


**Bayesian approaches to designing replication studies**

Samuel Pawel<sup>1</sup>, Guido Consonni<sup>2</sup>, and Leonhard Held<sup>1</sup>


<sup>1</sup>Department of Biostatistics, Center for Reproducible Science, University of Zurich

<sup>2</sup>Dipartimento di Scienze Statistiche, Università Cattolica del Sacro Cuore

### Author Note

Samuel Pawel  <https://orcid.org/0000-0003-2779-320X>

Guido Consonni  <https://orcid.org/0000-0002-1252-5926>

Leonhard Held  <https://orcid.org/0000-0002-8686-5325>

Preprint version May 9, 2023. Licensed under CC-BY 4.0  
(<https://creativecommons.org/licenses/by/4.0/>).

This research has not been preregistered. A preprint has been previously published on arXiv (<https://doi.org/10.48550/arXiv.2211.02552>) and included in the PhD thesis of Samuel Pawel. We declare that we have no conflicts of interest.

This work was supported by the Swiss National Science Foundation (#189295). The funder had no role in study design, data collection, data analysis, data interpretation, decision to publish, or preparation of the manuscript.

All our analyses were conducted in the R programming language version 4.3.0 (R Core Team, 2023). Code to reproduce this manuscript is available at <https://github.com/SamCH93/BAtDRS>. A snapshot of the Git repository at the time of writing is archived at <https://doi.org/10.5281/zenodo.7291076>. Methods for Bayesian design of replication studies are implemented in the R package BayesRepDesign which is available at <https://CRAN.R-project.org/package=BayesRepDesign>.

We thank Charlotte Micheloud and Angelika Stefan for helpful comments on drafts of the manuscript. We thank the anonymous reviewer for constructive comments. Our acknowledgment of these individuals does not imply their endorsement of this article.

We thank Protzko et al. (2020) for publicly sharing their data. Their CC-BY 4.0 licensed data were downloaded from <https://osf.io/42ef9/>. The R markdown script “Decline effects main analysis.Rmd” was executed and the relevant variables from the objects “ES\_experiments” and “decline\_effects” were saved.

Correspondence concerning this article should be addressed to Samuel Pawel, Epidemiology, Biostatistics and Prevention Institute, University of Zurich, Hirschengraben 84, 8001 Zurich, Switzerland. E-mail: [samuel.pawel@uzh.ch](mailto:samuel.pawel@uzh.ch)

### **Abstract**

Replication studies are essential for assessing the credibility of claims from original studies. A critical aspect of designing replication studies is determining their sample size; a too small sample size may lead to inconclusive studies whereas a too large sample size may waste resources that could be allocated better in other studies. Here we show how Bayesian approaches can be used for tackling this problem. The Bayesian framework allows researchers to combine the original data and external knowledge in a design prior distribution for the underlying parameters. Based on a design prior, predictions about the replication data can be made, and the replication sample size can be chosen to ensure a sufficiently high probability of replication success. Replication success may be defined by Bayesian or non-Bayesian criteria, and different criteria may also be combined to meet distinct stakeholders and enable conclusive inferences based on multiple analysis approaches. We investigate sample size determination in the normal-normal hierarchical model where analytical results are available and traditional sample size determination is a special case where the uncertainty on parameter values is not accounted for. An application to data from a multisite replication project of social-behavioral experiments illustrates how Bayesian approaches help to design informative and cost-effective replication studies. Our methods can be used through the R package BayesRepDesign.

*Keywords:* Bayesian design, design prior, multisite replication, sample size determination

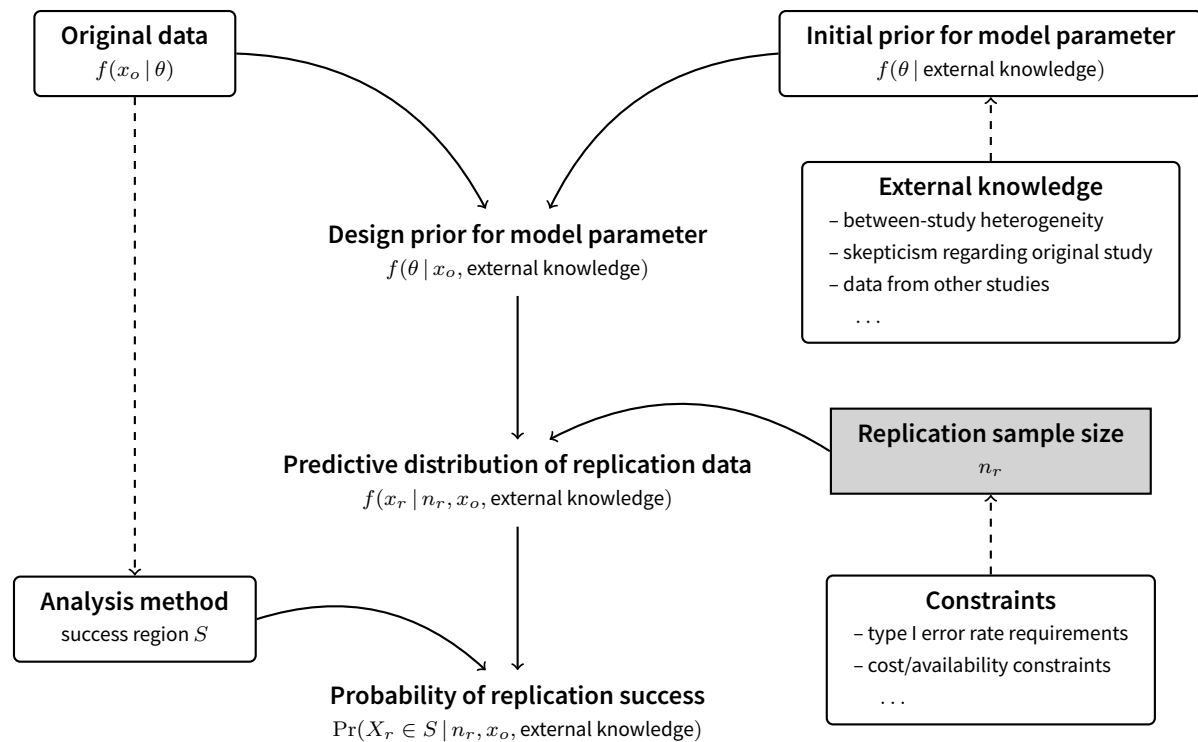
## Bayesian approaches to designing replication studies

### Introduction

The replicability of research findings is a cornerstone for the credibility of science. However, there is growing evidence that the replicability of many scientific findings is lower than expected (Camerer et al., 2018; Errington et al., 2021; Open Science Collaboration, 2015). This “replication crisis” has led to methodological reforms in various fields of science, one of which is an increased conduct of replication studies (Munafò et al., 2017). Statistical methodology plays a key role in the evaluation of replication studies, and various methods have been proposed for quantifying how “successful” a replication study was in replicating the original finding (Anderson & Maxwell, 2016; Bayarri & Mayoral, 2002; Bonett, 2020; Etz & Vandekerckhove, 2016; Harms, 2019; Hedges & Schauer, 2019; Held, 2020; Held, Micheloud, & Pawel, 2022; Johnson, Payne, Wang, Asher, & Mandal, 2016; Ly, Etz, Marsman, & Wagenmakers, 2018; Mathur & VanderWeele, 2020; Patil, Peng, & Leek, 2016; Pawel & Held, 2020, 2022; Simonsohn, 2015; van Aert & van Assen, 2017; Verhagen & Wagenmakers, 2014, among others). Yet, as with ordinary studies, statistical methodology is not only important for analyzing replication studies but also for designing them, in particular for their *sample size determination* (SSD). Optimal SSD is important since too small sample sizes may lead to inconclusive studies, whereas too large sample sizes may waste resources which could have been allocated better in other research projects.

SSD for replication studies comes with unique opportunities and challenges; the data from the original study can be used to inform SSD, at the same time the analysis of replication success based on original and replication study is typically different from an analysis of a single study for which traditional SSD methodology was developed. Since analysis and design of replication studies should be in accordance, a relatively small literature has emerged which specifically deals with replication study power calculations and SSD (Anderson & Kelley, 2022; Anderson & Maxwell, 2017; Bayarri & Mayoral, 2002; Goodman, 1992; Hedges & Schauer, 2021; Held, 2020; Micheloud & Held, 2022; Pawel & Held, 2022; Senn, 2002; van Zwet & Goodman, 2022). However, most of these articles only deal with selected analysis methods and data models. An exception is the excellent article by Anderson and Kelley (2022)

which discusses more general principles of replication SSD in the context of psychological research, mostly from a frequentist perspective. As they state “the literature on Bayesian sample size planning is still nascent, particularly with respect to Bayes Factors (Schönbrodt & Wagenmakers, 2018), and has not yet been clearly optimized for the context of most replication goals” (Anderson & Kelley, 2022, p. 18). Our goal is therefore to complement their article by developing a unified framework of replication SSD (schematically illustrated in Figure 1) based on principles from Bayesian design approaches (De Santis, 2004; Gelfand & Wang, 2002; Grieve, 2022; Kunzmann et al., 2021; O’Hagan & Stevens, 2001; Park & Pek, 2022; Pek & Park, 2019; Schönbrodt & Wagenmakers, 2018; Spiegelhalter, Abrams, & Myles, 2004; Spiegelhalter & Freedman, 1986; Spiegelhalter, Freedman, & Blackburn, 1986; Weiss, 1997). We aim to provide both a theoretical basis for methodologists developing new methods for design and analysis methods of replication studies, and also to illustrate how Bayesian design approaches can practically be used by researchers planning a replication study.



**Figure 1**

*Schematic illustration of Bayesian sample size determination for replication studies. The original and replication data are denoted by  $x_o$  and  $x_r$ , respectively. Both are assumed to come from a distribution with density/probability mass function denoted by  $f(x_i | \theta)$  for  $i \in \{o, r\}$ . An initial prior with density function  $f(\theta | \text{external knowledge})$  is assigned to the model parameter  $\theta$ .*

The design of replication studies is a natural candidate for Bayesian knowledge updating as it allows to combine uncertain information from different sources—for instance, the data from the original study and/or expert knowledge—in a so-called *design prior* distribution for the underlying model parameters. If the analysis of the replication data is also Bayesian, the design prior may be different from the so-called *analysis prior* which, unlike the design prior, is usually desired to be objective or “uninformative” (O’Hagan & Stevens, 2001). Based on the design prior, predictions about the replication data can then be made, and the sample size can be chosen such that the probability of replication success becomes sufficiently high. Importantly, Bayesian design approaches can also be used if the planned analysis of the replication study is non-Bayesian, which is the more common situation in practice. Bayesian design based on a frequentist analysis is known under various names, such as “hybrid classical-Bayesian design” (Spiegelhalter et al., 2004) or “Bayesian assurance” (O’Hagan, Stevens, & Campbell, 2005), and has also been used before for psychological applications (Park & Pek, 2022; Pek & Park, 2019) and replication studies (Anderson & Maxwell, 2017; Micheloud & Held, 2022).

This paper is structured as follows: We start with presenting a general framework for Bayesian SSD of replication studies which applies to any kind of data model and analysis method. We then investigate design priors and SSD in the normal-normal hierarchical model framework which provides sufficient flexibility for incorporating the original data and external knowledge in replication design. No advanced computational methods, such as (Markov Chain) Monte Carlo sampling, are required for conducting Bayesian SSD in this framework, and in many cases there are even simple formulae which generalize classical power and sample size calculations. We illustrate the methodology for several Bayesian and non-Bayesian analysis methods, and for both singlesite and multisite replication studies. Since multisite replication studies are becoming increasingly popular in psychology (e.g., Klein et al., 2018), we also discuss how to choose the optimum allocation of samples within and between sites from a Bayesian design point of view. As a running example we use data from a multisite replication project of social-behavioral experiments (Protzko et al., 2020). Finally, we close with concluding remarks, limitations, and open questions.

### General framework

Suppose an original study has been conducted and resulted in a data set  $x_o$ . These data are assumed to come from a distribution characterized by an unknown parameter  $\theta$  and with density function  $f(x_o | \theta)$ . To assess the replicability of a claim from the original study, an independent and identically designed (apart from the sample size) replication study is conducted, and the goal of the design stage is to determine its sample size  $n_r$ .

