

1 Reliable inference: benefits of open raw data may
2 be universal in meta-analysis

3 Danijela Žanko¹, Sunčana Geček¹, Azra Tafro², Livia Puljak³, Antica Culina^{1,*}

4 ¹Laboratory for Informatics and Environmental Modelling, Division for Ma-
5 rine and Environmental Research, Ruđer Bošković Institute, Bijenička Cesta
6 54, Zagreb 10000, Croatia

7 ²Faculty of Forestry and Wood Technology, University of Zagreb, Svetošimun-
8 ska cesta 23, Zagreb 10000, Croatia

9 ³Center for Evidence-Based Medicine, Catholic University of Croatia, Ilica
10 244, Zagreb 10000, Croatia

11 *Corresponding author: Antica Culina

12 **Keywords:** meta-analysis, publication bias, raw data, simulation, effect size
13 estimation, open science

14 **Abstract**

15 While the benefits of open data are often discussed, they are rarely
16 quantified. Here, we provide the first evidence of the potential gains
17 from using raw data for evidence synthesis and introduce a tool that
18 helps researchers determine when this approach is most beneficial.
19 Classical meta-analysis (CMA) relies on published results, making it
20 vulnerable to publication bias and p-hacking. We developed a simula-
21 tion framework comparing CMA with raw data meta-analysis (RDMA)
22 across varying true effect sizes, heterogeneity, and bias levels. Us-
23 ing ecology as a case study, we demonstrate that RDMA outperforms
24 CMA in most scenarios. When true effects are small and bias is severe,
25 RDMA reduces relative mean absolute error by 56–76%. Under moder-
26 ate bias, reductions reach 31–62%. For medium true effects, reductions
27 were 50–71% and 3–38%, respectively. RDMA maintained reliable confi-
28 dence interval coverage (>90%) across all scenarios, whereas CMA
29 failed to do so. Crucially, RDMA’s errors reflect natural sampling
30 variation, while CMA’s reflect systematic bias that persists regardless
31 of the sample size. We provide a decision tool for meta-analysts across
32 disciplines to calculate RDMA’s benefits. Our results offer the first
33 quantitative evidence that open data improves meta-analytic accuracy,
34 strengthening the case for open science.

35

1 Introduction

36

37

38

39

40

41

42

43

44

45

46

47

48

49

Meta-analysis is the main method for quantitatively synthesising existing evidence and is often used to inform future research, interventions, and policies [1–3]. Importantly, the statistical power of individual studies is often inadequate to detect true effects and can thus lead to imprecise effect estimates that limit the detection of meaningful effects [4–8]. Meta-analysis aggregates results from independent studies, increasing statistical power and precision [5]. Critically, this approach assumes that the included studies represent an unbiased sample of all research conducted on a given research question and that the included effects capture the true effects well. However, this assumption is often violated because selective publication and p-hacking can lead to overestimated effects and misguided decisions when meta-analysis relies on published results alone. Such selective publication and p-hacking are common in different scientific disciplines [9–13].

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

Publication bias [14], the preferential publication of statistically significant results, systematically distorts the published literature [12, 15–18]. Because significant results are more likely to be published, researchers might manipulate their analyses to achieve statistical significance (p-hacking) [19]. Data across disciplines demonstrate that these issues could be very widespread. For example, 94% of p-values reported in display items of the top multidisciplinary journals are significant [20], 96% of biomedical abstracts and full-text articles reporting p-values contain at least one significant result [21], while in ecology 98.9% of published articles report at least one significant result (computed from [22, 23]). Further, surveys of researchers show that p-hacking and other questionable research practices may be common: in psychology, self-admission rates for practices such as selective reporting and optional stopping range from 16–63% [24], with comparable rates confirmed across countries [25, 26] and in ecology and evolution [27], engineering [28], and across Dutch academia more broadly [29]. These findings convincingly indicate that researchers might selectively find and report significant effects, while null findings and non-significant results remain inaccessible, leading to a substantial loss of scientific knowledge [30].

70

71

72

73

74

75

76

77

78

Classical meta-analysis (CMA) commonly synthesises effect sizes calculated from the results reported in published studies, where each study contributes one or more effect size. Thus, the biases discussed above can lead to overstated meta-analytical means. For example, [31] simulations showed that p-hacking and publication bias interact to severely distort meta-analytic effect size estimates, particularly when true effects are small. Empirically, [15] found that correcting for publication bias across 87 ecological and evolutionary meta-analyses reduced effect sizes by at least 0.12 standard deviations on average, with 33 out

79 of 50 initially significant meta-analytic means becoming non-significant
80 after correction. These biases can also lead to confidence intervals with
81 sub-nominal coverage (i.e., intervals labelled as 95% may contain the
82 true effect less than 95% of the time) [15]. Statistical corrections for
83 publication bias exist, but often perform poorly in the presence of effect
84 size heterogeneity [16], which is widespread in different research fields
85 [32–34]. Consequently, even well-designed CMAs may systematically
86 overestimate true effect sizes when relying solely on published results.

87 Raw data meta-analysis (RDMA) provides an alternative approach
88 by calculating effect sizes directly from the original datasets of individ-
89 ual studies rather than extracting summary statistics from published
90 articles. This approach is becoming increasingly viable across the range
91 of disciplines, as open data becomes the norm [35]. Open data is often
92 publicly available in repositories, allowing researchers to recompute the
93 effect size for each study using raw observations. These recalculated
94 study-level effects are then synthesised using standard meta-analytic
95 models. This differs from the individual participant data (IPD) meta-
96 analysis often used in clinical research, where participant-level data
97 from multiple studies are combined into a single dataset and analysed
98 jointly. In RDMA, raw data are used only to recompute the effect sizes
99 within each study, and the meta-analysis still combines study-level es-
100 timates rather than pooling observations across studies.

101 However, while the growing availability of datasets through data-
102 sharing mandates and changes in research culture [36–38] makes RDMA
103 increasingly feasible, when and to what extent it outperforms CMA
104 is unknown. Thus, it remains unclear when researchers should aim to
105 conduct RDMA (or incorporate raw data in their meta-analysis) rather
106 than CMA. Although more data are becoming open [39, 40], open data
107 are often time-consuming to find [39, 40], interpret, and use [41, 42].
108 Here, we addressed this gap using a simulation-based approach, where
109 the main conditions under which meta-analysis operates can be set
110 to levels that reflect those present in a given research field. These
111 conditions include the magnitude of true effect size, heterogeneity be-
112 tween studies, typical sample size, and the severity of publication bias.
113 To parametrise the simulation and understand the potential benefits
114 of RDMA, we used ecology as a case example. However, the frame-
115 work was designed so that researchers can evaluate the performance of
116 RDMA versus CMA in their own domain by modifying the input pa-
117 rameters (effect size magnitude, heterogeneity level, sample size, and
118 bias severity). We provide worked examples using parameters from
119 ecology, complete documentation, and annotated scripts accessible in
120 the online repository [43].

121

2 Methodology

122

123

124

The study was preregistered on the Open Science Framework [44]. Deviations are explained in Supplementary Material S1. All code and data are openly available [43].

125

2.1 The overarching approach

126

127

128

129

130

131

132

133

134

135

136

137

138

139

Our framework simulates RDMA and CMA across different scenarios of true effect size magnitude and heterogeneity levels, and under different levels of presence of significant results in the published literature. For each scenario, the algorithm generates a population of N conducted studies, representing all studies that would be published if no selective publication occurred. Each study is represented by an effect size (Hedges' g) drawn from a distribution with mean θ (the true effect size) and variance τ^2 (the heterogeneity between studies). This population of studies is then subjected to a publication selection process designed to achieve target rates of statistical significance among published studies. We chose the Hedges' g as the effect size measure because standardised mean difference represents one of the two most common effect size measures in meta-analysis across disciplines [2, 3, 45, 46].

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

To model realistic conditions, the framework incorporates a simulated publication selection process that combines two mechanisms: (1) selective publication that favours significant results and (2) p-hacking that enables achieving the desired significance rates (Figure 1). The relative contribution of each mechanism to the final rates of significant results in the published literature is recorded across simulation scenarios (see Section 3.2). Some primary studies achieve statistical significance naturally (i.e., when the true effect size and the sample size allow for this). When the number of such naturally significant studies is insufficient to meet the target significance rates in published studies, p-hacking is applied. Here, the smallest needed number of non-significant studies, whose natural significance is $0.05 < p \leq 0.30$, is manipulated to reach significance. This approach ensures that p-hacking is applied conservatively—only when needed to reach empirically observed proportions of published significant results.

155

156

157

158

159

160

161

162

CMA then samples $k = N/2$ studies from the sample of studies that have gone through the simulated publication selection process. We have chosen to sample 50% of conducted studies as this reflects the publication rates detected across disciplines [18, 30, 47]. RDMA randomly samples $k = N/2$ studies from the population of all conducted studies (before the publication selection process). Performance of CMA and RDMA is then compared based on the accuracy of the effect size estimate, confidence interval coverage, and heterogeneity es-

163 timation (see Performance Metrics below for details). Compared to
164 the logic of IPD meta-analysis in clinical research, our process is more
165 relevant to other research areas (i.e. ecology, social sciences, environ-
166 mental sciences) where effect sizes are calculated per study, rather than
167 calculated for pooled data across studies.

168 The process is replicated 10,000 times for each combination of effect
169 size, heterogeneity level, and bias scenario ($3 \times 3 \times 2 = 18$ combina-
170 tions, of which 15 were valid in our case; see Section 2.2.3).

171 2.2 Empirical parameter estimation

172 To understand how different conditions affect the performance of CMA
173 and RDMA, we use ecology as a case study. However, many of the
174 parameters we use are realistic across scientific disciplines (see Discus-
175 sion). To obtain realistic simulation parameters, we used two ecology
176 literature datasets: (1) Costello and Fox [45, 46], which compiled 467
177 ecological meta-analyses with over 111,000 effect sizes, and (2) Kim-
178 mel et al. [22, 23], which contains 26,747 statistical estimates from 350
179 ecology papers published in Ecology, Ecology Letters, Journal of Ecol-
180 ogy, Science and Nature between January 2018 and May 2020. Based
181 on the Costello and Fox [45, 46] dataset, we calculated the typical ef-
182 fect size magnitude, between-study heterogeneity, and the number of
183 studies used in meta-analysis. Based on both datasets [22, 23, 45, 46],
184 we calculated the prevalence of statistically significant results in pub-
185 lished studies. Sample sizes of primary studies were derived from Fox
186 [48] and Kimmel et al. [22, 23] (see Section 2.2.2 for details).

187 2.2.1 Effect sizes and heterogeneity

188 Based on Costello and Fox [45, 46], ecological meta-analyses included
189 absolute Hedges' g values distributed with quartiles of 0.22 (Q_1), 0.55
190 (median), and 1.1 (Q_3). These empirical values align closely with Co-
191 hen's (1988) conventions for small ($g = 0.2$), medium ($g = 0.5$), and
192 large ($g = 0.8$) effect sizes. We therefore adopted Cohen's standardised
193 thresholds for our simulations.

