

# Statistical Assessment of Replicability via Bayesian Model Criticism

Yi Zhao\* and Xiaoquan Wen†

Department of Biostatistics, University of Michigan, Ann Arbor, MI, USA

## Abstract

Assessment of replicability is critical to ensure the quality and rigor of scientific research. In this paper, we discuss inference and modeling principles for replicability assessment. Targeting distinct application scenarios, we propose two types of Bayesian model criticism approaches to identify potentially irreproducible results in scientific experiments. They are motivated by established Bayesian prior and posterior predictive model-checking procedures and generalize many existing replicability assessment methods. Finally, we discuss the statistical properties of the proposed replicability assessment approaches and illustrate their usages by simulations and examples of real data analysis, including the data from the Reproducibility Project: Psychology and a systematic review of impacts of pre-existing cardiovascular disease on COVID-19 outcomes.

---

\*zhayi@umich.edu

†xwen@umich.edu

# 1 Introduction

Reproducibility is a hallmark of scientific research. It is a necessary characteristic to ensure correctness and rigor for scientific discoveries [1, 2, 3, 4, 5, 6, 7]. Despite the rising awareness across all scientific disciplines, there is a general lack of theoretical foundation and methodological work focusing on the statistical assessment of reproducibility/replicability in scientific practice. Consequently, inconsistent specifications of *replication successes* and misuses of subjective measures for assessing reproducibility are common in scientific practice [8, 9, 10]. For example, in reporting the findings from Open Science Collaboration’s Reproducibility Project: Psychology (RP:P) [11], media reports focused on one particular metric that more than half of the replications fail to reproduce statistically significant results in the expected directions and made a sensational statement that “more than half of the scientific results are not reproducible”. However, using repeated statistical significance to define replication success in such contexts is considered deeply flawed [12, 13, 14, 15, 16]. Because variability is at the core of reproducibility [17, 18, 19], statistics is uniquely positioned for providing sound and rigorous solutions. In a recent review, He and Lin list “Statistical Methods for Reproducibility and Replicability” as one of the ten immediate challenges and opportunities in statistics and data science [20].

In this paper, we focus on a particular mode of reproducibility, known as *replicability* or *results reproducibility*, according to the lexicon of reproducibility by Goodman *et al.* [5]. Roughly speaking, it refers to the consistency of the results from analyzing *different* data collected to address the same underlying scientific questions. For simplicity, we use the terminology *replicability* and *reproducibility* interchangeably throughout this paper.

There is a large volume of existing methodological works from relevant fields (e.g., systematic review and meta-analysis) related to replicability assessment [21, 22, 23]. Some are motivated by investigating effects of irreproducibility due to specific factors, e.g., publication bias/winner’s curse [24, 25, 26, 27, 28, 15] and batch effects [29, 30, 31], in specific settings. They utilize

different inference strategies, e.g., model comparison [32, 19], predictive interval checking [14, 33], and deals with distinct application scenarios. This diverse body of work lays a solid foundation for us to summarize important statistical principles in replication assessment. In this work, we take a Bayesian model criticism strategy for replicability assessment. The rich literature [34, 35, 36, 37, 38, 39] and the well-developed computational infrastructure on Bayesian model criticism make it feasible and attractive for our applications of interest.

In the rest of the paper, we first discuss statistical considerations for replicability assessment and argue that Bayesian model criticism is appropriate and effective for this purpose. We then propose two different model criticism approaches, targeting different application scenarios. Finally, we illustrate their usages through simulation studies and examples of real data analysis. We conclude the paper by discussing important technical issues in the proposed approaches and their connections to alternative statistical strategies.

The software implementation of the proposed statistical methods (R package, PRP) and the relevant code to reproduce all the simulation and data analysis results can be downloaded from the GitHub repository: <https://github.com/ArtemisZhao/PRP>.

## 2 Statistical Considerations for Assessing Replicability

This section first lays out general modeling and inference principles in defining replication successes and assessing replicability. Based on these guiding principles, we will further consider different application scenarios and discuss potential statistical solutions.

## 2.1 Modeling and Inference Principles

Consider a typical setup of replicability assessment where different experiments are conducted to address the same scientific questions. Two distinct sources contribute to the variability of analysis results from different experimental datasets. First, random noise is intrinsic to each dataset (and generally considered independent across datasets). Second, natural variations of underlying true effects of interest are expected across different experiments. The observed variations among analysis results reflect their combined effects. In assessing replicability, we argue that evaluating variability from the second source should be the sole focus. This principle is consistent with the common practice in quantifying and controlling heterogeneity in meta-analysis and systematic review literature [21, 26, 27].

Even though experiment-specific random noise is considered a nuisance in replicability assessment, the determination of the true effects is inevitably confounded by its existence. The following question best characterizes the role of random noise: should an extremely noisy replication experiment be considered as evidence against the replicability? The logical answer seems to be no. As accurate estimation of true underlying effects becomes infeasible in the presence of a high level of random experimental errors, the replication data are essentially *non-informative*, which is fundamentally distinct from evidence against reproducibility. This line of reasoning establishes an asymptotic behavior for desired replicability assessment approaches, which we refer to as the non-informativeness principle of noisy replications henceforth. An important corollary from the principle is the requirement of explicit uncertainty specification of observed effects as integral input information for replicability assessment. In comparison, characterizing signal-to-noise ratios as a compound quantity (e.g.,  $p$ -values derived from pivotal  $z$ -statistics) may be insufficient and violate the very principle (as the noise level can not be separately recovered).

Establishing an extent of heterogeneity for reproducible effects is another critical aspect for defining replication success. It is overly unrealistic and unnecessarily restrictive to assume re-

producible results must have *identical* underlying effects [21, 27]. Ultimately, defining acceptable heterogeneity in successful replications should be context-specific. Nevertheless, we argue that it is plausible to specify an *a priori* minimum requirement of the maximum tolerable heterogeneity generally, which is analogous to applying the minimax principle in decision theory. Recently, we have proposed such a standard, known as the directional consistency (DC) criterion, which emphasizes that reproducible effects are expected to show concordant (positive or negative) signs across replication experiments with a high probability [19]. Similar ideas have been successfully adopted in general hypothesis testing, meta-analysis, system reviews, and the statistical analysis of qualitative interactions [40, 41, 42, 43, 23]. It is worth noting that the implementation of DC the criterion in [19] is not *solely* focusing on the variability on a particular direction (i.e., the sign-flipping direction). It utilizes a variance parameter to define a range of acceptable heterogeneity centering around a true latent effect. This paper illustrates our proposed statistical methods by applying the DC criterion but note that other suitable definitions for acceptable replication variations are also applicable.

## 2.2 Application Scenarios

In practice, replicability assessment is applied in various distinct scenarios, representing some unique characteristics and requiring context-specific statistical treatments.

In one of the common application scenarios, a designated replication experiment is conducted to validate the result from an original study. The original and replication labels are typically assigned following the chronicle order of data generation. The RP:P and Experimental economics replication project [44] studies are representative examples of such kind. The setup can generally be extended to two groups of experiments (i.e., the original group and the replication group). The scientific question of interest is whether the results derived from the replication group confirm the finding from the original group. For such a design, the replication assessment should be

specific to the assigned labels, i.e., we may expect qualitatively different results if the original and replication labels are switched in the analysis. We will henceforth refer to this scenario as the two-group scenario.

Another application scenario of replicability assessment often arises from systematic review and meta-analysis, where a single group of multiple experiments is gathered. The relevant scientific question here is to assess the overall concordance of all experiments and identify potential outlying results. Notably, the chronicle order of data generation usually plays no role in the data analysis, and switching chronicle labels of the participating experiments should not yield quantitatively different assessment results. We will refer to this particular scenario as the (chronically) exchangeable group scenario henceforth.

From a statistical perspective, the two distinct scenarios above suggest different statistical treatments in replicability assessment: the analysis in the two-group scenario calls for explicit conditioning of the original experiment, whereas the chronically exchangeable structure among experiments should be preserved in the second scenario.

There are other distinguishing factors in applications of replicability assessment. For example, most experiments in physical and social sciences contain a single analysis unit. However, in biological science, high-throughput experiments, where tens of thousands of biological units (e.g., genes) are simultaneously measured, have become increasingly common. Thus, the unit- and study-level assessments are no longer synonymous and require specialized statistical treatments [32, 19]. This paper will focus on the former case, and we further address their connections and distinctions in the Discussion section.

## 2.3 Model Criticism Strategy

In this paper, we propose to assess replicability by applying a model criticism strategy. Specifically, we formulate a parametric model representing the expected characteristics from reproducible results. Subsequently, we fit such a model with observed data and evaluate the goodness-of-fit. Specifically, poor-fitting prompts to reject the notation of the observed data are likely reproducible. Because our employed parametric models are (Bayesian) hierarchical, we apply the principled framework of Bayesian model checking approaches for replicability assessment. Particularly, we utilize both prior checking [34] and posterior checking [35, 37] techniques to deal with different application scenarios.

In our proposed methods,  $p$ -value is one of the primary statistical instruments to summarize and quantify the results from model criticism procedures. Our usage of  $p$ -values in this context closely relates to the Fisherian significance testing and directly follows the applications in Bayesian model checking [34, 35]. Specifically,  $p$ -values quantify the discrepancy between the assumed reproducible model and the observed data, as in the original Fisherian design. We also adopt Fisher’s disjunction to interpret a small  $p$ -value, which indicates “either an exceptionally rare chance has occurred or the theory is not true” [45]. Here, the specific “theory” refers to the built-in reproducible assumptions. Additionally, our use and subsequent interpretations of  $p$ -values in the context of replicability assessment strictly follow the ASA’s statement on  $p$ -values [46].

Critically, our use of model criticism approaches and  $p$ -values should not be confused with the binary decision procedures characterized by the Neyman-Pearson hypothesis testing [47] and the hybrid null hypothesis significance testing (NHST). Our goal, in this paper, is *not* to classify the observed data into reproducible and irreproducible categories. Hence, the  $\alpha$  and  $\beta$  levels for type I and type II errors are irrelevant in our proposed statistical framework. On the other hand, we need to elucidate the statistical properties of the proposed procedures in dealing with truly replicable and irreproducible results. Specifically, the calibration of the proposed  $p$ -values

under the reproducible model and the sensitivity of the procedures under various irreproducible scenarios [47]. We will use dedicated terminologies in such cases to avoid confusion. The distinction between Fisherian significance testing and Neyman-Pearson style NHST has profound implications in our treatments of some important statistical issues, e.g., the simultaneous examinations of multiple replication experiments. The relevant points will be fully illustrated in the subsequent sections.

Notwithstanding the importance of  $p$ -values in quantifying evidence, our model criticism strategy for replication assessment also utilizes other suitable statistical instruments and visualization tools. We emphasize that replication assessment via model criticism is a comprehensive data analysis process and should not be simplified into a  $p$ -value-generating procedure.

## 3 Two Model Criticism Approaches

### 3.1 Model and Notation

Consider  $m$  experiments to address the same scientific question. The complete collection of the experiment data are denoted by  $(\mathbf{X}_1, \dots, \mathbf{X}_m)$ . Let  $\hat{\beta}_i$  and  $\sigma_i$  denote the measurement and its corresponding standard error of the underlying effect of interest from the  $i$ -th experiment, respectively. We assume, without loss of generality,  $(\hat{\beta}_i, \sigma_i)$  is the sufficient statistic derived from  $\mathbf{X}_i$ . Let  $\beta_i$  denote the true latent effect for the corresponding experiment and assume the grand effect in the population all (hypothetical) experiments is represented by  $\bar{\beta}$ . Abstracting away from the details of the experiment design and data analysis, we consider the following simple



probabilistic generative model:

$$\begin{aligned}
 \bar{\beta} &\sim \text{N}(0, \omega^2) \\
 \beta_i &= \bar{\beta} + \eta_i, \quad \eta_i \sim \text{N}(0, \phi^2) \\
 \hat{\beta}_i &= \beta_i + e_i, \quad e_i \sim \text{N}(0, \sigma_i^2)
 \end{aligned} \tag{1}$$

Both  $\eta_i$ 's and  $e_i$ 's are assumed mutually independent across experiments.

To formulate a reproducible model, we apply the DC criterion to constrain the variability of  $\beta_i$ . Specifically, we parametrize  $\gamma = \frac{\phi^2}{\phi^2 + \omega^2}$  and note that

$$p(\gamma) = \Pr(\text{sgn}(\beta_i) = \text{sgn}(\bar{\beta}) \mid \gamma) = \sqrt{\frac{2}{\pi}} \int_0^{+\infty} \Phi\left(\sqrt{\frac{1-\gamma}{\gamma}} \xi\right) e^{-\frac{\xi^2}{2}} d\xi, \tag{2}$$

where  $\Phi(\cdot)$  denotes the CDF of the standard normal distribution. Thus, we can define a level of acceptable variability through the sign-consistency probability,  $p(\gamma)$ , under the DC principle [19].