As the observed original data  $x_o$ , the yet unobserved replication data  $X_r$  are assumed to come from a distribution depending on the parameter  $\theta$ . The parameter  $\theta$  thus provides a link between the two studies, and the knowledge obtained from the original study can be used to make predictions about the replication. The central quantity for doing so is the so-called *design prior* of the parameter  $\theta$ , which we write as the posterior distribution of  $\theta$  based on the original data and an *initial prior* for  $\theta$

$$f(\theta | x_o, \text{external knowledge}) = \frac{f(x_o | \theta) f(\theta | \text{external knowledge})}{f(x_o | \text{external knowledge})}. \quad (1)$$

The initial prior of  $\theta$  may depend on external knowledge (e.g., data from other studies) and it represents the uncertainty about  $\theta$  before observing the original data. We will discuss common types of external knowledge in the replication setting in the next Section. The design prior (1) hence represents the state of knowledge and uncertainty about the parameter  $\theta$  before the replication is conducted, and, along with an assumed replication sample size  $n_r$ , it can be used to compute a predictive distribution for the replication data

$$f(x_r | n_r, x_o, \text{external knowledge}) = \int f(x_r | n_r, \theta) f(\theta | x_o, \text{external knowledge}) d\theta. \quad (2)$$

After completion of the replication, the observed data  $x_r$  will be analyzed in some way to quantify to what extent the original result could be replicated. The analysis may involve the original data (for example, a meta-analysis of the two data sets) or it may only use the replication data. Typically, there is a *success region*  $S$  which implies that if the replication data are contained within it ( $x_r \in S$ ), the replication is successful. The *probability of replication success* can thus be computed by integrating the predictive density (2) over  $S$ . To

ensure a sufficiently conclusive replication design, the sample size  $n_r$  is determined such that the probability of replication success is at least as high as a desired target probability of success, here and henceforth denoted by  $1 - \beta$ . The required sample size  $n_r^*$  is then the smallest sample size which leads to a probability of replication success of at least  $1 - \beta$ , i.e.,

$$n_r^* = \inf \{n_r : \Pr(X_r \in S \mid n_r, x_o, \text{external knowledge}) \geq 1 - \beta\}. \quad (3)$$

Often, replication studies are analyzed using several methods which quantify different aspects of replicability, and which have different success regions (e.g., one method for quantifying parameter compatibility and another for quantifying evidence against a null hypothesis). In this case, the sample size may be chosen such that the probability of replication success is as high as desired for all planned analysis methods.

There may sometimes be certain constraints which the replication sample size needs to satisfy. For instance, in most cases there is an upper limit on the sample size due to limited resources and/or availability of samples. Moreover, funders and regulators may also require methods to be *calibrated* (Grieve, 2016), that is, to have appropriate type I error rate control. The sample size  $n_r^*$  may thus also need to satisfy a type I error rate not higher than some required level.

### **Sample size determination in the normal-normal hierarchical model**

We will now illustrate the general methodology from the previous section in the *normal-normal hierarchical model* where predictive distributions and the probability of replication success can often be expressed in closed-form, permitting further insight. It is pragmatic to adopt a meta-analytic perspective and use only study level summary statistics instead of the raw study data since the raw data from the original study are not always available to the replicators. Typically, the underlying parameter  $\theta$  is a univariate effect size quantifying the effect on the outcome variable (e.g., a mean difference, a log odds ratio, or a log hazard ratio). The original and replication study can then be summarized through an effect estimate  $\hat{\theta}$ , possibly the maximum likelihood estimate, and a corresponding standard error  $\sigma$ , i.e.,  $x_o = \{\hat{\theta}_o, \sigma_o\}$  and  $x_r = \{\hat{\theta}_r, \sigma_r\}$ . Effect estimates and standard errors are routinely



reported in research articles or can, under some assumptions, be computed from  $p$ -values and confidence intervals. As in the conventional meta-analytic framework (Sutton & Abrams, 2001), we further assume that for study  $k \in \{o, r\}$  the (suitably transformed) effect estimate  $\hat{\theta}_k$  is approximately normally distributed around a study specific effect size  $\theta_k$  and with (known) variance equal to its squared standard error  $\sigma_k^2$ , here and henceforth denoted by  $\hat{\theta}_k | \theta_k \sim N(\theta_k, \sigma_k^2)$ . The standard error  $\sigma_k$  is typically of the form  $\sigma_k = \lambda/\sqrt{n_k}$  with  $\lambda^2$  some unit variance and  $n_k$  the sample size. The ratio of the original to the replication variance is thus the ratio of the replication to the original sample size

$$c = \sigma_o^2/\sigma_r^2 = n_r/n_o,$$

which is often the main focus of SSD as it quantifies how much the replication sample  $n_r$  size needs to be changed compared to the original sample size  $n_o$ . Depending on the effect size type, this framework might require slight modifications (see e.g., Spiegelhalter et al., 2004, section 2.4).

Assuming a normal sampling model for the effect estimates (4a), as described previously, and specifying an initial hierarchical normal prior for the study specific effect sizes (4b) and the effect size (4c), leads then to the normal-normal hierarchical model

$$\hat{\theta}_k | \theta_k \sim N(\theta_k, \sigma_k^2) \tag{4a}$$

$$\theta_k | \theta \sim N(\theta, \tau^2) \tag{4b}$$

$$\theta \sim N(\mu_\theta, \sigma_\theta^2). \tag{4c}$$

By marginalizing over the study specific effects sizes, the model (4) can alternatively be expressed as

$$\hat{\theta}_k | \theta \sim N(\theta, \sigma_k^2 + \tau^2) \tag{5a}$$

$$\theta \sim N(\mu_\theta, \sigma_\theta^2) \tag{5b}$$

which is often more useful for derivations and computations. In the following we will explain

how the normal-normal hierarchical model can be used for SSD of the replication study.

### Design prior and predictive distribution

The observed original data  $x_o = \{\hat{\theta}_o, \sigma_o\}$  can be combined with the initial prior (5b) by standard Bayesian theory for normal prior and likelihood (Spiegelhalter et al., 2004, section 3.7) to obtain a posterior distribution for the effect size  $\theta$

$$\theta | \hat{\theta}_o, \sigma_o^2 \sim N \left( \frac{\hat{\theta}_o}{1 + 1/g} + \frac{\mu_\theta}{1 + g}, \frac{\sigma_o^2 + \tau^2}{1 + 1/g} \right) \quad (6)$$

where  $g = \sigma_\theta^2 / (\sigma_o^2 + \tau^2)$  is the *relative prior variance*. This posterior serves then as the design prior for predicting the replication data.

It is interesting to contrast the design prior (6) to the “conditional” design prior (Micheloud & Held, 2022), that is, to assume that the unknown effect size  $\theta$  corresponds to the original effect estimate  $\hat{\theta}_o$ . This is a standard approach in practice, for instance, Open Science Collaboration (2015) determined the sample sizes of its 100 replications under this assumption. In our framework it implies that the normal design prior (6) becomes a point mass at the original effect estimate  $\hat{\theta}_o$ , which can either be achieved through overwhelmingly informative original data ( $\sigma_o^2 \downarrow 0$ ) along with no heterogeneity ( $\tau^2 = 0$ ), or through an overwhelmingly informative initial prior ( $g \downarrow 0$ ) centered around the original effect estimate ( $\mu_\theta = \hat{\theta}_o$ ). Both cases show that from a Bayesian perspective the standard approach is unnatural as it either corresponds to making the standard error  $\sigma_o$  smaller than it actually was, or to cherry-picking the prior based on the data.

Based on the design prior (6), a predictive distribution for the replication effect estimate  $\hat{\theta}_r$  can be computed. Specifically, assuming a replication standard error  $\sigma_r$  and integrating the marginal density of the replication effect estimate (5a) with respect to the prior density leads to

$$\hat{\theta}_r | \hat{\theta}_o, \sigma_o^2, \sigma_r^2 \sim N \left( \mu_{\hat{\theta}_r} = \frac{\hat{\theta}_o}{1 + 1/g} + \frac{\mu_\theta}{1 + g}, \sigma_{\hat{\theta}_r}^2 = \sigma_r^2 + \tau^2 + \frac{\sigma_o^2 + \tau^2}{1 + 1/g} \right), \quad (7)$$

which can again be shown using standard Bayesian theory (Spiegelhalter et al., 2004, section 3.13.3). The design prior (6) and the resulting predictive distribution (7) depend on the

parameters of the initial prior  $(\tau^2, \mu_\theta, \sigma_\theta^2)$ . We will now explain how these parameters can be specified based on external knowledge.

### **Incorporating external knowledge in the initial prior**

At least three common types of external knowledge can be distinguished in the replication setting: (i) expected heterogeneity between original and replication study due to differences in study design, execution, and population, (ii) prior knowledge about the effect size either from theory or from related studies, (iii) skepticism regarding the original study due to the possibility of exaggerated results.

#### ***Between-study heterogeneity***

The expected degree of between-study heterogeneity can be incorporated via the variance  $\tau^2$  in (4b). As  $\tau^2$  decreases, the study specific effect sizes become more similar, whereas for increasing  $\tau^2$  they become more unrelated. If the replicators do not expect any heterogeneity they can thus set  $\tau^2 = 0$  which will lead to the model collapsing to a common effect model.

If heterogeneity is expected, there are different approaches for specifying  $\tau^2$ . A domain expert may subjectively assess how much heterogeneity is to be expected due to the change in laboratory, study population, and other factors. An alternative is to take an estimate from the literature, e.g., from multisite replication projects or from systematic reviews. Finally, one can also specify an upper limit of “tolerable heterogeneity”. This approach is similar to specifying a minimal clinically relevant difference in classical power analysis in the sense that a true replication effect size which is intolerably heterogeneous from the original effect size is not relevant to be detected. An absolute (Spiegelhalter et al., 2004, section 5.7.3) and a relative approach (Held & Pawel, 2020) can be considered. In the absolute approach, a value of  $\tau^2$  is chosen such that a suitable range of study-specific effect sizes is not larger than an effect size difference considered negligible. For example, when 95% of the study specific effect sizes should not vary more than a small effect size e.g.,  $d = 0.2$  on standardized mean difference scale based on the Cohen (1992) effect size classification, this would lead to  $\tau = d/(2 \cdot 1.96) \approx 0.05$ . In the relative approach,  $\tau^2$  is specified relative to the variance of the original estimate  $\sigma_o^2$  using field conventions for tolerable relative heterogeneity. For example,

in the Cochrane guidelines for systematic reviews (Deeks, Higgins, & Altman, 2019) a value of  $I^2 = \tau^2 / (\tau^2 + \sigma_o^2) = 40\%$  is classified as “negligible”, which translates to  $\tau^2 = \sigma_o^2 / (1/I^2 - 1) = (2\sigma_o^2)/3$ .

We note that one can also assign a prior distribution to  $\tau^2$ . For an overview of prior distributions for heterogeneity variances in the normal-normal hierarchical model see Röver et al. (2021). In this case there is no closed-form expression for the predictive distribution of the replication effect estimate but numerical or Monte Carlo integration need to be used. We illustrate in the supplement how the probability of replication success can be computed in this case. The derived closed-form expressions conditional on  $\tau^2$  are still useful as they enable computation of the predictive distribution up to a one-dimensional integral which can be computed numerically.

### ***Knowledge about the effect size***

Prior knowledge about the effect size  $\theta$  can be incorporated via the prior mean  $\mu_\theta$  and the prior variance  $\sigma_\theta^2$  in (4c). For instance, the parameters may be specified based on a meta-analysis of related studies (McKinney, Stefan, & Gronau, 2021) or based on expert elicitation (O’Hagan, 2019). The resulting design prior will then contain more information than what was provided by the original data alone, leading to potentially more efficient designs. If there is no prior knowledge available, a standard approach is to specify an (improper) flat prior by letting the variance go to infinity ( $\sigma_\theta^2 \rightarrow \infty$ ). The resulting design prior will then only contain the information from the original study.

### ***Exaggerated original results***

Potentially exaggerated original results can be counteracted by setting  $\mu_\theta = 0$  which shrinks the design prior towards smaller effect sizes (in absolute value) than the observed effect estimate  $\hat{\theta}_o$ . For instance, replicators could believe that the results from the original study are exaggerated because there is no preregistered study protocol available. Even without such beliefs, weakly informative shrinkage priors may also be motivated from a “regularization” point of view as they can correct for statistical biases (Copas, 1983; Firth, 1993) or prevent unreasonable parameter values from taking over the posterior in settings with uninformative data (Gelman, 2009).