194 We computed between-study heterogeneity (τ^2) from the ecological
195 meta-analyses compiled by Costello and Fox [45, 46], finding a distri-
196 bution with $Q_1 = 0.26$, median = 0.53, and $Q_3 = 0.88$. Critically, τ^2
197 showed a modest positive correlation with the mean meta-analytical
198 effect size magnitude ($r = 0.41$, $p = 0.005$). To capture this empirical
199 pattern, rather than using fixed τ^2 values across all effect sizes, we
200 scaled heterogeneity proportionally to the effect size magnitude:

- 201 • Low heterogeneity: $\tau^2 = 0.40 \times \text{effect size}$
- 202 • Medium heterogeneity: $\tau^2 = 1.00 \times \text{effect size}$

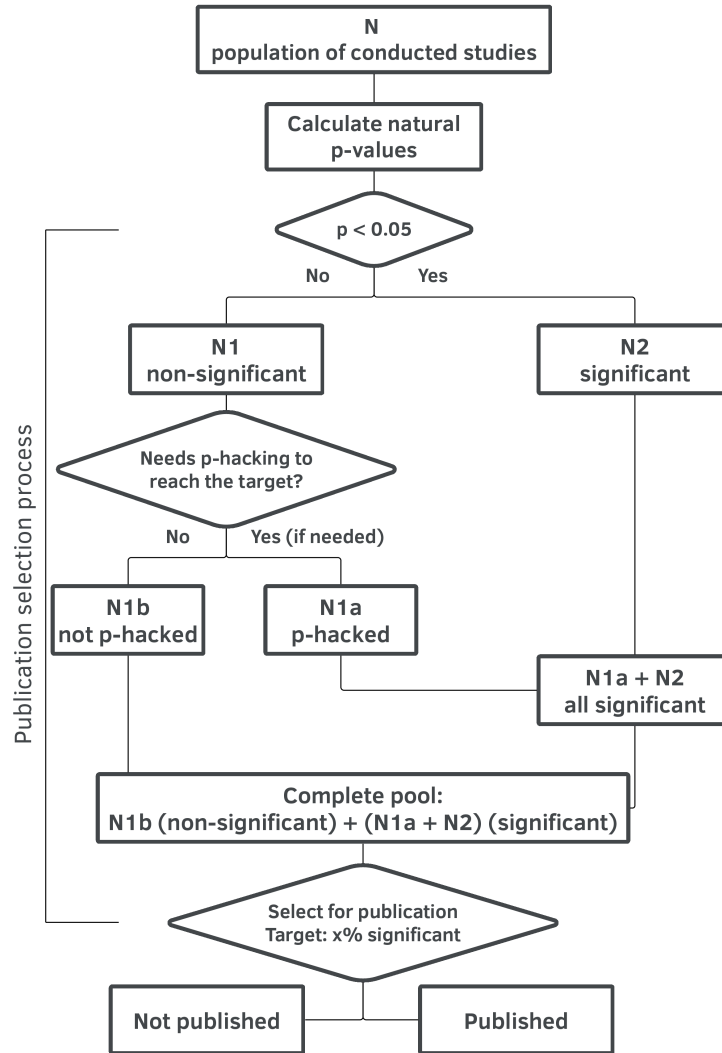


Figure 1: Flowchart illustrating the Classical meta-analysis (CMA) simulation process. The process begins with setting parameters: true effect size magnitude (θ), between-study heterogeneity (τ^2), sample sizes of primary studies (n_i), and population size of conducted studies (N). We then generate a population of N studies with known effect sizes and calculate natural p-values. Studies are then subjected to publication selection process with the aim to reach a target rate $x\%$ of $k = N/2$ published studies being significant. This process acts through selective publication and p-hacking.

203

- High heterogeneity: $\tau^2 = 1.50 \times$ effect size

204

205

206

207

208

This parametrisation produced τ^2 values that span from below the empirical Q_1 (0.08 for small effects with low heterogeneity) to above Q_3 (1.20 for large effects with high heterogeneity), covering the full empirical range while maintaining consistent relative heterogeneity across effect size magnitudes.

209

2.2.2 Meta-analysis and primary studies sample size

210

211

212

213

The median ecological meta-analysis comprises 64 effect sizes from 23 studies [45, 46]. Given that our simulation assumes one effect size per study, each meta-analysis in our simulation included 23 published studies (and thus effects).

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

For determining the realistic sample size of the primary studies included in the meta-analysis, we used two existing large datasets, both of which have some limitations. The first dataset, Fox [48], provides sample sizes only for studies that report correlation coefficients (median $n = 30$), and not for standardised mean differences (Hedges' g), which we use in our simulation. The second dataset, [22, 23] reports sample size for various effect sizes, including group comparisons suitable for Hedges' g (median $n = 79$). However, these studies come from the general primary literature rather than specifically from studies included in meta-analyses. To reconcile these sources, we used a sampling approach that reflects the variation in the sample sizes of ecological studies: 25% of studies with $n = 20 - 30$ (small samples), 50% with $n = 30 - 70$ (moderate samples), and 25% with $n = 70 - 100$ (larger samples). This stratified sampling approach was applied to all $N = 46$ conducted primary studies in each simulation iteration, with sample sizes randomly drawn from the specified ranges. This yielded a median of around $n = 40$, a value that is intermediate between the median sample sizes of the primary studies included in the ecological meta-analysis ($n = 30$, [48]) and the median of the primary literature ($n = 79$, [22, 23]). The total sample size was split equally between the treatment and control groups.

235

236

237

238

239

240

241

This approach captures the low statistical power typical of many research fields. For example, [4] found only 13–16% power to detect small effects and 40–47% for medium effects in behavioural ecology, reflecting the prevalent small sample sizes. Similarly, [6] reported mean power of just 48% for medium effects in psychology, [7] estimated a median power of 21% across neuroscience, and [8] found a median power of 18% or less across 159 areas of economics research.

242

2.2.3 Publication selection process rates

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

For the purpose of this study, we define bias as the percentage of published studies that contain a significant result achieved through both selective publication and p-hacking. Because each study contributes one effect size in our simulation, the percentage of significant studies equals the percentage of significant effects. Based on Kimmel et al. [22, 23] data, we calculated that 98.9% of published ecological studies contain at least one p-value below 0.05. From Costello and Fox [45, 46], we calculated that 64% of studies included in meta-analyses contain at least one significant result (amongst the results included in the meta-analysis). Consequently, we simulated two scenarios of bias severity, i.e., the percentage of statistically significant effects included in meta-analyses: 95% (near primary literature conditions, representing severe distortion) and 65% (matching typical meta-analytic data, representing moderate distortion). Throughout this paper, we refer to these scenarios as 'severe bias' and 'moderate bias'. For large effect sizes ($g = 0.8$), the natural significance rate exceeded 65% due to high statistical power (studies naturally achieved significance without bias), making the 65% bias scenario inapplicable. This resulted in 15 valid simulation scenarios.

262

2.3 Data generation process

263

264

265

266

267

268

269

270

271

272

273

274

Step 1: Generate population of all conducted studies. We generated 46 studies representing all conducted research before the publication selection process. This number is twice the number of studies used in meta-analysis in ecology ($k=23$, i.e. published studies in our case) and reflects empirical data on how much of conducted research is published in ecology (55.3% [30]), but also in other disciplines (approximately 50% of clinical trials [47], 25–29% of doctoral dissertations in psychology [49]) and selective non-publication of null results is pervasive across the social sciences (only 10 out of 48 null results were published, compared to 56 out of 91 strong results; [18]). This 2-to-1 ratio of conducted vs published studies determines the need for p-hacking.

275

276

277

278

279

280

281

282

283

For each study, we simultaneously: (1) drew one true effect size from $\theta_i = \mu_\theta + \delta_i$ where $\delta_i \sim N(0, \tau^2)$, using $\mu_\theta \in \{0.2, 0.5, 0.8\}$ and scaling τ^2 proportionally as described above; and (2) assigned a total sample size using stratified sampling (25% small: 20–30, 50% medium: 30–70, 25% large: 70–100), which were then split equally between the treatment and the control group.

Step 2: Simulate data and calculate effect sizes. For each study, we simulated individual-level data for treatment and control groups to generate the observed effect size (Hedges' g), its sampling

284 variance, and p-value, following standard simulation approaches for
285 meta-analysis [50, 51]. In this process, each study is also assigned its
286 natural significance.

287 **Step 3: Create published literature through selective pub-**
288 **lication process.** We constructed a published studies pool targeting
289 23 studies (our target meta-analysis sample size) where the target per-
290 centage (95% or 65%) of studies contained significant results. Some
291 studies from the population of 46 conducted studies achieved statisti-
292 cal significance ($p < 0.05$) naturally in Step 2, based on their true
293 effects and sample sizes. When the number of such naturally significant
294 studies was insufficient to meet the target significance rate, we applied
295 p-hacking to the minimum number of non-significant studies necessary
296 to reach the target. We used only those studies with the original signif-
297 icance $0.05 < p \leq 0.30$, and ordered them by statistical significance. If
298 the study was amongst the first 50% of studies (those with the lowest
299 p-values), it underwent mild manipulation: the p-value was replaced
300 with a random value drawn from a uniform distribution $U(0.01, 0.049)$
301 while retaining the original effect size. If it was in the remaining 50%
302 (with higher p-values), it underwent aggressive manipulation: a new
303 p-value was drawn from $U(0.01, 0.049)$, a corresponding t-statistic was
304 back-calculated from the new p-value and the study’s degrees of free-
305 dom, and the effect size was recalculated as $g_i = t_i \sqrt{v_i}$, where v_i
306 is the original sampling variance. Starting from those originally closest
307 to statistical significance, only the minimum number of studies needed
308 to reach the target significance rate were manipulated; all remaining
309 studies retained their original p-values and effect sizes. This approach
310 reflects realistic questionable research practices ranging from rounding
311 of p-values to optional stopping or manipulation of data [27]. Recent
312 evidence suggests that more than half of statistically significant find-
313 ings in environmental sciences may result from such practices [52]. The
314 resulting published literature of 23 studies combined naturally signif-
315 icant studies (randomly sampled from available naturally significant
316 studies), p-hacked studies, and a small proportion of non-significant
317 studies (randomly sampled from the remaining non-significant studies)
318 that were not excluded by the simulated publication selection process.

319 **Step 4: Sample for meta-analysis.** CMA used all 23 studies
320 from the published (biased) literature created in Step 3. RDMA ran-
321 domly sampled 23 studies from the original unbiased population of 46
322 conducted studies (Step 1, before p-hacking and publication selection).
323 This choice makes two simplified assumptions: that CMA includes all
324 available published studies (which is likely not true) and that raw data
325 meta-analysis has access to the same number of studies. However, de-
326 pending on the research question, RDMA might have access to more
327 studies (i.e. where data are published, though the study was not) or to
328 fewer studies (i.e. where datasets are not available). Our choice reflects

329 a state between the ideal one (all raw data are available) and the worst
330 one (no data available), while maintaining a balanced comparison by
331 keeping the number of studies equal between the two methods.

332 **Step 5: Meta-analytical estimates.** Both CMA and RDMA
333 used random-effects meta-analysis with REML estimation to obtain
334 effect size estimates ($\hat{\theta}$), heterogeneity estimates ($\hat{\tau}^2$), and 95% confi-
335 dence intervals.

336 A complete mathematical formalisation of the simulation algorithm—
337 including the exact conditioning structure, the construction of the pub-
338 lished sample, and the expectation and bias of the CMA estimator—is
339 provided in Supplementary Material S2.

340 2.4 Performance metrics

341 We evaluated the performance of CMA and RDMA using:

342 **Mean absolute error (MAE):** $\frac{1}{n} \sum_{i=1}^n |\hat{\theta}_i - \mu_{\theta}|$ is the average
343 magnitude of the estimation error across all iterations, where $n =$
344 10,000 is the number of simulation iterations.

345 **Relative mean absolute error (RelMAE):** $\frac{1}{n} \sum_{i=1}^n \frac{|\hat{\theta}_i - \mu_{\theta}|}{\mu_{\theta}}$ is
346 the estimation error expressed as the proportion of the true effect size,
347 enabling comparison across effect size magnitudes. Throughout this
348 manuscript, when reporting percentage improvements or reductions
349 between CMA and RDMA, we use this metrics to enable fair compar-
350 ison across different effect size magnitudes.

351 **Confidence interval coverage:** is the proportion of 95% CIs that
352 contain the true effect size. Well-calibrated methods should achieve
353 95% or higher coverage.

354 **Heterogeneity estimation accuracy:** the mean absolute devia-
355 tion $E[|\hat{\tau} - \tau|]$ and relative deviation $E[(\hat{\tau} - \tau)/\tau]$, measuring how well
356 the between-study standard deviation ($\tau = \sqrt{\tau^2}$) is estimated.

357 2.5 Software

358 All analyses were conducted in R version 4.4.2 [53], using metafor (v4.6-
359 0) [54], dplyr (v1.1.4) [55], tidyr (v1.3.1) [56], ggplot2 (v4.0.1) [57], and
360 patchwork (v1.2.0) [58].