To complete a reference reproducible model covering diverse practical cases, we consider a grid of  $(\omega^2, \gamma)$  from a Cartesian product  $\Omega \times \Gamma$ . Each grid value is assigned equal prior probability,  $\pi = 1/(|\Omega| \cdot |\Gamma|)$ . As the heterogeneity parameter  $\gamma$  is of primary interest for inference and the effect size parameter  $\omega$  is regarded as a nuisance parameter, we consider  $\Gamma = \{1.00, 0.99, 0.975, 0.95\}$  and construct  $\Omega$  adaptively according to observed data and various application scenarios, by default. (See Appendix A for details). The default construction of  $\Gamma$  specifies a lower bound for the sign-consistency probability (i.e., 0.95) but also allows more stringent criteria. The discrete construction for the overall hyperparameter space results in a finite mixture model as the reference reproducible model. Thus, it provides not only modeling flexibility and comprehensiveness but also computational convenience for inference procedures.

For inference,  $\hat{\boldsymbol{\beta}} = (\hat{\beta}_1, \dots, \hat{\beta}_m)$  are assumed observed. Without loss of generality in demonstrating

the proposed methods, we also consider the true standard errors  $\boldsymbol{\sigma} = (\sigma_1, \dots, \sigma_m)$  are known or accurately measured. Additionally, we illustrate the applications in the two-group scenario by considering two experiments ( $m = 2$ ), which are explicitly labeled as original and replication, respectively.

### 3.2 Prior Predictive Checking (for Two-group Scenario)

The replication assessment in the two-group scenario requires explicit conditioning on the original experiment result. It naturally leads to a predictive checking procedure in the Bayesian framework. Specifically, a Bayesian predictive distribution is computed conditional on the original experiment and a reference reproducible model (denoted by  $M_R$ ), namely,  $P(\mathbf{X} \mid \mathbf{X}_{\text{orig}}, M_R)$ . A  $p$ -value is subsequently evaluated to quantify the probability of obtaining results at least as extreme as the observed replication data, assuming the predictive distribution. A small  $p$ -value indicates poor model fitting by the replication data to the predictive distribution constructed from the original data and raises attention to potential irreproducibility. Henceforth, we call the  $p$ -values derived from this procedure as the *prior-predictive replication p-values* (prior-PRPs). Let  $T(\mathbf{X})$  denote a pre-defined test statistic, the one-sided version of the prior-predictive replication  $p$ -value is given by

$$p_{\text{prior}} := \Pr(T(\mathbf{X}) \geq T(\mathbf{X}_{\text{rep}}) \mid \mathbf{X}_{\text{orig}}, \mathbf{X}_{\text{rep}}, M_R); \quad (3)$$

and correspondingly, the two-sided version is defined by

$$p_{\text{prior}} := 2 \min \left\{ \Pr(T(\mathbf{X}) \geq T(\mathbf{X}_{\text{rep}}) \mid \mathbf{X}_{\text{orig}}, \mathbf{X}_{\text{rep}}, M_R), \right. \\ \left. \Pr(T(\mathbf{X}) \leq T(\mathbf{X}_{\text{rep}}) \mid \mathbf{X}_{\text{orig}}, \mathbf{X}_{\text{rep}}, M_R) \right\} \quad (4)$$

In the proposed reference reproducible model applying the DC criterion,  $T(\mathbf{X}_{\text{orig}}) = \hat{\beta}_{\text{orig}}$  and

$T(\mathbf{X}_{\text{rep}}) = \hat{\beta}_{\text{rep}}$  are natural test statistics. Under the mixture prior on  $(\omega, \gamma)$ , the required posterior and predictive distributions based on the original data are analytically available for evaluating  $p$ -values (Appendix D).

The proposed procedure is considered as Bayesian prior predictive checking regarding the analysis of the replication data. In the setting of the two-group scenario and following the principle of sequential Bayesian updating, the natural prior for estimating  $\bar{\beta}$  (or  $\beta_{\text{rep}}$ ) from the replication data is  $P(\bar{\beta} \mid \hat{\beta}_{\text{orig}}, M_{\text{R}})$ , which is the posterior from analyzing the original experiment. It is important to note that, under  $M_{\text{R}}$ , the original experiment can be highly informative on the magnitude of  $\bar{\beta}$  and  $\omega$ , but it lacks information to alter the prior assumption on the heterogeneity parameter  $\gamma$ . From this perspective, the replication data is used to examine the prior reproducibility assumption implementing the DC criterion. To improve the interpretability of the proposed approach, we explicitly set the grid values of  $\omega$  adaptive to the original data (Appendix A).

The prior-predictive replication  $p$ -values are  $u$ -values [38, 39] defined by the following proposition.

**Proposition 1.** *Prior-predictive replication  $p$ -values are uniformly distributed on  $[0, 1]$  under the assumed reference reproducible model.*

*Proof.* Appendix B □

*Remark.* The uniform distribution is established under the repeated sampling of both  $\mathbf{X}_{\text{orig}}$  and  $\mathbf{X}_{\text{rep}}$ . The exact uniformity is attained only under the *true* reference model.

The theoretical uniformity of the prior-predictive replication  $p$ -value is most helpful for benchmarking calibrations and evaluating sensitivity under a simulation setting when ground truths are known. In practice, it provides a useful reference for simultaneous inspections of multiple sets of replication experiments. For example, in RP:P,  $\sim 100$  pairs of original and replication experiments are selected to survey the replicability in psychology research. If there is no *system-*

*atic* irreproducibility, an empirical distribution of prior-predictive replication  $p$ -values pooling all studies should closely resemble a discrete uniform distribution. However, if the empirical distribution severely deviates from the uniformity and indicates clear positive skewness, the overall quality of replicability in the investigated set is questionable. More generally, particular patterns of deviations from uniformity are informative on the average heterogeneity in the examined set comparing to the reference model.

The two-sided prior-predictive replication  $p$ -values are connected to the predictive interval checking procedures in replication assessment [14, 33]. Specifically, the interval-checking procedures construct a  $(1 - \alpha)\%$  predictive interval based on an assumed reproducible model and the original experiment. They then examine if the observed replication data fall into the predictive intervals. Let  $q_\alpha(T(\mathbf{X}) \mid \mathbf{X}_{\text{orig}}, M_{\text{R}})$  denote the quantile function of the predictive distribution of  $T(\mathbf{X})$  based on the original experimental data, such that  $\Pr(T(\mathbf{X}) \leq q_\alpha(T(\mathbf{X}) \mid \mathbf{X}_{\text{orig}}, M_{\text{R}})) = \alpha$ . The predictive interval,  $(q_{\alpha/2}(T(\mathbf{X}) \mid \mathbf{X}_{\text{orig}}, M_{\text{R}}), q_{(1-\alpha/2)}(T(\mathbf{X}) \mid \mathbf{X}_{\text{orig}}, M_{\text{R}}))$ , is known as a  $(1 - \alpha)$  central interval [48]. The following proposition represents the equivalence of this specific form of interval-checking procedure and the model criticism approach using the two-sided prior-predictive  $p$ -values derived from the same predictive distribution.

**Proposition 2.**  $T(\mathbf{X}_{\text{rep}}) \notin (q_{\alpha/2}(T(\mathbf{X}) \mid \mathbf{X}_{\text{orig}}, M_{\text{R}}), q_{(1-\alpha/2)}(T(\mathbf{X}) \mid \mathbf{X}_{\text{orig}}, M_{\text{R}}))$ , if and only if the corresponding two-sided  $p_{\text{prior}} < \alpha$ .

*Proof.* Appendix C □

*Remark.* The similar equivalence can be generalized to one-sided prior-predictive replication  $p$ -values and the corresponding one-sided predictive intervals.

The connection described by Proposition 2 provides not only an alternative interpretation of the proposed  $p$ -values but also a way of visualizing replication assessment. Additionally, it highlights the distinction in the usage of  $p$ -values between our proposed model criticism framework and

the NHST, especially in a setting of multiple comparisons. As the type I error is irrelevant in our context, there is no need to adjust the predefined significance threshold (e.g., by Bonferroni correction) when multiple sets of experiments are simultaneously examined. Similarly, it seems illogical to modify a well-defined predictive interval (derived from the corresponding original experiment) just because of the co-existence of other sets of the experiments.

It is worth pointing out that, from the model criticism perspective, replication data falling within the corresponding predictive interval is not evidence *in favor of* replication success. This argument is based on Proposition 2 and the interpretation of relatively large  $p$ -values (Principle 6 of the ASA statement on  $p$ -values, [46]). Only the rejection region, defined by the complement of a predictive interval in Proposition 2, is informative for evidence *against* the reference replication model. A useful corollary is that the length of a rejection region can be used to measure the (relative) informativeness of the procedure.

By default, we use  $T(\mathbf{X}_{\text{orig}}) = \hat{\beta}_{\text{orig}}$ ,  $T(\mathbf{X}_{\text{rep}}) = \hat{\beta}_{\text{rep}}$  and compute two-sided  $p$ -values for replicability assessment under the two-group scenario. Its predictive distribution based on the reference model,  $M_R$ , can be analytically derived (Appendix D). The derivation shows that replication noise, characterized by  $\sigma_{\text{rep}}$ , does not change the location information of the prior predictive distribution. However, the length of a predictive interval monotonically increases with respect to  $\sigma_{\text{rep}}$ . Correspondingly, the length of the are of resulting rejection region is monotonically decreasing. That is, the procedure becomes less informative as replication noise increases, which satisfies the non-informativeness principle.

### 3.3 Posterior Predictive Checking (for Exchangeable-group Scenario)

In the exchangeable-group scenario, the replicability assessment is conditional on all observed experiments. To this end, we fit the reference model using the data from all experiments and compute a Bayesian  $p$ -value [35, 36, 37] based on the corresponding posterior predictive distri-

bution,  $P(\mathbf{X} \mid \mathbf{X}_1, \dots, \mathbf{X}_m, M_R)$ . We refer to the Bayesian  $p$ -values derived from this procedure as the *posterior-predictive replication  $p$ -values* (posterior-PRPs).

A unique property of the posterior predictive checking is its usage of test quantities (also known as discrepancy variables). A test quantity, denoted by  $T(\mathbf{X}, \boldsymbol{\theta})$ , is a function of both observed data ( $\mathbf{X}$ ) and latent variables/hyperparameters ( $\boldsymbol{\theta}$ ) defined by a reference model. In comparison, traditional test statistics are functions of only observed data ( $\mathbf{X}$ ). A one-sided posterior-predictive replication  $p$ -value is defined and computed by

$$\begin{aligned}
 p_{\text{posterior}} &:= \Pr \left( T(\mathbf{X}'_1, \dots, \mathbf{X}'_m, \boldsymbol{\theta}) \geq T(\mathbf{X}_1, \dots, \mathbf{X}_m, \boldsymbol{\theta}) \mid \mathbf{X}_1, \dots, \mathbf{X}_m, M_R \right) \\
 &= \int_{(\boldsymbol{\theta}, \mathbf{X}'_1, \dots, \mathbf{X}'_m)} \mathbb{1}(T(\mathbf{X}'_1, \dots, \mathbf{X}'_m, \boldsymbol{\theta}) \geq T(\mathbf{X}_1, \dots, \mathbf{X}_m, \boldsymbol{\theta})) \\
 &\quad \cdot \left( \prod_{i=1}^m P(\mathbf{X}'_i \mid \boldsymbol{\theta}) \right) P(\boldsymbol{\theta} \mid \mathbf{X}_1, \dots, \mathbf{X}_m) d\boldsymbol{\theta} d\mathbf{X}'_1 \dots d\mathbf{X}'_m.
 \end{aligned} \tag{5}$$

We use a posterior sampling scheme, outlined in Algorithm 1, to approximate the exact  $p$ -values. The relevant mathematical details in the algorithm are provided in Appendix E.

