_{sig}(-\infty)\phi(-\infty) = 0$ since $p_{sig}(\cdot)$ is bounded between zero and one. Substituting this into equation (45) gives

$$\mathbb{E}(X|X \leq -1.96\sigma) = \theta - \frac{\sigma}{\eta_1} p_{sig}(-1.96) \phi\left(-1.96 - \frac{\theta}{\sigma}\right) + \frac{\sigma}{\eta_1} \int_{-\infty}^{-1.96 - \frac{\theta}{\sigma}} p'_{sig}\left(z + \frac{\theta}{\sigma}\right) \phi(z) dz \quad (47)$$

Next, note that

$$\mathbb{E}(X|X \geq 2\theta + 1.96\sigma) = \theta + \sigma \mathbb{E}\left(Z|Z \leq \frac{\theta}{\sigma} + 1.96\right) \quad (48)$$

where

$$\mathbb{E}\left(Z|Z \leq \frac{\theta}{\sigma} + 1.96\right) = \frac{1}{\eta_2} \int_{1.96 + \frac{\theta}{\sigma}}^{\infty} z p_{sig}\left(z + \frac{\theta}{\sigma}\right) \phi(z) dz \geq \frac{1}{\eta_2} \int_{1.96 + \frac{\theta}{\sigma}}^{\infty} z p_{sig}\left(z - \frac{\theta}{\sigma}\right) \phi(z) dz \quad (49)$$

since $p_{sig}(z + \theta/\sigma) \geq p_{sig}(z - \theta/\sigma)$ for all $z \in (1.96 + \theta/\sigma, \infty)$ because $p_{sig}(t)$ is weakly increasing over $t > 1.96$. For the right-hand-side of this equation, we can apply similar arguments used to derive equation (46). Substituting the result into equation (48) gives

$$\mathbb{E}(X|X \geq 2\theta + 1.96\sigma) \geq \theta + \frac{\sigma}{\eta_2} p_{sig}(1.96) \phi\left(1.96 + \frac{\theta}{\sigma}\right) + \frac{\sigma}{\eta_2} \int_{1.96 + \frac{\theta}{\sigma}}^{\infty} p'_{sig}\left(z - \frac{\theta}{\sigma}\right) \phi(z) dz \quad (50)$$

Equations (47) and (50) imply

$$\begin{aligned} & \eta_1 \mathbb{E}(X|X \leq -1.96\sigma) + \eta_2 \mathbb{E}(X|X \geq 2\theta + 1.96\sigma) \\ & \geq (\eta_1 + \eta_2)\theta + \sigma \left[p_{sig}(1.96) \phi\left(1.96 + \frac{\theta}{\sigma}\right) - p_{sig}(-1.96) \phi\left(-1.96 - \frac{\theta}{\sigma}\right) \right] \\ & + \sigma \left[\int_{-\infty}^{-1.96 - \frac{\theta}{\sigma}} p'_{sig}\left(z + \frac{\theta}{\sigma}\right) \phi(z) dz + \int_{1.96 + \frac{\theta}{\sigma}}^{\infty} p'_{sig}\left(z - \frac{\theta}{\sigma}\right) \phi(z) dz \right] = (\eta_1 + \eta_2)\theta \quad (51) \end{aligned}$$

In the second line, the second term in the sum equals zero because symmetry of $p_{sig}()$ and $\phi()$ about zero implies that both terms in the brackets are equal. To see why the third term in the sum equals zero, note that

$$\int_{-\infty}^{-1.96 - \frac{\theta}{\sigma}} p'_{sig}\left(z + \frac{\theta}{\sigma}\right) \phi(z) dz = \int_{1.96 + \frac{\theta}{\sigma}}^{\infty} p'_{sig}\left(-u + \frac{\theta}{\sigma}\right) \phi(u) du = - \int_{1.96 + \frac{\theta}{\sigma}}^{\infty} p'_{sig}\left(u - \frac{\theta}{\sigma}\right) \phi(u) du \quad (52)$$

The first equality follows from both changing the order of the integral limits and applying the substitution $u = -x$; it also uses the symmetry of $\phi()$. The final equality holds because symmetry of $p_{sig}()$ about zero implies that for any $t > 1.96$, $p'_{sig}(t) = -p'_{sig}(-t)$.

Case 2.

Consider the second case where $\theta \geq 1.96\sigma$. For a given (θ, σ) , we have

$$\mathbb{E}(X|1.96\sigma \leq |X|) = \frac{1}{\bar{\eta}} \left(\eta_1 \mathbb{E}(X|X \leq -1.96\sigma) + \eta_2 \mathbb{E}(X|X \geq 2\theta + 1.96\sigma) \right)$$

$$\begin{aligned} & \left. \eta_3 \mathbb{E}(X|1.96\sigma \leq X \leq 2\theta - 1.96\sigma) + (\bar{\eta} - \eta_1 - \eta_2 - \eta_3) \mathbb{E}(X|2\theta - 1.96\sigma \leq X \leq 2\theta + 1.96\sigma) \right) \\ & > \frac{1}{\bar{\eta}} \left(\theta(\eta_1 + \eta_2) + (\bar{\eta} - \eta_1 - \eta_2 - \eta_3)\theta + \eta_3 \mathbb{E}(X|1.96\sigma \leq X \leq 2\theta - 1.96\sigma) \right) \end{aligned} \quad (53)$$

The inequality follows from two facts. First, the inequality proved in the first case: $\eta_1 \mathbb{E}(X|X \leq -1.96\sigma) + \eta_2 \mathbb{E}(X|X \geq 2\theta + 1.96\sigma) \geq (\eta_1 + \eta_2)\theta$. Second, the expectation in the third term of the sum satisfies $\mathbb{E}(X|2\theta - 1.96\sigma \leq X \leq 2\theta + 1.96\sigma) > \theta$ because $\theta \geq 1.96\sigma \iff 2\theta - 1.96\sigma \geq \theta$ and we assume that $p_{sig}() > 0$.

It remains to show that $\mathbb{E}(X|1.96\sigma \leq X \leq 2\theta - 1.96\sigma) \geq \theta$. Then it follows that $\mathbb{E}(X|1.96\sigma \leq |X|) > \theta$, which is what we want to show. First, note that

$$\mathbb{E}(X|1.96\sigma \leq X \leq 2\theta - 1.96\sigma) = \theta + \sigma \mathbb{E}\left(Z \middle| 1.96 - \frac{\theta}{\sigma} \leq Z \leq -1.96 + \frac{\theta}{\sigma}\right) \quad (54)$$

It is therefore sufficient to show that $\mathbb{E}\left(Z \middle| 1.96 - \frac{\theta}{\sigma} \leq Z \leq -1.96 + \frac{\theta}{\sigma}\right) \geq 0$. Writing out the expectation in full gives

$$\begin{aligned} \mathbb{E}\left(Z \middle| 1.96 - \frac{\theta}{\sigma} \leq Z \leq -1.96 + \frac{\theta}{\sigma}\right) &= \frac{1}{\eta_3} \left(\int_{1.96 - \frac{\theta}{\sigma}}^0 z p_{sig}\left(z + \frac{\theta}{\sigma}\right) \phi(z) dz + \int_0^{\frac{\theta}{\sigma} - 1.96} z p_{sig}\left(z + \frac{\theta}{\sigma}\right) \phi(z) dz \right) \\ &= \frac{1}{\eta_3} \left(\int_0^{\frac{\theta}{\sigma} - 1.96} z \left[p_{sig}\left(z + \frac{\theta}{\sigma}\right) - p_{sig}\left(-z + \frac{\theta}{\sigma}\right) \right] \phi(z) dz \right) \geq 0 \end{aligned} \quad (55)$$

The second equality follows because

$$\int_{1.96 - \frac{\theta}{\sigma}}^0 z p_{sig}\left(z + \frac{\theta}{\sigma}\right) \phi(z) dz = - \int_0^{1.96 - \frac{\theta}{\sigma}} z p_{sig}\left(z + \frac{\theta}{\sigma}\right) \phi(z) dz = - \int_0^{\frac{\theta}{\sigma} - 1.96} u p_{sig}\left(-u + \frac{\theta}{\sigma}\right) \phi(u) du \quad (56)$$

which uses the substitution $u = -x$ and the symmetry of $\phi()$. The weak inequality in equation (55) follows because $p_{sig}()$ is assumed to be weakly increasing over positive values. Thus, $z - \theta/\sigma > -z + \theta/\sigma$ for all $z \in (0, \theta/\sigma - 1.96)$ implies $p_{sig}(z + \theta/\sigma) - p_{sig}(-z + \theta/\sigma) \geq 0$.

This covers all cases and proves the proposition. \square

Proposition B3. *Under the fractional power rule which sets the replication standard error according to $\sigma_r(X, \beta, \psi) = \frac{\psi \cdot |X|}{1.96 - \Phi^{-1}(\beta)}$ with $\psi < 1$, the expected replication rate can range between 0.025 and $1 - \Phi\left[1.96 - \frac{1}{\psi}(1.96 - \Phi^{-1}(\beta))\right] > 1 - \beta$.*

Proof. Under the fractional power rule, the expected replication rate conditional on (θ, σ) is given by

$$\begin{aligned} & \mathbb{E}[RP(X, \Theta, \sigma_r(X, \beta, \psi)|\Theta = \theta, \Sigma = \sigma)] \\ &= \int \left[1 - \Phi\left(1.96 - \text{sign}(x) \frac{\theta}{\psi \cdot |x|} (1.96 - \Phi^{-1}(\beta))\right) \right] \frac{1}{\sigma} \phi\left(\frac{x - \theta}{\sigma}\right) dx \end{aligned} \quad (57)$$

If $\theta = 0$, then this equals 0.025. Next, suppose wlog that $\theta > 0$ and consider the case where $\sigma \rightarrow 0$ such that power in original studies approaches one. See that the integrand is bounded above by one and converges pointwise as $\sigma \rightarrow 0$ to

$$1 - \Phi\left(1.96 - \text{sign}(x) \frac{\theta}{\psi \cdot |x|} (1.96 - \Phi^{-1}(\beta))\right) \mathbb{1}\{x = \theta\} \quad (58)$$

since the normal distribution converges to a degenerate distribution when the variance goes to zero. Thus, by the dominated convergence theorem (and the fact that $\theta > 0$), we have that

$$\lim_{\sigma \rightarrow 0} \mathbb{E}[RP(X, \Theta, \sigma_r(X, \beta, \psi)|\Theta = \theta, \Sigma = \sigma)] = 1 - \Phi\left(1.96 - \frac{1}{\psi} (1.96 - \Phi^{-1}(\beta))\right) \quad (59)$$

When $\psi = 1$, this equals $1 - \beta$. Since equation (59) is strictly decreasing in ψ , it follows that equation (59) is strictly above $1 - \beta$ when $\psi < 1$.