361 3 Results

362 3.1 Overview

363 Table 1 presents performance metrics across simulation scenarios. Over-
364 all, under severe bias (95%), RDMA substantially outperformed CMA

365 across the performance metrics. RDMA’s advantages were most pro-
366 nounced under low to medium heterogeneity, and increased with de-
367 creasing sample size of primary studies used in meta-analysis. Un-
368 der moderate bias (65%), RDMA maintained advantages primarily for
369 small to medium effects with low heterogeneity. Conceptually, these
370 gains reflect the benefit of bypassing the publication filter, an advan-
371 tage also relevant to IPD meta-analysis in clinical research.

372 **3.2 P-hacking patterns**

373 The simulation revealed that the contribution of p-hacking to bias
374 considerably varied with the effect size magnitude. For small effects
375 (ES=0.2) under severe bias (95%), p-hacking rates also varied sub-
376 stantially across heterogeneity levels: 52.6% under low heterogeneity,
377 34.4% under medium heterogeneity, and 23.2% under high heterogen-
378 eity. For medium effects (ES=0.5), p-hacking was considerably lower
379 (9.8%, 2.7%, and 0.9% for low, medium, and high heterogeneity, re-
380 spectively), and for large effects (ES=0.8), essentially no p-hacking
381 was needed (0.1%, 0.1%, and 0.0%), suggesting that selective non-
382 publication dominated over questionable research practices in these
383 scenarios. Under moderate bias (65%), p-hacking was substantially
384 reduced: small effects required 25.1%, 8.4%, and 2.7% p-hacking (low,
385 medium, and high heterogeneity), while medium effects required min-
386 imal manipulation (0.3%, 0.0%, and 0.0%).

387 **3.3 Effect size estimation across the scenarios**

388 RDMA’s estimates of the mean effect size are consistently centred on
389 the true effect size (dashed lines), with distribution widths reflecting
390 the inherent sampling variability from $k = 23$ studies and heterogeneity
391 (Figure 2). In contrast, CMA’s estimates are systematically shifted
392 upward under severe bias (Figure 2, Panel B), with the degree of the
393 shift varying by effect size magnitude and heterogeneity level. RelMAE
394 values (Figure 2, Panels C and D) demonstrate how these patterns
395 translate to proportional estimation error.

396 The magnitude of CMA’s overestimation directionally varied with
397 the p-hacking rates (Section 3.2, Table 1). For small effects with low
398 heterogeneity, and under severe bias, where p-hacking was the high-
399 est (52.6%), CMA deviated the most from the true effect, and RDMA
400 achieved a 76% reduction in mean absolute error (MAE). This mecha-
401 nistic link explained a counter-intuitive pattern visible in Table 1: for
402 small effects, CMA’s MAE decreases as heterogeneity increases (0.29
403 \rightarrow 0.25 \rightarrow 0.23). This occurred because increased heterogeneity de-
404 creased the need for p-hacking (52.6% \rightarrow 34.4% \rightarrow 23.2%), thus reduc-
405 ing CMA’s systematic bias, more than it increased the random error.

Table 1: Performance Metrics

Heterogeneity Level	Method	65% Bias				95% Bias					
		ME est.	MAE	ReIMAE	CI Cov.	ME est.	MAE	ReIMAE	CI Cov.	Mean p-hack %	
True ES = 0.2											
Low	CMA	0.37	0.17	0.86	0.57	25.14	0.49	0.29	1.43	0.20	52.63
	RDMA	0.20	0.07	0.33	0.94	—	0.21	0.07	0.34	0.94	—
Medium	CMA	0.33	0.15	0.77	0.81	8.36	0.44	0.25	1.23	0.56	34.42
	RDMA	0.20	0.09	0.45	0.93	—	0.20	0.09	0.45	0.94	—
High	CMA	0.30	0.15	0.74	0.88	2.73	0.42	0.23	1.16	0.70	23.23
	RDMA	0.19	0.10	0.51	0.94	—	0.20	0.10	0.51	0.94	—
True ES = 0.5											
Low	CMA	0.64	0.14	0.29	0.78	0.32	0.81	0.31	0.62	0.20	9.85
	RDMA	0.49	0.09	0.18	0.94	—	0.49	0.09	0.18	0.94	—
Medium	CMA	0.59	0.14	0.29	0.92	0.02	0.81	0.31	0.63	0.56	2.73
	RDMA	0.49	0.13	0.25	0.93	—	0.49	0.13	0.25	0.94	—
High	CMA	0.56	0.16	0.31	0.94	0.00	0.78	0.30	0.60	0.72	0.92
	RDMA	0.49	0.15	0.30	0.94	—	0.49	0.15	0.30	0.94	—
True ES = 0.8											
Low	CMA	—	—	—	—	—	1.07	0.27	0.34	0.31	0.12
	RDMA	—	—	—	—	—	0.78	0.11	0.13	0.93	—
Medium	CMA	—	—	—	—	—	1.10	0.32	0.40	0.63	0.07
	RDMA	—	—	—	—	—	0.78	0.16	0.20	0.94	—
High	CMA	—	—	—	—	—	1.08	0.32	0.39	0.76	0.01
	RDMA	—	—	—	—	—	0.78	0.19	0.24	0.94	—

Note: True ES = true effect size (Hedges' g). ME est. = Mean estimated effect size across 10,000 iterations. MAE = Mean Absolute Error. ReIMAE = Relative Mean Absolute Error (MAE/true ES). CI Cov. = Confidence Interval Coverage (proportion of CIs containing the true ES; target: 0.95). Mean p-hack % = Mean percentage of studies requiring p-hacking to achieve target publication bias across iterations. Heterogeneity levels are defined relative to effect size ($\tau^2 = 0.40\times, 1.00\times, 1.50\times$ ES), so absolute τ^2 values differ across effect size scenarios (see Section 2.2.1)

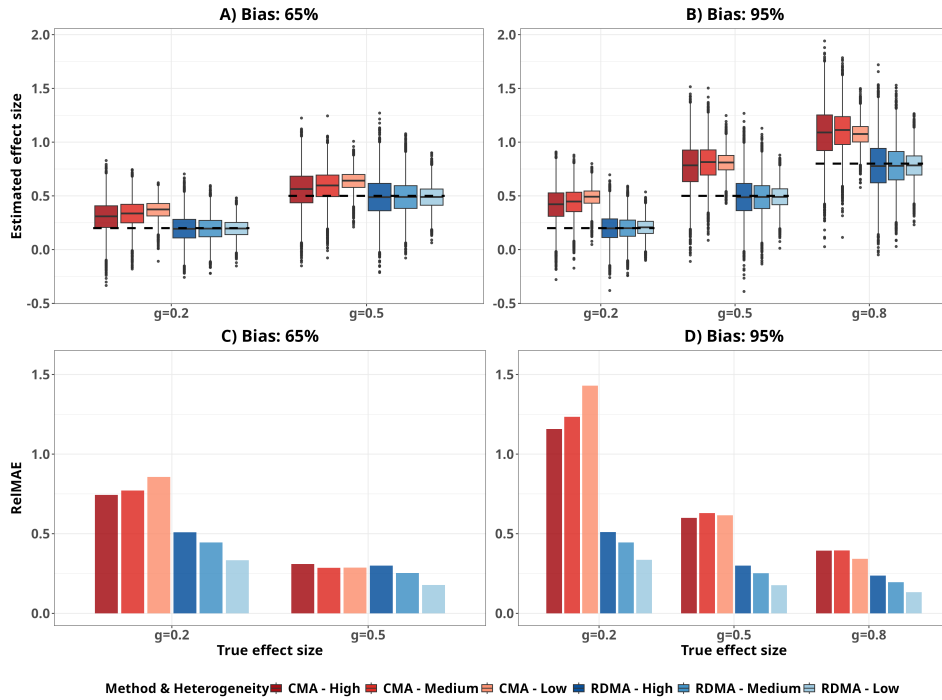


Figure 2: Effect size estimates and RelMAE comparing CMA and RDMA across scenarios. CMA is shown in red, RDMA in blue. The intensity of each colour represents the level of heterogeneity among the true effect sizes (darker colours represent higher heterogeneity). Panels A and B: Boxplots of estimated effect sizes under moderate (65%) and severe (95%) bias, summarising 10,000 iterations with dashed horizontal lines indicating true effect sizes. Panels C and D: RelMAE under moderate and severe bias. *Note: RDMA distributions are identical in panels A and B (and C and D) as this method samples from the unbiased population; RDMA is displayed in both bias scenarios to facilitate direct visual comparison with CMA*

406 For large effects with minimal p-hacking (0.1%), CMA’s MAE followed
407 the expected pattern, increasing with heterogeneity (0.27 → 0.32). In
408 contrast, RDMA’s MAE increased monotonically with heterogeneity,
409 confirming that errors reflect natural sampling variance rather than
410 bias. For medium and large effects with minimal p-hacking (0.9–9.7%
411 and 0.0–0.1%), CMA’s overestimation remained substantial, demon-
412 strating that selective publication alone distorts estimates.

413 Under moderate bias (Figure 2, Panel A), CMA’s estimates are
414 closer to true values than under severe bias, yet systematic overesti-
415 mation persisted, particularly for small effects with low heterogeneity,
416 where p-hacking remained at 25.1%. For medium to large effects with
417 high heterogeneity, where p-hacking was minimal ($\leq 0.3\%$), CMA and
418 RDMA distributions substantially overlap, indicating that CMA per-
419 forms reasonably well in these specific conditions.

420 3.4 Confidence interval coverage

421 RDMA maintained near-nominal CI coverage (93–94%) across all sim-
422 ulation scenarios, closely approaching the 95% target expected of a
423 well-calibrated method. CMA’s CI coverage, in contrast, was severely
424 degraded under severe bias, with the extent of coverage depending criti-
425 cally on the effect size magnitude and heterogeneity (Figure 3, Panel
426 B).

427 Under severe bias (95%), CMA’s coverage was worst under low het-
428 erogeneity: only 20% coverage for small and medium effects (i.e. 80%
429 of CIs failed to capture the true effect), and 31% for large effects. CMA
430 performed slightly better under medium heterogeneity (56% coverage
431 for small to medium effects, 63% for large effects), while high het-
432 erogeneity brought coverage closer to the acceptable levels (70–76%),
433 though still substantially below RDMA’s consistent $>92\%$. This pat-
434 tern is largely mechanical: higher heterogeneity inflates the estimated
435 between-study variance, widening the confidence intervals and increas-
436 ing the probability of capturing the true effect despite systematic bias
437 in point estimates.

438 Under moderate bias (65%), CMA approached acceptable perfor-
439 mance for medium effects with medium-to-high heterogeneity (92–94%
440 coverage), but still showed a substantial undercoverage for small effects
441 with low heterogeneity (57%).

442 3.5 Relative mean absolute error

443 RDMA exhibited minimal systematic bias (the mean estimates within
444 2–5% of true values across all scenarios). RDMA’s RelMAE ranged
445 from 13–51% depending on effect size and heterogeneity, primarily re-
446 flecting a natural sampling variation from $k = 23$ studies rather than

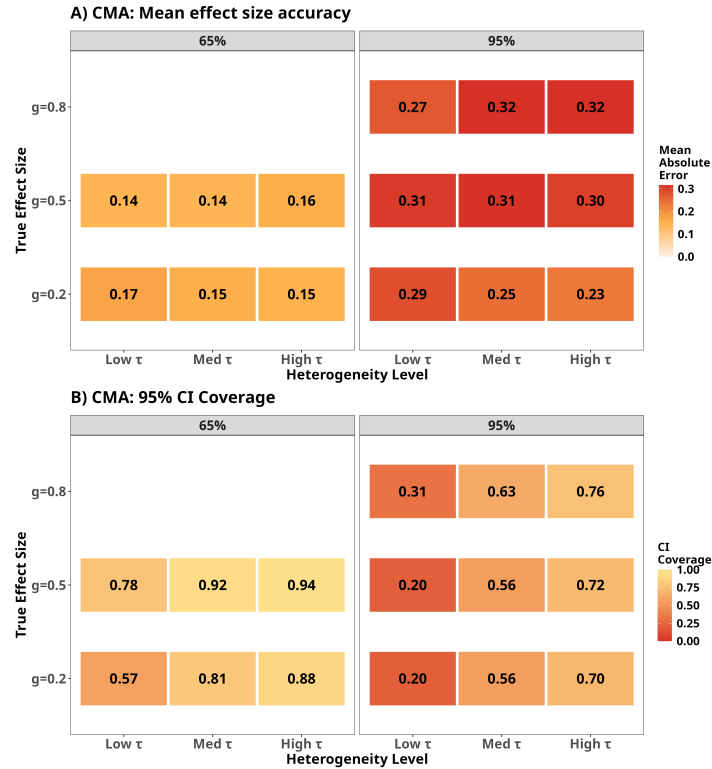


Figure 3: Performance heatmaps for CMA across simulation conditions. Panel A: Mean absolute error (MAE) in effect size estimates (lower values indicate larger accuracy). Panel B: 95% confidence interval coverage (target = 0.95). Rows represent true effect sizes ($g = 0.2, 0.5, 0.8$), columns represent heterogeneity levels (low, medium, high τ), and facets show bias scenarios (65% and 95% significant results amongst published results)

447 a systematic over- or underestimation. In some scenarios, RDMA
448 showed slight underestimation (e.g., mean estimate of 0.2 for true
449 $ES = 0.2$ under high heterogeneity), but RelMAE remained driven by
450 the heterogeneity-induced variance. As heterogeneity increased within
451 each magnitude of the effect size, RDMA’s RelMAE increased propor-
452 tionally (e.g., $ES = 0.2$: 34% \rightarrow 45% \rightarrow 51% across low, medium, and
453 high heterogeneity), confirming that errors likely stem from sampling
454 variance in heterogeneous populations rather than from bias.