---

**Algorithm 1** Computing posterior-predictive replication  $p$ -values

---

```

1 procedure APPROXIMATING  $p$ -VALUE BY POSTERIOR SAMPLING
2   compute posterior distribution  $P(\boldsymbol{\theta} \mid \mathbf{X}_1, \dots, \mathbf{X}_m, M_R)$ 
3   initialize counter  $l \leftarrow 0$ 
4   for  $k = 1$  to  $L$  do
5     sample  $\boldsymbol{\theta} \sim P(\boldsymbol{\theta} \mid \mathbf{X}_1, \dots, \mathbf{X}_m, M_R)$ 
6     independently sample  $\mathbf{X}'_i \sim P(\mathbf{X} \mid \boldsymbol{\theta}, M_R)$  for  $i = 1, \dots, m$ 
7     evaluate  $T(\mathbf{X}_1, \dots, \mathbf{X}_m, \boldsymbol{\theta})$  and  $T(\mathbf{X}'_1, \dots, \mathbf{X}'_m, \boldsymbol{\theta})$ 
8     set  $l \leftarrow l + 1$  if  $T(\mathbf{X}'_1, \dots, \mathbf{X}'_m, \boldsymbol{\theta}) \geq T(\mathbf{X}_1, \dots, \mathbf{X}_m, \boldsymbol{\theta})$ 
9   end for
10  return  $p$ -value =  $l/L$ 
11 end procedure

```

---

The interpretation of the posterior-PRPs is not different from the prior-PRPs: a small value indicates poor model fitting and suggests that the observed data are unlikely under the assumed reference model. In our implementation, we set the hyperparameter,  $\omega$ , adaptively to

the observed data (Appendix A). Hence, the lack of model fit implies potential violations of the replicability assumption defined by the DC criterion.

The posterior-predictive replication  $p$ -values, as a special case of Bayesian  $p$ -values, are not necessarily uniformly distributed on  $[0, 1]$  under the assumed reference model and repeated sampling of  $(\mathbf{X}_1, \dots, \mathbf{X}_m)$ . Meng [36] shows that under the true reference model, a Bayesian  $p$ -value can be stochastically less variable than a uniform distribution on  $[0, 1]$  but with the same mean (1/2). This implies that extremely small Bayesian  $p$ -values are less likely under the true reference model than a frequentist (or prior-predictive replication)  $p$ -value. As we have emphasized that the proposed model criticism approaches fundamentally differ from the Neyman-Pearson hypothesis testing and NHST, the non-uniformity property has little impact on interpreting results from evaluating an individual set of experiments. This is because the Bayesian  $p$ -values are valid posterior probabilities, as Gelman points out in [38, 39], and their usage in model criticism approaches is well-justified. On the other hand, it is relevant to establish a benchmark if the goal is to jointly evaluate multiple sets of experiments (as we have discussed in the case of prior-PRPs). In such a scenario, generating empirical distribution of the posterior-predictive replication  $p$ -values under the well-defined reference model for corresponding test quantities should suffice.

Motivated by Cochran's  $Q$  statistic, we use a test quantity,  $T(\mathbf{X}, \boldsymbol{\theta}) = T_Q(\hat{\boldsymbol{\beta}}, (\bar{\beta}, \phi))$ , by default in our proposed reference reproducible model, i.e.,

$$T_Q(\hat{\boldsymbol{\beta}}, (\bar{\beta}, \phi)) = \sum_{i=1}^m w_i (\hat{\beta}_i - \bar{\beta})^2, \text{ where } w_i = \frac{1}{\sigma_i^2 + \phi^2} \quad (6)$$

Similar to Cochran's  $Q$  statistic, this quantity measures the spread of the observed effects and quantifies the heterogeneity across  $m$  participating experiments. But unlike the  $Q$  statistic, which is computed under a fixed-effect assumption,  $\phi^2$  is no longer constrained to be 0.

## 4 Constructing Test Statistics and Quantities

Although the proposed prior- and posterior-predictive checking procedures using the default test statistics/quantities serve the general purpose of replicability assessment, some applications attempt to investigate specific mechanisms of irreproducibility that lead to particular patterns in observed data. To this end, the most effective way is to design and incorporate special-purpose test statistics/quantities into the proposed model criticism framework.

Take detection of publication bias as an illustrating example in two-group scenarios. Because publication bias often leads to overestimating true effects in original experiments, the re-estimated effects in replication experiments tend to shrink noticeably towards 0. To capture such particular directional shift of effect estimates, we consider the following test statistic,

$$T_{\text{pb}} = \frac{\hat{\beta}_{\text{rep}}}{\hat{\beta}_{\text{orig}}}. \quad (7)$$

Stronger shrinkage effects lead  $T_{\text{pb}}$  further away from 1, regardless of the sign of  $\hat{\beta}_{\text{orig}}$ . Thus, a one-sided prior-predictive replication  $p$ -value, i.e.,

$$p_{\text{prior, pb}} := \Pr \left( T_{\text{pb}} \leq \frac{\hat{\beta}_{\text{rep}}}{\hat{\beta}_{\text{orig}}} \mid \hat{\beta}_{\text{orig}}, \hat{\beta}_{\text{rep}}, \text{MR} \right) \quad (8)$$

can be applied to detect patterns of publication bias. The resulting procedure show much improved sensitivity to publication bias than the default statistic.

In the exchangeable-group scenario, detecting publication bias is a long-standing problem in meta-analysis and systematic review. The patterns of publication bias are traditionally investigated by detecting the asymmetry of funnel plots. The inspection can be analytically carried out by Egger regression or other types of meta-regression procedures. These approaches examine the correlation between the estimated effect,  $\hat{\beta}_i$ , and the corresponding standard error,  $\sigma_i$ ,



which characterizes a potential linear trend in funnel plots [49]. The implementation by Egger regression estimates the correlation in a weighted linear regression model, whose test statistic can be trivially replicated in the framework of posterior-PRPs. In the presence of genuine between-experiment heterogeneity,  $\phi^2$ , the theory requires to compute the correlation between  $\hat{\beta}_i$  and  $\sqrt{\sigma_i^2 + \phi^2}$ . The prevailing approach is to plug in an unconstrained point estimate of  $\phi^2$ . However, the accuracy of the point estimate is questionable, especially when the number of available experiments ( $m$ ) is limited, and the solution also ignores its uncertainty. The proposed Bayesian model criticism approach naturally addresses this issue by explicitly sampling reference model-defined  $\phi^2$  from the posterior distribution in computing the test quantity.

In summary, constructing test statistics for computing prior-PRPs is not fundamentally different from the common practice of significance testing, whereas the test quantities in posterior-predicting replication  $p$ -value computation present some unique methodological and practical advantages. First, Bayesian  $p$ -values allow direct incorporation of hyperparameters and latent variables (i.e.,  $\theta$ ), which are typically considered nuisance parameters in constructing traditional test statistics. The Bayesian procedure provides a principled way to account for their uncertainty without restricting the test quantities to be pivotal. This feature allows for adopting more rational reference models permitting reasonable between-experiment heterogeneity in the context of replicability assessment. Second, the “null” distribution of the test quantity does not need to be pre-defined or pre-computed. This flexibility allows practitioners to emphasize context-dependent information in constructing relevant test quantities rather than focus on their statistical properties.

## 5 Numerical Illustrations

We use simulation studies to mimic two common phenomena that can lead to irreproducible results, namely, batch effects and publication bias. We examine the behaviors of the proposed

statistical methods in different application scenarios.

## 5.1 Simulations of Batch Effects

In association analysis, batch effects refer to unaccounted experimental factors causing spurious correlations. We consider a balanced case-control analysis of a quantitative trait with 200 individuals per experiment. For each experiment contaminated by batch effects, we simulate binary batch labels for all samples by permuting their case-control labels. On average, the correlation between the batch labels and the case-control status is 0.7. The continuous outcome of interest is generated from a linear regression model.

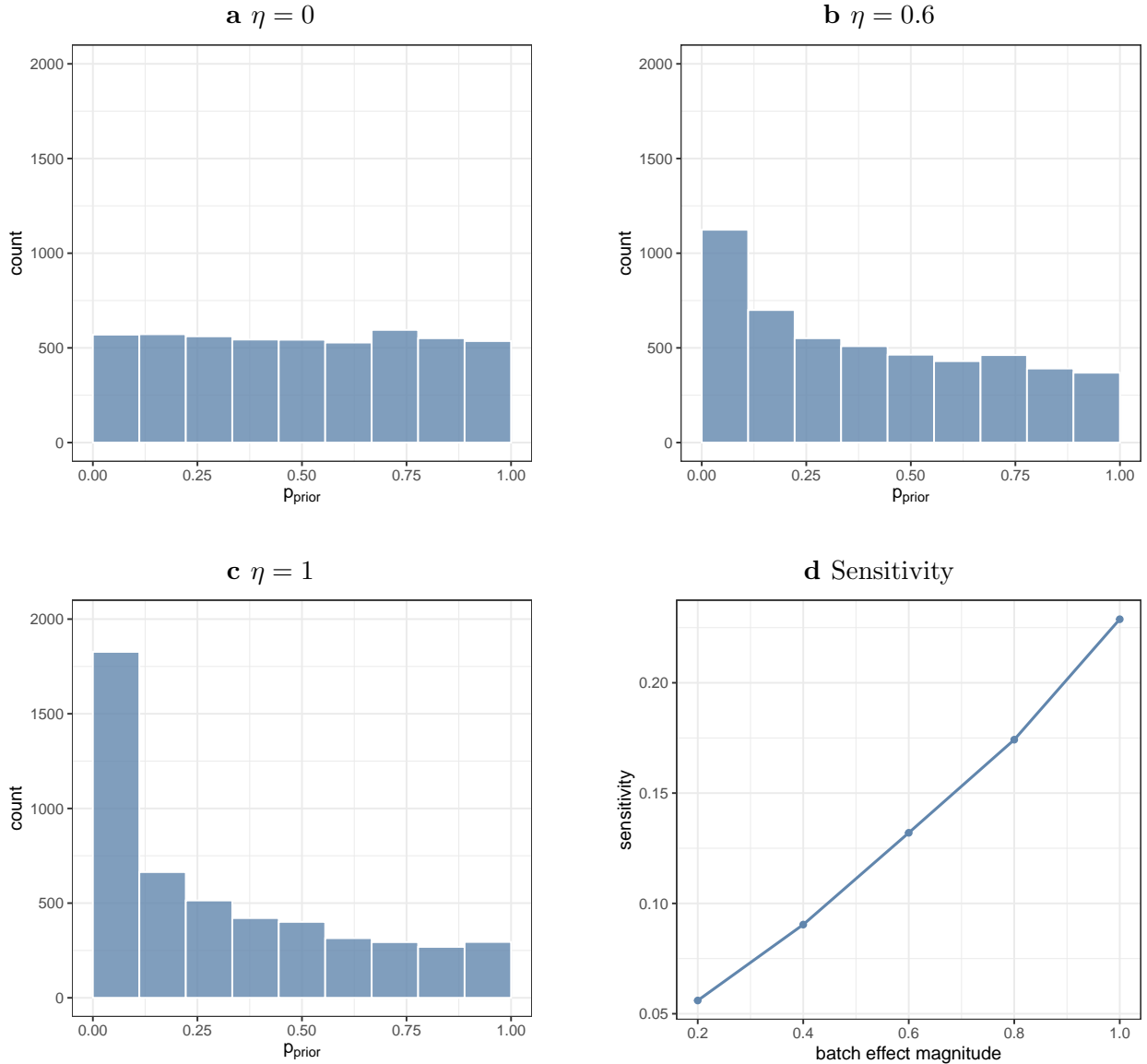
### 5.1.1 Batch Effect Detection in Two-group Scenario

To create a two-group scenario, we simulate two experiments following the general descriptions above. More specifically, the data from the original experiment is batch contaminated, where the batch effect is drawn from the distribution,  $N(0, \eta^2)$ . The replication experiment is generated free of batch effect. The true case-control associations are fixed at 0.5 for both original and replication experiments, and the residual errors for each individual are independently sampled from the standard normal distribution. We vary  $\eta$  from 0 to 1 to create different levels of contamination. For each experiment, regardless of its contamination status, we fit a simple linear model by regressing the simulated outcome on the corresponding case-control status. The resulting association effect estimate and the corresponding standard error are used as summary statistics for input. We repeat the data generation and analysis procedures 5,000 times for each experimented  $\eta$  value.

We compute the prior-PRP for each simulated dataset, assuming the default reproducible model. A simulated dataset is flagged if its corresponding prior-predictive  $p$ -value is smaller than the

commonly used significance threshold 0.05. We define *sensitivity* as the proportion of flagged datasets for each experimented non-zero  $\eta$  value.