This shows that the expected replication of an individual study can range between 0.025 and $1 - \Phi[1.96 - \frac{1}{\psi} (1.96 - \Phi^{-1}(\beta))] > 1 - \beta$. Integrating over the distribution of latent studies gives the desired result. \square

Proposition B4. *For any function $g(X, \Sigma, X_r, \beta)$, $\mathbb{E}[g(X, \Sigma, X_r, \beta)|D = 1, R = 1, S_X = 1]$ does not depend on $p_{insig}()$.*

Proof. We can write $\mathbb{E}[g(X, \Sigma, X_r, \beta)|D = 1, R = 1, S_X = 1]$ as

$$\begin{aligned} & \int g(x, \sigma, x_r, \beta) f_{X^*, \Sigma^*, \Theta^*, X_r | D, R, S_X^*}(x, \sigma, \theta, x_r | D = 1, R = 1, S_X^* = 1) dx d\sigma d\theta dx_r \\ &= \int_{x, \sigma, \theta} \left(\int_{x_r} g(x, \sigma, x_r, \beta) f_{X_r | X^*, \Sigma^*, \Theta^*}(x_r | \theta, \sigma_r(x, \sigma, \beta)) dx_r \right) f_{X^*, \Sigma^*, \Theta^* | D, R, S_X^*}(x, \sigma, \theta | D = 1, R = 1, S_X^* = 1) dx d\sigma d\theta \end{aligned} \quad (60)$$

The equality uses the Law of Iterated Expectations and $f_{X_r | X^*, \Sigma^*, \Theta^*, D, R, S_X^*}(x_r | \theta, \sigma_r(x, \sigma, \beta)) = f_{X_r | X^*, \Sigma^*, \Theta^*}(x_r | \theta, \sigma_r(x, \sigma, \beta))$. Replication estimates are not subject to selective publication, which implies this is a normal density that

does not depend on $p()$. Hence, the term in parentheses can only be affected by $p()$ indirectly through $f_{X^*, \Sigma^*, \Theta^* | D, R, S_X^*}$, which is the joint distribution of original studies conditional on being published, chosen for replication, and statistically significant at the 5% level. However, this distribution does not depend on the probability of publishing insignificant findings. To see this, apply Bayes rule twice to get

$$\begin{aligned}
& f_{X^*, \Sigma^*, \Theta^* | D, R, S_X^*}(x, \sigma, \theta | D = 1, R = 1, S_X^* = 1) \\
= & \frac{\mathbb{P}(D = 1 | X^* = x, \Sigma^* = \sigma, \Theta^* = \theta, R = 1, S_X^* = 1)}{\mathbb{P}(D = 1 | R = 1, S_X^* = 1)} \times \frac{\mathbb{P}(R = 1 | X^* = x, \Sigma^* = \sigma, \Theta^* = \theta, S_X^* = 1)}{\mathbb{P}(R = 1 | S_X^* = 1)} \\
& \times f_{X^*, \Theta, \Sigma^* | S_X^*}(x, \theta, \sigma | S_X^* = 1) \\
= & \frac{p_{sig}(x/\sigma)}{\mathbb{E}(p_{sig}(X^*/\Sigma^*) | S_X^* = 1)} \cdot \frac{r_{sig}(x/\sigma)}{\mathbb{E}(r_{sig}(X^*/\Sigma^*) | S_X^* = 1)} \cdot f_{X^*, \Sigma^*, \Theta^* | S_X^*}(\theta, x, \sigma | S_X^* = 1) \quad (61)
\end{aligned}$$

In the final line, the first factor in the product includes only $p_{sig}()$; the denominator does not condition on R because replication selection is assumed to be random for significant findings. The second factor equals one because replication selection for significant results is assumed to be random. The final factor in the product is the density of latent studies conditional on significance, which is not affected by selective publication. \square

C. Replication Rate Gap Decomposition

This appendix derives a decomposition that measures the relative importance of non-linearities as compared to distortions from selection on significance and issues that arise when original estimates have. It then applies this decomposition to economics experiments and psychology.

C.1. Derivation

The decomposition is based on the expected replication rate under two regimes, which differ according to which studies are published and chosen for replication.

1. **Actual:** this regime is based on the actual expected replication rate, that is, where the studies chosen for replication depends both on the degree of selective publication and how replicators actually choose which results to replication e.g. only choosing statistically significant results to replicate.
2. **Counterfactual:** this regime considers a counterfactual scenario where all results are published and replication is random. This implies the distribution of published, replicated studies coincides with the distribution of latent studies.

Let the subscript ‘Ac’ denote the actual regime and ‘Cf’ the counterfactual regime. Formally, the expectation operators under both regimes are defined by:

$$\mathbb{E}_{Ac}[RP(X, \Theta, \sigma_r(X, \beta))] = \int RP(x, \theta, \sigma_r(x, \beta)) f_{X^*, \Theta^* | D, R, S_X^*}(x, \theta | D = 1, R = 1, S_X^* = 1) dx d\theta \quad (62)$$

$$\mathbb{E}_{Cf}[RP(X, \Theta, \sigma_r(X, \beta))] = \int RP(x, \theta, \sigma_r(x, \beta)) f_{X^*, \Theta^*}(x, \theta) dx d\theta \quad (63)$$

See that the expressions differ based on the distribution of studies over which we integrate. Under the ‘actual’ regime, we integrate over the distribution of latent studies conditional on being selected for publication ($D = 1$) and replication ($R = 1$). By contrast, for the ‘counterfactual’ regime, we integrate over the unconditional distribution of latent studies (since there is no distortive selection in this regime).

Using these expressions, we have the following decomposition:

$$\begin{aligned} \underbrace{(1 - \beta) - \mathbb{E}_{Ac}[RP(X, \Theta, \sigma_r(X, \beta))]}_{\text{replication rate gap}} &= \underbrace{(1 - \beta) - \mathbb{E}_{Cf}[RP(X, \Theta, \sigma_r(X, \beta)) | X \geq 0]}_{\text{(i) concavity gap}} \\ &+ \underbrace{\mathbb{P}_{Ac}(X < 0) \left(\mathbb{E}_{Ac}[RP(X, \Theta, \sigma_r(X, \beta)) | X \geq 0] - \mathbb{E}_{Ac}[RP(X, \Theta, \sigma_r(X, \beta)) | X < 0] \right)}_{\text{(ii) wrong-sign gap}} \\ &+ \underbrace{\mathbb{E}_{Cf}[RP(X, \Theta, \sigma_r(X, \beta)) | X \geq 0] - \mathbb{E}_{Ac}[RP(X, \Theta, \sigma_r(X, \beta)) | X \geq 0]}_{\text{(iii) selection-on-significance gap}} \end{aligned} \quad (64)$$

Proof. Write the expected replication probability under model 1 as

$$\begin{aligned} \mathbb{E}_{Ac}[RP(X, \Theta, \sigma_r(X, \beta))] &= \mathbb{E}_{Ac}[RP(X, \Theta, \sigma_r(X, \beta)) | X \geq 0] \\ &+ \mathbb{P}_{Ac}(X < 0) \left(\mathbb{E}_{Ac}[RP(X, \Theta, \sigma_r(X, \beta)) | X < 0] \right) - \mathbb{E}_{Ac}[RP(X, \Theta, \sigma_r(X, \beta)) | X \geq 0] \end{aligned} \quad (65)$$

To arrive at equation (64), substitute equation (65) into the replication rate gap; add and subtract $\mathbb{E}_{Cf}[RP(X, \Theta, \sigma_r(X, \beta)) | X \geq 0]$; and rearrange the terms. \square

Equation (64) states that we can express the difference between the nominal power target $1 - \beta$ and the expected replication rate – i.e. the replication rate gap – into the sum of three components: the concavity gap; the wrong-sign gap; and the selection-on-significance gap.

The concavity gap measures how far we move from the nominal target $1 - \beta$ once we account for the concavity of the $RP(\cdot)$ function. That is, how much does the replication rate decrease once we taken the expectation of the function (rather than evaluation the function at its expectation)? Note that this portion of the gap isolates the impact of the concavity of the function by abstracting from the two other issues. First, it abstracts from distortions due to selection-on-significance since it is based on the counterfactual regime with no selection.

Second, it abstracts from the issue original estimate having the wrong sign by conditioning of $X > 0$.

The wrong-sign gap examines the difference in the expected replication rate between estimates with the correct sign $X > 0$ and the incorrect sign $X < 0$. In that latter case, replication probabilities are very low.

Finally, the selection on significance gap examines the difference in the expected replication rate between the actual regime, which may be distorted by publication and replication selection, and the counterfactual regime, which is not. Like the concavity gap, it abstracts away from the issue of original estimates with the wrong sign.