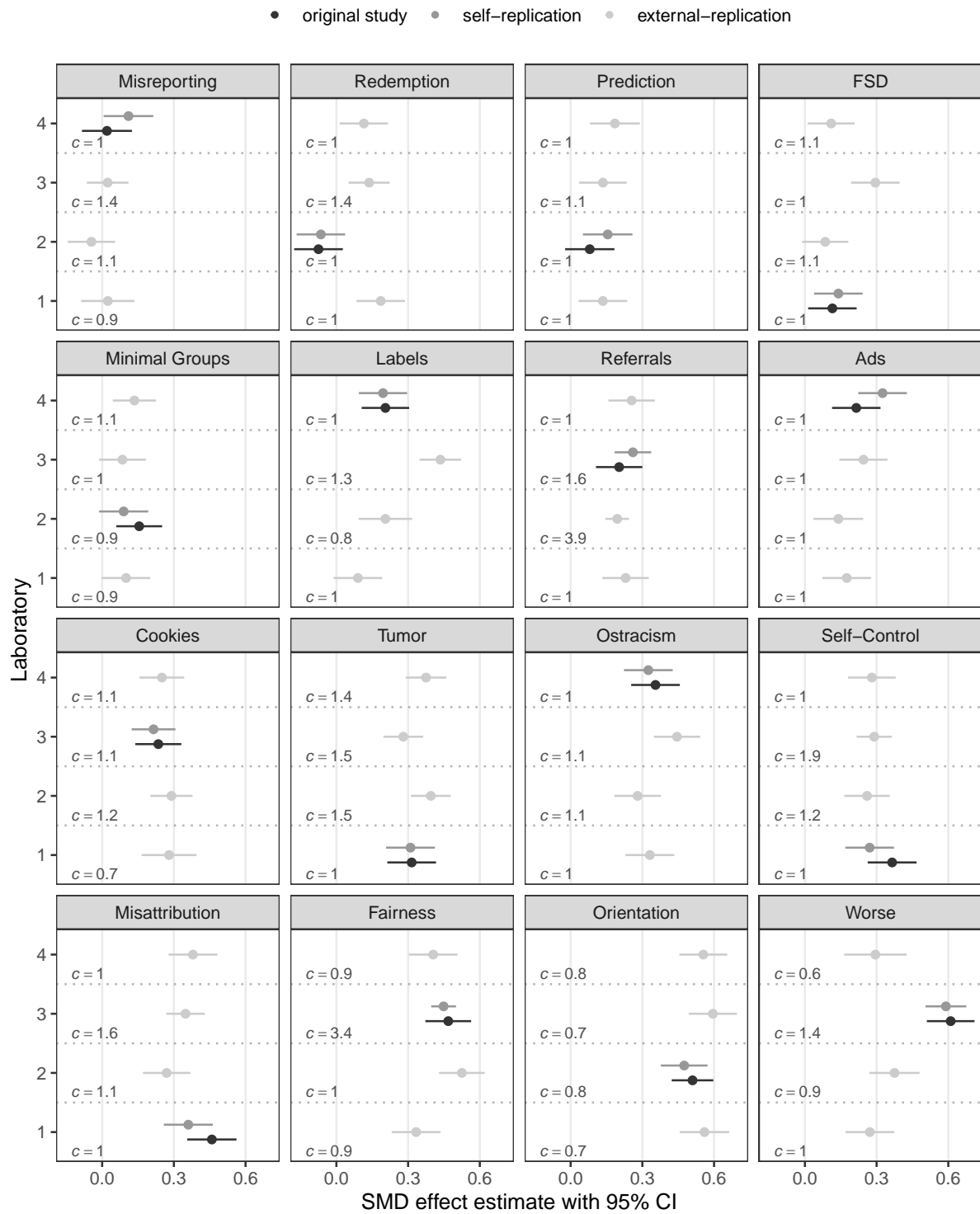
The amount of shrinkage is determined via the prior variance  $\sigma_\theta^2$ . A flat prior ( $\sigma_\theta^2 \rightarrow \infty$ ) will lead to no shrinkage, while a highly concentrated prior ( $\sigma_\theta^2 \downarrow 0$ ) will completely shrink the design prior to a point mass on zero. One option for specifying  $\sigma_\theta^2$  is to use an estimate from a corpus of related studies. For instance, van Zwet, Schwab, and Senn (2021) used the Cochrane library of systematic reviews to specify design priors for hypothetical replication studies of RCTs. If no corpus is available, a pragmatic alternative is to use the empirical Bayes estimate based on the original data

$$\hat{\sigma}_\theta^2 = \max\{(\hat{\theta}_o - \mu_\theta)^2 - \tau^2 - \sigma_o^2, 0\}. \quad (8)$$

The estimate (8) will lead to adaptive shrinkage (Pawel & Held, 2020) in the sense that shrinkage is large for unconvincing original studies (those with small effect estimates in absolute value  $|\hat{\theta}_o|$  and/or large standard errors  $\sigma_o$ ), but disappears as the data become more convincing (through larger effect estimates in absolute value  $|\hat{\theta}_o|$  and/or smaller standard errors  $\sigma_o$ ).

### **Example: Cross-laboratory replication project**

We will now illustrate the construction of design priors based on data from a recently conducted replication project (Protzko et al., 2020), see Figure 2 for a summary of the data. The data were collected in four laboratories which, over the course of five years, conducted their typical social-behavioral experiments on topics such as psychology, communication, or political science. From the experiments conducted in this period, each lab submitted four original findings to be replicated. For instance, the original finding from the “Labels” experiment was: “When a researcher uses a label to describe people who hold a certain opinion, he or she is interpreted as disagreeing with those attributes when a negative label is used and agreeing with those attributes when a positive label is used” (Protzko et al., 2020, p. 17), which was based on an effect estimate  $\hat{\theta}_o = 0.205$  with 95% confidence interval from 0.11 to 0.3. For each submitted original finding, four replication studies were then carried out, one by the same lab (a *self-replication*) and three by the other three labs (three *external-replications*).

**Figure 2**

Data from cross-laboratory replication project by Protzko et al. (2020). Shown are standardized mean difference (SMD) effect estimates with 95% confidence intervals stratified by experiment and laboratory. For each replication study the relative sample size  $c = n_r/n_o$  is shown.

Most studies used simple between-subject designs with two groups and a continuous outcome so that for a study  $i \in \{o, r\}$  the standardized mean difference (SMD) effect estimate  $\hat{\theta}_i$  can be computed from the group means  $\bar{y}_{i1}, \bar{y}_{i2}$ , group standard deviations  $s_{i1}, s_{i2}$ , and group sample sizes  $n_{i1}, n_{i2}$  by

$$\hat{\theta}_i = \frac{\bar{y}_{i1} - \bar{y}_{i2}}{s_i}$$

with  $s_i^2 = \{(n_{i1} - 1)s_{i1}^2 + (n_{i2} - 1)s_{i2}^2\} / (n_{i1} + n_{i2} - 2)$  the pooled sample variance. Under a normal sampling model and assuming equal variances in both groups, the approximate variance of  $\hat{\theta}_i$  is

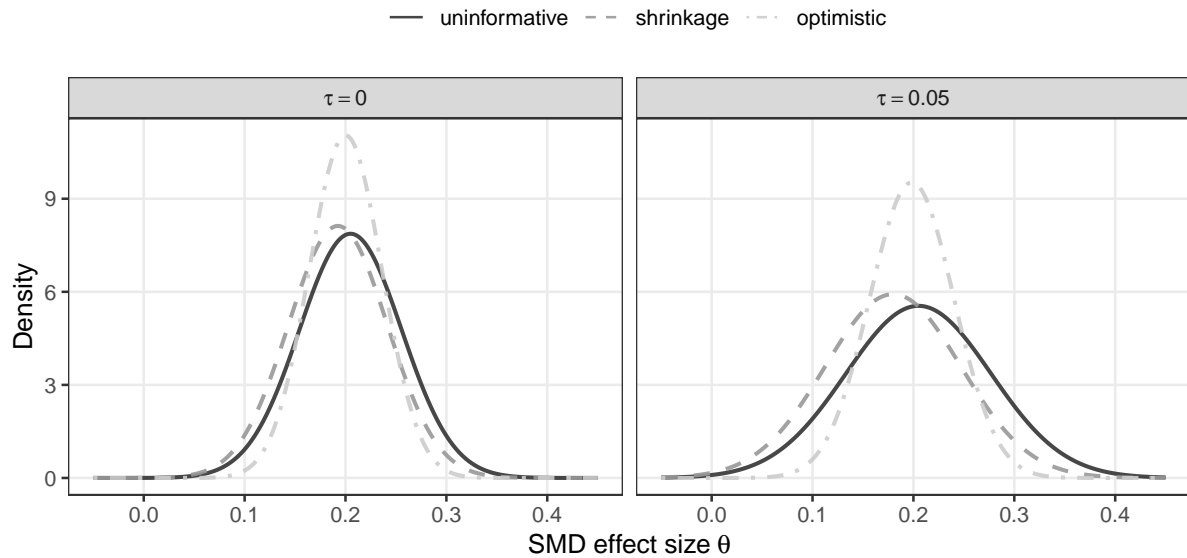
$$\sigma_i^2 = \frac{n_{i1} + n_{i2}}{n_{i1}n_{i2}} + \frac{\hat{\theta}_i^2}{2(n_{i1} + n_{i2})} \quad (9)$$

(Hedges, 1981). A cruder, but for SSD more useful, approximation  $\sigma_i^2 \approx 4/n_i$  is obtained by assuming the same sample size in both groups  $n_{i1} = n_{i2} = n_i/2$ , with  $n_i$  the total sample size, and neglecting the second term in (9) which will be close to zero for small effect estimates and/or large sample sizes (Hedges & Schauer, 2021). We thus have the approximate unit variance  $\lambda^2 = 4$  and the relative variance  $c = \sigma_o^2 / \sigma_r^2 = n_r / n_o$ , which can be interpreted as the ratio of the replication to the original sample size.

Suppose now the original studies have been finished, and we want to conduct SSD for the not yet conducted replication studies. We start by specifying the design priors (one for each replication). Since the original studies have been preregistered, we do not expect an exaggeration of their effect estimates due to selective reporting or other questionable research practices. Therefore, we choose a flat initial prior for  $\theta$ , which leads to design prior and predictive distribution both centered around the original effect estimate  $\hat{\theta}_o$ .

For specifying the between-study heterogeneity variance  $\tau^2$ , a distinction needs to be made between self-replications and external-replications. For self-replications it is reasonable to set  $\tau^2 = 0$  because we would expect no between-study heterogeneity as the experimental conditions will be nearly identical in both studies. In contrast, one would expect some between-study heterogeneity for external-replications as the experimental conditions may

slightly differ between the labs. In the following, we will use  $\tau^2 = 0.05^2$  elicited via the “absolute” approach as discussed previously, so that the range between the 2.5% and the 97.5% quantile of the study specific effect size distribution is equal to a small effect size  $d = 0.2$ .



**Figure 3**

*Design priors for the SMD effect size  $\theta$  in the “Labels” experiment based on the original effect estimate  $\hat{\theta}_o = 0.205$  with standard error  $\sigma_o = 0.051$ . Shown are different choices for the between-study heterogeneity  $\tau$  and the initial prior for the effect size  $\theta$ , “uninformative” corresponds to a flat prior, “shrinkage” corresponds to a zero-mean normal prior with empirical Bayes variance estimate (8), and “optimistic” corresponds to a flat prior updated by the data from a pilot study with effect estimate  $\hat{\theta}_p = 0.195$  and standard error  $\sigma_p = 0.052$ .*

Taken together, we obtain the design prior  $\theta \mid \hat{\theta}_o, \sigma_o^2 \sim N(\hat{\theta}_o, \sigma_o^2)$  for self-replications and the design prior  $\theta \mid \hat{\theta}_o, \sigma_o^2 \sim N(\hat{\theta}_o, \sigma_o^2 + \tau^2)$  with  $\tau^2 = 0.05^2$  for external-replications. For the “Labels” experiment the design prior would be centered around the original effect estimate  $\hat{\theta}_o = 0.205$  with variance  $\sigma_o^2 = 0.05^2$  for a self-replication, and with variance  $\sigma_o^2 + \tau^2 = 0.05^2 + 0.05^2 \approx 0.07^2$  for an external-replication. Figure 3 (dark-gray solid lines) shows the densities of the two priors.

While these two priors seem sensible for the Protzko et al. (2020) data, it is interesting to think about alternative scenarios. If there had been reasons to believe that the original result might be exaggerated, we could have specified an initial shrinkage prior. For instance, the



empirical Bayes estimate for the prior variance  $\sigma_\theta^2$  from (8) leads to a prior whose mean and variance are shrunk towards zero by 12% (medium-gray dashed lines in Figure 3). In contrast, if we had prior knowledge about the effect size  $\theta$  from another study, we could have specified an initial “optimistic” prior. For example, if the self-replication of the “Labels” experiment had been a pilot study and we used its effect estimate  $\hat{\theta}_p = 0.195$  and standard error  $\sigma_p = 0.05$  to specify the initial prior, this would lead to a design prior centered around the weighted mean of original and pilot study, and a prior precision equal to the sum of the precision of both estimates (light-gray dot-dashed lines in Figure 3). Due to the inclusion of the external data, this design prior is much more concentrated than the other two.

### Probability of replication success and required sample size

To compute the probability of replication success one needs to select an analysis method and integrate the predictive distribution (7) over the associated success region  $S$ . There is no universally accepted method for quantifying replicability and here we do not intend to contribute to the debate about the most appropriate method. We will simply show the success regions of different methods, and how the replication sample size can be computed from them. Some methods depend on the direction of the original effect estimate  $\hat{\theta}_o$  and throughout we will assume that it was positive ( $\hat{\theta}_o > 0$ ). Functions for computing the probability of replication success and the required sample size are implemented in the R package BayesRepDesign (see Appendix) for all analysis methods discussed in the following.

#### *The two-trials rule*

The most common approach for the analysis of replication studies is to declare replication success when both the original and replication study lead to a  $p$ -value for testing the null hypothesis  $H_0: \theta = 0$  smaller than a pre-specified threshold  $\alpha$ , usually  $\alpha = 5\%$  for two-sided tests and  $\alpha = 2.5\%$  for one-sided tests. This procedure is known as the *two-trials rule* in drug regulation (Senn, 2008, section 12.2.8).