455 In contrast, CMA showed both systematic overestimation and high
456 RelMAE. For small effects with low heterogeneity under severe bias,
457 CMA’s mean estimate represented 143% overestimation with 143%
458 RelMAE; nearly identical values indicating a systematic bias rather
459 than a random variation. Standardising absolute error by effect size
460 revealed different relative impacts: a similar MAE represents 143%
461 error for $ES = 0.2$ but only 34% error for $ES = 0.8$, demonstrating
462 why RelMAE is essential for comparing accuracy across effect sizes.

463 Under moderate bias (65%), CMA’s RelMAE improved substan-
464 tially (29–86% across scenarios), approaching RDMA’s performance
465 for medium effects with high heterogeneity where p-hacking was min-
466 imal (Table 1, Figure 2).

467 3.6 Heterogeneity estimation

468 Both methods estimated heterogeneity reasonably well, with RDMA
469 performing slightly better than the CMA (Figure 4). RDMA’s MAE
470 for heterogeneity estimation (τ) was consistently lower than CMA’s
471 across scenarios, with the largest advantage under severe bias (95%):
472 RDMA MAE 0.088 vs. CMA MAE 0.187, a difference of 0.099 units
473 ($ES = 0.2$, high heterogeneity), and near-zero differences in the least
474 biased scenarios. Complete heterogeneity estimation results for all the
475 iterations are available in the accompanying data files.

476 Under moderate bias (65%), the performance gap nearly disap-
477 peared for medium and large effects with medium-to-high heterogene-
478 ity, although RDMA maintained its advantage for small effects and low
479 heterogeneity scenarios (Table 1, Figure 4).

480 4 Discussion and Conclusions

481 We provide the first estimate of the benefits that raw (open) data can
482 bring to meta-analysis by mitigating the effects of publication bias
483 and p-hacking. Our results reveal that raw data meta-analysis yields
484 widespread advantages over the classical meta-analysis (when the sam-
485 ple number of included primary studies and their sample size is the
486 same between the two). These advantages depend on the extent of

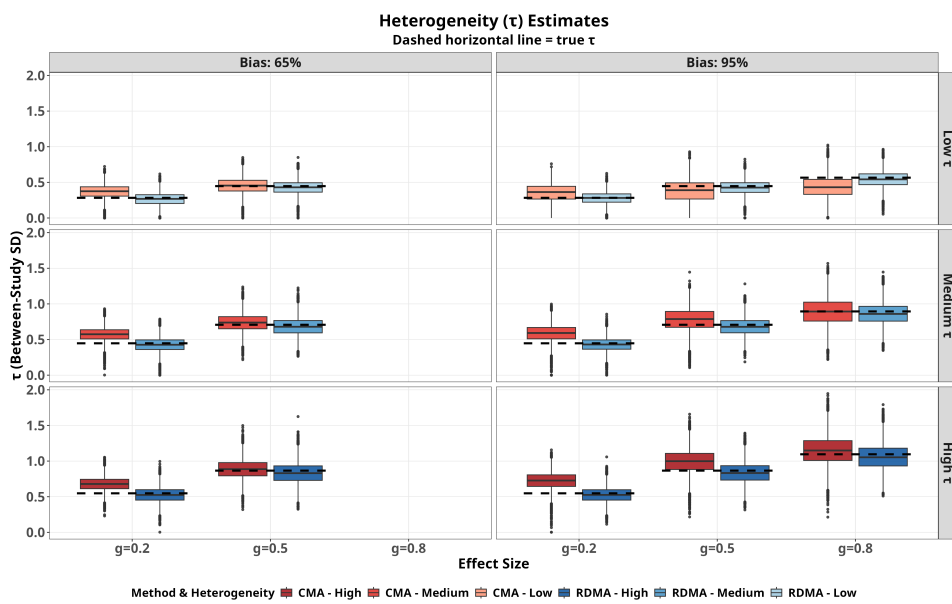


Figure 4: Heterogeneity estimates (τ , between-study standard deviation) from CMA (red) and RDMA (blue). Box-plots show distributions from 10,000 iterations. Dashed lines indicate true τ values under three scenarios of the true effect size values.

487 publication bias, with the magnitude of the true effect size and the
488 heterogeneity level both playing important moderating roles. Further,
489 our simulations also revealed under what conditions p-hacking is sup-
490 ported. We discuss our findings and their implications below.

491 **4.1 Conditions favouring p-hacking**

492 Our simulation revealed that the contribution of p-hacking to overall
493 bias (which we define as the percentage of significant results in pub-
494 lished literature) largely depends on three factors: the magnitude of the
495 true effect size, the heterogeneity between the studies, and the severity
496 of bias. The need for p-hacking is substantial when the true effects are
497 small ($ES = 0.2$) and when bias is large (95% of published results are
498 significant), peaking at 52.6% under low heterogeneity. When the het-
499 erogeneity of small effects is low, true effect sizes cluster tightly around
500 a small mean value. Given that sample sizes typical for ecology (median
501 $n \approx 40$) rarely achieve sufficient power to reach significance for such
502 effects, p-hacking becomes necessary to reach publication thresholds
503 of significant p-values. As heterogeneity of small effects increases, the
504 distribution of true effect sizes widens, so some studies have larger true
505 effects that achieve natural significance (statistical significance under
506 the data-generating model), reducing the need for p-hacking. This ex-
507 plains why, paradoxically, CMA's mean absolute error decreases with
508 increasing heterogeneity for $ES = 0.2$: higher heterogeneity reduces
509 the proportion of studies requiring p-hacking, thereby reducing the
510 severity of bias.

511 In contrast, under severe bias, minimal p-hacking was needed for
512 medium effects (mean 4.5%) and essentially none for large effects (0.07%),
513 where natural significance rates approach or exceed target levels due to
514 adequate statistical power. Here, selective publication dominated over
515 p-hacking. Under moderate bias (65%), p-hacking requirements were
516 substantially reduced but not eliminated: small effects still required
517 considerable p-hacking (25.1% under low heterogeneity), while medium
518 effects required minimal manipulation ($\leq 0.3\%$). These patterns reflect
519 our modelling choice of an approximate 50% publication rate, which is
520 likely realistic, and consistent with empirical evidence across fields [18,
521 30, 47, 49]. This ratio makes our estimates conservative: lower pub-
522 lication rates would increase p-hacking requirements. Further, effects
523 in ecology are typically small to medium (median $|g| = 0.55$; [45, 46]),
524 and low statistical power is prevalent across many research fields [4,
525 7, 8]. Thus, our results demonstrate that the current incentives and
526 publishing models, where significant results are favoured, are likely the
527 key to pushing scientists to p-hack.

4.2 When to use raw data meta-analysis

529

530

531

532

533

534

535

536

537

538

539

Overall, the benefits of RDMA for accurately estimating the mean effect increase as true effect size and heterogeneity are smaller, yet remain substantial in most scenarios. The benefits are also larger under severe bias (95% of published effects are statistically significant) compared to moderate bias. These findings suggest that RDMA provides the greatest value for research questions with relatively homogeneous effects, such as controlled experiments with standardised protocols, well-defined interventions, or phenomena governed by consistent mechanisms. It provides somewhat less value when studying highly context-dependent phenomena with inherent large heterogeneity, although it still yields meaningful improvements over CMA in these contexts.

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

RDMA's advantages are most pronounced for small effects with low heterogeneity (76% error reduction) because RDMA circumvents both selective publication and severe p-hacking in these scenarios. For larger effects, RDMA primarily addresses selective publication alone (as p-hacking rates in this scenario are small), yielding smaller but still substantial advantages (38–62% error reduction). Higher heterogeneity reduces RDMA's relative advantage across all effect sizes because true variation among studies becomes large relative to the distortion introduced by bias. Critically, the type of error differs between the two methods. RDMA's errors stem from natural sampling variation: individual meta-analyses vary around the true effect, as expected when randomly sampling $k = 23$ studies from a larger population. This produces unbiased estimates with moderate precision. CMA's errors, in contrast, reflect systematic bias, with estimates consistently shifted away from true values in a predictable direction. This distinction is fundamental: random error decreases with more studies, while systematic bias persists regardless of sample size. RDMA therefore provides not just more accurate estimates, but estimates that improve predictably as data accumulate.

559

560

561

562

563

564

565

566

567

568

569

570

571

Our simulation suggests that RDMA's primary advantage lies in accessing the unbiased population of conducted studies, bypassing the publication selection filter that distorts CMA estimates. In our framework, both RDMA and CMA apply the same analytical method (random-effects meta-analysis of study-level effect sizes); the difference is solely in the sampling frame and the integrity of the effect estimates entering the meta-analysis. CMA uses summary statistics as reported in publications, which may be distorted by selective publication and p-hacking, whereas RDMA accesses the original unmanipulated data. Our results therefore quantify the access channel through which raw data benefits meta-analysis. This advantage may be shared by other approaches that bypass the publication filter and use original datasets, such as IPD meta-analysis in clinical research, though the mecha-

572 nisms differ. For simple treatment–control comparisons with continu-
573 ous outcomes—the design that our simulation captures—published evi-
574 dence suggests that one-stage participant-level modelling and two-stage
575 study-level approaches yield practically equivalent summary treatment
576 effects [59–61], though our simulation does not test this directly. The
577 additional analytical advantages of IPD meta-analysis—such as consis-
578 tent covariate adjustment, exploration of participant-level treatment-
579 effect modifiers, and appropriate handling of time-to-event outcomes—
580 become important for more complex research questions that our frame-
581 work does not address.

582 4.3 Limitations and future directions

583 We make several simplifying assumptions. First, we conservatively
584 modelled p-hacking only when needed to reach target bias rates, mak-
585 ing our RDMA benefit estimates conservative; in reality, questionable
586 practices are more diverse. Second, we assumed RDMA randomly
587 samples all conducted studies with complete data. This design implies
588 that RDMA has access to all conducted studies, including unpublished
589 studies, and that it can use original unmanipulated data, randomly
590 sampling the entire research population. In practice, selective data
591 sharing and quality issues [41] could reduce RDMA’s advantages, so
592 results represent an upper bound. Similar constraints are well recog-
593 nised in IPD meta-analysis, where data availability and access decisions
594 can shape the final evidence base. However, RDMA can provide larger
595 sample sizes than CMA when published literature is limited [40]. On
596 the other hand, in an ideal scenario where open data are the norm,
597 RDMA might have access to a larger pool of studies than CMA, likely
598 increasing the benefits beyond those reported here.