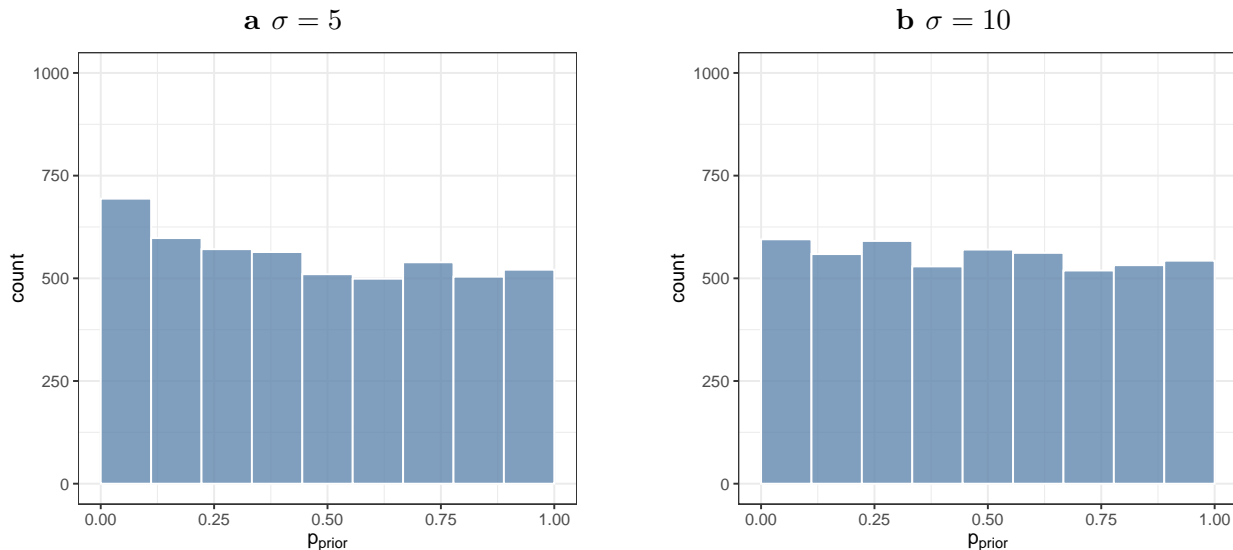
The simulation results are summarized in Figure 1. The empirical distribution of the prior-predictive  $p$ -values under  $\eta = 0$  resembles a discrete uniform distribution, suggesting the desired calibration under the reference model. As expected, the sensitivity in detecting irreproducibility increases as the magnitude of batch effects increase. The empirical  $p$ -value distributions also show increased positive skewness as  $\eta$  increases.



**Figure 1: Batch effect detection in two-group scenario.** Panels **a**, **b**, and **c** present the histograms of the two-sided prior-PRPs from different simulation settings, where the magnitude of batch effects in the original experiments varies from  $\eta = 0$  (i.e., no batch effect),  $\eta = 0.6$ , and to  $\eta = 1.0$ . Panel **d** displays the relationship between the sensitivity (percentage of contaminated experiments detected by the replicability analysis) and the experimented batch effect magnitude.

To demonstrate the non-informativeness principle, we simulate additional datasets with  $\eta = 1$  but increased residual errors ( $\sigma = 5, 10$ ) in replication experiments. Note that  $\sigma_{\text{rep}}$  is proportional to  $\sigma$ . The resulting empirical distributions of prior-predictive replication  $p$ -values display a clear

trend of reduced skewness (Figure 2), suggesting decreased sensitivity.



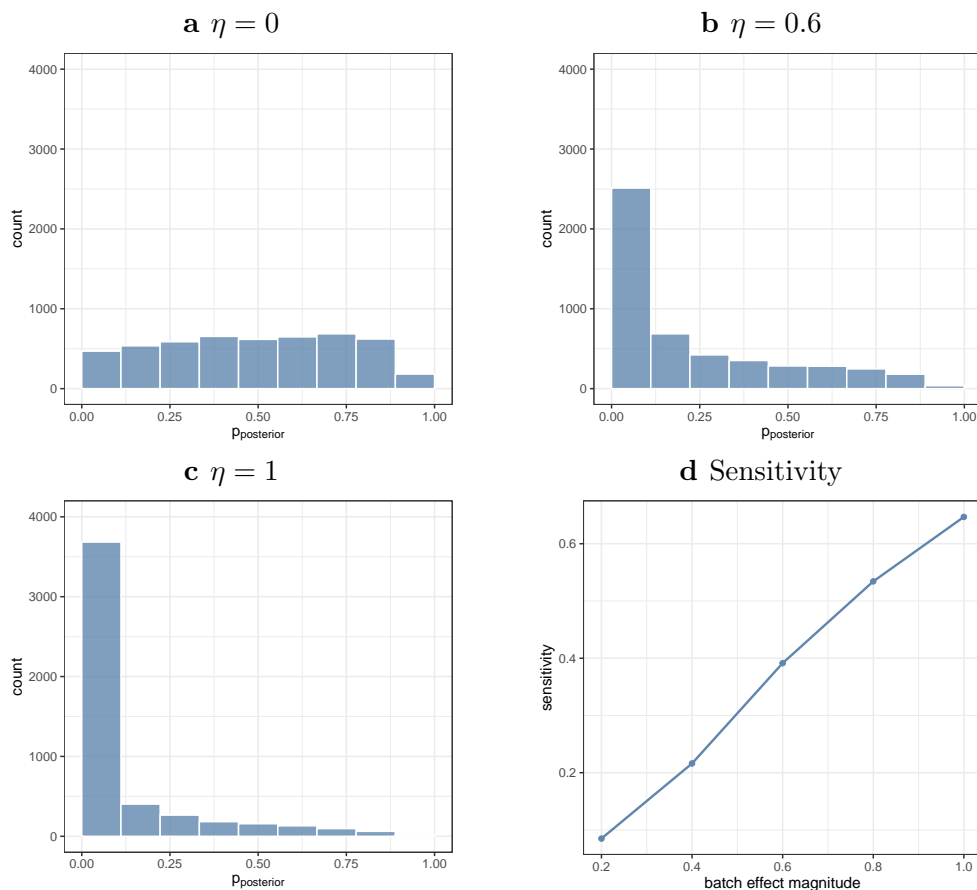
**Figure 2: Impacts of noisy replications on prior-PRPs** The original experiments are simulated with batch effect magnitude  $\eta = 1$ . The data from replication experiments are not batch contaminated but generated with increasing levels of residual errors, i.e.,  $\sigma = 5$  (Panel **a**) and  $\sigma = 10$  (Panel **b**). Note that  $\sigma_{\text{rep}}$  is proportional to  $\sigma$ , and the baseline comparison is the panel **c** of Figure 1, where  $\eta = 1$  and  $\sigma = 1$ .

### 5.1.2 Batch Effect Detection in Exchangeable-group Scenario

We adopt a similar scheme to simulate batch effects and generate observed summary statistics for the exchangeable-group scenario. For each dataset, we consider 5 experiments with 2 out of 5 experiments batch contaminated. Additionally, we introduce a low-level heterogeneity for experiment-specific association effects, such that the sign-consistency probability  $\approx 0.96$  across experiments.

We compute a posterior-predictive replication  $p$ -value for each simulated dataset using the default Q test quantity and assuming the default mixture reference model. Same as in the two-group scenario, increased sensitivity in detecting batch effects is observed as the magnitude of the batch effects ( $\eta$ ) increases (Figure 3). Although the distributional properties of posterior-PRPs

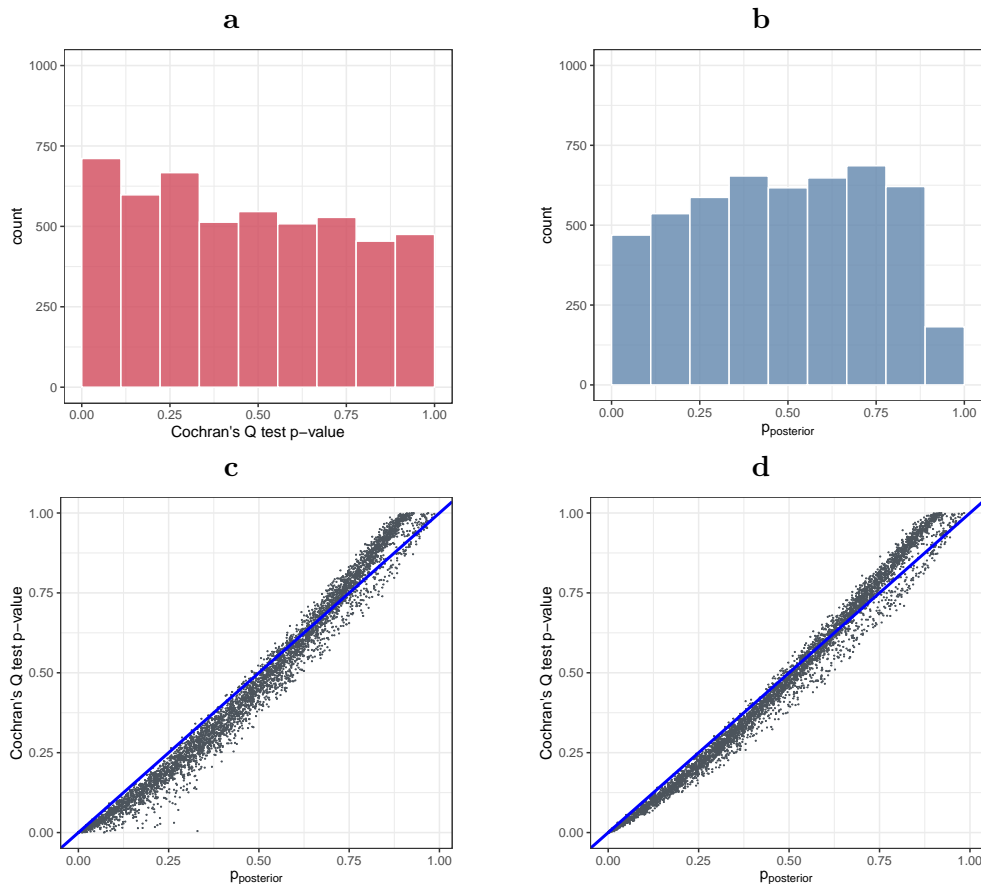
are not completely understood under the reference model, their empirical distributions under the simulation setting  $\eta \neq 0$  all show clear positive skewness, indicating excessive small  $p$ -values.



**Figure 3: Batch effect detection in exchangeable-group scenario** In each dataset, we simulate five experiments and two are contaminated by batch effects with magnitude  $\eta = 0$  (Panel **a**, no batch effect),  $\eta = 0.6$  (Panel **b**), and  $\eta = 1$  (Panel **c**). Panel **d** shows increased sensitivity of the posterior PRPs with respect to the increased batch effect magnitude.

We also take the opportunity to empirically compare the Bayesian  $p$ -values with its frequentist counterpart, the  $p$ -values derived from Cochran’s  $Q$  statistic. Our focus is on their behaviors when  $\eta = 0$ , which relates to specificity. Cochran’s  $Q$  statistics are computed assuming a fixed-effect meta-analysis model, under which they follow a pivotal  $\chi^2$  distribution with  $m - 1$  degrees of freedom. Figure 4 shows that the Bayesian and the frequentist  $p$ -values are largely aligned. However, the frequentist  $p$ -values are over-sensitive to the intrinsic heterogeneity within the

data, likely because of its stringent fixed-effect assumption. To further confirm this point, we re-compute the posterior-PRPs by modifying the reference model to include only the fixed effect assumption (i.e.,  $\gamma = 0$ ). The resulting Bayesian  $p$ -values become much more concordant to the frequentist  $p$ -values, especially for extremely small values (e.g.,  $p < 0.05$ ).



**Figure 4: Comparison of posterior-PRPs and classic  $p$ -values for testing heterogeneity under the reference reproducible model** Each dataset consists of five experiments without batch effect contamination ( $\eta = 0$ ) but is generated with an intrinsic low level of heterogeneity. Panels **a** and **b** compare the histograms of classic  $p$ -values, obtained from Cochran’s test of heterogeneity, and the proposed posterior-PRPs using the  $Q$  quantity. Cochran’s  $Q$  test assumes a fixed-effect model and is shown sensitive to the low level of heterogeneity. The posterior-checking procedure, based on the default reference model, tolerates such a level of heterogeneity acceptable. Panels **c** and **d** compare the two types of  $p$ -values for each simulated dataset. Overall, they show good agreement. In Panel **c**, the default reference model is used to generate posterior-PRPs. In Panel **d**, the reference model is modified to allow only the fixed-effect, and the results indicate increased agreement with the classic  $p$ -values, especially in the range of small values.