C.2. Estimation

The decomposition can be calculated for economics and psychology using the model estimates in Table 1 in the main text. I do not calculate the decomposition for social science experiments because the fractional power rule does not provide a nominal power target, and hence there is no gap that can be decomposed.

Expectations are calculated by simulation. For example, to calculate the expected replication rate in the counterfactual regime, we randomly draw from the estimated distribution of latent true effects and standard errors in Table 1 and then, without censoring any draws due to selective publication or replication, we calculate their average replication probability. Expectations in the actual regime are calculated by drawing from the same latent distribution of true effects and standard errors, but in this case, some studies are now censored due to selective publication and replication, as determined by the estimated selection parameters $\beta_{\mathbf{p}}$ in Table 1. This is a similar procedure as in Section 4 in the main text, except that we condition here on the sign of original estimates for the decomposition.

Note that the distribution of latent studies is the same across both regimes, but the key difference is that studies are subject to selection in the ‘actual’ regime while they are not in the ‘counterfactual’ regime.

C.3. Results

Table C1 presents the results. Panel A reproduces the results in the main text and Panel B presents the decomposition results. The empirical results for the decomposition show that failing to account for the concavity of the replication power function explains the overwhelming majority of the explained replication rate gap in both economics and psychology. The selection-on-significance gap is small, explaining only 3.1% of the gap in economics, while actually *decreasing* the expected replication rate in psychology.

TABLE C1 – REPLICATION RATE PREDICTIONS AND DECOMPOSITION RESULTS

	Economics experiments	Psychology	Social sciences
<i>A. Replication rate predictions</i>			
Nominal target (intended power)	0.92	0.92	–
Observed replication rate	0.611	0.348	0.571
Predicted replication rate	0.600	0.545	0.543
<i>B. Decomposition of explained gap</i>			
Predicted replication rate gap	0.320 (100%)	0.375 (100%)	–
Concavity gap	0.292 (91.16%)	0.364 (97.16%)	–
Wrong-sign gap	0.018 (5.72%)	0.030 (8.03%)	–
Selection-on-significance gap	0.010 (3.12%)	-0.019 (-5.18%)	–

Notes: Economics experiments refers to [Camerer et al. \(2016\)](#), psychology experiments to [Open Science Collaboration \(2015\)](#) and social sciences to [Camerer et al. \(2018\)](#). The replication rate is defined as the share of original estimate whose replications have statistically significant findings of the same sign. Figures in the first row report the mean intended power reported in both applications. The second row shows observed replication rates. The third row reports the predicted replication rate in equation (9) calculated using parameter estimates Table 1. In social sciences, power is set to detect three-quarters of the original effect size with 90% power. This approach does not have a fixed nominal target for the replication rate.

The small to negative impact of selection on significance on the replication rate is perhaps surprising. This is because the selection gap is the net outcome of two offsetting effects. First, conditioning on significance leads to exaggerated estimates and hence lower replication rates *conditional* on any given true effect θ . However, selection on significance also changes the distribution of true effects that make it into the published literature. In particular, the average true effect size increases with selection since larger true effects are more likely to produce statistically significant results. Moreover, larger true effects have higher replication probabilities as compared with smaller true effects, leading to an increase in the probability of replication. These two opposing forces lead to a small decline in replication rates in economics and a slight *increase* in replication rates in psychology. In the subsections below, I provide additional evidence for the intuition underpinning these results.

Concavity gap.—Figure C1 presents normal simulations showing that the non-linearity gap is largest for standardized true effects $\omega \equiv \theta/\sigma$ which are close to 0, and remains above 0.2 for $\omega \leq 1$. It decreases monotonically as the true effect size ω increases and approaches zero in the limit (Lemma A1.5 in Appendix A). It follows that the size of the non-linearity gap depends on the distribution of ω . The first row of graphs in Figure C2 plot the distribution of latent studies (conditional on having the ‘correct’ sign to match the decomposition). We see that a high fraction of latent studies have $\omega < 1$, which explains why the non-linearity gap explains such a large portion of the replication rate gap.

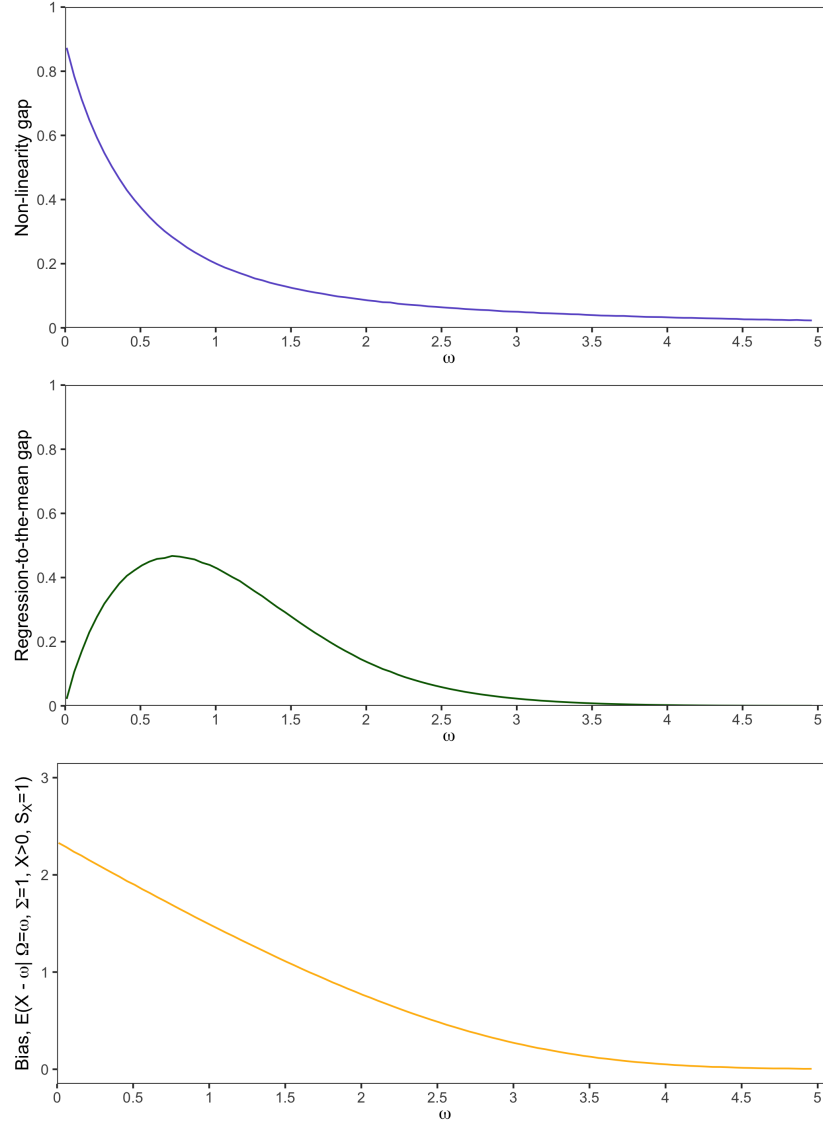


FIGURE C1. REPLICATION RATE GAP DECOMPOSITION: MONTE CARLO SIMULATIONS

Notes: Plots are based on simulating studies from an $N(\omega, 1)$ distribution, for different values of ω . Replication estimates are drawn from a $N(\omega, \sigma_r(x, \beta)^2)$, where $\sigma_r(x, \beta)$ is set based on the common power rule to detect the original effect x with $1 - \beta = 0.92$ intended power. The non-linearity gap and regression-to-the-mean gap are based on equation (64) and calculated using Monte Carlo methods.

Wrong-sign gap.—Random sampling variation means that original estimates will occasionally have the ‘wrong’ sign. When this occurs, the replication probability is bounded above by 0.025. The extent to which this issue contributes to low replication rates therefore depends on the share of studies that have the wrong sign among significant studies. Since the probability of an original study having the wrong sign and being statistically significant is fairly low (3% in economics and 5% in psychology), the contribution of the wrong-sign gap is relatively small.