599 Third, our simulation examined only a single effect size per study
600 from simple treatment-control comparisons and therefore did not cap-
601 ture potential within-study heterogeneity in effect sizes for studies that
602 report on more than one effect. Related to this, our representation of
603 the publication selection process enforced a target prevalence of sta-
604 tistically significant results among published studies, with one effect
605 size per study. In practice, selection mechanisms are more complex
606 across disciplines: entire studies may remain unpublished (file drawer
607 effect), researchers may selectively report significant results while omit-
608 ting non-significant ones from the same study, and published stud-
609 ies typically contain multiple effect sizes rather than one. Clinical
610 research may exhibit additional structured selection mechanisms, in-
611 cluding time-lag bias and sponsor-dependent dissemination patterns.
612 Fourth, we tested only random-effects models with REML estimation;
613 other meta-analytic approaches might show different patterns. Finally,
614 we did not examine meta-regression or subgroup analyses, which are

615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659

common in practice.

Future work should thus examine scenarios with multiple effects per study to understand how within-study selection interacts with publication bias. Future extensions could also implement explicit selection probability models, in which the probability of publication is modelled as a function of statistical significance (e.g., high probability for $p < 0.05$ and lower probability otherwise) and study characteristics, including sponsor status. Such models are widely used in clinical selection-model literature and would allow direct comparison with aggregate-data correction approaches developed for medical meta-analysis. Incorporating sponsor-dependent or outcome-level selection would be particularly relevant for extending our framework to broader bias scenarios encountered in regulatory or guideline-setting contexts.

Finally, our choice of 50% publication rate was based on empirical evidence from ecology and other fields [18, 30, 47, 49], and represents a conservative scenario. Our severe bias scenario (95% significant results) may be less realistic for ecology, where meta-analyses typically include studies that report multiple effect sizes, and only 64% of such studies contain at least one significant result among the effects included in the meta-analysis [45, 46]. However, this scenario may be more realistic in fields where primary studies commonly report a single effect size, as the proportion of significant results among published studies tends to be higher in such contexts. Publication rates might also vary across fields and research questions: if true publication rates in a specific domain are lower (e.g., 30%, requiring $N \approx 77$ conducted studies for $k = 23$ published), p-hacking requirements would be substantially higher, and RDMA's advantages would exceed those reported here. Conversely, if publication rates are higher (e.g., 70%, requiring only $N \approx 33$ conducted studies), RDMA's benefits would be smaller. Future work could explore sensitivity to this parameter by systematically varying the N:k ratio to establish boundary conditions for RDMA's benefits across different publication rate scenarios. This is easily done by adjusting the corresponding part of our simulation code. Such analyses would help researchers assess whether RDMA is worthwhile in their specific field based on field-specific publication rates.

Our estimates cannot at the moment be evaluated against real data as raw data meta-analysis remains very rare. However, if more meta-analyses begin using raw data, following existing guidelines [62], it might be possible to evaluate the assumptions of our model and its estimates. Hybrid approaches that combine published effect sizes with available raw data merit investigation, as they may provide a pragmatic middle ground. Empirical studies quantifying the time and resource costs of data acquisition would help researchers make informed decisions about when to pursue RDMA. Such extensions would strengthen understanding of when RDMA justifies the additional time investment

660 it requires.

661 **4.4 Tool for meta-analysts**

662 Our simulation framework is designed to be field-agnostic. Researchers
663 can assess RDMA’s potential benefits in their domain (or a particular
664 research area) by replacing our ecological parameters with values rele-
665 vant to their scenario. Required inputs include: (1) expected (typical)
666 effect size magnitudes (Hedges’ g or other standardised effect sizes), (2)
667 between-study heterogeneity estimates (τ^2), (3) typical sample sizes of
668 primary studies (e.g., treatment and control group sizes), (4) the typ-
669 ical number of published studies per meta-analysis ($k = 23$ in our
670 simulations, representing the median in ecology), (5) the estimated
671 publication rate in the field (approximately 50% in ecology, mean-
672 ing $N = 46$ conducted studies are needed to yield $k = 23$ published
673 studies), and (6) bias severity (proportion of significant results in the
674 published literature). Both k and the publication rate are necessary
675 because together they determine N , the population of conducted stud-
676 ies from which RDMA samples, and consequently how much p-hacking
677 is required to achieve target bias rates: a lower publication rate (re-
678 quiring a larger N of conducted studies to yield k published studies)
679 increases p-hacking requirements and amplifies RDMA’s advantages.

680 The code, which can be accessed at [43], automatically generates
681 performance metrics (mean absolute error, confidence interval cover-
682 age, heterogeneity estimation accuracy) across user-specified scenar-
683 ios. For researchers who want a quick preliminary estimate of ex-
684 pected CMA bias without running the full simulation, we also pro-
685 vide a closed-form approximation based on truncated normal theory
686 (Supplementary Material S3) with an accompanying R script (S3.R)
687 that returns an approximate bias estimate for a given combination of
688 effect size, heterogeneity, sample sizes, and publication bias severity.
689 Although our simulation used Hedges’ g as the effect size measure,
690 the framework can be adapted to other common meta-analytic met-
691 rics such as the log response ratio (lnRR) or Fisher’s z -transformed
692 correlation coefficient. The overall architecture—including the publi-
693 cation bias mechanism, the RDMA versus CMA sampling design, and
694 all performance metrics—is effect size agnostic. Adaptation requires
695 replacing the data-generating function with one appropriate for the
696 targeted effect size, adjusting the p-hacking mechanism to match the
697 corresponding test statistic, and recalibrating input parameters to re-
698 flect empirical distributions for the chosen metric. Detailed guidance
699 is provided in the code repository [43].

700

4.5 Conclusions

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

RDMA is overall superior to CMA; however, its advantages depend on the interaction of bias severity, the true effect size magnitude, and the heterogeneity level. Under severe bias (95% significant results published) with low to medium heterogeneity, RDMA substantially reduces mean absolute error compared to CMA: by 56–76% for small effects ($g = 0.2$), 50–71% for medium effects ($g = 0.5$), and 38–62% for large effects ($g = 0.8$). RDMA consistently maintained 93–94% confidence interval coverage across all scenarios, while CMA’s coverage ranged from only 20–31% under low heterogeneity to 70–76% under high heterogeneity. Under moderate bias (65%) with high heterogeneity and medium effects, the performance gap narrows substantially, with CMA performing adequately.

Given that raw data meta-analysis provides substantially more accurate effect size estimates across the scenarios most plausible for a range of research fields, researchers should more often seek to incorporate raw data into meta-analyses. This can be done by either using raw data from published studies or by identifying raw data independently [40, 62]. However, raw data are often difficult to find and use, as their quality is often very low [41, 42]. The research community, including funders and publishers, should thus seek to increase data quality and thereby increase the reliability of estimated effects.

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

Author contributions: Conceptualisation: D.Ž., A.C.; Methodology: D.Ž., S.G., A.T., A.C.; Software: D.Ž.; Validation: S.G., A.T.; Formal analysis: D.Ž.; Visualisation: D.Ž.; Funding acquisition: A.C.; Writing—original draft: D.Ž., A.C.; Writing—review and editing: S.G., A.T., A.C., L.P.

Competing interests: The authors declare that no competing interests exist.

Funding: This work was supported by the Croatian Science Foundation under the project number HRZZ-IP-2022-10-2872.

Data availability: All simulation code, parameter estimation scripts, input data, and complete results are available in the OSF repository: <https://doi.org/10.17605/OSF.IO/VEBYM> [43].

Acknowledgements: We are grateful to Dr Alfredo Sánchez-Tójar for providing feedback on aspects of the methodology. AI tools (Claude, Anthropic) were used to assist with manuscript formatting and code documentation. All content was reviewed and verified by the authors.

Ethics: This study is based entirely on simulated data. No ethical approval was required.

740

References

- 741 [1] S. Nakagawa et al. “Quantitative evidence synthesis: A
742 practical guide on meta-analysis, meta-regression, and pub-
743 lication bias tests for environmental sciences”. In: *Environ-*
744 *mental Evidence* 12.1 (2023), p. 8. DOI: [10.1186/s13750-](https://doi.org/10.1186/s13750-023-00301-6)
745 [023-00301-6](https://doi.org/10.1186/s13750-023-00301-6).
- 746 [2] M. Borenstein et al. *Introduction to Meta-Analysis*. 2nd.
747 Wiley, 2021. ISBN: 978-1-119-55835-4.
- 748 [3] J.P.T. Higgins et al. *Cochrane Handbook for Systematic Re-*
749 *views of Interventions*. 2nd. Wiley-Blackwell, 2019. ISBN:
750 978-1-119-53662-8.
- 751 [4] M.D. Jennions and A.P. Møller. “A survey of the statistical
752 power of research in behavioral ecology and animal behav-
753 ior”. In: *Behavioral Ecology* 14.3 (2003), pp. 438–445. DOI:
754 [10.1093/beheco/14.3.438](https://doi.org/10.1093/beheco/14.3.438).
- 755 [5] Y. Yang et al. “Low statistical power and overestimated
756 anthropogenic impacts, exacerbated by publication bias,
757 dominate field studies in global change biology”. In: *Global*
758 *Change Biology* 28.3 (2022), pp. 969–989. DOI: [10.1111/](https://doi.org/10.1111/gcb.15972)
759 [gcb.15972](https://doi.org/10.1111/gcb.15972).
- 760 [6] J. Cohen. “The statistical power of abnormal-social psy-
761 chological research: A review”. In: *The Journal of Abnor-*
762 *mal and Social Psychology* 65.3 (1962), pp. 145–153. DOI:
763 [10.1037/h0045186](https://doi.org/10.1037/h0045186).
- 764 [7] K.S. Button et al. “Power failure: why small sample size
765 undermines the reliability of neuroscience”. In: *Nature Re-*
766 *views Neuroscience* 14.5 (2013), pp. 365–376. DOI: [10 .](https://doi.org/10.1038/nrn3475)
767 [1038/nrn3475](https://doi.org/10.1038/nrn3475).
- 768 [8] J.P.A. Ioannidis, T.D. Stanley, and H. Doucouliagos. “The
769 power of bias in economics research”. In: *The Economic*
770 *Journal* 127.605 (2017), F236–F265. DOI: [10.1111/ecoj.](https://doi.org/10.1111/ecoj.12461)
771 [12461](https://doi.org/10.1111/ecoj.12461).
- 772 [9] M.L. Head et al. “The extent and consequences of p-hacking
773 in science”. In: *PLOS Biology* 13.3 (2015), e1002106. DOI:
774 [10.1371/journal.pbio.1002106](https://doi.org/10.1371/journal.pbio.1002106).

- 775 [10] D. Fanelli. ““Positive” results increase down the hierarchy
776 of the sciences”. In: *PLOS ONE* 5.4 (2010), e10068. DOI:
777 [10.1371/journal.pone.0010068](https://doi.org/10.1371/journal.pone.0010068).
- 778 [11] D. Fanelli. “Negative results are disappearing from most
779 disciplines and countries”. In: *Scientometrics* 90.3 (2012),
780 pp. 891–904. DOI: [10.1007/s11192-011-0494-7](https://doi.org/10.1007/s11192-011-0494-7).
- 781 [12] A. Brodeur, N. Cook, and A. Heyes. “Methods matter:
782 p-hacking and publication bias in causal analysis in eco-
783 nomics”. In: *American Economic Review* 110.11 (2020),
784 pp. 3634–3660. DOI: [10.1257/aer.20190687](https://doi.org/10.1257/aer.20190687).
- 785 [13] S.B. Bruns and J.P.A. Ioannidis. “p-Curve and p-hacking
786 in observational research”. In: *PLOS ONE* 11.2 (2016),
787 e0149144. DOI: [10.1371/journal.pone.0149144](https://doi.org/10.1371/journal.pone.0149144).
- 788 [14] T.D. Sterling. “Publication decisions and their possible ef-
789 fects on inferences drawn from tests of significance—or vice
790 versa”. In: *Journal of the American Statistical Association*
791 54.285 (1959), pp. 30–34. DOI: [10.1080/01621459.1959.
792 10501497](https://doi.org/10.1080/01621459.1959.10501497).
- 793 [15] Y. Yang et al. “Publication bias impacts on effect size, sta-
794 tistical power, and magnitude (Type M) and sign (Type S)
795 errors in ecology and evolutionary biology”. In: *BMC Biol-*
796 *ogy* 21.1 (2023), p. 71. DOI: [10.1186/s12915-022-01485-
797 y](https://doi.org/10.1186/s12915-022-01485-y).
- 798 [16] S. Nakagawa et al. “Methods for testing publication bias
799 in ecological and evolutionary meta-analyses”. In: *Methods*
800 *in Ecology and Evolution* 13.1 (2022), pp. 4–21. DOI: [10.
801 1111/2041-210X.13724](https://doi.org/10.1111/2041-210X.13724).
- 802 [17] E.H. Turner et al. “Selective publication of antidepressant
803 trials and its influence on apparent efficacy”. In: *The New*
804 *England Journal of Medicine* 358.3 (2008), pp. 252–260.
805 DOI: [10.1056/NEJMs065779](https://doi.org/10.1056/NEJMs065779).
- 806 [18] A. Franco, N. Malhotra, and G. Simonovits. “Publication
807 bias in the social sciences: unlocking the file drawer”. In:
808 *Science* 345.6203 (2014), pp. 1502–1505. DOI: [10.1126/
809 science.1255484](https://doi.org/10.1126/science.1255484).