We now assume that the one-sided original  $p$ -value was significant at some level  $\alpha$ , i.e.,  $p_o = 1 - \Phi(\hat{\theta}_o/\sigma_o) \leq \alpha$ . Replication success at level  $\alpha$  is then achieved if the replication

$p$ -value is also significant, i.e.,  $p_r = 1 - \Phi(\hat{\theta}_r/\sigma_r) \leq \alpha$ , which implies a success region

$$S_{2\text{TR}} = [z_\alpha \sigma_r, \infty),$$

where  $z_\alpha$  is the  $1 - \alpha$  quantile of the standard normal distribution. The probability of replication success is thus given by

$$\Pr(\hat{\theta}_r \in S_{2\text{TR}} | \hat{\theta}_o, \sigma_o, \sigma_r) = \Phi\left(\frac{\mu_{\hat{\theta}_r} - z_\alpha \sigma_r}{\sigma_{\hat{\theta}_r}}\right) \quad (10)$$

with  $\Phi(\cdot)$  the standard normal cumulative distribution function and  $\mu_{\hat{\theta}_r}$  and  $\sigma_{\hat{\theta}_r}$  the mean and standard deviation of the predictive distribution (7). Importantly, by decreasing the standard error  $\sigma_r$  (through increasing the sample size  $n_r$ ), the probability of replication success (10) cannot become arbitrarily high but is bounded from above by

$$\lim \Pr_{2\text{TR}} = \Phi\left(\frac{\mu_{\hat{\theta}_r}}{\sqrt{\tau^2 + (\sigma_o^2 + \tau^2)/(1 + 1/g)}}\right). \quad (11)$$

The required replication standard error  $\sigma_r^*$  to achieve a target probability of replication success  $1 - \beta < \lim \Pr_{2\text{TR}}$  can now be obtained by equating (10) to  $1 - \beta$  and solving for  $\sigma_r$ . This leads to

$$\sigma_r^* = \frac{\mu_{\hat{\theta}_r} z_\alpha - z_\beta \sqrt{(z_\alpha^2 - z_\beta^2) \{\tau^2 + (\sigma_o^2 + \tau^2)/(1 + 1/g)\} + \mu_{\hat{\theta}_r}^2}}{z_\alpha^2 - z_\beta^2} \quad (12)$$

for  $\alpha < \beta$ . The standard error  $\sigma_r^*$  can subsequently be translated in a sample size. The translation depends on the type of effect size, for instance, for SMD effect sizes we can use the approximation  $n_r^* \approx 4/(\sigma_r^*)^2$  from earlier. Moreover, by assuming a standard error of the form  $\sigma_r = \lambda/\sqrt{n_r}$  and plugging in the parameters of the “conditional” design prior ( $\tau^2 = 0$ ,  $\mu_\theta = \hat{\theta}_o$ ,  $g \downarrow 0$ ), we obtain the well-known sample size formula (Matthews, 2006, section 3.3)

$$n_r^* = \frac{(z_\alpha + z_\beta)^2}{(\hat{\theta}_o/\lambda)^2}$$

for a one-sided significance test at level  $\alpha$  with power  $1 - \beta$  to detect the original effect

estimate  $\hat{\theta}_o$ . The formula (12) thus generalizes standard sample size calculation to take into account the uncertainty of the original estimate, between-study heterogeneity and other types of external knowledge.

### *Fixed effect meta-analysis*

The data from the original and replication studies are sometimes pooled via fixed effect meta-analysis. The pooled effect estimate  $\hat{\theta}_m$  and standard error  $\sigma_m$  are then given by

$$\hat{\theta}_m = (\hat{\theta}_o/\sigma_o^2 + \hat{\theta}_r/\sigma_r^2) \sigma_m^2 \quad \text{and} \quad \sigma_m = (1/\sigma_o^2 + 1/\sigma_r^2)^{-1/2},$$

and they are also equivalent to the mean and standard deviation of a posterior distribution for the effect size  $\theta$  based on the data from both studies and a flat initial prior for  $\theta$ . The success region

$$S_{\text{MA}} = \left[ \sigma_r z_\alpha \sqrt{1 + \sigma_r^2/\sigma_o^2} - (\hat{\theta}_o \sigma_r^2)/\sigma_o^2, \infty \right) \quad (13)$$

then corresponds to both replication success defined via a one-sided meta-analytic  $p$ -value being smaller than level  $\alpha$ , i.e.,  $p_m = 1 - \Phi(\hat{\theta}_m/\sigma_m) \leq \alpha$ , or to replication success defined via a Bayesian posterior probability  $\Pr(\theta > 0 \mid \hat{\theta}_o, \hat{\theta}_r, \sigma_o, \sigma_r) \geq 1 - \alpha$ . Based on the success region (13) and an assumed standard error  $\sigma_r$ , the probability of replication success can be computed by

$$\Pr(\hat{\theta}_r \in S_{\text{MA}} \mid \hat{\theta}_o, \sigma_o, \sigma_r) = \Phi \left( \frac{\mu_{\hat{\theta}_r} - \sigma_r z_\alpha \sqrt{1 + \sigma_r^2/\sigma_o^2} + (\hat{\theta}_o \sigma_r^2)/\sigma_o^2}{\sigma_{\hat{\theta}_r}} \right). \quad (14)$$

As for the two-trials rule, the probability (14) cannot be made arbitrarily high by decreasing the standard error  $\sigma_r$  but is bounded from above by  $\lim \Pr_{2\text{TR}}$  defined in (11). The required standard error  $\sigma_r^*$  to achieve a target probability of replication success  $1 - \beta < \lim \Pr_{2\text{TR}}$  can be computed numerically using root finding algorithms.

### *Effect size equivalence test*

Anderson and Maxwell (2016) proposed a method for quantifying replicability based on effect size equivalence. Under normality, replication success at level  $\alpha$  is achieved if the

$(1 - \alpha)$  confidence interval for the effect size difference  $\theta_r - \theta_o$

$$\hat{\theta}_r - \hat{\theta}_o \pm z_{\alpha/2} \sqrt{\sigma_r^2 + \sigma_o^2}$$

is fully inside an equivalence region  $[-\Delta, \Delta]$  defined via the pre-specified margin  $\Delta > 0$ .

This procedure corresponds to rejecting the null hypothesis  $H_0: |\theta_r - \theta_o| > \Delta$  in an equivalence test, and it implies a success region for the replication effect estimate  $\hat{\theta}_r$  given by

$$S_E = \left[ \hat{\theta}_o - \Delta + z_{\alpha/2} \sqrt{\sigma_o^2 + \sigma_r^2}, \hat{\theta}_o + \Delta - z_{\alpha/2} \sqrt{\sigma_o^2 + \sigma_r^2} \right] \quad (15)$$

for  $\Delta \geq z_{\alpha/2} \sqrt{\sigma_o^2 + \sigma_r^2}$ . For too small margins ( $\Delta < z_{\alpha/2} \sqrt{\sigma_o^2 + \sigma_r^2}$ ), the success region (15) becomes the empty set meaning that replication success is impossible. Assuming now that the margin is large enough, the probability of replication success can be computed by

$$\begin{aligned} \Pr(\hat{\theta}_r \in S_E | \hat{\theta}_o, \sigma_o, \sigma_r) &= \Phi \left( \frac{\hat{\theta}_o + \Delta - z_{\alpha/2} \sqrt{\sigma_o^2 + \sigma_r^2} - \mu_{\hat{\theta}_r}}{\sigma_{\hat{\theta}_r}} \right) \\ &\quad - \Phi \left( \frac{\hat{\theta}_o - \Delta + z_{\alpha/2} \sqrt{\sigma_o^2 + \sigma_r^2} - \mu_{\hat{\theta}_r}}{\sigma_{\hat{\theta}_r}} \right). \end{aligned} \quad (16)$$

As with the previous methods, the probability (16) cannot be made arbitrarily high by decreasing the replication standard error  $\sigma_r$ , but is bounded by

$$\lim \Pr_E = \Phi \left( \frac{\hat{\theta}_o + \Delta - z_{\alpha/2} \sigma_o - \mu_{\hat{\theta}_r}}{\sqrt{\tau^2 + (\sigma_o^2 + \tau^2)/(1 + 1/g)}} \right) - \Phi \left( \frac{\hat{\theta}_o - \Delta + z_{\alpha/2} \sigma_o - \mu_{\hat{\theta}_r}}{\sqrt{\tau^2 + (\sigma_o^2 + \tau^2)/(1 + 1/g)}} \right).$$

The required replication standard error  $\sigma_r^*$  to achieve a target probability of replication success  $1 - \beta < \lim \Pr_E$  can again be computed numerically.

### ***The replication Bayes factor***

A Bayesian hypothesis testing approach for assessing replication success was proposed by Verhagen and Wagenmakers (2014) and further developed by Ly et al. (2018).

They define a “replication Bayes factor”

$$\text{BF}_R = \frac{f(x_r | H_0)}{f(x_r | H_1)}$$

which is the ratio of the marginal likelihood of the replication data  $x_r$  under the null hypothesis  $H_0: \theta = 0$  to the marginal likelihood of  $x_r$  under the alternative hypothesis  $H_1: \theta \sim f(\theta | x_o)$ , that is, the posterior of the effect size  $\theta$  based on the original data  $x_o$ . If the original study provides evidence against the null hypothesis, replication Bayes factor values  $\text{BF}_R < 1$  indicate replication success, and the smaller the value the higher the degree of success.

Under normality and assuming no heterogeneity, the success region for achieving  $\text{BF}_R \leq \gamma$  is given by

$$S_{\text{BF}_R} = \left( -\infty, -\sqrt{A} - (\hat{\theta}_o \sigma_r^2) / \sigma_o^2 \right] \cup \left[ \sqrt{A} - (\hat{\theta}_o \sigma_r^2) / \sigma_o^2, \infty \right) \quad (17)$$

with  $A = \sigma_r^2(1 + \sigma_r^2/\sigma_o^2)\{\hat{\theta}_o^2/\sigma_o^2 - 2\log \gamma + \log(1 + \sigma_o^2/\sigma_r^2)\}$ . Details of this calculation are given in the supplement. The fact that the success region (17) is defined on both sides around zero shows that replication success is also possible if the replication effect estimate goes in opposite direction of the original one, which is known as the “replication paradox” (Ly et al., 2018). The paradox can be avoided using a modified version of the replication Bayes factor but the success region is no longer available in closed-form (Pawel & Held, 2022, Appendix D). Based on the success region (17), the probability of replication success can be computed by

$$\Pr(\hat{\theta}_r \in S_{\text{BF}_R} | \hat{\theta}_o, \sigma_o, \sigma_r) = \Phi \left( \frac{\mu_{\hat{\theta}_r} - \sqrt{A} + (\hat{\theta}_o \sigma_r^2) / \sigma_o^2}{\sigma_{\hat{\theta}_r}} \right) + \Phi \left( \frac{-\sqrt{A} - (\hat{\theta}_o \sigma_r^2) / \sigma_o^2 - \mu_{\hat{\theta}_r}}{\sigma_{\hat{\theta}_r}} \right). \quad (18)$$

To avoid powering the replication study for the replication paradox, one may want to compute the probability of replication success only for the part of the success region with the same sign as the original effect estimate. As for the other methods, the probability (18) is bounded from above by a constant  $\lim \Pr_{\text{BF}_R} = \lim_{\sigma_r \downarrow 0} \Pr(\hat{\theta}_r \in S_{\text{BF}_R} | \hat{\theta}_o, \sigma_o, \sigma_r)$ , and root finding algorithms

can be used to numerically determine the required standard error  $\sigma_r^*$  for achieving a target probability of replication success  $1 - \beta < \lim\Pr_{\text{BF}_R}$ .

### *The skeptical $p$ -value*

Held (2020) proposed a reverse-Bayes approach for quantifying replication success. The main idea is to determine the variance of a “skeptical” zero-mean normal prior for the effect size  $\theta$  such that its posterior distribution based on the original study no longer indicates evidence for a genuine effect. Replication success is then achieved if the replication data are in conflict with the skeptical prior. The procedure can be summarized by a “skeptical  $p$ -value”  $p_s$ , and the lower the  $p$ -value the higher the degree of replication success. Held, Micheloud, and Pawel (2022, sec. 2.1) showed that the success region for replication success defined by  $p_s \leq \alpha$  is given by

$$S_{p_s} = \left[ z_\alpha \sqrt{\sigma_r^2 + \frac{\sigma_o^2}{(z_o^2/z_\alpha^2) - 1}}, \infty \right). \quad (19)$$

From the success region (19) the probability of replication success at level  $\alpha$  is

$$\Pr(\hat{\theta}_r \in S_{p_s} \mid \hat{\theta}_o, \sigma_o, \sigma_r) = \Phi \left( \frac{\mu_{\hat{\theta}_r} - z_\alpha \sqrt{\sigma_r^2 + \sigma_o^2 / \{(z_o^2/z_\alpha^2) - 1\}}}{\sigma_{\hat{\theta}_r}} \right),$$

and also bounded from above by a constant  $\lim\Pr_{p_s} = \lim_{\sigma_r \downarrow 0} \Pr(\hat{\theta}_r \in S_{p_s} \mid \hat{\theta}_o, \sigma_o, \sigma_r)$ . As for the two-trials rule, the required standard error  $\sigma_r^*$  to achieve a probability of replication success  $1 - \beta < \lim\Pr_{p_s}$  can be computed analytically for  $\alpha < \beta$ :

$$\sigma_r^* = \sqrt{x^2 - \frac{\sigma_o^2}{(z_o/z_\alpha)^2 - 1}}$$

with

$$x = \frac{z_\alpha \mu_{\hat{\theta}_r} - z_\beta \sqrt{\mu_{\hat{\theta}_r}^2 - (z_\alpha^2 - z_\beta^2)[\tau^2 + (\sigma_o^2 + \tau^2)/(1 + 1/g) - \sigma_o^2 / \{(z_o/z_\alpha)^2 - 1\}]}}{z_\alpha^2 - z_\beta^2}.$$