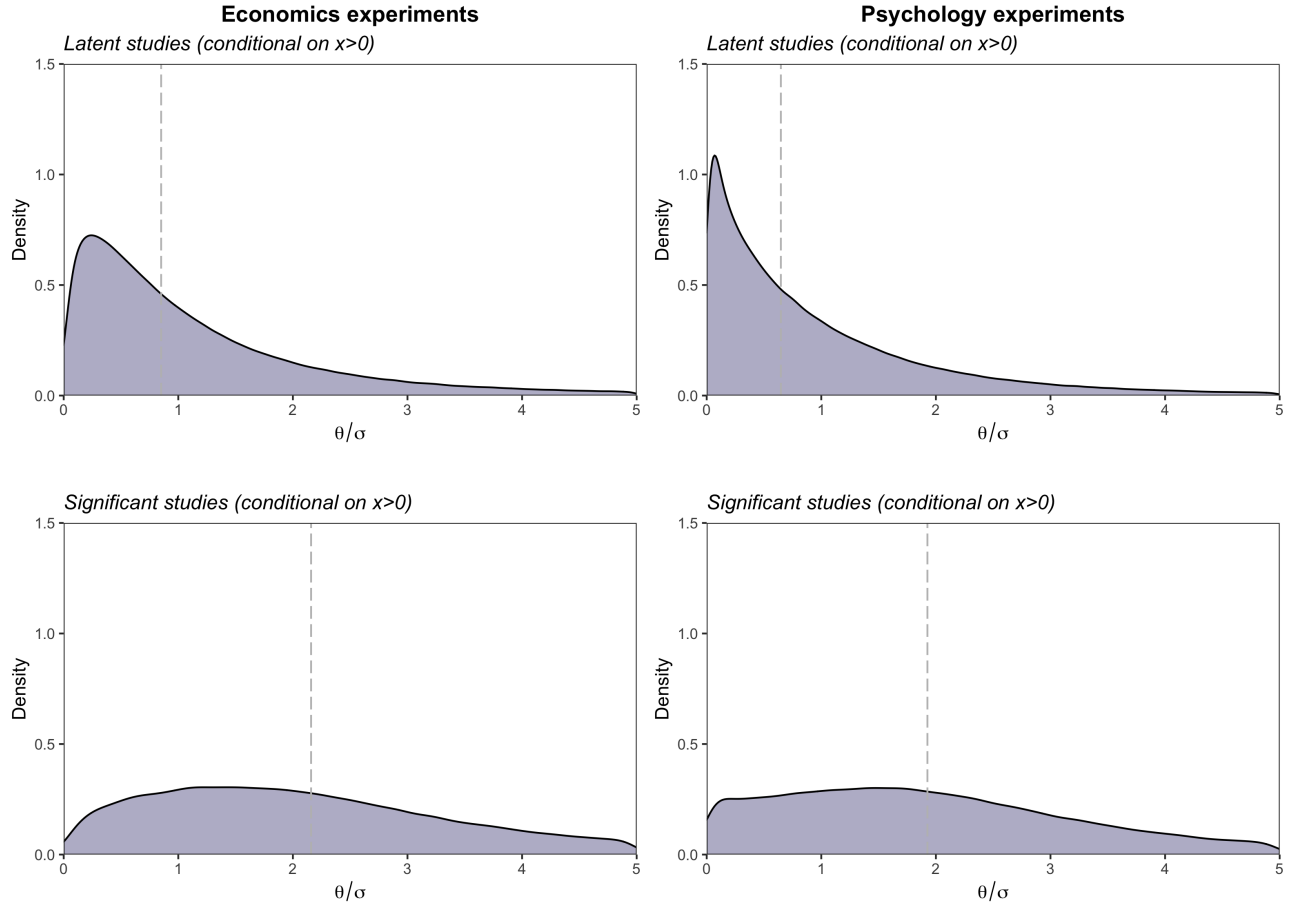


FIGURE C2. DISTRIBUTION OF NORMALIZED TRUE EFFECTS: LATENT STUDIES AND SIGNIFICANT STUDIES

Notes: Economics experiments refers to [Camerer et al. \(2016\)](#) and psychology experiments to [Open Science Collaboration \(2015\)](#). Densities are based on simulated draws from the estimated distribution of latent studies in Table 1 in the main text. Dashed vertical lines show the median of the distribution.

Selection-on-significance gap.—The selection-on-significance gap is 1% in economics and slightly negative for psychology (i.e. conditioning on statistical significance increases the replication rate compared to when there is no conditioning). The sign of this gap is ambiguous because of two opposing effects from conditioning on statistical significance, as discussed above. To better understand these two effects, consider the figures in Table C2 which are based on the estimated empirical models. For the first effect, note that conditioning on significant findings increases mean bias in both applications.²⁴ This makes replication more difficult for any given true effect and standard error. For the second effect, note that conditioning also tends to select studies with larger standardized true effects ω , which have higher replication probabilities.²⁵

²⁴Note that bias is positive for latent studies because these statistics condition on original estimates X^* to have the same sign as true effects.

²⁵The impact of conditioning on the full distribution of ω can be seen in Figure C2.

TABLE C2 – TRUE EFFECT SIZES AND BIAS FOR STUDIES WITH THE ‘CORRECT’ SIGN

	Economics experiments		Psychology experiments	
	Latent	Published & significant	Latent	Published & significant
Mean standardized bias	1.522	3.352	1.323	2.950
Mean standardized true effect	1.415	2.915	1.084	2.367

Notes: Economics experiments refers to [Camerer et al. \(2016\)](#) and psychology experiments to [Open Science Collaboration \(2015\)](#). Figures are based on simulated draws from the estimated distribution of latent studies from Table 1 in the main text. The mean of the standardized true effect is equal to $\mathbb{E}[\Omega^*|S_X^*, X^* > 0, D]$. Mean standardized bias is equal to $\mathbb{E}[X^*/\Sigma^* - \Omega^*|S_X^*, X^* > 0, D]$. ‘Published & significant studies’ set $S_X^* = 1$ and $D = 1$. ‘Latent studies’ do not condition on significance or publication.

Higher replication probabilities arise because for more highly powered studies: non-linearity effects are less severe (Panel 1, Figure C1); and bias is smaller (Panel 3, Figure C1). Bias is smaller because censoring insignificant original estimates has little ‘bite’ when the true effect is very large, since the probability of drawing an insignificant estimate is very small.

D. Alternative Measures of Selective Publication

Proposition 1 shows that the replication rate is unresponsive to the most salient form of selective publication. For journals and policymakers seeking to change current norms, this highlights the need for more informative measures. In this section, I conduct policy simulations using the estimated model to show how three alternative measures respond to changes in the selective publication of null results:

1. **Replication CI:** This measure counts a replication as ‘successful’ if its 95% confidence interval covers the original estimate: $\mathbb{1}[X \in (X_r - 1.96\Sigma_r, X_r + 1.96\Sigma_r)]$.
2. **Meta-analysis:** The standard criterion of replication with the same sign and significance is applied to a fixed-effect meta-analytic estimate combining the original and replication estimate (uncorrected for selective publication): $\mathbb{1}[|X_m| \geq 1.96\Sigma_m, \text{sign}(X_m) = \text{sign}(X)]$ where X_m and Σ_m are the meta-analytic estimate and standard error, respectively.²⁶
3. **Prediction interval:** Original and replication estimates are counted as ‘consistent’ under this approach if their difference is not statistically different from zero at the 5% level ([Patil et al., 2016](#)). This is equivalent to estimating a 95% ‘prediction interval’ for the

²⁶The fixed-effects meta-analytic estimate is a weighted average of original and replication estimates: $X_m = (\omega_o X + \omega_r X_r) / (\omega_o + \omega_r)$, where the weights are equal to the precision of each estimate i.e. $(\omega_o, \omega_r) = (\Sigma^{-2}, \Sigma^{-2})$. These weights minimize the mean-squared error of X_m ([Laird and Mosteller, 1990](#)). The variance of this estimator is given by $\Sigma_m^2 = 1 / (\omega_o + \omega_r)$.

original estimate and then determining if it covers the replication estimate: $\mathbb{1}[X_r \in (X - 1.96\sqrt{\Sigma^2 + \Sigma_r^2}, X + 1.96\sqrt{\Sigma^2 + \Sigma_r^2})]$.²⁷

These alternative replication measures are frequently reported in large-scale replication studies (Open Science Collaboration, 2015; Camerer et al., 2016, 2018). In simulations, I calculate these measures over significant and insignificant published results, since conditioning on statistical significance makes them unresponsive to selective publication on null results (Proposition B4).

Simulations assume that all results significant at the 5% level are published, and that results insignificant at the 5% level are published with probability β_p . I then calculate how the various measures change with β_p to see how well they capture changes in selective publication (e.g. because of policy changes that reduce selective publication). Policymakers' successful efforts to increase the probability of publishing null results lead to an increase in the policy variable, β_p . Note that while model estimation assumes multiple cutoffs, policy simulations are performed assuming policymakers influence publication probabilities at a single cutoff (1.96) for simplicity (i.e. in the policy simulations I set $\beta_p = \beta_{p1} = \beta_{p2}$ and $\beta_{p3} = 1$ in social science).

Figure D1 shows the results. In line with Proposition 1, the replication rate is completely unresponsive to changes in the probability of publishing null results, making it a poor measure to evaluate efforts to reduce selective publication. Turning to alternative measures, note that the replication CI and meta-analysis measures actually *worsen* when more null results are published ($\beta_p \rightarrow 1$). This is because less selective publication leads to more small effects being selected for replication, which have relatively low replication probabilities under these approaches. By contrast, the prediction interval measure is low when selective publication is high, and approaches close to 95% as the probability of publishing null results approach one.²⁸ The prediction interval measure performs well because it explicitly accounts for the decline in original power as more small effects are selected for replication. Noisy low-powered original studies contain limited information about true effects, which implies that a large range of replication estimates are statistically consistent with them.

Overall, for the purpose of evaluating efforts to reduce selective publication, these results suggest that calculating the prediction interval measure over a random sample of all published results could provide a useful alternative to the replication rate.

²⁷This approach assumes that original and replication estimates share the same true effect and are statistically independent. For more details, see the Supplementary Materials for Patil et al. (2016).

²⁸When $\beta_p = 1$, the prediction interval measure is slightly higher than 95% in all applications. This is because it assumes that the original estimate X and the replication estimate X_r are uncorrelated. In practice, the replication standard error is a function of the original estimate via the common power rule, which generates some correlation between X and X_r .