- 810 [19] J.P. Simmons, L.D. Nelson, and U. Simonsohn. “False-positive
811 psychology: Undisclosed flexibility in data collection and
812 analysis allows presenting anything as significant”. In: *Psy-*
813 *chological Science* 22.11 (2011), pp. 1359–1366. DOI: [10 .
814 1177/0956797611417632](https://doi.org/10.1177/0956797611417632).
- 815 [20] I.A. Cristea and J.P.A. Ioannidis. “P values in display items
816 are ubiquitous and almost invariably significant: A sur-
817 vey of top science journals”. In: *PLOS ONE* 13.5 (2018),
818 e0197440. DOI: [10.1371/journal.pone.0197440](https://doi.org/10.1371/journal.pone.0197440).
- 819 [21] D. Chavalarias et al. “Evolution of reporting p values in the
820 biomedical literature, 1990-2015”. In: *JAMA* 315.11 (2016),
821 pp. 1141–1148. DOI: [10.1001/jama.2016.1952](https://doi.org/10.1001/jama.2016.1952).
- 822 [22] K. Kimmel, M.L. Avolio, and P.J. Ferraro. “Empirical evi-
823 dence of widespread exaggeration bias and selective report-
824 ing in ecology”. In: *Nature Ecology & Evolution* 7.9 (2023),
825 pp. 1525–1536. DOI: [10.1038/s41559-023-02144-3](https://doi.org/10.1038/s41559-023-02144-3).
- 826 [23] K. Kimmel and P.J. Ferraro. *Replicability in ecology*. Dataset,
827 OSF. Accessed 15 December 2025. 2023. URL: [https://
828 osf.io/9yd2b](https://osf.io/9yd2b).
- 829 [24] L.K. John, G. Loewenstein, and D. Prelec. “Measuring the
830 prevalence of questionable research practices with incen-
831 tives for truth telling”. In: *Psychological Science* 23.5 (2012),
832 pp. 524–532. DOI: [10.1177/0956797611430953](https://doi.org/10.1177/0956797611430953).
- 833 [25] F. Agnoli et al. “Questionable research practices among
834 Italian research psychologists”. In: *PLOS ONE* 12.3 (2017),
835 e0172792. DOI: [10.1371/journal.pone.0172792](https://doi.org/10.1371/journal.pone.0172792).
- 836 [26] K. Fiedler and N. Schwarz. “Questionable research prac-
837 tices revisited”. In: *Social Psychological and Personality*
838 *Science* 7.1 (2016), pp. 45–52. DOI: [10.1177/1948550615612150](https://doi.org/10.1177/1948550615612150).
- 839 [27] H. Fraser et al. “Questionable research practices in ecology
840 and evolution”. In: *PLOS ONE* 13.7 (2018), e0200303. DOI:
841 [10.1371/journal.pone.0200303](https://doi.org/10.1371/journal.pone.0200303).
- 842 [28] A.A. Baptista and F. Pereira. “Questionable research prac-
843 tices in engineering research”. In: *International Journal of*
844 *Information Systems and Project Management* 14.1 (2026),
845 pp. 1–24. DOI: [10.12821/ijispm140101](https://doi.org/10.12821/ijispm140101).

- 846 [29] G. Gopalakrishna et al. “Prevalence of questionable re-
847 search practices, research misconduct and their potential
848 explanatory factors: a survey among academic researchers
849 in The Netherlands”. In: *PLOS ONE* 17.2 (2022), e0263023.
850 DOI: [10.1371/journal.pone.0263023](https://doi.org/10.1371/journal.pone.0263023).
- 851 [30] M. Purgar, T. Klanjscek, and A. Culina. “Quantifying re-
852 search waste in ecology”. In: *Nature Ecology & Evolution*
853 6.9 (2022), pp. 1390–1397. DOI: [10.1038/s41559-022-01820-0](https://doi.org/10.1038/s41559-022-01820-0).
- 854 [31] M. Friese and J. Frankenbach. “p-Hacking and publication
855 bias interact to distort meta-analytic effect size estimates”.
856 In: *Psychological Methods* 25.4 (2020), pp. 456–471. DOI:
857 [10.1037/met0000246](https://doi.org/10.1037/met0000246).
- 858 [32] A.M. Senior et al. “Heterogeneity in ecological and evolu-
859 tionary meta-analyses: its magnitude and implications”. In:
860 *Ecology* 97.12 (2016), pp. 3293–3299. DOI: [10.1002/ecy.1591](https://doi.org/10.1002/ecy.1591).
- 861 [33] S. van Erp et al. “Estimates of between-study heterogeneity
862 for 705 meta-analyses reported in Psychological Bulletin
863 from 1990–2013”. In: *Journal of Open Psychology Data* 5
864 (2017), p. 4. DOI: [10.5334/jopd.33](https://doi.org/10.5334/jopd.33).
- 865 [34] J.P. Higgins et al. “Measuring inconsistency in meta-analyses”.
866 In: *BMJ* 327.7414 (2003), pp. 557–560. DOI: [10.1136/bmj.327.7414.557](https://doi.org/10.1136/bmj.327.7414.557).
- 867 [35] Digital Science et al. *The State of Open Data 2023*. 2023.
868 DOI: [10.6084/m9.figshare.24428194.v1](https://doi.org/10.6084/m9.figshare.24428194.v1). URL: [https://digitalscience.figshare.com/articles/report/
870 The_State_of_Open_Data_2023/24428194](https://digitalscience.figshare.com/articles/report/The_State_of_Open_Data_2023/24428194).
- 871 [36] B.L. Houtkoop et al. “Data sharing in psychology: A survey
872 on barriers and preconditions”. In: *Advances in Methods
873 and Practices in Psychological Science* 1.1 (2018), pp. 70–
874 85. DOI: [10.1177/2515245917751886](https://doi.org/10.1177/2515245917751886).
- 875 [37] A. Pepe et al. “How do astronomers share data? Reliabil-
876 ity and persistence of datasets linked in AAS publications
877 and a qualitative study of data practices among US as-
878 tronomers”. In: *PLOS ONE* 9.8 (2014), e104798. DOI: [10.1371/journal.pone.0104798](https://doi.org/10.1371/journal.pone.0104798).
- 879
880
881
882

- 883 [38] G. Colavizza et al. “The citation advantage of linking pub-
884 lications to research data”. In: *PLOS ONE* 15.4 (2020),
885 e0230416. DOI: [10.1371/journal.pone.0230416](https://doi.org/10.1371/journal.pone.0230416).
- 886 [39] L. Tedersoo et al. “Data sharing practices and data avail-
887 ability upon request differ across scientific disciplines”. In:
888 *Scientific Data* 8 (2021), p. 192. DOI: [10.1038/s41597-](https://doi.org/10.1038/s41597-021-00981-0)
889 [021-00981-0](https://doi.org/10.1038/s41597-021-00981-0).
- 890 [40] A. Culina et al. “Navigating the unfolding open data land-
891 scape in ecology and evolution”. In: *Nature Ecology & Evo-*
892 *lution* 2.3 (2018), pp. 420–426. DOI: [10.1038/s41559-017-](https://doi.org/10.1038/s41559-017-0458-2)
893 [0458-2](https://doi.org/10.1038/s41559-017-0458-2).
- 894 [41] D.G. Roche et al. “Public data archiving in ecology and
895 evolution: how well are we doing?” In: *PLOS Biology* 13.11
896 (2015), e1002295. DOI: [10.1371/journal.pbio.1002295](https://doi.org/10.1371/journal.pbio.1002295).
- 897 [42] D.G. Roche et al. “Slow improvement to the archiving qual-
898 ity of open datasets shared by researchers in ecology and
899 evolution”. In: *Proceedings of the Royal Society B: Biologi-*
900 *cal Sciences* 289.1975 (2022), p. 20212780. DOI: [10.1098/](https://doi.org/10.1098/rspb.2021.2780)
901 [rspb.2021.2780](https://doi.org/10.1098/rspb.2021.2780).
- 902 [43] D. Žanko. *Reliable inference: benefits of open raw data may*
903 *be universal in meta-analysis - data and code*. OSF. Ac-
904 cessed 27 March 2026. 2026. URL: [https://doi.org/10.](https://doi.org/10.17605/OSF.IO/VEBYM)
905 [17605/OSF.IO/VEBYM](https://doi.org/10.17605/OSF.IO/VEBYM).
- 906 [44] D. Žanko et al. *Performance of open data versus classical*
907 *meta-analysis in ecology - a simulation study*. OSF Preprints.
908 Accessed 15 December 2025. 2025. URL: [https://doi.org/](https://doi.org/10.17605/OSF.IO/KAH42)
909 [10.17605/OSF.IO/KAH42](https://doi.org/10.17605/OSF.IO/KAH42).
- 910 [45] L. Costello and J.W. Fox. “Decline effects are rare in ecol-
911 ogy”. In: *Ecology* 103.6 (2022), e3680. DOI: [10.1002/ecy.](https://doi.org/10.1002/ecy.3680)
912 [3680](https://doi.org/10.1002/ecy.3680).
- 913 [46] J.W. Fox and L. Costello. *Decline effects are rare in ecology:*
914 *a meta-meta-analysis*. Accessed 15 December 2025. Dryad,
915 2022. URL: <https://doi.org/10.5061/dryad.zkh1893b7>.
- 916 [47] I. Chalmers and P. Glasziou. “Avoidable waste in the pro-
917 duction and reporting of research evidence”. In: *Lancet* 374.9683
918 (2009), pp. 86–89. DOI: [10.1016/S0140-6736\(09\)60329-9](https://doi.org/10.1016/S0140-6736(09)60329-9).

- 919 [48] J.W. Fox. “The sample size of the typical ecological corre-
920 relation coefficient is small and slowly declining”. In: *Oikos*
921 2025.12 (2025), e11430. DOI: [10.1002/oik.11430](https://doi.org/10.1002/oik.11430).
- 922 [49] S.C. Evans et al. ““Are you gonna publish that?” Peer-
923 reviewed publication outcomes of doctoral dissertations in
924 psychology”. In: *PLOS ONE* 13.2 (2018), e0192219. DOI:
925 [10.1371/journal.pone.0192219](https://doi.org/10.1371/journal.pone.0192219).
- 926 [50] F. Gambarota and G. Altoè. “Understanding meta-analysis
927 through data simulation with applications to power analy-
928 sis”. In: *Advances in Methods and Practices in Psychological*
929 *Science* 7.1 (2024). DOI: [10.1177/25152459231209330](https://doi.org/10.1177/25152459231209330).
- 930 [51] L.V. Hedges. “Distribution theory for Glass’s estimator of
931 effect size and related estimators”. In: *Journal of Educa-*
932 *tional Statistics* 6.2 (1981), pp. 107–128. DOI: [10.3102/
933 10769986006002107](https://doi.org/10.3102/10769986006002107).
- 934 [52] T. Deressa et al. *More Than Half of Statistically Significant*
935 *Research Findings in the Environmental Sciences are Ac-*
936 *tually Not*. Accessed 15 December 2025. 2023. URL: [https:
937 //doi.org/10.32942/x24g6z](https://doi.org/10.32942/x24g6z).
- 938 [53] R Core Team. *R: A language and environment for statis-*
939 *tical computing*. Accessed 15 December 2025. R Founda-
940 tion for Statistical Computing. Vienna, Austria, 2024. URL:
941 <https://www.R-project.org/>.
- 942 [54] W. Viechtbauer. “Conducting meta-analyses in R with the
943 metafor package”. In: *Journal of Statistical Software* 36.3
944 (2010), pp. 1–48. DOI: [10.18637/jss.v036.i03](https://doi.org/10.18637/jss.v036.i03).
- 945 [55] H. Wickham et al. *dplyr: A grammar of data manipulation*.
946 R package version 1.1.4. Accessed 15 December 2025. 2023.
947 URL: <https://CRAN.R-project.org/package=dplyr>.
- 948 [56] H. Wickham, D. Vaughan, and M. Girlich. *tidyr: Tidy messy*
949 *data*. R package version 1.3.1. Accessed 15 December 2025.
950 2024. URL: [https://CRAN.R-project.org/package=
951 tidyr](https://CRAN.R-project.org/package=tidyr).
- 952 [57] H. Wickham. *ggplot2: Elegant graphics for data analysis*.
953 New York: Springer-Verlag, 2016. ISBN: 978-3-319-24277-4.