## 5.2 Simulations of Publication Bias

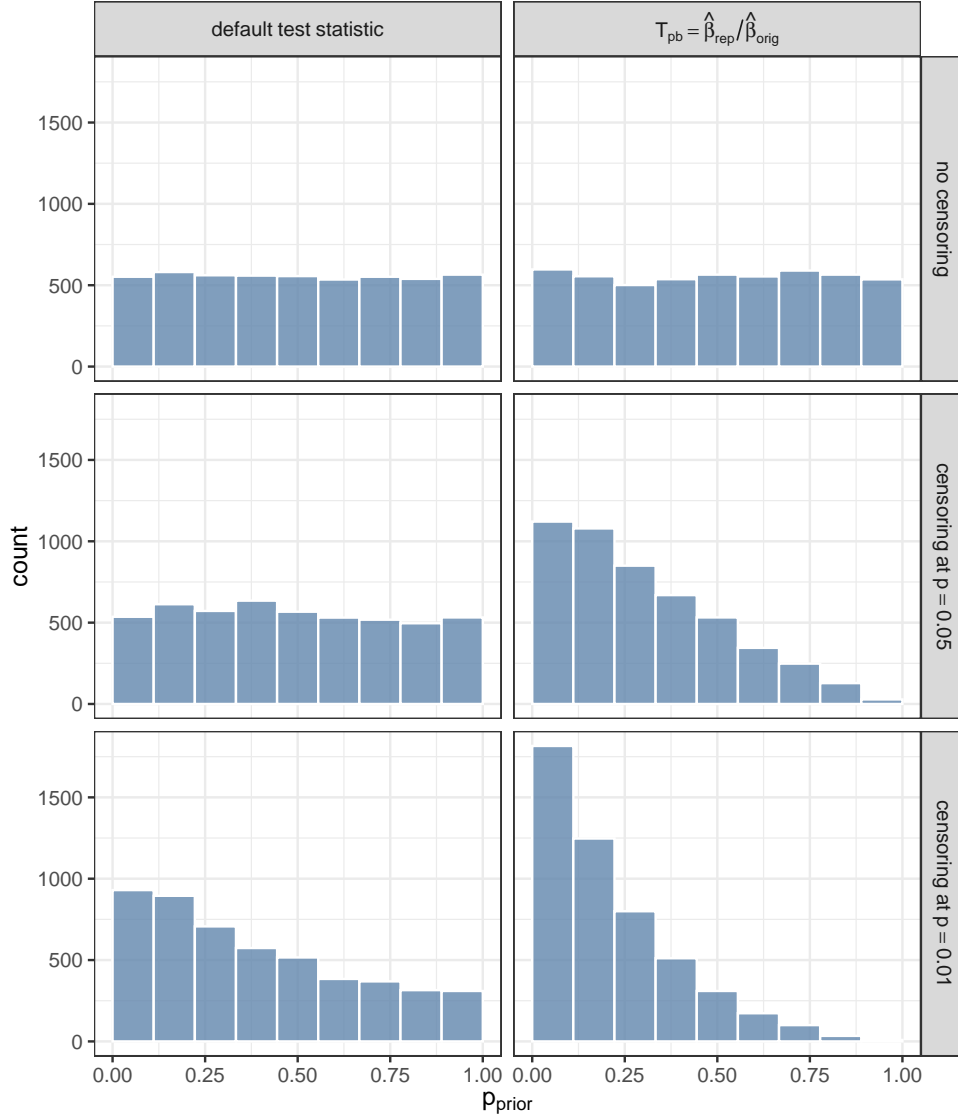
Our general simulation scheme considers a case-control study of a binary outcome. To introduce selection bias, we impose censoring mechanisms into otherwise normal data generative procedures. In all cases, following [24], the censoring schemes utilize the  $p$ -values derived from NHSTs.

### 5.2.1 Publication Bias Detection in Two-group Scenario

In the two-group scenario, we simulate two experiments per dataset. Both the original and replication experiments use a balanced case-control design with 200 samples. The true association effects are fixed at odds ratio,  $2/3$ , in both experiments. The outcomes are generated by a binomial sampling procedure, and the summary statistics (in the form of estimated log-odds ratios and their corresponding standard errors) are obtained by fitting standard logistic regression models. The censoring scheme only applies to the original experiment. Specifically, a simulated original experiment data is retained only if its association testing  $p$ -value is smaller than the pre-defined threshold,  $p_t$ . We experiment with  $p_t = 0.01, 0.05$ , and 1 (i.e., no censoring). For each threshold value, 5,000 datasets are generated.

For each simulated dataset, we compute both the default two-sided prior-PRPs and the one-sided values specifically designed for detecting publication bias (i.e., using  $T_{\text{pb}}$ ). The results are summarized in Figure 5.





**Figure 5: Publication bias detection in two-group scenario** The original experiment data are simulated via a censoring mechanism mimicking publication bias. The replication experiment data are generated without censoring. The left column represents the histograms of the default two-sided prior-PRPs, and the right column represents histograms of the one-sided prior-PRPs designed to detect publication bias. The different rows represent different censoring strengths. Both types of prior-PRPs are well-behaved when publication bias is absent. However, the one-sided  $p$ -values are clearly more sensitive in detecting publication bias.

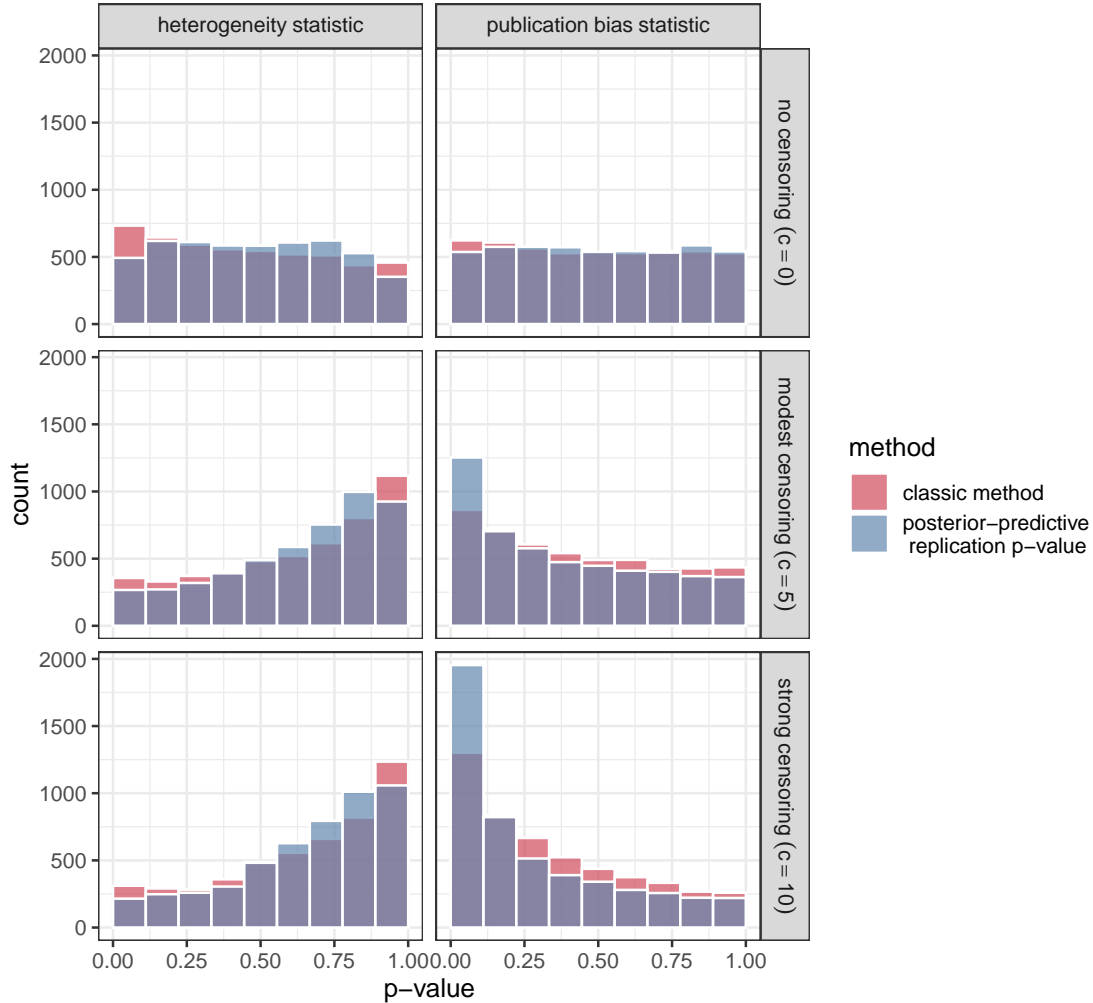
Both types of  $p$ -values are seemingly calibrated when there is no selection bias (i.e.,  $p_t = 1$ ), and both show positively skewed empirical distributions when the censoring mechanism is in place. Nevertheless, the special-purpose replication  $p$ -values exhibit superior sensitivity over

their general-purpose counterpart in the presence of publication bias.

### 5.2.2 Publication Bias Detection in Exchangeable-group Scenario

For the exchangeable-group scenario, our simulations mimic practical settings commonly encountered in meta-analysis and systematic review. We consider a dataset consisting of 10 experiments with various sample sizes (5 with 200 samples, 3 with 500 samples, and 2 with 1,000 samples) and balanced case-control designs. The true association effects are centered at the odds ratio = 2/3 with a low level of heterogeneity (the sign consistency probability = 0.96). Following [24], we apply a “soft” censoring scheme to each participating experiment. Specifically, the simulated data for an experiment is retained with probability,  $\exp(-cp^{3/2})$ , where  $p$  represents the association testing  $p$ -value and parameter  $c$  characterizes the censoring strength. We examine strong, modest, and no censoring effects by setting  $c = 10, 5$ , and  $0$ , respectively. For each  $c$  value, 5,000 datasets are generated. Given the unequal sample sizes in the 10 participating experiments, this scheme preferentially censors data from small studies and introduces selection bias when  $c \neq 0$ .

We analyze each simulated dataset by computing the posterior-predictive replication  $p$ -values derived from both the  $Q$  quantity and Egger regression statistic. The results are summarized in Figure 6.



**Figure 6: Publication bias detection in exchangeable-group scenario** The right column shows the histograms of posterior-PRPs obtained using Egger’s test quantity, overlaid with the histograms of  $p$ -values from classic Egger regressions. The left column represents the histograms of posterior-PRPs derived from the  $Q$  quantity, overlaid with the histograms of  $p$ -values from Cochran’s  $Q$  tests. The strength of the censoring mechanism increases from the top row to the bottom row. While both heterogeneity tests indicate a lack of heterogeneity, the posterior-PRPs designed for detecting publication bias show good sensitivity compared to the classic approach.

When selection bias is absent ( $c = 0$ ), there are no excessive small Bayesian  $p$ -values for either type of test statistic. The average  $p$ -values from both test statistics over 5,000 datasets are close to  $1/2$ , as expected by the theory. In the presence of publication bias, the posterior-predictive replication  $p$ -values based on Egger test quantity show increased sensitivity as the censoring strength ( $c$ ) increases. Interestingly, the  $p$ -value distributions based on the  $Q$  quantity

indicate reduced heterogeneity (comparing to the reference model specification), and the  $p$ -values derived from Cochran’s Q test show similar patterns. This observation demonstrates that a single diagnosis may not address all issues in replicability. Hence, the assessment of replicability should be multifaceted. Comparing to the existing frequentist model criticism approaches in detecting funnel plot asymmetry, the proposed approach (with Egger statistic) shows slightly better (but overall similar) sensitivity as the classic Egger regression test. In this particular simulation setting, it also outperforms the meta-analytic regression test implemented in the R package `metafor`, which similarly tests the linear trend between  $\hat{\beta}$  and  $\sqrt{\sigma^2 + \phi^2}$  (by plugging in a point estimate of  $\phi^2$ ).