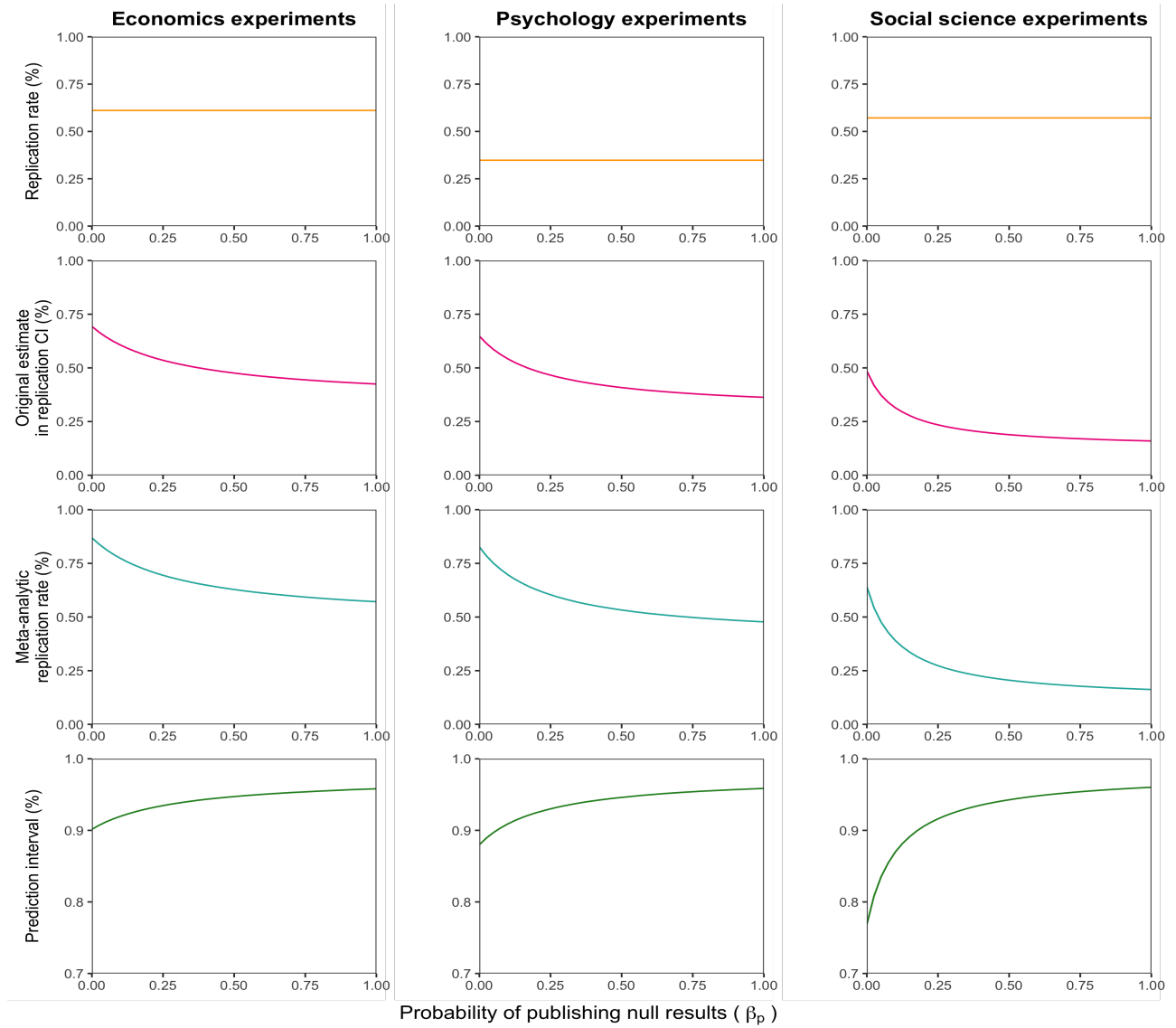


FIGURE D1. POLICY SIMULATIONS: ALTERNATIVE MEASURES OF REPLICATION AND SELECTIVE PUBLICATION

Notes: Details of each measure are provided in the main text. All measures except for the replication rate are calculated over significant and insignificant published results. Simulations use model estimates of the latent distribution of studies from Table 1 and set different levels of selective publication β_p . The first column reproduces replication rate predictions in Table 2.

E. Likelihood and Robustness to Alternative Parametric Specifications

This appendix provides details on the empirical model used in the applications. The first subsection provides details on the likelihood; the second discusses sign normalization in estimation; and the third presents robustness results for alternative parametric specifications.

Likelihood.—In each empirical application, the model is estimated via maximum likelihood. The marginal likelihood of published estimates and standard errors, (X, Σ) , is given by

$$f_{X,\Sigma}(x, \sigma) = \frac{p\left(\frac{x}{\sigma}\right)r\left(\frac{x}{\sigma}\right) \int_{\theta} \frac{1}{\sigma} \phi\left(\frac{x-\theta}{\sigma}\right) dF_{\Theta}(\theta|\kappa_{\theta}, \lambda_{\theta}) \cdot g_{\Sigma}(\sigma|\kappa_{\sigma}, \lambda_{\sigma})}{\int_{x'} \int_{\sigma'} p\left(\frac{x'}{\sigma'}\right)r\left(\frac{x'}{\sigma'}\right) \int_{\theta} \frac{1}{\sigma'} \phi\left(\frac{x'-\theta}{\sigma'}\right) dF_{\Theta}(\theta|\kappa_{\theta}, \lambda_{\theta}) g_{\Sigma}(\sigma'|\kappa_{\sigma}, \lambda_{\sigma}) dx' d\sigma'} \quad (66)$$

where $\phi()$ is the standard normal density, F_{Θ} is the distribution function for (normalized) latent true effects $|\Theta^*|$, and g_{Σ} is the density function for latent standard errors Σ^* . Assuming independence across studies, the log-likelihood of the data $\{x_i, \sigma_i\}_i$ is $\ell(\kappa_{\theta}, \lambda_{\theta}, \kappa_{\sigma}, \lambda_{\sigma}, \beta_{\mathbf{p}}) = \sum_i \log f_{X,\Sigma}(x_i, \sigma_i)$, where $\beta_{\mathbf{p}}$ is a vector of the parameters of the publication probability function.

Sign Normalization.—Following [Andrews and Kasy \(2019\)](#)²⁹, we normalize initial estimates to be positive and conduct estimation based on $W \equiv |X|$ and Σ using the marginal likelihood $f_{W,\Sigma}(w, \sigma) = f_{X,\Sigma}(w, \sigma) + f_{X,\Sigma}(-w, \sigma)$.

Robustness to Alternative Parametric Specifications.—The model-based replication rate predictions in [Table 2](#) are based on the maximum likelihood model estimates in [Table 1](#). These model estimates are based on certain parametric assumptions about the distribution of latent (normalized) true effects $|\Theta^*|$ and latent standard error Σ^* , namely, that both follow gamma distributions.

[Table E1](#) presents replication rate predictions using three alternative sets of parametric assumptions. For reference, the results of the baseline gamma-gamma model from the main text are reprinted in the first row under ‘model predictions’. The other models consider different combinations of gamma and log-normal assumptions for the distribution of $|\Theta^*|$ and Σ^* . For example, (Gamma, Log-normal) refers to a model where $|\Theta^*|$ follows a gamma distribution and Σ^* follows a log-normal distribution.

Overall, the accuracy of the replication rate predictions is robust to alternative parametric specifications. For economics experiments and psychology, the replication rate predictions are very similar across all models. For social sciences, the replication rate predictions are somewhat higher in models where we assume that the latent distribution of true effects follows a gamma

²⁹[Andrews and Kasy \(2019\)](#) analyze the same three application in their paper: [Camerer et al. \(2016\)](#) and [Open Science Collaboration \(2015\)](#) in their main paper and [Camerer et al. \(2018\)](#) in the Online Appendix.

distribution (62% v. 54%). However, the accuracy of the models, as measured by the distance of their predictions from the observed replication rate, is similar (5 percentage points vs. 3 percentage points).

	Economics experiments	Psychology	Social science
Nominal power (intended target)	0.92	0.92	–
Observed replication rate	0.611	0.348	0.571
<i>Model predictions</i>			
Baseline model: (Gamma, Gamma)	0.600	0.545	0.543
(Log-normal, Gamma)	0.602	0.547	0.621
(Gamma, Log-normal)	0.606	0.543	0.542
(Log-normal, Log-normal)	0.607	0.545	0.623

TABLE E1 – ROBUSTNESS: REPLICATION RATE PREDICTIONS

Notes: Model predictions are based on models with different parametric assumptions. Each model is denoted by a tuple with two distributions. The first refers to the distribution of normalized latent true effects $|\Theta^*|$ and the second to the distortion of latent standard errors Σ^* . For example, (Gamma, Log-normal) refers to a model where $|\Theta^*|$ follows a gamma distribution and Σ^* follows a log-normal distribution.

F. Replication Selection in Empirical Applications

Replication selection is a multi-step mechanism that first selects studies, and then selects results within those studies to replicate (since studies typically report multiple results). It consists of three steps:

1. **Eligibility:** define the set of eligible studies (e.g. journals, time-frame, study designs).
2. **Study selection:** on the set of eligible studies, a mechanism that select which studies will be included in the replication study.
3. **Within-study replication selection:** for selected studies, a mechanism for selecting which result(s) to replicate.

These three features of the replication selection mechanism influence the interpretation of the selection parameters $(\beta_{p1}, \beta_{p2}, \beta_{p3})$.

Economics experiments.—Consider these three steps in [Camerer et al. \(2016\)](#):

1. **Eligibility:** Between-study laboratory experiments in *American Economic Review* and *Quarterly Journal of Economics* published between 2011 and 2014.

2. **Study selection:** [Camerer et al. \(2016\)](#) select for replication all eligible studies that had ‘at least one significant between subject treatment effect that was referred to as statistically significant in the paper.’ [Andrews and Kasy \(2019\)](#) review eligible studies and conclude that no studies were excluded by this restriction. Thus, the complete set of eligible studies was selected for replication.
3. **Within-study replication selection:** the most important *statistically significant* result within a study, as emphasized by the authors, was chosen for replication. Further details are in the supplementary materials in [Camerer et al. \(2016\)](#). Of the 18 replication studies, 16 were significant at the 5% level and two had p -values slightly above 0.05 but were treated as ‘positive’ results for replication and included in the replication rate calculation.

I assume replication selection is random with respect to the t -ratio for results whose p -values are below or only slightly above 0.05. This implies that β_{p2} measures the relative probability of being published and chosen for replication for a result whose p -value is slightly above 0.05, compared to if it were strictly below 0.05. Overall, the empirical results are valid for the population of ‘most important’ significant (or ‘almost significant’) results, as emphasized by authors, in experimental economics papers published in top economics journals between 2011 and 2014.