- 954 [58] T.L. Pedersen. *patchwork: The composer of plots*. R pack-
955 age version 1.2.0. Accessed 15 December 2025. 2024. URL:
956 <https://CRAN.R-project.org/package=patchwork>.
- 957 [59] D.L. Burke, J. Ensor, and R.D. Riley. “Meta-analysis using
958 individual participant data: one-stage and two-stage ap-
959 proaches, and why they may differ”. In: *Statistics in Medicine*
960 36.5 (2017), pp. 855–875. DOI: [10.1002/sim.7141](https://doi.org/10.1002/sim.7141).
- 961 [60] R.D. Riley et al. “Two-stage or not two-stage? That is the
962 question for IPD meta-analysis projects”. In: *Research Syn-
963 thesis Methods* 14.6 (2023), pp. 903–910. DOI: [10.1002/
964 jrsm.1661](https://doi.org/10.1002/jrsm.1661).
- 965 [61] T.P. Morris et al. “Meta-analysis of Gaussian individual
966 patient data: two-stage or not two-stage?” In: *Statistics in
967 Medicine* 37.9 (2018), pp. 1419–1438. DOI: [10.1002/sim.
968 7589](https://doi.org/10.1002/sim.7589).
- 969 [62] A. Culina et al. “How to do meta-analysis of open datasets”.
970 In: *Nature Ecology & Evolution* 2.7 (2018), pp. 1053–1056.
971 DOI: [10.1038/s41559-018-0579-2](https://doi.org/10.1038/s41559-018-0579-2).

972 **Supplementary Material S1: Deviations from**
973 **Preregistration**

974 Our simulation approach evolved during empirical data analysis to bet-
975 ter reflect realistic conditions in ecological meta-analysis. Key mod-
976 ifications align simulation parameters with empirical evidence while
977 maintaining core objectives.

978 **Number of studies per meta-analysis**

979 We initially planned $k = 70$ studies based on the median number of ef-
980 fect sizes per meta-analysis. However, examining Fox and Costello[46]
981 more carefully revealed that the median meta-analysis comprises 64 ef-
982 fect sizes from only 23 unique studies. Since our simulation generates
983 one effect size per study, $k = 23$ better reflects typical primary studies
984 in ecological meta-analyses.

985 **Sample size allocation**

986 Our preregistration specified 50% of studies with sample sizes $n =$
987 $20 - 30$, 25% with $n = 30 - 70$, and 25% with $n = 70 - 100$. We revised
988 this to 25% with $n = 20 - 30$, 50% with $n = 30 - 70$, and 25% with
989 $n = 70 - 100$ to better represent the range of sample sizes in ecological
990 research while maintaining prevalence of small samples.

991 **Heterogeneity parametrisation**

992 We initially specified heterogeneity as $\tau = 10 - 30\%$ of effect size (τ be-
993 ing standard deviation). Empirical analysis revealed modest positive
994 correlation between absolute effect size and τ^2 (variance) ($r = 0.41$,
995 $p = 0.005$), leading us to express heterogeneity as proportions of ef-
996 fect size magnitude: $\tau^2 = 0.40 \times$ effect size (low), $1.00 \times$ effect size
997 (medium), and $1.50 \times$ effect size (high). This produces values span-
998 ning from below empirical Q_1 (0.26) to above Q_3 (0.88) observed by
999 Fox and Costello[46].

1000 **Additional bias scenario**

1001 While our preregistration focused on severe bias (95%), we added a
1002 moderate bias scenario (65%) based on Fox and Costello[46], which
1003 showed 64% of studies in ecological meta-analyses contain significant
1004 results. These refinements improve our simulation's ecological real-

1005 ism without changing fundamental research questions or analytical ap-
1006 proach.

1007

Supplementary Material S2: Mathematical formalisation of the simulation algorithm

1008

1009

1010

1011

1012

Here we formalize the publication-bias and p-hacking mechanism exactly as implemented in the simulation code. We explicitly state the conditions under which the algorithm is successful, and calculate the expected total reported effect under those conditions.

1013

Exact conditioning structure induced by the simulation algorithm

1014

1015

1016

1017

1018

1019

1020

1021

1022

Let k denote the target number of published studies entering the classical meta-analysis, and let K denote the size of the conducted-study population generated in each attempt ($k = 23$, $K = 46$). The code repeats generation attempts until a valid published literature is obtained, so all expectations relevant to the simulation are conditional on the event that the generation algorithm succeeds.

For each conducted study $i = 1, \dots, K$, a total sample size N_i is generated, then split as

$$n_{Ti} = \left\lceil \frac{N_i}{2} \right\rceil, \quad n_{Ci} = \left\lfloor \frac{N_i}{2} \right\rfloor, \quad df_i = N_i - 2.$$

1023

A true study-specific effect is then generated as

$$\Delta_i = \mu + U_i, \quad U_i \sim \mathcal{N}(0, \tau^2).$$

1024

Conditional on (N_i, Δ_i) , the study-level raw data are generated as

$$Y_{Ci1}, \dots, Y_{Cin_{Ci}} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1), \quad Y_{Ti1}, \dots, Y_{Tin_{Ti}} \stackrel{\text{iid}}{\sim} \mathcal{N}(\Delta_i, 1).$$

1025

Let

$$\bar{Y}_{Ci} = \frac{1}{n_{Ci}} \sum_{j=1}^{n_{Ci}} Y_{Cij}, \quad \bar{Y}_{Ti} = \frac{1}{n_{Ti}} \sum_{j=1}^{n_{Ti}} Y_{Tij},$$

1026

and let the pooled standard deviation be

$$S_{pi} = \left[\frac{(n_{Ti} - 1)S_{Ti}^2 + (n_{Ci} - 1)S_{Ci}^2}{N_i - 2} \right]^{1/2}.$$

1027

The simulated Cohen's d and the Hedges correction factor are

$$d_i = \frac{\bar{Y}_{Ti} - \bar{Y}_{Ci}}{S_{pi}}, \quad J_i = 1 - \frac{3}{4(N_i - 2) - 1}.$$

1028

The observed study estimate is therefore

$$Y_i = J_i d_i,$$

1029 and the within-study sampling variance is

$$V_i = \left(\frac{N_i}{n_{T_i} n_{C_i}} + \frac{Y_i^2}{2N_i} \right) J_i^2.$$

1030 The associated test statistic and two-sided p -value are

$$T_i = \frac{Y_i}{\sqrt{V_i}}, \quad P_i = 2 \left(1 - F_{t,df_i}(|T_i|) \right),$$

1031 where $F_{t,df}$ is the cumulative distribution function of the Student t
 1032 distribution with df degrees of freedom.

1033 The simulation partitions the conducted studies into three sets:

$$A = \{i : P_i \leq 0.05\}, \quad B = \{i : 0.05 < P_i \leq 0.30\}, \quad C = \{i : P_i > 0.30\}.$$

1034 Set A contains the originally significant studies, B contains the non-
 1035 significant studies eligible for p-hacking, and C contains the remaining
 1036 non-significant studies. Let

$$a = |A|, \quad b = |B|, \quad c = |C|, \quad M = a + b + c.$$

1037 Note that we retain only studies with finite Y_i , finite V_i , and $V_i > 0$; all
 1038 formulas below are conditional on the retained conducted literature.

1039 Let

$$q \in \{0.65, 0.95\}$$

1040 denote the target proportion of published studies that must be signif-
 1041 icant. We have

$$s^* = \text{round}(qk), \quad \ell^* = k - s^*.$$

1042 so that s^* is the target number of significant published studies and ℓ^*
 1043 is the target number of non-significant published studies.

1044 If $a \geq s^*$, then no p-hacking is needed. If $a < s^*$, the number
 1045 of additional significant studies that must be created by p-hacking is
 1046 $h = s^* - a$.

1047 In all cases define

$$s = \min(a, s^*).$$

1048 which is the number of originally significant studies that are actually
 1049 published.

1050 A generation attempt is accepted only if the algorithm can con-
 1051 struct the published sample and the resulting sample size falls within
 1052 a tolerance window of the target. Let $\delta = \max(3, \text{round}(0.15k))$. The
 1053 exact success event is

$$\mathcal{E} = \{M \geq K\} \cap \{b \geq h\} \cap \{M - a - h \geq \ell^*\} \cap \{|k_{\text{actual}} - k| \leq \delta\},$$

1054 where $k_{\text{actual}} = s+h+\ell_{\text{actual}}^*$ is the realised number of published studies.
1055 The first condition is the code's requirement that at least $K = 2k$ valid
1056 studies remain after filtering. The second condition ensures that there
1057 are enough eligible studies for p-hacking. The third condition ensures
1058 that after removing the hacked studies from the non-significant pool
1059 there are still enough remaining non-significant studies to fill the ℓ^*
1060 non-significant publication slots. The fourth condition is the code's
1061 acceptance tolerance, which in practice is rarely binding ($\delta = 3$ for
1062 $k = 23$).

1063 All expectations corresponding to the actual simulation are there-
1064 fore conditional on \mathcal{E} .

1065 Exact construction of the published sample

1066 Conditional on the realized conducted literature and on \mathcal{E} , we construct
1067 the published literature as follows.

1068 If $s > 0$, the set of originally significant published studies, denoted
1069 S , is drawn uniformly without replacement from A so that $|S| = s$.
1070 Hence

$$\Pr(S = \mathcal{S} \mid A, \mathcal{E}) = \binom{a}{s}^{-1}$$

1071 for every $\mathcal{S} \subseteq A$ such that $|\mathcal{S}| = s$.

1072 If $h > 0$, eligible p-hackable studies in B are ordered by increasing
1073 p -value:

$$P_{(1)}^B \leq P_{(2)}^B \leq \dots \leq P_{(b)}^B.$$

1074 The hacked set is then chosen deterministically as the h studies closest
1075 to significance:

$$H = \{i_{(1)}, \dots, i_{(h)}\}.$$

1076 Let

$$r = \left\lceil \frac{h}{2} \right\rceil.$$

1077 The first r hacked studies, namely $i_{(1)}, \dots, i_{(r)}$, undergo "mild" p-
1078 hacking, while the remaining $h-r$ hacked studies, namely $i_{(r+1)}, \dots, i_{(h)}$,
1079 undergo "aggressive" p-hacking.

1080 For the mild subgroup, the reported effect size is unchanged:

$$Y_{i_{(j)}}^* = Y_{i_{(j)}}, \quad j = 1, \dots, r,$$

1081 and only the reported p -value is replaced by an independent draw

$$P_{i_{(j)}}^* \sim \text{Unif}(0.01, 0.049).$$

1082 For the aggressive subgroup, we generate

$$U_{i_{(j)}} \stackrel{\text{iid}}{\sim} \text{Unif}(0.01, 0.049), \quad j = r+1, \dots, h,$$

1083 then define

$$T_{i(j)}^* = F_{t, \text{df}_{i(j)}}^{-1} \left(1 - \frac{U_{i(j)}}{2} \right).$$

1084 The observed effect size becomes

$$Y_{i(j)}^* = T_{i(j)}^* \sqrt{V_{i(j)}}.$$

1085 The reported p -value is set equal to

$$P_{i(j)}^* = U_{i(j)}.$$

1086 Note the within-study variance is not altered, so for all hacked studies

$$1087 V_{i(j)}^* = V_{i(j)}.$$

1088 The remaining non-significant studies are

$$R = (B \cup C) \setminus H.$$

1089 The published non-significant set N is sampled uniformly without re-
1090 placement from R so that $|N| = \ell^*$. Hence

$$\Pr(N = \mathcal{N} \mid R, \mathcal{E}) = \binom{|R|}{\ell^*}^{-1}$$

1091 for every $\mathcal{N} \subseteq R$ such that $|\mathcal{N}| = \ell^*$.