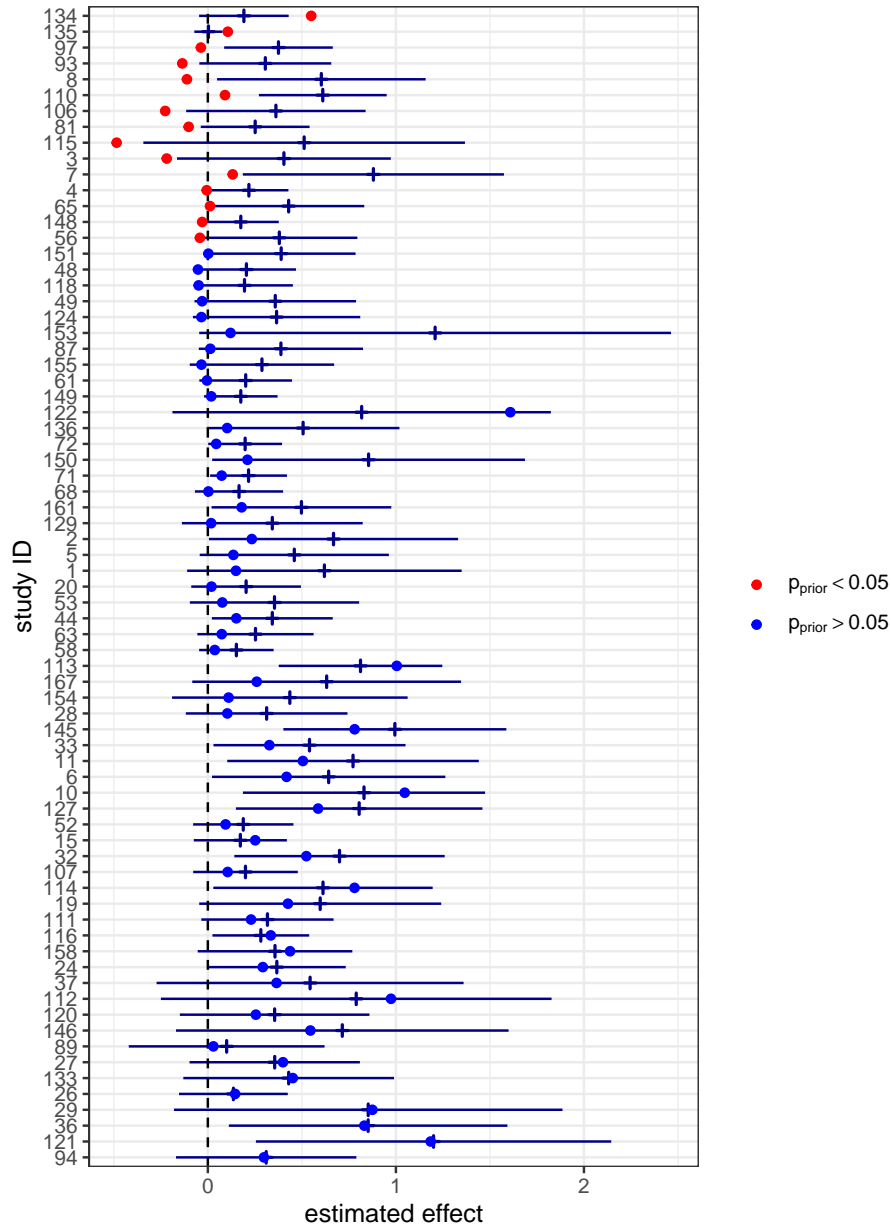
## 6 Real Data Applications

### 6.1 Re-analysis of RP:P Data

We apply the prior-predictive  $p$ -values to re-analyze the RP:P data, whose design falls in the category of the two-group scenario. The RP:P project attempts to replicate 100 psychology studies published in three top psychology journals during 2008. The replication experiment of each study is designed to match the design, sample size, and analysis methods in the corresponding original experiment. Following [14, 33, 15], we focus our re-analysis on a subset of 73 studies, for which the corresponding effect sizes and standard errors are computed via Fisher’s  $z$ -transformation.

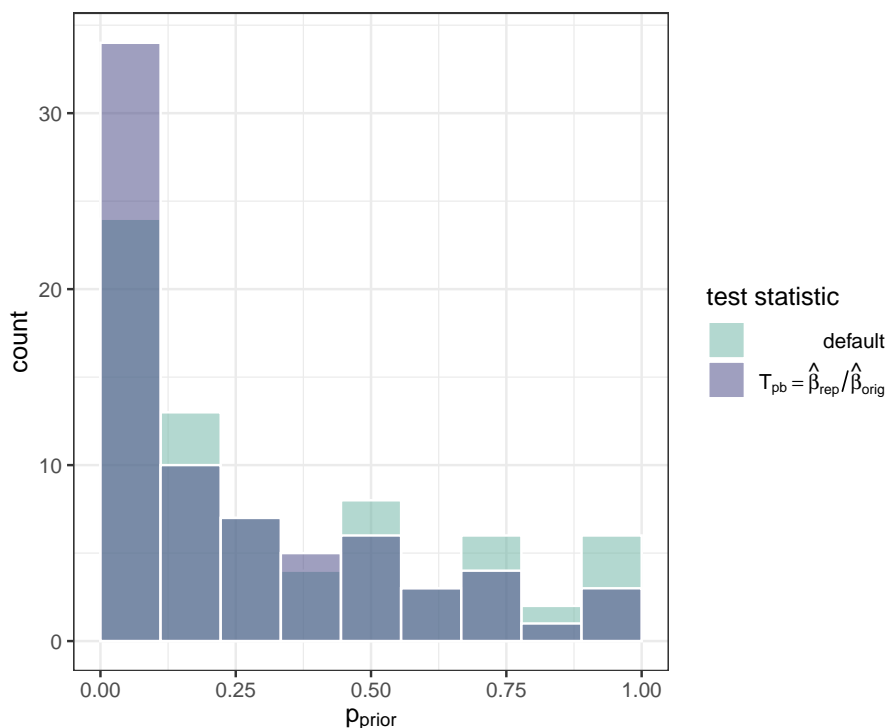
We apply the default reference model and compute the two-sided prior-predictive  $p$ -values for all 73 studies. The result is summarized in Figure 7, where we highlight the consistency between the prior-predictive  $p$ -values and the Bayesian predictive intervals indicated by Proposition 2. 15 out of 73 studies are flagged at the 5% significance level. In comparison, Patil *et al.* identify 22 studies by constructing 95% predictive intervals assuming a fixed-effect meta-analysis model. We are able to reproduce their result by resetting our default mixture reference model to a more

restrictive special case (i.e.,  $\omega^2 \in (10, \infty)$  and  $\gamma = 0$ ).



**Figure 7: Analysis of RP:P data by the prior-checking procedure** The plot summarizes the analysis result of all 73 studies in the RP:P dataset. The original estimate of effect is represented by a vertical tick centered at the corresponding 95% predictive interval (horizontal line) for each study. The filled dots are estimates from the replication studies. The red dots are from the studies with the default two-sided prior-PRP  $< 0.05$  and do not fall in the corresponding intervals. The studies are sorted by the increasing prior-PRP values. For most participating studies showing small prior-PRPs, the estimated effects tend to shrink toward 0, compared to the original estimates.

Figure 7 also reveals a striking pattern in the RP:P data: the effect estimates in the replications are predominantly shrunk towards 0 in flagged studies. This observation prompts computing the one-sided prior-predictive  $p$ -values based on (8) for detecting publication bias. This new analysis flags 22 studies at the 5% significance level, with 13 overlapping with those previously flagged by the two-sided prior-PRPs. This result suggests that publication bias can be responsible for a large proportion of irreproducible results in psychology research, a conclusion similarly drawn by [15].



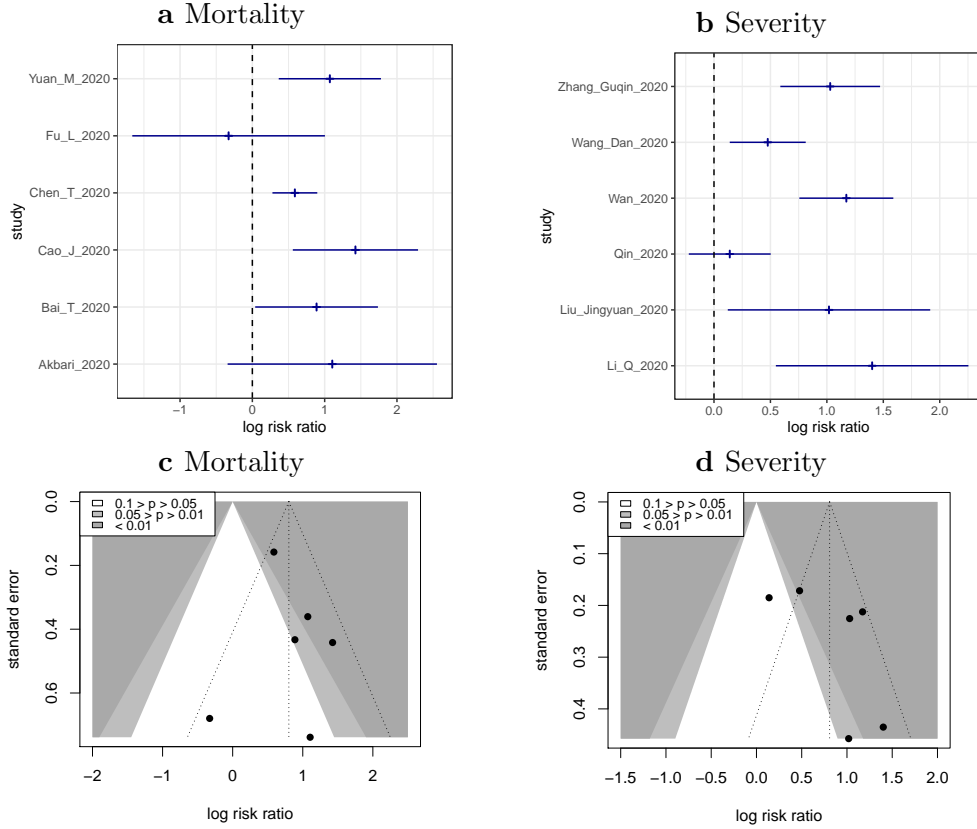
**Figure 8: Histograms of two types of prior-PRPs in assessing replicability of the RP:P data** The comparison is between general-purpose and default two-sided prior-PRPs (teal) and the one-sided prior-PRPs specifically designed for detecting publication bias (purple). The two histograms are overlaid, and both show clear positive-skewness. The prior-PRPs based on  $T_{pb}$  show more excessive small values and increased skewness.

Figure 8 shows the histograms of two types of prior-predictive replication  $p$ -values in analyzing RP:P data. Both empirical distributions show *severe* positive-skewness, indicating a worrisome trend in psychology research. Such a finding of a global pattern highlights the unique value of

the systematic efforts in simultaneously investigating multiple studies by the design of RP:P.

## **6.2 Systematic Review of Impact of Cardiovascular Diseases on Mortality and Severity of COVID-19**

Coronavirus disease 2019 (COVID-19), caused by the infection of the SARS-Cov-2 virus, has been linked to a worse prognosis in patients with pre-existing cardiovascular diseases. Pranta *et al.* [50] recently perform a systematic review on association evidence between pre-existing cardiovascular conditions and COVID-19 outcomes. Six studies investigating the mortality and another six studies investigating the severity of COVID-19 are included in the review. All 12 studies are observational and retrospective. They also have diverse sample sizes, ranging from 24 to 441 (median = 156). Estimates of log-risk ratios and the corresponding standard errors are extracted from each study and are shown in Figure 9.



**Figure 9: Forest and funnel plots of the COVID-19 datasets** Panels **a** and **b** show the forest plots for the estimated effects in the six different experiments on the cardiovascular disease impact the COVID-19 mortality, and the COVID-19 severe case rate, respectively. Panel **c** and **d** show their contour-enhanced funnel plots.

We compute the posterior-predictive replication  $p$ -values using  $Q$  and Egger statistics for both datasets and compare the results to the available traditional analyses described in [50]. The results are summarized in Table 1.

For heterogeneity assessment, Cochran’s tests and the Bayesian  $p$ -values are qualitatively similar in both mortality and severity studies. As expected, the Bayesian  $p$ -values are more conservative by allowing reasonable heterogeneity in the reference model. In all test cases, only the severity study is flagged by the posterior-predictive replication  $p$ -value based on  $Q$  quantity at 5% significance level (Bayesian  $p$ -value = 0.010, Cochran’s  $\chi^2$  test  $p$ -value = 0.001). Upon close inspection, the larger-than-expected heterogeneity is seemingly driven by a single study with



COVID-19 Outcome	Test of Heterogeneity		Test of Publication Bias	
	classic	posterior-PRP	classic	posterior-PRP
Mortality	0.221	0.278	0.429	0.598
Severity	0.001	0.010	0.010	0.240

**Table 1: Systematic review of pre-existing cardiovascular disease on COVID-19 outcomes** The table shows the comparisons between the classic approaches and the proposed posterior-predictive replication  $p$ -values in assessing excessive heterogeneity and potential publication bias. The classic  $p$ -values are obtained from Cochran’s  $Q$  test for testing heterogeneity, and the posterior-PRPs are computed using the  $Q$  quantity. For detecting publication bias, the classic  $p$ -values are obtained from Egger regression test implemented in R package `metafor`, and the posterior-PRPs are computed using Egger statistic with sampled hyperparameters.

a large sample size (Qin 2020) in the analysis, which under-estimates the effect size than the remaining studies. (The posterior-PRP based on the remaining five studies is 0.157.) The only disagreement between the classic and the proposed assessment approaches appears in detecting publication bias for the severity outcome. The Bayesian PRP based on Egger regression quantity does not deem the data to provide strong enough evidence against reproducibility. Figure 9 does not seem to show a clear linear trend between the estimates and the corresponding standard errors upon close examination. As a rule of thumb, it is commonly recommended that the classic Egger regression test should be applied to  $\geq 10$  studies, partially because of the calibration of the null distribution and the considerations for lack of sensitivity [51]. (Pranata *et al.* do not formally list the classic Egger regression result as numerical evidence but briefly comment on potential funnel plot asymmetry.) In this particular case, the Bayesian  $p$ -value solution shows some advantages without requiring a pre-defined null distribution. However, the second concern remains, and we caution the interpretation that the publication bias is *absent*. We should conclude that the potential publication bias may not be strong enough to be flagged by the given experiments using the proposed method.