Psychology.—Next, consider replication selection in [Open Science Collaboration \(2015\)](#):

1. **Eligibility:** Studies published in 2008 in one of the following journals: *Psychological Science*, *Journal of Personality and Social Psychology*, and *Journal of Experimental Psychology: Learning, Memory, and Cognition*.
2. **Study selection:** [Open Science Collaboration \(2015\)](#) write: ‘The first replication teams could select from a pool of the first 20 articles from each journal, starting with the first article published in the first 2008 issue. Project coordinators facilitated matching articles with replication teams by interests and expertise until the remaining articles were difficult to match. If there were still interested teams, then another 10 articles from one or more of the three journals were made available from the sampling frame.’ Importantly, the most common reason why an article was not matched was due to feasibility constraints (e.g. time, resources, instrumentation, dependence on historical events, or hard-to-access samples).
3. **Within-study replication selection:** the last experiment reported in each article was chosen for replication. [Open Science Collaboration \(2015\)](#) write that, ‘Deviations from

selecting the last experiment were made occasionally on the basis of feasibility or recommendations of the original authors.’ A small number of results had p -values just above 0.05 but were treated as ‘positive’ results for replication, as in [Camerer et al. \(2016\)](#).

This selection mechanism implies that the empirical results are valid for the distribution of last experiments in the set of eligible journals. Since neither studies nor results were selected based on statistical significance, it is reasonable to treat the ‘last experiment’ rule as effectively random. In this case, we can interpret the results are being valid for all results in the eligible set of journals.

Social science experiments.—Finally, consider replication selection in [Camerer et al. \(2018\)](#):

1. **Eligibility:** Experimental studies in the social sciences published in *Nature* or *Science* between 2010 and 2015.
2. **Study selection:** [Camerer et al. \(2018\)](#) include all studies that: ‘(1) test for an experimental treatment effect between or within subjects, (2) test at least one clear hypothesis with a statistically significant finding, and (3) were performed on students or other accessible subject pools. Twenty-one studies were identified to meet these criteria.’
3. **Within-study replication selection:** [Camerer et al. \(2018\)](#) write, ‘We used the following three criteria in descending order to determine which treatment effect to replicate within each of these 21 papers: (a) select the first study reporting a significant treatment effect for papers reporting more than one study, (b) from that study, select the statistically significant result identified in the original study as the most important result among all within- and between-subject treatment comparisons, and (c) if there was more than one equally central result, randomly select one of them for replication.’ All results selected for replication had p -values strictly below 0.05.

This selection mechanism implies that the empirical results are valid for the population of statistically significant between- or within-subject treatment comparisons in experimental social science, which were identified by authors as the most ‘important’ and published in *Nature* or *Science* between 2010 and 2015.

G. Predicted Replication Rates Under Alternative Power Calculations

This appendix presents several extensions to the main empirical results on predicting replication rates in experimental economics, psychology and social science. The first extension

allows for variation in the application of the common power rule around mean intended power. Results are similar to those in the main text, which assume no variability in the application of the common power rule. The second extension generates replication rate predictions under the rule of setting replication power equal to original power. This delivers lower replication rates than the common power rule.

Alternative power calculation rules.—Consider first the rule used for calculating replication power in the main text, and then two additional approaches. For concreteness, suppose we want to calculate the replication standard error for a simulated original study $(x^{sim}, \sigma^{sim}, \theta^{sim})$.

1. **Common power rule (mean):** This is the rule reported in the results in the main text. It assumes no variability in the application of the common power rule, such that all replications have mean intended power $1 - \beta$. This rule implies

$$\sigma_r^{sim}(x^{sim}, \beta) = \frac{|x^{sim}|}{1.96 - \Phi^{-1}(\beta)} \quad (67)$$

2. **Common power rule (realized):** Intended power for individual replications varied around mean intended power for at least two reasons. First, replication teams were instructed to meet minimum levels of statistical power, and encouraged to obtain higher power if feasible. Second, a number of replication in [Open Science Collaboration \(2015\)](#) did not meet this requirement. Figure G1 shows the distribution of realized intended power in replications for experimental economics and psychology. Realized intended power is right-skewed for psychology. In experimental economics and social science, realized intended power is distributed more tightly around mean.

To capture variability in the application of the common power rule, take a random draw from the empirical distribution of $|x|/\sigma_r$ and denote it $1.96 - \hat{\beta}^n$. Then realized intended power for simulated study $(x^{sim}, \sigma^{sim}, \theta^{sim})$ is equal to

$$\sigma_r^{sim}(x^{sim}, \hat{\beta}^n) = \frac{|x^{sim}|}{1.96 - \Phi^{-1}(\hat{\beta}^n)} \quad (68)$$

3. **Same power:** Set replication power equal to the power in the original study:

$$\sigma_r^{sim}(\sigma^{sim}) = \sigma^{sim} \quad (69)$$

This rule has been proposed as a straightforward, intuitive approach for designing replication studies. In a review of replication studies by [Anderson and Maxwell \(2017\)](#), 19 of

108 studies used this approach.

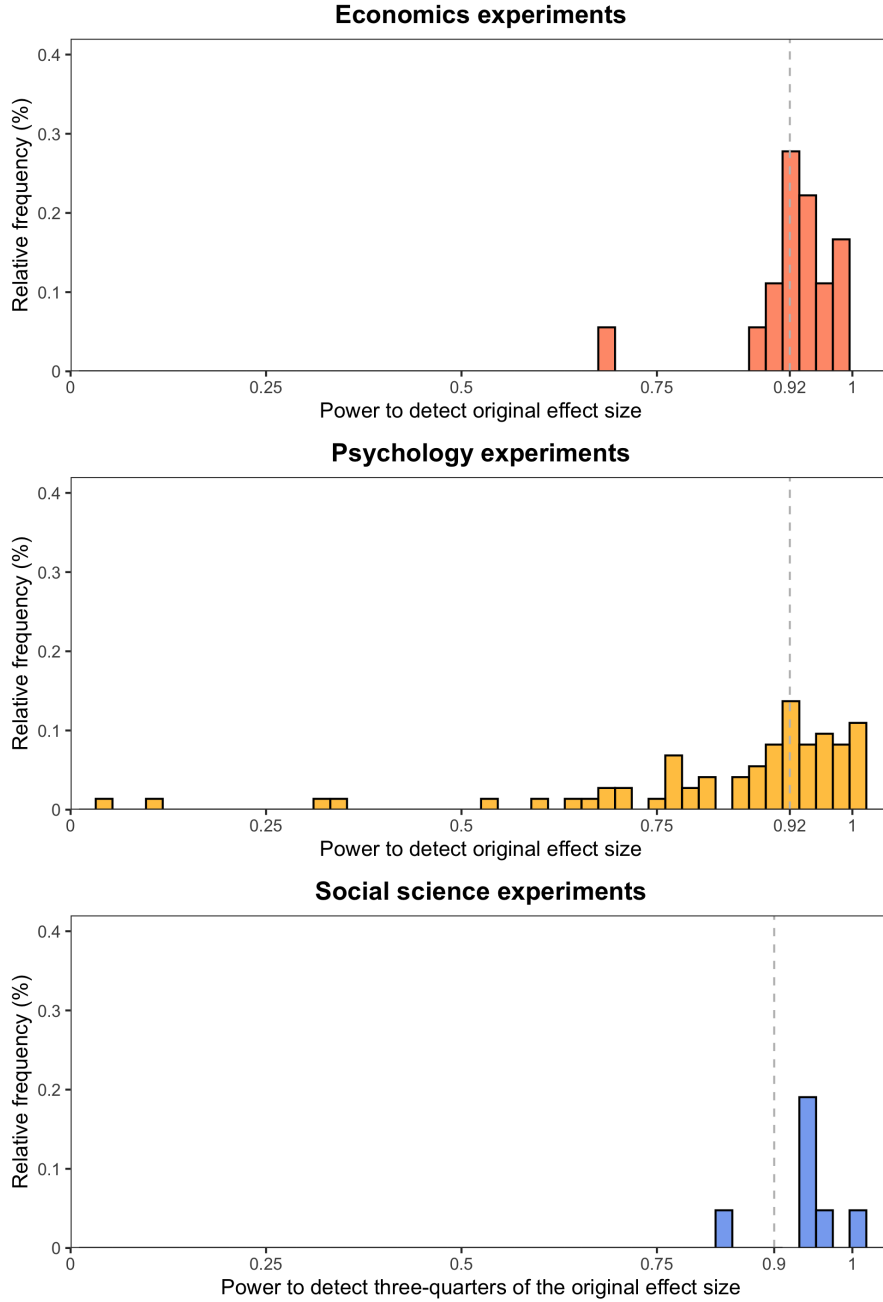


FIGURE G1. Realized Intended Power

Notes: Data are from [Camerer et al. \(2016\)](#), [Open Science Collaboration \(2015\)](#), and [Camerer et al. \(2018\)](#), respectively. Realized intended power is defined as $1 - \Phi(1.96 - \psi \cdot \frac{x}{\sigma_r})$ with $\psi = 1$ in economics and psychology and $\psi = 3/4$ in social science. The horizontal dashed line is reported mean power in each application. In economics and psychology, this is 92% to detect the original effect size. In social science, this is 90% to detect three quarters of the effect size.

Results.—Table G1 presents the results for all three applications. Panel A shows that

allowing intended power to vary across replications (‘Realized power’) yields similar replication rate prediction to assuming all replications have intended power equal to the report mean (‘92% on X ’). In fact, in all three applications, the accuracy improves very slightly under the realized power rule. The biggest differences is in psychology, because the realized power rule accounts for the fact that the distribution of intended power is right skewed.