### *The skeptical Bayes factor*

Pawel and Held (2022) modified the previously described reverse-Bayes assessment of

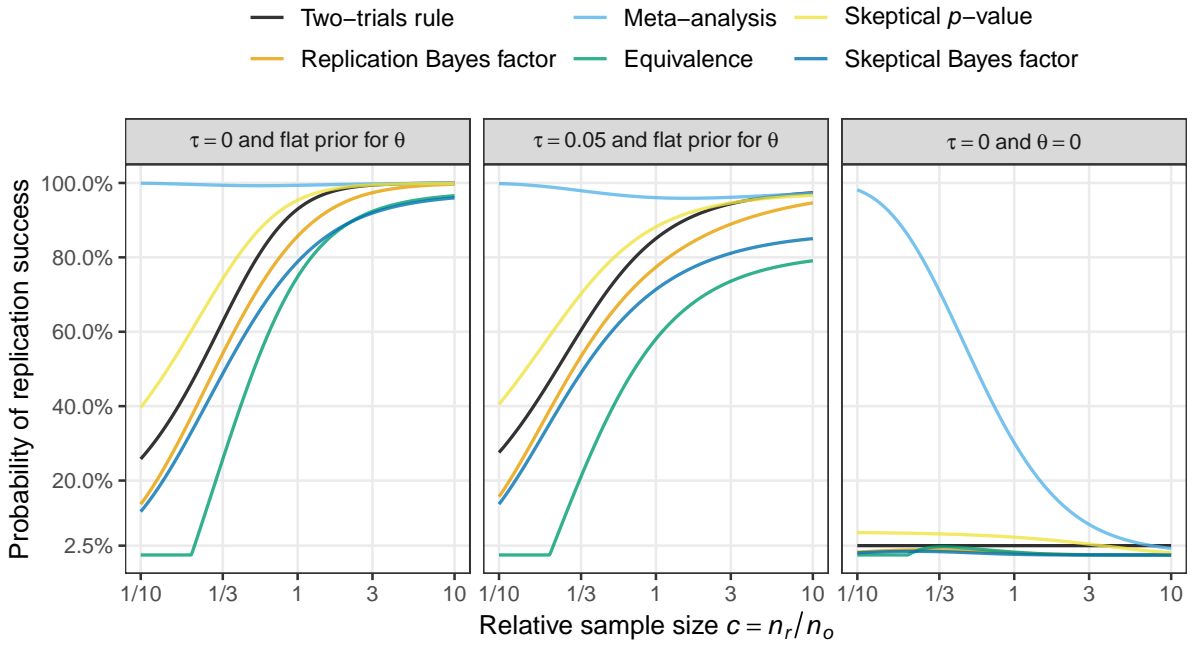
replication success from Held (2020) to use Bayes factors instead of tail probabilities as measures of evidence. Again, the procedure can be summarized in a single quantity termed the “skeptical Bayes factor”  $BF_s$ , with lower values of  $BF_s$  pointing to higher degrees of replication success. Also for this method, the success region and the probability of replication success can be expressed in closed-form but the derivations are more involved than for the other methods. For this reason, they are only given in the supplement.

### **Example: Cross-laboratory replication project (continued)**

We will now revisit the “Labels” experiment and compute the probability of replication success. The parameters of the analysis methods are specified as follows: For the two-trials rule we use the conventional one-sided significance level  $\alpha = 0.025$ , while for meta-analysis we use the more stringent level  $\alpha = 0.025^2$  as the method is based on two data sets rather than one. We use a  $1 - \alpha = 90\%$  confidence interval which is conventionally used in equivalence testing, along with a margin  $\Delta = 0.2$  corresponding to a small SMD effect size according to the classification from Cohen (1992). For the skeptical  $p$ -value we use the recommended “golden” level  $\alpha = 0.062$  as it guarantees that for original studies which were just significant at  $\alpha = 0.025$  replication success is only possible if the replication effect estimate is larger than the original one (Held, Micheloud, & Pawel, 2022). Finally, for the replication Bayes factor and the skeptical Bayes factor we use the “strong evidence” level  $\gamma = 1/10$  from Jeffreys (1961).

Figure 4 shows the probability of replication success as a function of the relative sample size  $c = n_r/n_o$  and for different initial priors. The left and middle plot are based on a flat initial prior for the effect size without heterogeneity ( $\tau = 0$ ) and with heterogeneity ( $\tau = 0.05$ ), respectively. The right plot shows the prior corresponding to the “fixed effect null hypothesis”  $H_0: \theta = 0, \tau^2 = 0$ , so that the probability of replication success is the type I error rate which some stakeholders might require to be “controlled” at some adequate level.

We see from the left and middle plots that increasing the relative sample size monotonically increases the probability of replication success for all methods but meta-analysis (light blue). Meta-analysis shows a non-monotone behavior because the original study was already highly significant so that the pooled effect estimate is significant even for

**Figure 4**

Probability of replication success as a function of relative sample size  $c = n_r/n_o$  for the “Labels” experiment with original effect estimate  $\hat{\theta}_o = 0.205$  and standard error  $\sigma_o = 0.051$  under different initial prior distributions. Replication success is defined by the two-trials rule at level  $\alpha = 0.025$ , the replication Bayes factor at level  $\gamma = 1/10$ , fixed effect meta-analysis at level  $\alpha = 0.025^2$ , effect size equivalence based on 90% confidence interval and with margin  $\Delta = 0.2$ , skeptical  $p$ -value at level  $\alpha = 0.062$ , and skeptical Bayes factor at level  $\gamma = 1/10$ .

replication studies with very small sample size (Micheloud & Held, 2022). The uncertainty regarding the replication effect estimate  $\hat{\theta}_r$  may therefore even reduce the probability of replication success for meta-analysis if the sample size is increased. If heterogeneity is taken into account (middle plot) the probability of replication success becomes closer to 50% for all methods except the equivalence test, reflecting the larger uncertainty about the effect size  $\theta$ . To achieve 80% probability of replication success the fewest samples are required with meta-analysis, followed by the skeptical  $p$ -value, the two-trials rule, the replication Bayes factor, the skeptical Bayes factor, and lastly the equivalence test. If the sample size should guarantee a sufficiently conclusive replication study with all these methods, the replication sample size has to be slightly larger than the original one if no heterogeneity is assumed ( $\tau = 0$ ), while it has to be increased more than ten-fold if heterogeneity is assumed ( $\tau = 0.05$ ). However, this is mostly due to the equivalence test which requires by far the most samples. If

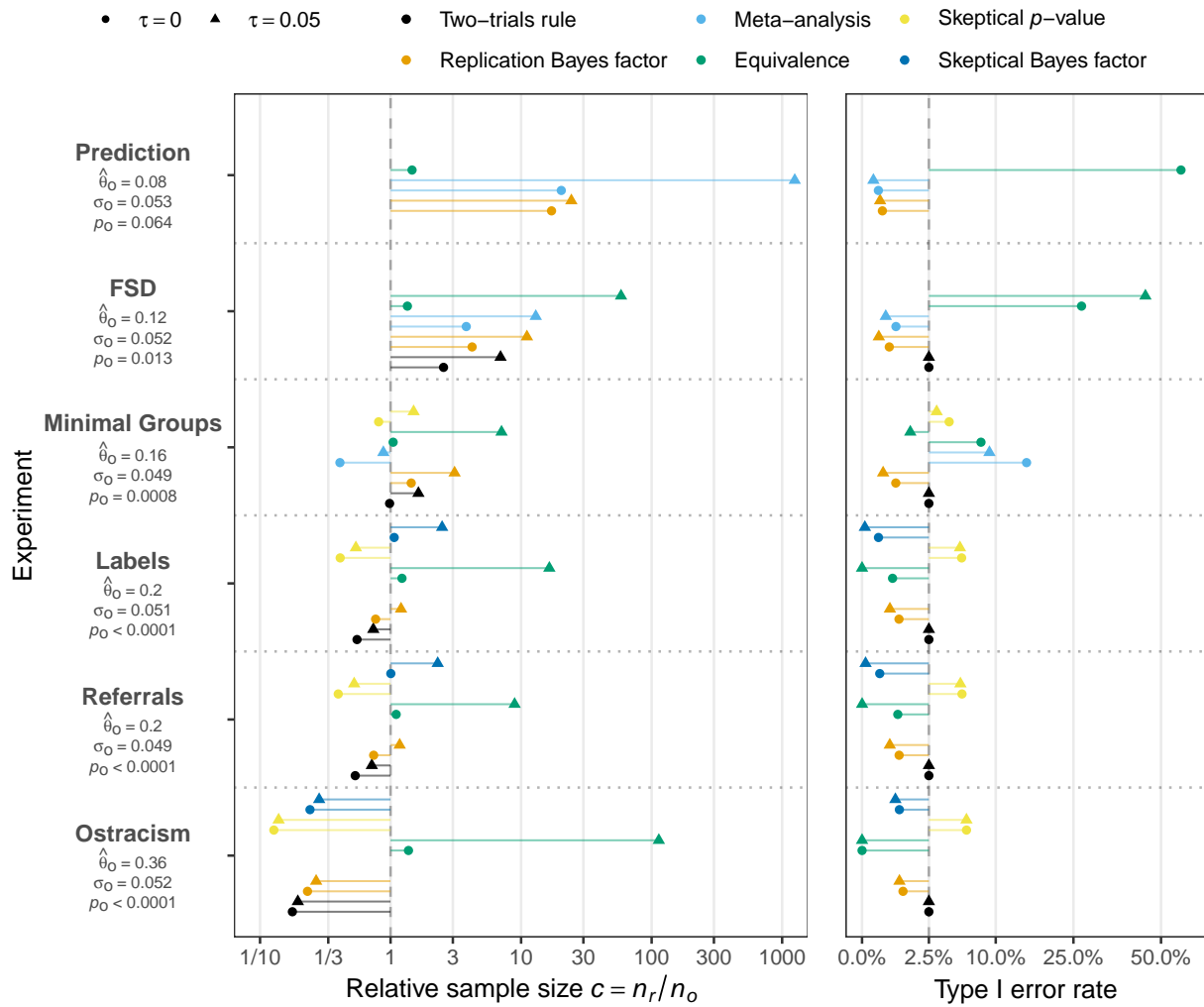


the equivalence test sample size is ignored, the relative sample size  $c = 2.5$  ensures at least 80% probability of replication success under heterogeneity with the remaining methods.

The right plot in Figure 4 shows that the type I error rate of the two-trials rule (black) stays constant at  $\alpha = 0.025$ , as expected by definition of the method. In contrast, the type I error rates of the other methods vary with the relative sample size  $c$  but most of them stay below  $\alpha = 0.025$  for all  $c$  with the exception of meta-analysis and the skeptical  $p$ -value. Meta-analysis (light blue) has an extremely high type I error rate as the pooling with the highly significant original data leads to replication success if the replication sample size is not drastically increased. The type I error rate of the skeptical  $p$ -value (yellow) is only slightly higher than  $\alpha = 0.025$  which is expected since the level  $\alpha = 0.062$  is used for declaring replication success with the skeptical  $p$ -value, and its type I error rate is always smaller than the level for thresholding it (Held, 2020). The type I error rate of the skeptical  $p$ -value decreases to values smaller than  $\alpha = 0.025$  of the two-trials rule at approximately  $c = 3$ .

We now perform SSD for an illustrative subset of studies from the Protzko et al. (2020) replication project. Figure 5 shows the required relative sample size and the associated type I error rates if a sample size can be computed for a target probability of replication success of  $1 - \beta = 80\%$ . If there is no sample size for which a probability of 80% can be achieved, the space is left blank. For example, in the application of the meta-analysis method to the “Labels” experiment the probability remains above 80% for any relative sample size, and therefore no sample size is shown.

We see that for all methods except the equivalence test, the required relative sample size  $c$  decreases as the original  $p$ -value  $p_o$  decreases and original studies with very small  $p$ -values require much fewer samples in the replication study. For example, in the “Ostracism” experiment with  $p_o < 0.0001$  the required sample size for all methods except the equivalence test is at most one-third the size of the original. For the equivalence test the required sample size depends instead on the size of the original standard error  $\sigma_o$  and smaller standard errors lead to smaller required sample sizes in the replication. For example, the “Referrals” experiment with original standard error  $\sigma_o = 0.049$  requires fewer samples for the equivalence test than the “Ostracism” experiment with original standard error  $\sigma_o = 0.052$ .

**Figure 5**

The left plot shows the required relative sample size  $c = n_r/n_o$  to achieve a target probability of replication success of  $1 - \beta = 80\%$  (if possible). Replication success is defined through the two-trials rule at level  $\alpha = 0.025$ , replication Bayes factor at level  $\gamma = 1/10$ , fixed effect meta-analysis at level  $\alpha = 0.025^2$ , effect size equivalence at level  $\alpha = 0.1$  with margin  $\Delta = 0.2$ , skeptical p-value at level  $\alpha = 0.062$ , and skeptical Bayes factor at level  $\gamma = 1/10$  for an illustrative subset of studies from the Protzko et al. (2020) replication project, the supplement shows results for all studies. A flat initial prior is used for the effect size  $\theta$  either without ( $\tau = 0$ ) or with heterogeneity ( $\tau = 0.05$ ). The right plot shows the type I error rate associated with the required sample size. Experiments are ordered (top to bottom) by their original one-sided p-value.