## 7 Discussion

Applying model criticism approaches in assessing replicability is not new. The uses of traditional Cochran’s  $Q$  test and Egger regression in exchangeable-group scenarios all fall in this statistical framework. A common limitation of these existing (frequentist) methods is that the reference (i.e., the null) model is often too restrictive to be realistic (e.g., the fixed-effect assumption). We address this limitation by adopting more general Bayesian hierarchical models while following the same model criticism principles. The existing toolkit built for Bayesian model criticism has its unique, pragmatic flexibility to adopt the extended reference models and make inferences in the presence of nuisance parameters. More importantly, it is versatile in dealing with different application scenarios arising from replicability assessment.

A key element in replication assessment is to define replication success. We argue that this definition should be ultimately context-dependent. Nevertheless, we believe definitions based on the DC criterion can be useful if prior established quantitative standards for reproducibility lack in some scientific domains. In other words, replication definitions based on the DC criterion can serve as a reasonable starting point, and other domain-specific standards should gradually replace it with the accumulation of new data and knowledge. Additionally, our presented model is not the only way to set up the reference model for replicability assessment. Our choice to present this particular version of the reference model is due to its connections to the commonly-used random-effect meta-analysis model. Thus, it is more convenient to compare the proposed Bayesian model criticism approaches and the traditional approaches assuming similar models. The curved exponential family normal (CEF<sub>N</sub>) model [42, 19] is another parametric model capable of faithfully implementing the DC criterion.

A well-known drawback of model criticism approaches is that their conclusions are often misinterpreted and misused. The binary outcome from a model criticism procedure is either flagging/rejecting a reference model or not. If the corresponding reference model is flagged, the

correct interpretation is via Fisher’s disjunction. A common mistake is to misinterpret a lack of evidence for flagging as evidence for supporting the reference model, and such a mistake is not unique for replication assessment. This is why we wish to highlight the importance of the non-informativeness principle and its implications in our proposed procedures: it serves as a good reminder to avoid misinterpreting the proposed replication  $p$ -values in assessing replicability.

Model criticism is not the only statistical inference framework capable of assessing replicability. For example, inference procedures based on model comparison principles are also applicable. A model comparison approach aims to classify the observed data into mutually exclusive latent data-generative models, e.g., irreproducible vs. reproducible models. Our previous work [19] and the influential work by Li *et al.* [32] are both Bayesian model comparison procedures for replicability assessment in high-throughput experiments. Nevertheless, we note that prior quantification on data-generative models is critical for the success of model comparison procedures. In high-throughput settings, this information can be sufficiently “learned” from data by partial pooling, whereas in other classic settings of replicability assessment, justifying particular prior choices can be challenging. In comparison, the model criticism strategy presents some attractive and practical simplicity: it does not require a prior specification on generative models, nor does it requires a parametric specification of an irreproducible model. For many applications discussed in this paper, this operational simplicity is highly desired. In summary, we believe both model criticism and model comparison approaches can be useful for different application settings; they can be complementary in some cases. It is also important to realize that, while they represent different inference strategies, they share the same guiding principles for replication assessment (or even the underlying probabilistic generative models).

## References

- [1] Claerbout, J. F. & Karrenbach, M. Electronic documents give reproducible research a new meaning. In *SEG Technical Program Expanded Abstracts 1992*, 601–604 (Society of Exploration Geophysicists, 1992).
- [2] Peng, R. D. Reproducible research in computational science. *Science* **334**, 1226–1227 (2011).
- [3] Crook, S. M., Davison, A. P. & Plesser, H. E. Learning from the past: approaches for reproducibility in computational neuroscience. In *20 Years of Computational Neuroscience*, 73–102 (Springer, 2013).
- [4] Heller, R., Bogomolov, M. & Benjamini, Y. Deciding whether follow-up studies have replicated findings in a preliminary large-scale omics study. *Proceedings of the National Academy of Sciences* **111**, 16262–16267 (2014).
- [5] Goodman, S. N., Fanelli, D. & Ioannidis, J. P. What does research reproducibility mean? *Science translational medicine* **8**, 341ps12–341ps12 (2016).
- [6] Nichols, T. E. *et al.* Best practices in data analysis and sharing in neuroimaging using mri. *Nature neuroscience* **20**, 299–303 (2017).
- [7] Haihe-Kains, B. *et al.* Transparency and reproducibility in artificial intelligence. *Nature* **586**, E14–E16 (2020).
- [8] Schmidt, S. Shall we really do it again? the powerful concept of replication is neglected in the social sciences. *Review of general psychology* **13**, 90–100 (2009).
- [9] Plesser, H. E. Reproducibility vs. Replicability: A brief history of a confused terminology. *Frontiers in Neuroinformatics* **11** (2018).

- [10] Mikowski, M., Hensel, W. M. & Hohol, M. Replicability or reproducibility? on the replication crisis in computational neuroscience and sharing only relevant detail. *Journal of Computational Neuroscience* **45**, 163–172 (2018).
- [11] Open Science Collaboration. Estimating the reproducibility of psychological science. *Science* **349** (2015).
- [12] Goodman, S. N. A comment on replication, p-values and evidence. *Statistics in medicine* **11**, 875–879 (1992).
- [13] Gilbert, D. T., King, G., Pettigrew, S. & Wilson, T. D. Comment on “estimating the reproducibility of psychological science”. *Science* **351**, 1037–1037 (2016).
- [14] Patil, P., Peng, R. D. & Leek, J. T. What should researchers expect when they replicate studies? a statistical view of replicability in psychological science. *Perspectives on Psychological Science* **11**, 539–544 (2016).
- [15] Hung, K., Fithian, W. *et al.* Statistical methods for replicability assessment. *Annals of Applied Statistics* **14**, 1063–1087 (2020).
- [16] Gibson, E. W. The Role of p-Values in Judging the Strength of Evidence and Realistic Replication Expectations. *Statistics in Biopharmaceutical Research* (2020).
- [17] Karp, N. A. Reproducible preclinical research—is embracing variability the answer? *PLoS Biology* **16**, e2005413 (2018).
- [18] Kochunov, P., Thompson, P. M. & Hong, L. E. Toward high reproducibility and accountable heterogeneity in schizophrenia research. *JAMA psychiatry* **76**, 680–681 (2019).
- [19] Zhao, Y., Sampson, M. G. & Wen, X. Quantify and control reproducibility in high-throughput experiments. *Nature Methods* **17**, 1207–1213 (2020).

- [20] He, X. & Lin, X. Challenges and opportunities in statistics and data science: Ten research areas. *2.3* (2020).
- [21] Higgins, J. P. & Thompson, S. G. Quantifying heterogeneity in a meta-analysis. *Statistics in medicine* **21**, 1539–1558 (2002).
- [22] Magosi, L. E., Goel, A., Hopewell, J. C. & Farrall, M. Identifying systematic heterogeneity patterns in genetic association meta-analysis studies. *PLOS Genetics* **13**, e1006755 (2017).
- [23] Jaljuli, I., Benjamini, Y., Shenhav, L., Panagiotou, O. & Heller, R. Quantifying replicability and consistency in systematic reviews. *arXiv* 1907.06856 (2019).
- [24] Macaskill, P., Walter, S. D. & Irwig, L. A comparison of methods to detect publication bias in meta-analysis. *Statistics in medicine* **20**, 641–654 (2001).
- [25] Zöllner, S. & Pritchard, J. K. Overcoming the winner’s curse: estimating penetrance parameters from case-control data. *The American Journal of Human Genetics* **80**, 605–615 (2007).
- [26] Peters, J. L. *et al.* Assessing publication bias in meta-analyses in the presence of between-study heterogeneity. *Journal of the Royal Statistical Society: Series A (Statistics in Society)* **173**, 575–591 (2010).
- [27] Lin, L. & Chu, H. Quantifying publication bias in meta-analysis. *Biometrics* **74**, 785–794 (2018).
- [28] Palmer, C. & Pe’er, I. Statistical correction of the winner’s curse explains replication variability in quantitative trait genome-wide association studies. *PLoS genetics* **13**, e1006916 (2017).
- [29] Johnson, W. E., Li, C. & Rabinovic, A. Adjusting batch effects in microarray expression data using empirical bayes methods. *Biostatistics* **8**, 118–127 (2007).

- [30] Leek, J. T. *et al.* Tackling the widespread and critical impact of batch effects in high-throughput data. *Nature Reviews Genetics* **11**, 733 (2010).
- [31] Goh, W. W. B., Wang, W. & Wong, L. Why batch effects matter in omics data, and how to avoid them. *Trends in biotechnology* **35**, 498–507 (2017).
- [32] Li, Q., Brown, J. B., Huang, H., Bickel, P. J. *et al.* Measuring reproducibility of high-throughput experiments. *The annals of applied statistics* **5**, 1752–1779 (2011).
- [33] Pawel, S. & Held, L. Probabilistic forecasting of replication studies. *PLoS ONE* **15** (2020).
- [34] Box, G. E. Sampling and bayes' inference in scientific modelling and robustness. *Journal of the Royal Statistical Society: Series A (General)* **143**, 383–404 (1980).
- [35] Rubin, D. B. Bayesianly justifiable and relevant frequency calculations for the applies statistician. *The Annals of Statistics* 1151–1172 (1984).
- [36] Meng, X.-L. Posterior predictive p-values. *The Annals of Statistics* **22**, 1142–1160 (1994).
- [37] Gelman, A., Meng, X.-L. & Stern, H. Posterior predictive assessment of model fitness via realized discrepancies. *Statistica Sinica* **6**, 733–807 (1996).
- [38] Gelman, A. Comment: Fuzzy and bayesian p-values and u-values. *Statistical Science* **20** (2005).
- [39] Gelman, A. *et al.* Two simple examples for understanding posterior p-values whose distributions are far from uniform. *Electronic Journal of Statistics* **7**, 2595–2602 (2013).
- [40] Gelman, A. & Tuerlinckx, F. Type s error rates for classical and bayesian single and multiple comparison procedures. *Computational Statistics* **15**, 373–390 (2000).
- [41] Owen, A. B. *et al.* Karl pearson's meta-analysis revisited. *The annals of statistics* **37**, 3867–3892 (2009).

- [42] Wen, X. & Stephens, M. Bayesian methods for genetic association analysis with heterogeneous subgroups: from meta-analyses to gene-environment interactions. *The annals of applied statistics* **8**, 176 (2014).
- [43] Stephens, M. False discovery rates: a new deal. *Biostatistics* **18**, 275–294 (2016).
- [44] Camerer, C. F. *et al.* Evaluating replicability of laboratory experiments in economics. *Science* **351**, 1433–1436 (2016).
- [45] Fisher, R. A. *Statistical methods and scientific inference*. (Hafner Publishing Co., 1956).
- [46] Wasserstein, R. L. & Lazar, N. A. The asa statement on p-values: context, process, and purpose (2016).
- [47] Hubbard, R. & Bayarri, M. J. Confusion over measures of evidence (p's) versus errors ( $\alpha$ 's) in classical statistical testing. *The American Statistician* **57**, 171–178 (2003).
- [48] Thulin, M. Decision-theoretic justifications for bayesian hypothesis testing using credible sets. *Journal of Statistical Planning and Inference* **146**, 133–138 (2014).
- [49] Sterne, J. A. & Egger, M. Regression methods to detect publication and other bias in meta-analysis. *Publication bias in meta-analysis: Prevention, assessment and adjustments* 99–110 (2005).
- [50] Pranata, R., Huang, I., Lim, M. A., Wahjoepramono, E. J. & July, J. Impact of cerebrovascular and cardiovascular diseases on mortality and severity of covid-19 : systematic review, meta-analysis, and meta-regression. *Journal of Stroke and Cerebrovascular Diseases* (2020).
- [51] Higgins, J. P. *et al.* *Cochrane handbook for systematic reviews of interventions* (John Wiley & Sons, 2019).