Panel B examines the proposed rule of setting replication power equal to original power. In all three cases, the expected replication rate is lower than under the common power rule.

TABLE G1 – REPLICATION RATE PREDICTIONS UNDER ALTERNATIVE REPLICATION POWER RULES

	Economics	Psychology	Social science
<i>A. Replication rate predictions</i>			
Nominal target (intended power)	0.92	0.92	–
Observed replication rate	0.611	0.348	0.571
Mean power	0.600	0.545	0.543
Realized power	0.615	0.522	0.555
<i>B. Alternative rule</i>			
Same power	0.550	0.486	0.494

Notes: Economics experiments refer to [Camerer et al. \(2016\)](#), psychology experiments to [Open Science Collaboration \(2015\)](#), and social science experiments to [Camerer et al. \(2018\)](#). The replication rate is defined as the share of original estimate whose replications have statistically significant findings of the same sign. Figures in the first row are observed outcomes from large-scale replication studies. Remaining rows report predicted replication rates using parameter estimates Table 1 in the main text and assuming different rules for calculating replication power.

H. Relative Effect Size Predictions

The main focus of this article is the binary measure of replication based on the statistical significance criterion. This is because of its status as the primary replication indicator in the large-scale replication studies.³⁰ However, complementary measures are frequently presented alongside the replication rate. Perhaps the most common is the relative effect size, a continuous measure of replication defined as the ratio of replication effect size and original effect size. Relative effect sizes typically range between 0.35 and 0.7. Below, I include a brief theoretical discussion of the relative effect size and then present predictions of this measure using the estimated models.

Theoretical discussion.—The relative effect size for individual studies may be informative about biases affecting original studies, especially when original studies are well-powered. However, as an *aggregate* measure of reproducibility, the relative effect size measure may be subject

³⁰Power calculations in replications are themselves typically designed to measure a binary notation of replication ‘success’ or ‘failure’.

to similar issues to the replication rate, at least in the case where it is defined exclusively over significant findings.

First, if the relative effect size is defined over significant original results, then it will be largely uninformative about the ‘file-drawer’ problem (Proposition B4).³¹ Second, non-random sampling of significant results for replication mechanically induces inflationary bias in original estimates and regression to the mean in replication estimates, such that relative effect sizes are below one in expectation. Thus, similar to the replication rate, it has no natural benchmark against which to judge deviations, making it challenging to interpret. Relatedly, the average relative effect size is also very sensitive to power in original studies, which is unobserved. Figure H1 provides an illustration with intended power set to 0.9, which shows that the expected relative effect size for significant results is increasing in the power of original studies, and approaches one only as statistical power approaches 100%.

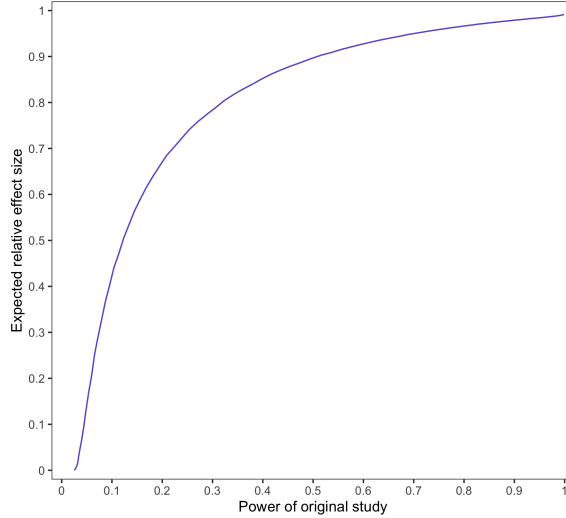


FIGURE H1. EXPECTED RELATIVE EFFECT SIZE OF SIGNIFICANT ORIGINAL STUDIES AND STATISTICAL POWER

Notes: Illustration for the relationship between original power and the expected relative effect size of significant findings under the common power rule are both functions of $\omega = \theta/\sigma$ (normalized to be positive). Original power to obtain a significant effect with the same sign as the true effect is equal to $1 - \Phi(1.96 - \omega)$. The expected relative effect size is calculated by taking 10^6 draws of Z from $N(\omega, 1)$ and then calculating $\frac{1}{M_{sig}} \sum_{i=1}^{M_{sig}} \rho_{i,r}^{sig} / \rho_i^{sig}$, where $\rho = \tanh z$ denotes the Pearson correlation coefficient obtained by transforming the Fisher-transformed correlation coefficient (Fisher, 1915); and M_{sig} is the number of significant latent studies. The superscript *sig* reflects the fact that only statistically significant original results at the 5% level and their replications are included in the calculation. Replication estimates $z_{i,r}$ are drawn from an $N(\omega, \sigma_{r,i}(z_i, \beta)^2)$ distribution. The replication standard error is calculated using the common power rule to detect original effect sizes with 90% power (i.e. $1 - \beta = 0.9$), which is given by $\sigma_r(z_i, \beta) = |z_i| / [1.96 - \Phi^{-1}(\beta)] = |z_i| / 3.242$.

³¹Defining it over null results may present its own difficulties. For a perfectly measured null effect, the denominator in the statistic is equal to zero and the statistic is not well defined. On the other hand, if it is close but not equal to zero, then the statistic is highly sensitive to the precision of replication estimates; this raises questions about how one should set replication power when replicating a null effect.

Empirical results.—The estimated models in Table 1 in the main text can be used to generate predictions of the average relative effect sizes. To procedure for simulating replications is identical to the procedure outlined in the main text for the replication rate case. Let $\{x_i, \sigma_i, x_{r,i}, \sigma_{r,i}\}_{i=1}^{M_{sig}}$ be the set of simulated original studies that are published and significant, and their corresponding replication results; M_{sig} is the size of the set. The predicted relative effect size is equal to

$$\frac{1}{M_{sig}} \sum_{i=1}^{M_{sig}} \frac{\rho_{i,r}^{sig}}{\rho_i^{sig}} \quad (70)$$

where $\rho = \tanh x$ denotes the Pearson correlation coefficient which is obtained by transforming the Fisher-transformed correlation coefficient (Fisher, 1915). I also present results for the median relative effect size. Results are presented in Table H1. The predicted average relative effect size is relatively close to observed average relative effect size in economics, somewhat further off in social science, and quite far off in psychology. In each case, the predicted average relative effect size is optimistic compared to the observed value. In economics and psychology, the difference in predicted and observed relative effect sizes is not statistically different from zero, while in psychology it is. Predictions for median relative effect sizes show qualitatively similar results.

TABLE H1 – AVERAGE RELATIVE EFFECT SIZE PREDICTIONS

	Economics	Psychology	Social Sciences
Observed relative effect size (mean)	0.657	0.374	0.443
Predicted relative effect size (mean)	0.703 (0.135)	0.637 (0.060)	0.533 (0.141)
Observed relative effect size (median)	0.691	0.292	0.527
Predicted relative effect size (median)	0.747 (0.129)	0.674 (0.063)	0.595 (0.240)

Notes: Economics experiments refers to Camerer et al. (2016), psychology experiments to Open Science Collaboration (2015) and social science experiments to Camerer et al. (2018). Observed relative effect sizes are based on data from large-scale replication studies. Predicted average relative effect sizes are calculated using equation (70) and the procedure outlined in the text. Standard errors are calculated using the delta method.

I. Extended Replication Rate Definition

This appendix analyzes a generalization of the replication rate definition that extends to insignificant results. It outlines a number of issues with this proposal.

The Generalized Replication Rate.—Suppose we extend the definition of the replication rate such that insignificant original results are counted as ‘successfully replicated’ if they are also insignificant in replications. Assume replication selection is a random sample of published results. Then we have the following definitions:

Definition I1 (Generalized Replication Probability of Individual Study). *The replication probability of a study (X, Σ, Θ) which is published ($D = 1$) and chosen for replication ($R = 1$) is*

$$\widetilde{RP}(X, \Theta, \sigma_r(X, \Sigma, \beta)) = \begin{cases} \mathbb{P}\left(\frac{|X_r|}{\sigma_r(X, \Sigma, \beta)} \geq 1.96, \text{sign}(X) = \text{sign}(X_r) \mid X, \Theta, \sigma_r(X, \Sigma, \beta)\right) & \text{if } 1.96.\Sigma \leq |X| \\ \mathbb{P}\left(\frac{|X_r|}{\sigma_r(X, \Sigma, \beta)} < 1.96 \mid X, \Theta, \sigma_r(X, \Sigma, \beta)\right) & \text{if } 1.96.\Sigma > |X| \end{cases} \quad (71)$$

Definition I2 (Expected Generalized Replication Probability). *The expected generalized replication probability equals*

$$\begin{aligned} \mathbb{E}\left[\widetilde{RP}(X, \Theta, \sigma_r(X, \Sigma, \beta))\right] &= \mathbb{P}(1.96.\Sigma \leq |X|) \mathbb{E}\left[\widetilde{RP}(X, \Theta, \sigma_r(X, \Sigma, \beta) \mid X, \Theta, \sigma_r(X, \Sigma, \beta), 1.96.\Sigma \leq |X|)\right] \\ &+ \left(1 - \mathbb{P}(1.96.\Sigma \leq |X|)\right) \mathbb{E}\left[\widetilde{RP}(X, \Theta, \sigma_r(X, \Sigma, \beta) \mid X, \Theta, \sigma_r(X, \Sigma, \beta), 1.96.\Sigma > |X|)\right] \end{aligned} \quad (72)$$

First, note that Definition I2 equals the standard replication rate definition when the expectation is taken only over significant studies because, in this case, $\mathbb{P}(|X| \leq 1.96.\Sigma) = 0$. Thus, the degree to which the expected generalized replication probability differs from the standard expected replication probability depends on two factors. First, the share of published results that are insignificant. Second, the expected probability that replications will be insignificant conditional on original estimates being insignificant.³²

Empirical Results.—To analyze the generalized replication rate, we can apply the empirical approach outlined in the main text, but using the generalized definition in place of the original definition. Recall that the original replication rate is invariant to publication bias against null results. The generalized replication rate, by contrast, does vary as the degree of selective publication against null results changes. Thus, two sets of results are presented for comparison. The first set assumes selective publication using estimated selection parameters in Table 1 in the main text. The second set assumes no selective publication (i.e. that all results are published with equal probability). We examine two rules for calculating replication power: the common power rule and the original power rule (where the replication standard error is set equal to the original standard error). For more details on different rules for calculating replication power, see Appendix G.