Figure 5 also shows that accounting for heterogeneity (triangles) increases the required sample size for all methods compared to ignoring it (points). Although more costly to the researcher, larger sample sizes also reduce the type I error rate for most methods (right plot). Comparing the type I error rates of the different methods, we see again the pattern that the

type I error rates of the equivalence test and the skeptical  $p$ -value are higher than the type I error rate of 2.5% of the two-trials rule. However, while the type I error rate of the skeptical  $p$ -value decreases when replication studies require larger samples sizes, the type I error rate of the equivalence test may also be high if the replication requires very large sample sizes (e.g., for the “FSD” experiment), since it depends on whether the original effect estimate  $\hat{\theta}_o$  is sufficiently different from zero. If the original effect estimate  $\hat{\theta}_o$  is close to zero, the type I error rate of the equivalence test increases drastically, since equivalence can be established even if the original and replication effect estimates are close to zero.

The supplement shows the same analysis for all studies in the Protzko et al. (2020) project. Most original studies were highly significant and therefore require fewer samples in the replication than in the original study to achieve a target probability of  $1 - \beta = 80\%$  for replication success with all methods except the equivalence test. Some original studies were less convincing and therefore require larger replication sample sizes. The additional samples needed for these studies could thus be reallocated from the studies that require fewer samples. The project would still use the same total sample size, but it would be more efficiently allocated. An exception to this conclusion is the equivalence test, which in most cases requires larger replication sample sizes. This is because the original standard errors of all studies are relatively large compared to the specified equivalence margin. Therefore, if one plans to analyze the original and replication pair with an equivalence test, this should already be taken into account at the design stage of the original study, since an imprecise original study will diminish the chances of replication success with this method.

### Sample size determination for multisite replication projects

So far we considered the situation where a pair of a single original and a single replication study are analyzed in isolation. However, if multiple replications per single original study are conducted (so-called *multisite* replication studies) the ensemble of replications can also be analyzed jointly. In this case, some adaptations of the SSD methodology are required.

The replication effect estimate and its standard error are now vectors

$\hat{\theta}_r = (\hat{\theta}_{r1}, \dots, \hat{\theta}_{rm})^\top$  and  $\sigma_r^2 = (\sigma_{r1}^2, \dots, \sigma_{rm}^2)^\top$  consisting of  $m$  replication effect estimates and their standard errors. The normal hierarchical model for the replication estimates  $\hat{\theta}_r$  then

becomes

$$\hat{\boldsymbol{\theta}}_r | \boldsymbol{\theta}_r \sim N_m \left\{ \boldsymbol{\theta}_r, \text{diag}(\boldsymbol{\sigma}_r^2) \right\} \quad (20a)$$

$$\boldsymbol{\theta}_r | \theta \sim N_m \left\{ \theta \mathbf{1}_m, \tau^2 \text{diag}(\mathbf{1}_m) \right\}, \quad (20b)$$

where  $\boldsymbol{\theta}_r$  is a vector of  $m$  study specific effect sizes,  $\mathbf{1}_m$  is a vector of  $m$  ones, and  $N_m(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  denotes the  $m$ -variate normal distribution with mean vector  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Sigma}$ . By marginalizing over the study specific effect size  $\boldsymbol{\theta}_k$ , the model can alternatively be expressed by

$$\hat{\boldsymbol{\theta}}_r | \theta \sim N_m \left\{ \theta \mathbf{1}_m, \text{diag}(\boldsymbol{\sigma}_r^2 + \tau^2 \mathbf{1}_m) \right\}, \quad (21)$$

so the predictive distribution of  $\hat{\boldsymbol{\theta}}_r$  based on the design prior (6) is given by

$$\hat{\boldsymbol{\theta}}_r | \hat{\theta}_o, \sigma_o^2, \boldsymbol{\sigma}_r^2 \sim N_m \left\{ \mu_{\hat{\theta}_r} \mathbf{1}_m, \text{diag}(\boldsymbol{\sigma}_r^2 + \tau^2 \mathbf{1}_m) + \left( \frac{\tau^2 + \sigma_o^2}{1 + 1/g} \right) \mathbf{1}_m \mathbf{1}_m^\top \right\} \quad (22)$$

with  $\mu_{\hat{\theta}_r}$  the mean of the predictive distribution of a single replication effect estimate from (7). Importantly, the replication effect estimates are correlated as the covariance matrix in (22) has  $(\tau^2 + \sigma_o^2)/(1 + 1/g)$  in the off-diagonal entries.

Often the assessment of replication success can be formulated in terms of a weighted average of the replication effect estimates  $\hat{\theta}_{r*} = (\sum_{i=1}^m w_i \hat{\theta}_{ri}) / (\sum_{i=1}^m w_i)$  with  $w_i$  the weight of replication  $i$ . For instance, several multisite replication projects (e. g., Klein et al., 2018) have defined replication success by the fixed or random effect(s) meta-analytic effect estimate of the replication effect estimates achieving statistical significance. Based on the predictive distribution of the replication effect estimate vector (22), the predictive distribution of the weighted average  $\hat{\theta}_{r*}$  is given by

$$\hat{\theta}_{r*} | \hat{\theta}_o, \sigma_o^2, \boldsymbol{\sigma}_r^2 \sim N \left\{ \mu_{\hat{\theta}_r}, \sigma_{\hat{\theta}_{r*}}^2 = \left( \sum_{i=1}^m w_i^2 \sigma_{\hat{\theta}_{ri}}^2 + \sum_{i=1}^m \sum_{\substack{j=1 \\ j \neq i}}^m w_i w_j \frac{\tau^2 + \sigma_o^2}{1 + 1/g} \right) / \left( \sum_{i=1}^m w_i \right)^2 \right\} \quad (23)$$

with  $\sigma_{\hat{\theta}_{ri}}^2$  the predictive variance of a single replication effect estimate with standard error  $\sigma_{ri}$

as in (7). In particular when the studies receive equal weights ( $w_i = w$  for  $i = 1, \dots, m$ ) and the standard errors of the replication effect estimates are equal ( $\sigma_{ri} = \sigma_r$  for  $i = 1, \dots, m$ ), the predictive variance becomes

$$\sigma_{\hat{\theta}_{r*}}^2 = \frac{\sigma_r^2 + \tau^2}{m} + \frac{\tau^2 + \sigma_o^2}{1 + 1/g}. \quad (24)$$

The probability of replication success can now be obtained by integrating (22), respectively (23), over the corresponding success region  $S$ . This may be more involved if the success region is defined in terms of the replication effect estimate vector  $\hat{\theta}_r$ , whereas it is as simple as in the singlesite replication case if the success region is formulated in terms of the weighted average  $\hat{\theta}_{r*}$ .

### ***Optimal allocation within and between sites***

A key challenge in SSD for multisite replication studies is the optimal allocation of samples within and between sites, that is, how many sites  $m$  and how many samples  $n_{ri}$  per site  $i$  should be used. A similar problem exists in SSD for cluster randomized trials and we can adapt the common solution based on cost functions (Raudenbush, 1997). That is, the optimal configuration is determined so that the probability of replication success is maximized subject to a constrained cost function which accounts for the (typically different) costs of additional samples and sites.

For example, assume a balanced design ( $n_{ri} = n_r$  for  $i = 1, \dots, m$ ) and that the standard errors of the replication effect estimates are inversely proportional to the square-root of the sample size  $\sigma_{ri} = \lambda/\sqrt{n_r}$  for some unit variance  $\lambda^2$ . Further assume that maximizing the probability of replication success corresponds to minimizing the variance of the weighted average  $\sigma_{\hat{\theta}_{r*}}^2$  in (24). Let  $K_s$  denote the cost of an additional site, and  $K_c$  the cost of an additional sample/case. The total cost of the project is then  $K = m(K_c n_r + K_s)$ , and constrained minimization of the predictive variance (24) leads to the optimal sample size per site

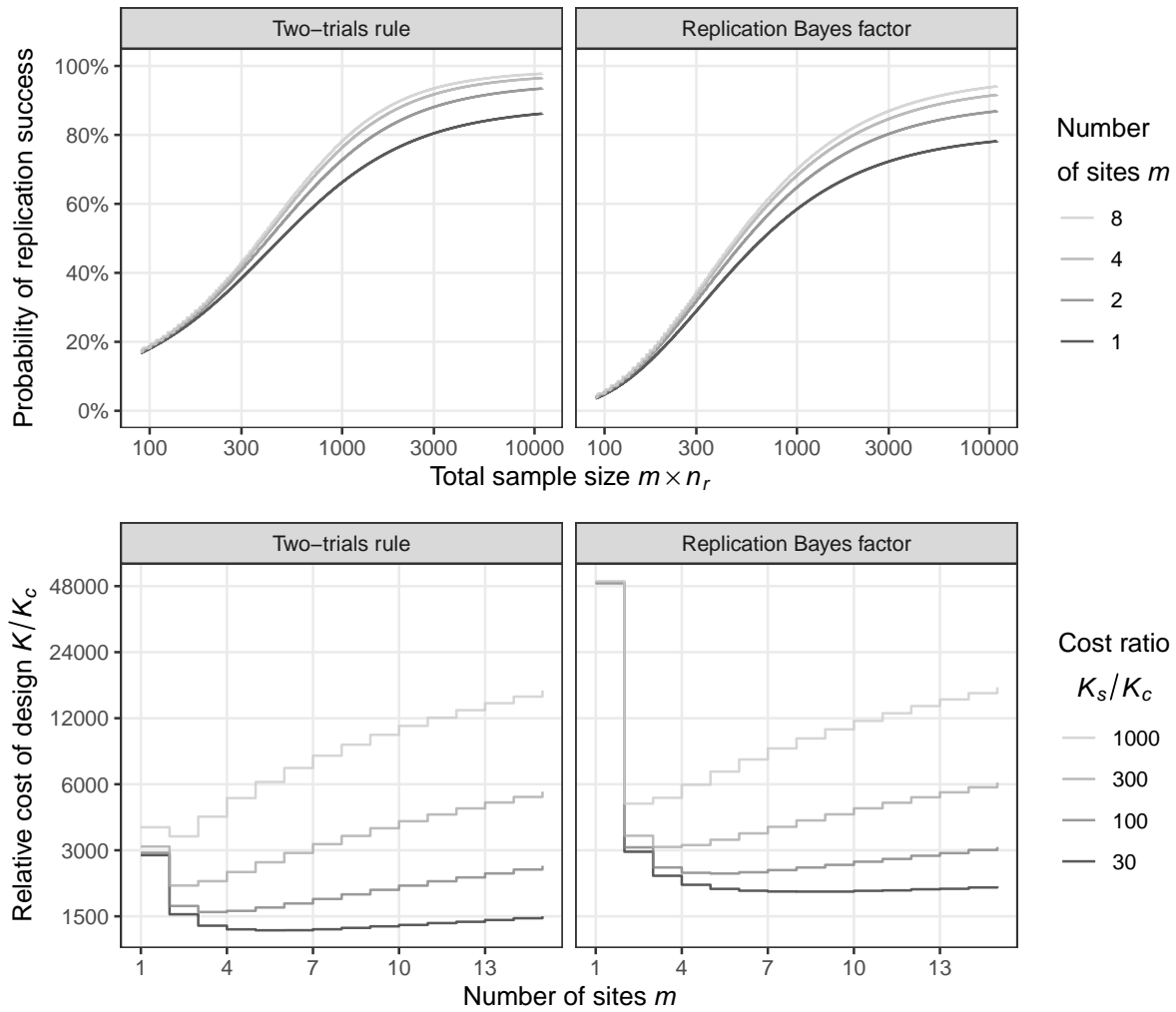
$$n_r^* = \frac{\lambda}{\tau} \sqrt{\frac{K_s}{K_c}}$$

which is equivalent to the optimal cluster sample size known from cluster randomized trials (Raudenbush & Liu, 2000). Note that the optimal sample size per site may be different for other analysis approaches where maximizing the probability of replication success does not correspond to minimizing the variance of the weighted average. Moreover, there are also practical considerations which affect the choice of how many sites should be included in a project. For instance, there may simply not be enough labs available with the required expertise to perform the replication experiments.

**Example: Cross-laboratory replication project (continued)**

Figure 6 illustrates multisite SSD for the “Labels” experiment from Protzko et al. (2020) for planned analyses based on the two-trials rule and the replication Bayes factor (see the supplement for details on the multisite extension of these two methods). As for singlesite SSD, we use the design prior based on a flat initial prior for the effect size and taking into account heterogeneity ( $\tau = 0.05$ ). The top plots show the probability of replication success as a function of the total sample size  $m \times n_r$  for different number of sites  $m$ . We see that for the same total sample size a larger number of sites increases the probability of replication success. For instance, a total sample size of roughly 3000 is required to achieve an 80% target probability with one site for the two-trials rule, whereas only approximately half as many samples are required for two sites.