# Appendix A Construction of Reference Reproducible Model

The main text has detailed the specification of set  $\Gamma$ , which defines various levels of tolerable heterogeneity for the reference model. Here we describe the construction of set  $\Omega$ , which specifies the prior effect size of  $\bar{\beta}$ . Intentionally, we set up the grid values of  $\omega$  to be compatible with the observed data, such that the potential poor fitting of the reference model can only be attributed to the incompatibility of the heterogeneity parameter,  $\gamma$ .

This general principle is applied to both prior and posterior checking procedures, but the implementations are slightly different. In both cases, instead of directly setting  $\Omega$ , we first specify a grid of values for  $\lambda^2 = (\phi^2 + \omega^2)$ , then compute  $\omega^2 = (1 - \gamma) \lambda^2$ .

## Two-group Scenario

In this scenario, we determine  $\Lambda$ , the set of  $\lambda^2$  values, by  $\hat{\beta}_{\text{orig}}$  and  $\sigma_{\text{orig}}$ . Let  $q_1, q_2, q_3$  denote the first, second, and the third quartiles from the  $\chi_1^2$  distribution. We regard  $\hat{\beta}_{\text{orig}}^2 + \sigma_{\text{orig}}^2$  as a proxy for an estimate of  $\lambda^2$  and reason that  $\frac{\hat{\beta}_{\text{orig}}^2 + \sigma_{\text{orig}}^2}{\lambda^2}$  should be a *plausible* draw from the  $\chi_1^2$  distribution. Thus, we construct

$$\Lambda = \left\{ \frac{\hat{\beta}_{\text{orig}}^2 + \sigma_{\text{orig}}^2}{q_3}, \frac{\hat{\beta}_{\text{orig}}^2 + \sigma_{\text{orig}}^2}{q_2}, \frac{\hat{\beta}_{\text{orig}}^2 + \sigma_{\text{orig}}^2}{q_1} \right\}$$

## Exchangeable-group Scenario

In this scenario, we first compute a fixed-effect estimate of  $\bar{\beta}$ , i.e.,

$$\hat{\bar{\beta}} = \frac{\sum_{i=1}^m w_i \hat{\beta}_i}{\sum_{i=1}^m w_i}, \quad \text{where } w_i = \frac{1}{\sigma_i^2},$$

and

$$\text{se}(\hat{\beta}) = \sqrt{\frac{1}{\sum_{i=1}^m w_i}}.$$

We construct set  $\Lambda$  by

$$\Lambda = \left\{ \left( \hat{\beta} - \text{se}(\hat{\beta}) \right)^2, \hat{\beta}^2, \left( \hat{\beta} + \text{se}(\hat{\beta}) \right)^2 \right\}.$$

## Appendix B Proof of Proposition 1

*Proof.* We introduce a shorthand notation,  $P \circ T(\mathbf{X}_{\text{orig}}, \mathbf{X}_{\text{rep}})$ , to represent prior-predictive replication  $p$ -value,  $\Pr(T(\mathbf{X}) \leq T(\mathbf{X}_{\text{rep}}) | \mathbf{X}_{\text{orig}}, \mathbf{X}_{\text{rep}}, \mathbf{M}_{\text{R}})$  and emphasize it as a function of random variables  $\mathbf{X}_{\text{orig}}$  and  $\mathbf{X}_{\text{rep}}$ .

It follows that

$$\begin{aligned} & \Pr [ P \circ T(\mathbf{X}_{\text{orig}}, \mathbf{X}_{\text{rep}}) \leq p \mid \mathbf{M}_{\text{R}} ] \\ &= \Pr [ \Pr ( P \circ T(\mathbf{X}_{\text{orig}}, \mathbf{X}_{\text{rep}}) \leq p \mid \mathbf{X}_{\text{orig}}, \mathbf{M}_{\text{R}} ) \mid \mathbf{M}_{\text{R}} ] \end{aligned}$$

Let  $F_{T | \mathbf{X}_{\text{orig}}, \mathbf{M}_{\text{R}}}$  denote the continuous cumulative distribution function of the predictive distribution,  $T(\mathbf{X}) | \mathbf{X}_{\text{orig}}, \mathbf{M}_{\text{R}}$ . Note that,

$$\begin{aligned} & \Pr ( P \circ T(\mathbf{X}_{\text{orig}}, \mathbf{X}_{\text{rep}}) \leq p \mid \mathbf{X}_{\text{orig}}, \mathbf{M}_{\text{R}} ) \\ &= \Pr \left( T(\mathbf{X}_{\text{rep}}) \leq F_{T | \mathbf{X}_{\text{orig}}, \mathbf{M}_{\text{R}}}^{-1}(p) \mid \mathbf{X}_{\text{orig}}, \mathbf{M}_{\text{R}} \right) \\ &= p \end{aligned}$$

Thus,

$$\Pr [ P \circ T(\mathbf{X}_{\text{orig}}, \mathbf{X}_{\text{rep}}) \leq p \mid M_{\text{R}} ] = p.$$

That is, the prior-predictive replication  $p$ -values are uniformly distributed under the resampling of  $\mathbf{X}_{\text{orig}}$  and  $\mathbf{X}_{\text{rep}}$ .

□

## Appendix C Proof of Proposition 2

*Proof.*

$$\begin{aligned} & T(\mathbf{X}_{\text{rep}}) \notin (q_{\alpha/2}(T(\mathbf{X}) \mid \mathbf{X}_{\text{orig}}, M_{\text{R}}), q_{(1-\alpha/2)}(T(\mathbf{X}) \mid \mathbf{X}_{\text{orig}}, M_{\text{R}})) \\ \iff & T(\mathbf{X}_{\text{rep}}) \leq q_{\alpha/2}(T(\mathbf{X}) \mid \mathbf{X}_{\text{orig}}, M_{\text{R}}) \quad \text{or} \quad T(\mathbf{X}_{\text{rep}}) \geq q_{(1-\alpha/2)}(T(\mathbf{X}) \mid \mathbf{X}_{\text{orig}}, M_{\text{R}}) \\ \iff & \min \left\{ \Pr(T(\mathbf{X}) \geq T(\mathbf{X}_{\text{rep}}) \mid \mathbf{X}_{\text{orig}}, \mathbf{X}_{\text{rep}}, M_{\text{R}}), \right. \\ & \left. \Pr(T(\mathbf{X}) \leq T(\mathbf{X}_{\text{rep}}) \mid \mathbf{X}_{\text{orig}}, \mathbf{X}_{\text{rep}}, M_{\text{R}}) \right\} \leq \alpha/2 \\ \iff & p_{\text{prior}} \leq \alpha \end{aligned}$$

□

## Appendix D Computation of Prior-predictive Replication $p$ -values

As described in the main text, we consider  $K$  different  $(\omega_k, \gamma_k)$  hyperparameter values with equal prior probabilities. Correspondingly,  $\phi_k^2 = \omega_k^2 \gamma / (1 - \gamma)$ .

Given  $(\omega_k, \phi_k)$ , the posterior distribution of  $\bar{\beta}$  is

$$\bar{\beta} | \hat{\beta}_{\text{orig}}, \omega_k, \phi_k \sim \text{N} \left( \left( \frac{1}{\sigma_{\text{orig}}^2 + \phi_k^2} + \frac{1}{\omega_k^2} \right)^{-1} \left( \frac{\hat{\beta}_{\text{orig}}}{\sigma_{\text{orig}}^2 + \phi_k^2} \right), \left( \frac{1}{\sigma_{\text{orig}}^2 + \phi_k^2} + \frac{1}{\omega_k^2} \right)^{-1} \right)$$

The desired posterior predictive distribution is a mixture normal, i.e.,

$$\hat{\beta} | \hat{\beta}_{\text{orig}}, \text{M}_R \sim \sum_k w_k \cdot \text{N} \left( \left( \frac{1}{\sigma_{\text{orig}}^2 + \phi_k^2} + \frac{1}{\omega_k^2} \right)^{-1} \left( \frac{\hat{\beta}_{\text{orig}}}{\sigma_{\text{orig}}^2 + \phi_k^2} \right), \left( \frac{1}{\sigma_{\text{orig}}^2 + \phi_k^2} + \frac{1}{\omega_k^2} \right)^{-1} + \phi_k^2 + \sigma_{\text{rep}}^2 \right),$$

where

$$w_k = \frac{\mathcal{N}(\hat{\beta}_{\text{orig}} | 0, \omega_k^2 + \phi_k^2 + \sigma_{\text{orig}}^2)}{\sum_{j=1}^K \mathcal{N}(\hat{\beta}_{\text{orig}} | 0, \omega_j^2 + \phi_j^2 + \sigma_{\text{orig}}^2)},$$

and  $\mathcal{N}(\hat{\beta}_{\text{orig}} | 0, \omega_k^2 + \phi_k^2 + \sigma_{\text{orig}}^2)$  denote a normal density function with mean 0 and variance  $\omega_k^2 + \phi_k^2 + \sigma_{\text{orig}}^2$  evaluated at  $\hat{\beta}_{\text{orig}}$ .

## Appendix E Computation of Posterior-predictive replication $p$ -values

Here we detail the sampling scheme described in Algorithm 1 for the reference reproducible model.

As described in the main text, we consider  $K$  different  $(\omega_k, \gamma_k)$  hyperparameter values with equal prior probabilities. Correspondingly,  $\phi_k^2 = \omega_k^2 \gamma / (1 - \gamma)$ . The marginal likelihood of observed

data for  $(\omega_k, \phi_k)$  is represented by  $\mathcal{N}(\hat{\boldsymbol{\beta}} \mid \mathbf{0}, \boldsymbol{\Sigma}_k)$ , where  $\boldsymbol{\Sigma}_k$  is an  $m \times m$  covariance matrix and

$$\boldsymbol{\Sigma}_k = \begin{bmatrix} \omega_k^2 + \phi_k^2 + \sigma_1^2 & \omega_k^2 & \dots & \omega_k^2 \\ \omega_k^2 & \omega_k^2 + \phi_k^2 + \sigma_2^2 & \dots & \omega_k^2 \\ \vdots & \vdots & \ddots & \vdots \\ \omega_k^2 & \omega_k^2 & \dots & \omega_k^2 + \phi_k^2 + \sigma_m^2 \end{bmatrix}$$

The key procedure in the sampling scheme is to draw from the posterior,  $P(\boldsymbol{\theta} \mid \hat{\boldsymbol{\beta}}, M_R)$ . Specifically,

1. Draw a single generative model  $(\omega_k, \phi_k)$  from the  $K$  candidates with probability

$$p_k = \frac{\mathcal{N}(\hat{\boldsymbol{\beta}} \mid \mathbf{0}, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \mathcal{N}(\hat{\boldsymbol{\beta}} \mid \mathbf{0}, \boldsymbol{\Sigma}_j)}$$

2. Conditional on  $(\omega_k, \phi_k)$  and  $\hat{\boldsymbol{\beta}}$ , draw  $\bar{\boldsymbol{\beta}}$  by

$$\bar{\boldsymbol{\beta}} \mid \hat{\boldsymbol{\beta}}, \omega_k, \phi_k \sim \mathcal{N} \left( \left( \frac{1}{\omega_k^2} + \sum_{j=1}^m \frac{1}{\sigma_j^2 + \phi_k^2} \right)^{-1} \left( \sum_{j=1}^m \frac{\hat{\beta}_j}{\sigma_j^2 + \phi_k^2} \right), \left( \frac{1}{\omega_k^2} + \sum_{j=1}^m \frac{1}{\sigma_j^2 + \phi_k^2} \right)^{-1} \right)$$

3. Draw  $\beta_j, j = 1, \dots, m$  by

$$\beta_j \mid \hat{\boldsymbol{\beta}}, \bar{\boldsymbol{\beta}}, \omega_k, \phi_k \sim \mathcal{N} \left( \left( \frac{1}{\phi_k^2} + \frac{1}{\sigma_j^2} \right)^{-1} \left( \frac{\bar{\beta}}{\phi_k^2} + \frac{\hat{\beta}_j}{\sigma_j^2} \right), \left( \frac{1}{\phi_k^2} + \frac{1}{\sigma_j^2} \right)^{-1} \right)$$

Finally, given all  $\beta_j$ 's, we re-sample the observed data by

$$\hat{\beta}'_j \mid \beta_j \sim \mathcal{N}(\beta_j, \sigma_j^2).$$