³²Additionally, note that this definition implies that if $\theta = 0$, then $\widetilde{RP}(X, \Theta, \sigma_r(X, \Sigma, \beta) \mid \Theta = 0) = 0.90375$. That is, the replication probability of null results is constant and independent of power in original studies and replication studies.

Table II reports the results for both applications. Under the common power rule, the simulated generalized replication rate remains below intended power in both publication regimes. Under the original power rule, it is relatively low when there is selective publication and around 80% when there is no selective publication.

These generalized replication rate predictions differs from the standard replication rate predictions for two reasons: (i) the share of insignificant results in the published literature and (ii) the replication probability when results are insignificant, which depends on the power rule used in replication studies. On the first point, moving from the selective publication regime to the no selective publication regime implies a dramatic increase in the share of insignificant published results; in both applications, null results change from a minority of published results to a majority. On the second point, the results show that the replication power rules considered here have some undesirable properties. First, note that the common power rule is designed to detect original estimates with high statistical power. This implies that low-powered, insignificant original results will be high-powered in replications, which increases the probability that they are significant and thus counted as replication ‘failures’ under the generalized definition. The original power rule has the reverse problem. On the one hand, low-powered, insignificant original studies are likely to be insignificant in replications, which counts as a ‘successful’ replication under the generalized definition. However, on the other hand, low-powered, significant original studies will have low replication probabilities when the same low-powered design is repeated in replications. The generalized replication rate therefore depends crucially on the share of significant and insignificant findings in the published literature, and the distribution of standard errors. Under the original power rule with no selective publication, the generalized replication rate is around 80% in both applications; however, with greater power in original studies, the replication rate would fall.

While the generalized replication rate changes as selective publication is reduced, the direction of this change depends on which replication power rule is used: with the original power rule the replication rate increases, while with the common power rule it decreases.

Overall, generalizing the replication rate with Definition H2 does not deliver replication rates close to intended power under the common power rule. For the original power rule, it is higher when there is no selective publication because replications repeat low-power designs for low-powered original studies with insignificant results. The generalized replication rate under this original power rule will therefore be sensitive to the distribution of power in original studies.

TABLE II – PREDICTED GENERALIZED REPLICATION RATE RESULTS

A Economics experiments	<i>Simulated statistics</i>	
	92% for X	Original power
<i>Selective publication</i>		
Generalized replication rate	0.600	0.553
$\mathbb{P}(\text{Replicated} S_X = 1)$	0.600	0.551
$\mathbb{P}(\text{Replicated} S_X = 0)$	0.574	0.789
$\mathbb{P}(S_X = 1)$	0.993	0.993
$\mathbb{P}(S_X = 0)$	0.007	0.007
<i>No selective publication</i>		
Generalized replication rate	0.432	0.773
$\mathbb{P}(\text{Replicated} S_X = 1)$	0.582	0.515
$\mathbb{P}(\text{Replicated} S_X = 0)$	0.378	0.867
$\mathbb{P}(S_X = 1)$	0.268	0.268
$\mathbb{P}(S_X = 0)$	0.732	0.732
B Psychology experiments		
<i>Selective Publication</i>		
Generalized replication rate	0.546	0.526
$\mathbb{P}(\text{Replicated} S_X = 1)$	0.544	0.487
$\mathbb{P}(\text{Replicated} S_X = 0)$	0.563	0.839
$\mathbb{P}(S_X = 1)$	0.890	0.890
$\mathbb{P}(S_X = 0)$	0.110	0.110
<i>No selective publication</i>		
Generalized replication rate	0.490	0.798
$\mathbb{P}(\text{Replicated} S_X = 1)$	0.535	0.469
$\mathbb{P}(\text{Replicated} S_X = 0)$	0.478	0.886
$\mathbb{P}(S_X = 1)$	0.209	0.209
$\mathbb{P}(S_X = 0)$	0.791	0.791

Notes: Economics experiments refer to [Camerer et al. \(2016\)](#) and psychology experiments to [Open Science Collaboration \(2015\)](#). The generalized replication rate is defined in the text. The indicator variable S_X equals one for significant results and zero otherwise. Economics experiments refers to [Camerer et al. \(2016\)](#) and psychology experiments to [Open Science Collaboration \(2015\)](#). Simulated statistics are based on parameter estimates in Table 1 in the main text. Different column represent different rules for calculating power in replications.

J. Prediction Interval Approach

This appendix provides details on the prediction interval approach from [Patil et al. \(2016\)](#) and provides some guidance on implementing it for replicators.

The prediction interval approach in [Patil et al. \(2016\)](#) tests the following null hypothesis:

$$H_0 : X \sim N(\Theta, \Sigma^2) \text{ and } X_r \sim N(\Theta, \Sigma_r^2)$$

This captures many of the primary concerns of replicators. For instance, this null is false if the original study is biased due to publication bias of p -hacking, or if there is treatment effect heterogeneity such that the true effect in the original and replication study differ.

Under this setup, [Patil et al. \(2016\)](#) test whether the difference in the original and replication

estimates, $X - X_r$ is statistically different from zero at the 5% level. Assuming $X \perp\!\!\!\perp X_r$, we have $X - X_r \sim N(0, \tilde{\Sigma}^2)$, with $\tilde{\Sigma} \equiv \sqrt{\Sigma^2 + \Sigma_r^2}$. Thus, under the null hypothesis, it follows that:

$$\mathbb{P}(-1.96\tilde{\Sigma} \leq X - X_r \leq 1.96\tilde{\Sigma}) = \mathbb{P}(X - 1.96\tilde{\Sigma} \leq X_r \leq X + 1.96\tilde{\Sigma}) = 0.95 \quad (73)$$

A replication estimate X_r is therefore deemed to be statistically consistent with original estimate X under the null if the replication estimate lies within the ‘unnormalized prediction interval’ $(X - 1.96\tilde{\Sigma}, X + 1.96\tilde{\Sigma})$. Crucially, the power to reject this test depends on the statistical power of both the original study and the replication study. Specifically, if either original or replication studies have low power, then it may be difficult to detect even large differences in original and replication estimates.

In the main text, I normalize the prediction interval to be in terms of relative effect sizes. The motivation behind this is simply that relative effect sizes are commonly reported in the replication literature, making them more convenient to interpret and compare across studies. To do this, simply divide the interval bounds by the original estimate X . This implies that a replication is deemed statistically consistent with the original study if the relative effect size lies within the (normalized) prediction interval i.e. $\frac{X_r}{X} \in (1 - 1.96\frac{\tilde{\Sigma}}{X}, 1 + 1.96\frac{\tilde{\Sigma}}{X})$. Thus, when the relative effect size deviates sufficiently far from one, we can reject the null hypothesis that both estimates come from a normal distribution centered at the same true effect.

Note that for the replication of null results, it may be more appropriate to use the unnormalized prediction interval approach, namely, to reported whether $X_r \in (X - 1.96\tilde{\Sigma}, X + 1.96\tilde{\Sigma})$. This is because dividing by an original estimate that is close to zero might lead to extreme values which are difficult to interpret.

J.1. Implementation

When implementing the prediction interval approach, care must be taken with which units are used. Large-scale replication studies commonly convert all findings across different studies into a standardized unit for comparing effect sizes (and relative effect sizes). For example, in all three empirical applications in the main text, effect sizes are reported in correlation coefficient units. However, inference is typically done by applying a Fisher transformation to correlation coefficients, since this transformation provides approximately normally distributed random variables for hypothesis testing (Fisher, 1915). The prediction interval approach is also based on normally distributed estimates. Thus, reporting prediction intervals in their original units requires a conversion back to the original units.

To illustrate how this could be done in practice, suppose the replicator observes original and replication outcomes (ρ, N, ρ_r, N_r) , where (ρ, ρ_r) denote the correlation coefficients of the

original and replication study, respectively; and (N, N_r) denote their respective sample sizes. Suppose we are interested in applying the prediction interval approach, as recommended in the main text. Doing so involves the following steps:

1. **Normalization:** normalize the original estimate to be positive. For example, if the unnormalized correlation coefficients in the original and replication study have the same sign, then both would be normalized as positive. If, on the other hand, they have opposite signs then the original coefficient would be positive and the replication coefficient negative.
2. **Fisher transformed variables:** calculate $x = \arctan \rho$; $\sigma = 1/\sqrt{N-3}$; $x_r = \arctan \rho_r$; and $\sigma_r = 1/\sqrt{N_r-3}$.³³ Define $\tilde{\sigma} \equiv \sqrt{\sigma^2 + \sigma_r^2}$.
3. **Calculate prediction interval:** finally, convert the prediction interval into the original correlation coefficient units using the following formula:³⁴

$$\left(\frac{\tanh(x - 1.96\tilde{\sigma})}{\tanh x}, \frac{\tanh(x + 1.96\tilde{\sigma})}{\tanh x} \right)$$

and test whether the relative effect size ρ_r/ρ falls within this range.

³³The formula for the standard deviation is from [Fisher \(1915\)](#).

³⁴Note that normalization of the original effect to be positive implies $\tanh x > 0$.