However, focusing only on the total sample size ignores the fact that the cost of an additional site is usually larger than the cost of an additional sample. The bottom plot shows the total cost  $K$  of a design (relative to the cost of one sample  $K_c$ ) whose sample size is determined for a target probability of replication success  $1 - \beta = 80\%$ . We see that if the cost of an additional site  $K_s$  is not much larger than the cost of an additional sample  $K_c$ , e.g.,  $K_s/K_c = 30$  the optimal number of sites is  $m = 5$  for the two-trials rule and  $m = 8$  for the replication Bayes factor. If an additional site is more costly the optimal number of sites is lower, e.g., if the cost ratio is  $K_s/K_c = 300$ , the optimal number of sites is  $m = 2$  for the two-trials rule and  $m = 3$  for the replication Bayes factor. This is similar to the actually used number of sites  $m = 3$  (counting only external-replications), respectively,  $m = 4$  (counting also the internal-replication) from Protzko et al. (2020).

**Figure 6**

The top plots show the probability of replication success based on the two-trials rule at level  $\alpha = 0.025$  (left) and the replication Bayes factor at level  $\gamma = 1/10$  (right) as a function of the total sample size and for different number of sites  $m$  for data from the “Labels” experiment. A design prior with heterogeneity  $\tau = 0.05$  and flat initial prior for the effect size  $\theta$  is used. The same heterogeneity value is assumed in the analysis of the replications. The bottom plot shows the total cost  $K$  of the design (relative to the cost of a single sample  $K_c$ ) as a function of the number of sites  $m$  and for different site costs  $K_s$ . The sample size of each design corresponds to a target probability of replication success  $1 - \beta = 80\%$ .

## Discussion

We showed how Bayesian approaches can be used to determine the sample size of replication studies based on all the available information and the associated uncertainty. A key strength of the approach is that it can be applied to any type of replication analysis method, Bayesian or non-Bayesian, as long as there is a well-defined success region for the replication

effect estimate. Methods for assessing replication success which have not yet been adapted to Bayesian design approaches in the normal-normal hierarchical model (or not even proposed) can thus benefit from our methodology. For instance, our methods could easily be applied to the “dual-criterion” from Rosenkranz (2021), which defines replication success via simultaneous statistical significance and practical relevance of the effect estimates from the original and replication studies.

There are some limitations and possible extensions: we have developed the methodology for “direct” replication studies (Simons, 2014), which attempt to replicate the conditions of the original study as closely as possible; however, SSD methodology is also needed for “conceptual” replication or “generalization” studies, which may have systematic deviations from the original study. While the heterogeneity variance in the design prior allows SSD to account for effect size heterogeneity to some extent, more research is needed to investigate how to account for systematic study variation. For the same reason, it is unclear how our Bayesian design approach can be applied to a “causal” replication framework (Steiner, Wong, & Anglin, 2019; Wong, Anglin, & Steiner, 2021), where the focus is on the ability of the original and replication studies to estimate the same causal estimand, rather than on similar study procedures. In addition, as in standard meta-analysis, we assumed that the variances of the effect estimates are known, which can sometimes be inadequate (Jackson & White, 2018). Specifying priors also for the variances could better reflect the available uncertainty but would come at the cost of reduced interpretability and increased computational complexity. We also did not consider designs in which the replication data are analyzed sequentially. Ideas from Bayesian sequential designs (Schönbrodt & Wagenmakers, 2018; Stefan, Gronau, & Wagenmakers, 2022) or from adaptive clinical trials (Bretz, Koenig, Brannath, Glimm, & Posch, 2009) could be adapted to the replication setting, as in Micheloud and Held (2022). A sequential analysis of the replication data could possibly increase the efficiency of the replication. An additional point is that we assumed that the original study has been completed when planning the replication study. One could also consider a scenario where both the original and the replication study are planned simultaneously and adopt a “project” perspective (Held, Micheloud, & Pawel, 2022; Maca, Gallo, Branson, & Maurer,



2002). In this case, however, no information from the original study is available and the design prior must be specified entirely based on external knowledge. Finally, researchers have limited resources and may not be able afford a large enough sample size to achieve their desired probability of replication success. In this situation, a reverse-Bayes approach (Held, Matthews, Ott, & Pawel, 2022) could be used to determine the prior for the effect size required to achieve the desired probability of replication success based on the maximally affordable sample size. Researchers can then judge whether or not such prior beliefs are scientifically sensible, and decide whether to conduct the replication study with their limited resources.

### Appendix: The BayesRepDesign R package

The R package BayesRepDesign can be installed from the Comprehensive R Archive Network (CRAN) by running the following command from an R console

```
install.packages("BayesRepDesign")
```

Once the package is installed, it can be loaded with

```
library("BayesRepDesign")
```

To see an overview of the functionality of the package, run

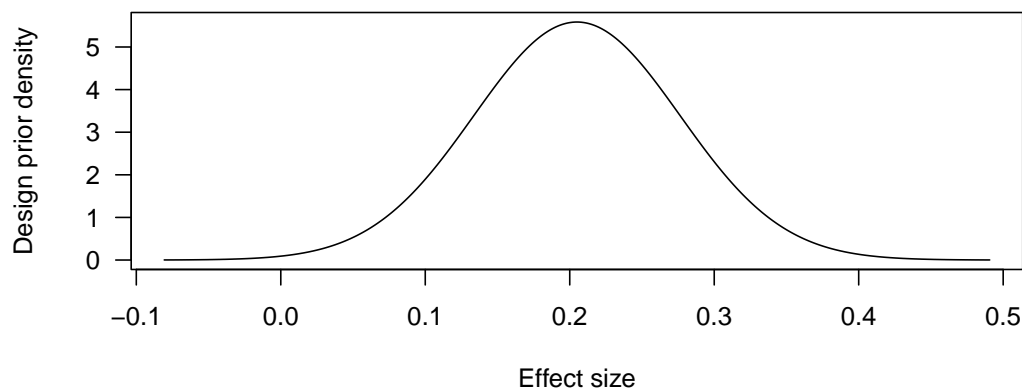
```
help(package = "BayesRepDesign")
```

The first step in Bayesian design of a replication study is to create a design prior for the effect size  $\theta$ . We use the original effect estimate  $\hat{\theta}_o = 0.205$  and standard error  $\sigma_o = 0.051$  from the “Labels” experiment along with a flat initial prior for  $\theta$  (the default) and a heterogeneity deviation of  $\tau = 0.05$  as inputs to the `designPrior` function

```
dp <- designPrior(to = 0.205, so = 0.051, tau = 0.05)
```

The resulting design prior object can be visualized with

```
plot(dp)
```



The design prior can now be used to compute the probability of replication success with the `pors` functions or to compute the replication standard error with the `ssd` functions. Each analysis method discussed in this paper has dedicated `pors` and `ssd` functions. For example, `porsSig` can be used to compute the probability of replication success defined by a significant replication  $p$ -value for a given replication standard error, while `ssdSig` can be used to compute the replication standard error required to achieve significance for a given target probability of replication success. In the following, we will compute the replication standard error for achieving replication success with a target probability of 80%.

```
(ssd1 <- ssdSig(level = 0.025, dprior = dp, power = 0.8))

##          Bayesian sample size calculation for replication studies
##          =====
##
## success criterion and computation
## -----
##   replication p-value <= 0.025 (exact computation)
##
## original data and initial prior for effect size
## -----
##   to = 0.2 : original effect estimate
##   so = 0.051 : standard error of original effect estimate
##   tau = 0.05 : assumed heterogeneity standard deviation
##   N(mean = 0, sd = Inf) : initial normal prior
```

```
##
## design prior for effect size
## -----
##   N(mean = 0.2, sd = 0.071) : normal design prior
##
## probability of replication success
## -----
##   PoRS = 0.8 : specified
##   PoRS = 0.8 : recomputed with sr
##
## required sample size
## -----
##   sr = 0.059 : required standard error of replication effect estimate
##   c = so^2/sr^2 ~= nr/no = 0.74 : required relative variance / sample size
```

The output shows the relative variance  $c = \sigma_o/\sigma_r$  which, assuming a standard error form  $\sigma_i = \lambda/\sqrt{n_i}$ , is equal to the relative sample size  $c = n_r/n_o$ . The parameter  $c$  thus quantifies by how much the replication sample size  $n_r$  must be increased/decreased compared to the original sample size  $n_o$ . The replication standard error can also be converted to an absolute sample size using

```
se2n(se = ssdl$sr, unitSD = 2)

## [1] 1137
```

By default, the function assumes a unit standard deviation of  $\lambda = 2$  for the conversion, which is a reasonable approximation of the unit standard deviation for standardized mean differences and log odds/hazard/rate ratios for balanced group designs (Spiegelhalter et al., 2004, section 2.4). However, more exact conversions can be obtained by considering the exact form of the standard error and solving for the sample size.

Finally, the BayesRepDesign package can be easily extended to replication analysis methods other than those for which dedicated functions are provided. To do so, users need to define a function that returns the success region for the replication effect estimate for a given

replication standard error. The function is then passed as an argument to the `ssd` function, which then numerically determines the required standard error. The following code illustrates how the significance method from earlier can be reimplemented in this way.

```
sregionfunSig <- function(sr, alpha = 0.025) {  
  za <- qnorm(p = 1 - alpha)  
  sregion <- successRegion(intervals = cbind(za*sr, Inf))  
  return(sregion)  
}  
ssd2 <- ssd(sregionfun = sregionfunSig, dprior = dp, power = 0.8)  
se2n(se = ssd2$sr, unitSD = 2)  
  
## [1] 1137
```

We see that this results in the same sample size as the `ssdSig` function (which uses a closed-form solution).

## References

- Anderson, S. F., & Kelley, K. (2022). Sample size planning for replication studies: The devil is in the design. *Psychological Methods*. doi: 10.1037/met0000520
- Anderson, S. F., & Maxwell, S. E. (2016). There's more than one way to conduct a replication study: Beyond statistical significance. *Psychological Methods*, 21(1), 1–12. doi: 10.1037/met0000051
- Anderson, S. F., & Maxwell, S. E. (2017). Addressing the “replication crisis”: Using original studies to design replication studies with appropriate statistical power. *Multivariate Behavioral Research*, 52(3), 305–324. doi: 10.1080/00273171.2017.1289361
- Bayarri, M. J., & Mayoral, A. M. (2002). Bayesian design of “successful” replications. *The American Statistician*, 56(3), 207–214. doi: 10.1198/000313002155
- Bonett, D. G. (2020). Design and analysis of replication studies. *Organizational Research Methods*, 24(3), 513–529. doi: 10.1177/1094428120911088
- Bretz, F., Koenig, F., Brannath, W., Glimm, E., & Posch, M. (2009). Adaptive designs for confirmatory clinical trials. *Statistics in Medicine*, 28(8), 1181–1217. doi: 10.1002/sim.3538
- Camerer, C. F., Dreber, A., Holzmeister, F., Ho, T., Huber, J., Johannesson, M., . . . others (2018). Evaluating the replicability of social science experiments in Nature and Science between 2010 and 2015. *Nature Human Behavior*, 2, 637–644. doi: 10.1038/s41562-018-0399-z
- Cohen, J. (1992). A power primer. *Psychological Bulletin*, 112(1), 155–159.
- Copas, J. B. (1983). Regression, prediction and shrinkage (with discussion). *Journal of the Royal Statistical Society: Series B (Methodological)*, 45, 311–354. doi: 10.1111/j.2517-6161.1983.tb01258.x
- Deeks, J. J., Higgins, J. P., & Altman, D. G. (2019). Analysing data and undertaking meta-analyses. In *Cochrane handbook for systematic reviews of interventions* (p. 241–284). John Wiley & Sons, Ltd. doi: 10.1002/9781119536604.ch10
- De Santis, F. (2004). Statistical evidence and sample size determination for Bayesian hypothesis testing. *Journal of Statistical Planning and Inference*, 124(1), 121–144. doi:

10.1016/s0378-3758(03)00198-8

- Errington, T. M., Mathur, M., Soderberg, C. K., Denis, A., Perfito, N., Iorns, E., & Nosek, B. A. (2021). Investigating the replicability of preclinical cancer biology. *eLife*, *10*. doi: 10.7554/elife.71601
- Etz, A., & Vandekerckhove, J. (2016). A Bayesian perspective on the reproducibility project: Psychology. *PLOS ONE*, *11*(2), e0149794. doi: 10.1371/journal.pone.0149794
- Firth, D. (1993). Bias reduction of maximum likelihood estimates. *Biometrika*, *80*(1), 27–38. doi: 10.1093/biomet/80.1.27
- Gelfand, A. E., & Wang, F. (2002). A simulation-based approach to Bayesian sample size determination for performance under a given model and for separating models. *Statistical Science*, *17*(2), 193–208. doi: 10.1214/ss/1030550861
- Gelman, A. (2009). Bayes, Jeffreys, prior distributions and the philosophy of statistics. *Statistical Science*, *24*(2), 76–178. doi: 10.1214/09-sts284d
- Goodman, S. N. (1992). A comment on replication, *p*-values and evidence. *Statistics in Medicine*, *11*(7), 875–879. doi: 10.1002/sim.4780110705
- Grieve, A. P. (2016). Idle thoughts of a ‘well-calibrated’ Bayesian in clinical drug development. *Pharmaceutical Statistics*, *15*(2), 96–108. doi: 10.1002/pst.1736
- Grieve, A. P. (2022). *Hybrid frequentist/Bayesian power and Bayesian power in planning clinical trials*. London, England: Taylor & Francis.
- Harms, C. (2019). A Bayes factor for replications of ANOVA results. *The American Statistician*, *73*(4), 327–339. doi: 10.1080/00031305.2018.1518787
- Hedges, L. V. (1981). Distribution theory for Glass's estimator of effect size and related estimators. *Journal of Educational Statistics*, *6*(2), 107–128. doi: 10.3102/10769986006002107
- Hedges, L. V., & Schauer, J. M. (2019). More than one replication study is needed for unambiguous tests of replication. *Journal of Educational and Behavioral Statistics*, *44*(5), 543–570. doi: 10.3102/1076998619852953
- Hedges, L. V., & Schauer, J. M. (2021). The design of replication studies. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, *184*(3), 868–886. doi:

10.1111/rssa.12688

- Held, L. (2020). A new standard for the analysis and design of replication studies (with discussion). *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 183(2), 431–448. doi: 10.1111/rssa.12493
- Held, L., Matthews, R., Ott, M., & Pawel, S. (2022). Reverse-Bayes methods for evidence assessment and research synthesis. *Research Synthesis Methods*, 13(3), 295–314. doi: 10.1002/jrsm.1538
- Held, L., Micheloud, C., & Pawel, S. (2022). The assessment of replication success based on relative effect size. *The Annals of Applied Statistics*, 16(2), 706–720. doi: 10.1214/21-aos1502
- Held, L., & Pawel, S. (2020). Comment on “the role of  $p$ -values in judging the strength of evidence and realistic replication expectations”. *Statistics in Biopharmaceutical Research*, 13(1), 46–48. doi: 10.1080/19466315.2020.1828161
- Jackson, D., & White, I. R. (2018). When should meta-analysis avoid making hidden normality assumptions? *Biometrical Journal*, 60(6), 1040–1058. doi: 10.1002/bimj.201800071
- Jeffreys, H. (1961). *Theory of probability* (third ed.). Oxford: Clarendon Press.
- Johnson, V. E., Payne, R. D., Wang, T., Asher, A., & Mandal, S. (2016). On the reproducibility of psychological science. *Journal of the American Statistical Association*, 112(517), 1–10. doi: 10.1080/01621459.2016.1240079
- Klein, R. A., Vianello, M., Hasselman, F., Adams, B. G., Reginald B. Adams, J., Alper, S., . . . others (2018). Many labs 2: Investigating variation in replicability across samples and settings. *Advances in Methods and Practices in Psychological Science*, 1(4), 443–490. doi: 10.1177/2515245918810225
- Kunzmann, K., Grayling, M. J., Lee, K. M., Robertson, D. S., Rufibach, K., & Wason, J. M. S. (2021). A review of Bayesian perspectives on sample size derivation for confirmatory trials. *The American Statistician*, 75(4), 424–432. doi: 10.1080/00031305.2021.1901782
- Ly, A., Etz, A., Marsman, M., & Wagenmakers, E.-J. (2018). Replication Bayes factors from

- evidence updating. *Behavior Research Methods*, 51(6), 2498–2508. doi: 10.3758/s13428-018-1092-x
- Maca, J., Gallo, P., Branson, M., & Maurer, W. (2002). Reconsidering some aspects of the two-trials paradigm. *Journal of Biopharmaceutical Statistics*, 12(2), 107–119. doi: 10.1081/bip-120006450
- Mathur, M. B., & VanderWeele, T. J. (2020). New statistical metrics for multisite replication projects. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 183(3), 1145–1166. doi: 10.1111/rssa.12572
- Matthews, J. N. (2006). *Introduction to randomized controlled clinical trials*. New York: Chapman and Hall/CRC. doi: 10.1201/9781420011302
- McKinney, K., Stefan, A., & Gronau, Q. F. (2021). Developing prior distributions for Bayesian meta-analyses. (Preprint) doi: 10.31234/osf.io/2v5bz
- Micheloud, C., & Held, L. (2022). Power calculations for replication studies. *Statistical Science*, 37(3), 369–379. doi: 10.1214/21-sts828
- Munafò, M. R., Nosek, B. A., Bishop, D. V. M., Button, K. S., Chambers, C. D., Sert, N. P., ... Ioannidis, J. P. A. (2017). A manifesto for reproducible science. *Nature Human Behaviour*, 1(0021). doi: 10.1038/s41562-016-0021
- O'Hagan, A. (2019). Expert knowledge elicitation: Subjective but scientific. *The American Statistician*, 73(sup1), 69–81. doi: 10.1080/00031305.2018.1518265
- O'Hagan, A., & Stevens, J. (2001). Bayesian assessment of sample size for clinical trials of cost-effectiveness. *Medical Decision Making*, 21(3), 219–230. doi: 10.1177/02729890122062514
- O'Hagan, A., Stevens, J. W., & Campbell, M. J. (2005). Assurance in clinical trial design. *Pharmaceutical Statistics*, 4(3), 187–201. doi: 10.1002/pst.175
- Open Science Collaboration. (2015). Estimating the reproducibility of psychological science. *Science*, 349(6251), aac4716. doi: 10.1126/science.aac4716
- Park, J., & Pek, J. (2022). Conducting Bayesian-classical hybrid power analysis with R package hybridpower. *Multivariate Behavioral Research*. doi:



10.1080/00273171.2022.2038056

- Patil, P., Peng, R. D., & Leek, J. T. (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. doi: 10.1177/1745691616646366
- Pawel, S., & Held, L. (2020). Probabilistic forecasting of replication studies. *PLOS ONE*, 15(4), e0231416. doi: 10.1371/journal.pone.0231416
- Pawel, S., & Held, L. (2022). The sceptical Bayes factor for the assessment of replication success. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 84(3), 879–911. doi: 10.1111/rssb.12491
- Pek, J., & Park, J. (2019). Complexities in power analysis: Quantifying uncertainties with a Bayesian-classical hybrid approach. *Psychological Methods*, 24(5), 590–605. doi: 10.1037/met0000208
- Protzko, J., Krosnick, J., Nelson, L. D., Nosek, B. A., Axt, J., Berent, M., . . . Schooler, J. (2020). High replicability of newly-discovered social-behavioral findings is achievable. (Preprint) doi: 10.31234/osf.io/n2a9x
- R Core Team. (2023). R: A language and environment for statistical computing [Computer software manual]. Vienna, Austria. Retrieved from <https://www.R-project.org/>
- Raudenbush, S. W. (1997). Statistical analysis and optimal design for cluster randomized trials. *Psychological Methods*, 2(2), 173–185. doi: 10.1037/1082-989x.2.2.173
- Raudenbush, S. W., & Liu, X. (2000). Statistical power and optimal design for multisite randomized trials. *Psychological Methods*, 5(2), 199–213. doi: 10.1037/1082-989x.5.2.199
- Rosenkranz, G. K. (2021). Replicability of studies following a dual-criterion design. *Statistics in Medicine*, 40(18), 4068–4076. doi: 10.1002/sim.9014
- Röver, C., Bender, R., Dias, S., Schmid, C. H., Schmidli, H., Sturtz, S., . . . Friede, T. (2021). On weakly informative prior distributions for the heterogeneity parameter in Bayesian random-effects meta-analysis. *Research Synthesis Methods*, 12(4), 448–474. doi:

10.1002/jrsm.1475

Schönbrodt, F. D., & Wagenmakers, E.-J. (2018). Bayes factor design analysis: Planning for compelling evidence. *Psychonomic Bulletin & Review*, 25(1), 128–142. doi:

10.3758/s13423-017-1230-y

Senn, S. (2002). Letter to the editor: A comment on replication,  $p$ -values and evidence by S. N. Goodman, *Statistics in Medicine* 1992; 11:875–879. *Statistics in Medicine*, 21(16), 2437–2444. doi: 10.1002/sim.1072

Senn, S. (2008). *Statistical issues in drug development* (Vol. 69). John Wiley & Sons.

Simons, D. J. (2014). The value of direct replication. *Perspectives on Psychological Science*, 9(1), 76–80. doi: 10.1177/1745691613514755

Simonsohn, U. (2015). Small telescopes: Detectability and the evaluation of replication results. *Psychological Science*, 26(5), 559–569. doi: 10.1177/0956797614567341

Spiegelhalter, D. J., Abrams, R., & Myles, J. P. (2004). *Bayesian approaches to clinical trials and health-care evaluation*. New York: Wiley.

Spiegelhalter, D. J., & Freedman, L. S. (1986). A predictive approach to selecting the size of a clinical trial, based on subjective clinical opinion. *Statistics in Medicine*, 5(1), 1–13. doi: 10.1002/sim.4780050103

Spiegelhalter, D. J., Freedman, L. S., & Blackburn, P. R. (1986). Monitoring clinical trials: Conditional or predictive power? *Controlled Clinical Trials*, 7(1), 8–17. doi: 10.1016/0197-2456(86)90003-6

Stefan, A., Gronau, Q. F., & Wagenmakers, E.-J. (2022). Interim design analysis using Bayes factor forecasts.

(Preprint) doi: 10.31234/osf.io/9sazk

Steiner, P. M., Wong, V. C., & Anglin, K. (2019). A causal replication framework for designing and assessing replication efforts. *Zeitschrift für Psychologie*, 227(4), 280–292. doi: 10.1027/2151-2604/a000385

Sutton, A. J., & Abrams, K. R. (2001). Bayesian methods in meta-analysis and evidence synthesis. *Statistical Methods in Medical Research*, 10(4), 277–303. doi:

10.1177/096228020101000404

- van Zwet, E. W., & Goodman, S. N. (2022). How large should the next study be? Predictive power and sample size requirements for replication studies. *Statistics in Medicine*, 41(16), 3090–3101. doi: 10.1002/sim.9406
- van Zwet, E. W., Schwab, S., & Senn, S. (2021). The statistical properties of RCTs and a proposal for shrinkage. *Statistics in Medicine*, 40(27), 6107–6117. doi: 10.1002/sim.9173
- van Aert, R. C. M., & van Assen, M. A. L. M. (2017). Bayesian evaluation of effect size after replicating an original study. *PLOS ONE*, 12(4), e0175302. doi: 10.1371/journal.pone.0175302
- Verhagen, J., & Wagenmakers, E. J. (2014). Bayesian tests to quantify the result of a replication attempt. *Journal of Experimental Psychology: General*, 143(4), 1457–1475. doi: 10.1037/a0036731
- Weiss, R. (1997). Bayesian sample size calculations for hypothesis testing. *Journal of the Royal Statistical Society: Series D (The Statistician)*, 46(2), 185–191. doi: 10.1111/1467-9884.00075
- Wong, V. C., Anglin, K., & Steiner, P. M. (2021). Design-based approaches to causal replication studies. *Prevention Science*, 23(5), 723–738. doi: 10.1007/s11121-021-01234-7

### Computational details

```
cat(paste(Sys.time(), Sys.timezone(), "\n"))

## 2023-05-11 13:12:25.336059 Europe/Zurich

sessionInfo()

## R version 4.3.0 (2023-04-21)
## Platform: x86_64-pc-linux-gnu (64-bit)
## Running under: Ubuntu 20.04.6 LTS
##
## Matrix products: default
## BLAS: /usr/lib/x86_64-linux-gnu/blas/libblas.so.3.9.0
## LAPACK: /usr/lib/x86_64-linux-gnu/lapack/liblapack.so.3.9.0
##
## locale:
##  [1] LC_CTYPE=en_US.UTF-8      LC_NUMERIC=C
##  [3] LC_TIME=de_CH.UTF-8      LC_COLLATE=en_US.UTF-8
##  [5] LC_MONETARY=de_CH.UTF-8  LC_MESSAGES=en_US.UTF-8
##  [7] LC_PAPER=de_CH.UTF-8     LC_NAME=C
##  [9] LC_ADDRESS=C             LC_TELEPHONE=C
## [11] LC_MEASUREMENT=de_CH.UTF-8 LC_IDENTIFICATION=C
##
## time zone: Europe/Zurich
## tzcode source: system (glibc)
##
## attached base packages:
## [1] stats      graphics  grDevices  utils      datasets  methods    base
##
## other attached packages:
##  [1] BayesRep_0.1              BayesRepDesign_0.42      ggpubr_0.6.0
##  [4] colorspace_2.1-0         scales_1.2.1             ggplot2_3.4.2
##  [7] tidyr_1.3.0              dplyr_1.1.2              ReplicationSuccess_1.3.1
## [10] knitr_1.42
##
## loaded via a namespace (and not attached):
```

```
## [1] gtable_0.3.3      highr_0.10      compiler_4.3.0  ggsignif_0.6.4
## [5] tidyselect_1.2.0  Rcpp_1.0.10     gridExtra_2.3   R6_2.5.1
## [9] labeling_0.4.2    generics_0.1.3  backports_1.4.1  tibble_3.2.1
## [13] car_3.1-2         munsell_0.5.0   pillar_1.9.0     rlang_1.1.0
## [17] utf8_1.2.3        broom_1.0.4     xfun_0.39        RcppParallel_5.1.7
## [21] cli_3.6.1         withr_2.5.0     magrittr_2.0.3   grid_4.3.0
## [25] cowplot_1.1.1     lifecycle_1.0.3 lamW_2.1.2        vctrs_0.6.2
## [29] rstatix_0.7.2     evaluate_0.20   glue_1.6.2       farver_2.1.1
## [33] abind_1.4-5       carData_3.0-5   fansi_1.0.4      purrr_1.0.1
## [37] tools_4.3.0       pkgconfig_2.0.3
```