1092 The final published set is

$$\mathcal{P} = S \cup H \cup N, \quad |\mathcal{P}| = k.$$

1093 **Expected effect among originally significant published** 1094 **studies**

1095 Because S is a simple random sample without replacement from A , the
1096 sample average of the originally significant published studies is unbi-
1097 ased for the finite-population average over A . Therefore, conditional
1098 on the realized conducted literature,

$$\mathbb{E} \left(\frac{1}{s} \sum_{i \in S} Y_i \mid A, \mathcal{E} \right) = \frac{1}{a} \sum_{i \in A} Y_i, \quad a > 0.$$

1099 Equivalently, the exact expected effect among originally significant
1100 published studies is

$$\mu_S^{\text{pub}}(D) = \frac{1}{a} \sum_{i \in A} Y_i, \quad a > 0,$$

1101 where D denotes the realized conducted literature.

1102 Passing from the realized literature to the simulation distribution,
 1103 the exact expected effect among originally significant published studies
 1104 is

$$\mu_S^{\text{pub}} = \mathbb{E} \left(\frac{1}{|A|} \sum_{i \in A} Y_i \mid |A| > 0, \mathcal{E} \right).$$

1105 By exchangeability of the K conducted studies, this can also be written
 1106 as

$$\mu_S^{\text{pub}} = \mathbb{E}(Y_1 \mid P_1 \leq 0.05, \mathcal{E}).$$

1107 **Expected reported effect among hacked studies**

1108 Let

$$\bar{Y}_H^* = \frac{1}{h} \sum_{i \in H} Y_i^*$$

1109 denote the average reported effect among hacked studies, defined only
 1110 when $h > 0$. Conditional on the realized conducted literature D and on
 1111 the success event \mathcal{E} , the hacked set H is deterministic, and randomness
 1112 arises only from the uniforms used in the severe hacking step. Hence

$$\mathbb{E}(\bar{Y}_H^* \mid D, \mathcal{E}, h > 0) = \frac{1}{h} \left[\sum_{j=1}^r Y_{i(j)} + \sum_{j=r+1}^h \mathbb{E}(Y_{i(j)}^* \mid D, \mathcal{E}) \right].$$

1113 For $j = r + 1, \dots, h$,

$$Y_{i(j)}^* = \sqrt{V_{i(j)}} F_{t, \text{df}_{i(j)}}^{-1} \left(1 - \frac{U_{i(j)}}{2} \right), \quad U_{i(j)} \sim \text{Unif}(0.01, 0.049),$$

1114 so

$$\mathbb{E}(Y_{i(j)}^* \mid D, \mathcal{E}) = \sqrt{V_{i(j)}} m(\text{df}_{i(j)}),$$

1115 where

$$m(d) = \frac{1}{0.039} \int_{0.01}^{0.049} F_{t,d}^{-1} \left(1 - \frac{u}{2} \right) du.$$

1116 Therefore the exact conditional expectation among hacked studies is

$$\mu_H(D) := \mathbb{E}(\bar{Y}_H^* \mid D, \mathcal{E}, h > 0) = \frac{1}{h} \left[\sum_{j=1}^r Y_{i(j)} + \sum_{j=r+1}^h \sqrt{V_{i(j)}} m(\text{df}_{i(j)}) \right].$$

1117 Averaging over the simulation law yields the exact unconditional ex-
 1118 pectation

$$\mu_H = \mathbb{E}(\mu_H(D) \mid h > 0, \mathcal{E}).$$

1119 Equivalently,

$$\mu_H = \mathbb{E} \left[\frac{1}{h} \left(\sum_{j=1}^{\lceil h/2 \rceil} Y_{i(j)} + \sum_{j=\lceil h/2 \rceil + 1}^h \sqrt{V_{i(j)}} m(\text{df}_{i(j)}) \right) \mid h > 0, \mathcal{E} \right].$$

1120 **Expected reported effect among all significant pub-**
 1121 **lished studies**

1122 The set of significant published studies is

$$G = S \cup H, \quad |G| = s + h = s^*.$$

1123 Its average reported effect is

$$\bar{Y}_G^* = \frac{1}{s^*} \left(\sum_{i \in S} Y_i + \sum_{i \in H} Y_i^* \right).$$

1124 Conditional on D and \mathcal{E} ,

$$\mathbb{E}(\bar{Y}_G^* \mid D, \mathcal{E}) = \frac{1}{s^*} \left[\frac{s}{a} \sum_{i \in A} Y_i + \sum_{j=1}^r Y_{i(j)} + \sum_{j=r+1}^h \sqrt{V_{i(j)}} m(\text{df}_{i(j)}) \right].$$

1125 **Exact expectation and bias of the CMA estimator**

1126 In the simulation, an intercept-only random-effects meta-analysis model
 1127 is fitted to the published sample by REML. Let

$$\mathbf{Y}^* = (Y_1^*, \dots, Y_k^*)^\top, \quad \mathbf{V}^* = (V_1^*, \dots, V_k^*)^\top$$

1128 denote the published effect sizes and within-study variances after the
 1129 publication and p-hacking algorithm has been applied. The published-
 1130 study ordering is immaterial. Define

$$\hat{\tau}_{\text{REML}}^2(\mathbf{y}, \mathbf{v})$$

1131 to be the REML estimate of the between-study variance in the intercept-
 1132 only model returned by the `rma` fit, and define

$$\psi_{\text{REML}}(\mathbf{y}, \mathbf{v}) = \frac{\sum_{j=1}^k \frac{y_j}{v_j + \hat{\tau}_{\text{REML}}^2(\mathbf{y}, \mathbf{v})}}{\sum_{j=1}^k \frac{1}{v_j + \hat{\tau}_{\text{REML}}^2(\mathbf{y}, \mathbf{v})}}.$$

1133 Then the CMA point estimator is exactly

$$\hat{\mu}_{\text{CMA}} = \psi_{\text{REML}}(\mathbf{Y}^*, \mathbf{V}^*).$$

1134 Conditional on the realized conducted literature D and on \mathcal{E} , the
 1135 only remaining randomness comes from the uniform subset S of origi-
 1136 nally significant studies, the uniform subset N of non-significant pub-
 1137 lished studies, and the uniforms used in the severe p-hacking step.
 1138 Hence the exact conditional expectation of the CMA estimator is

$$\mathbb{E}(\hat{\mu}_{\text{CMA}} \mid D, \mathcal{E}) = \frac{1}{\binom{a}{s} \binom{M-a-h}{\ell^*}} \sum_{\substack{S \subseteq A \\ |S|=s}} \sum_{\substack{N \subseteq (B \cup C) \setminus H \\ |N|=\ell^*}} \left[\frac{1}{0.039^{h-r}} \int_{[0.01, 0.049]^{h-r}} \psi_{\text{REML}}(\mathbf{y}^*, \mathbf{v}^*) d\mathbf{u} \right],$$

1139 where

$$\psi_{\text{REML}}(\mathbf{y}^*, \mathbf{v}^*) = \psi_{\text{REML}}(\mathbf{y}^*(D, \mathcal{S}, \mathcal{N}, \mathbf{u}), \mathbf{v}^*(D, \mathcal{S}, \mathcal{N})),$$

1140 $\mathbf{y}^*(D, \mathcal{S}, \mathcal{N}, \mathbf{u})$ is the k -vector of published effect sizes formed by con-
 1141 catenating

$$\{Y_i : i \in \mathcal{S}\}, \{Y_{i_{(j)}} : j = 1, \dots, r\}, \left\{ \sqrt{V_{i_{(j)}}} F_{t, \text{df}_{i_{(j)}}}^{-1} \left(1 - \frac{u_j}{2}\right) : j = r + 1, \dots, h \right\}, \{Y_i : i \in \mathcal{N}\},$$

1142 and $\mathbf{v}^*(D, \mathcal{S}, \mathcal{N})$ is the matching vector of within-study variances,
 1143 which is simply the corresponding collection of original V_i 's.

1144 The exact unconditional expectation of the CMA estimator under
 1145 the simulation is therefore

$$\mathbb{E}(\hat{\mu}_{\text{CMA}} | \mathcal{E}) = \mathbb{E}[\mathbb{E}(\hat{\mu}_{\text{CMA}} | D, \mathcal{E}) | \mathcal{E}].$$

1146 Accordingly, the exact bias of the CMA estimator induced by the
 1147 simulation algorithm is

$$\text{Bias}_{\text{CMA}}^{\text{REML}} = \mathbb{E}(\hat{\mu}_{\text{CMA}} | \mathcal{E}) - \mu.$$

1148 This expression is exact for finite sample representation as defined by
 1149 the algorithm. It is not reducible to a simpler elementary closed form
 1150 without some approximations and inexact simplifications.

1151

Supplementary Material S3: Closed-form approximation for CMA bias

1152

1153

1154

1155

1156

We derive a closed-form approximation for CMA bias under selective publication and p-hacking, using the truncated normal distribution. This avoids the combinatorial complexity of the exact expressions in S2 while closely matching the simulation results (Table 2).

1157

Setup

1158

1159

1160

1161

1162

1163

We use the notation of S2 throughout. Each conducted study i has true effect $\Delta_i \sim \mathcal{N}(\mu, \tau^2)$, sample size N_i , observed effect $Y_i = J_i d_i$ (Hedges' g), and within-study variance $V_i \approx (N_i/n_{Ti}n_{Ci}) J_i^2$. A study is significant when $|Y_i| > Y_{\text{crit},i}$, where $Y_{\text{crit},i} = t_{0.975, N_i-2} \sqrt{V_i}$. As in S2, $s^* = \text{round}(qk)$ studies in the published sample must be significant and $\ell^* = k - s^*$ are non-significant.

1164

Assumptions

1165

We make two simplifications. First, conditional on (Δ_i, N_i) :

$$Y_i \mid \Delta_i, N_i \sim \mathcal{N}(J_i \Delta_i, V_i).$$

1166

1167

1168

1169

Second, we approximate the REML-weighted estimator $\hat{\mu}_{\text{CMA}} = \psi_{\text{REML}}(\mathbf{Y}^*, \mathbf{V}^*)$ (S2) by a simple average of published effects. This is reasonable because the REML weights $1/(V_i + \hat{\tau}^2)$ become approximately uniform as $\hat{\tau}^2$ grows.

1170

Truncated normal expectations

1171

For a study with $\Delta_i = \delta$, $N_i = n$, let $\sigma_i = \sqrt{V_i}$, $\bar{Y}_i = J\delta$, and define

$$\alpha_u = \frac{Y_{\text{crit}} - \bar{Y}_i}{\sigma_i}, \quad \alpha_\ell = \frac{-Y_{\text{crit}} - \bar{Y}_i}{\sigma_i},$$

1172

1173

1174

where ϕ and Φ are the standard normal density and CDF. The expected effect among significant studies (right tail, which dominates for positive effects) is

$$\mathbb{E}[Y_i \mid Y_i > Y_{\text{crit}}, \delta, n] = \bar{Y}_i + \sigma_i \cdot \frac{\phi(\alpha_u)}{1 - \Phi(\alpha_u)},$$

1175

1176

with probability of significance $\pi(\delta, n) = 1 - \Phi(\alpha_u) + \Phi(\alpha_\ell)$. For non-significant studies (doubly-truncated normal):

$$\mathbb{E}[Y_i \mid |Y_i| < Y_{\text{crit}}, \delta, n] = \bar{Y}_i + \sigma_i \cdot \frac{\phi(\alpha_\ell) - \phi(\alpha_u)}{\Phi(\alpha_u) - \Phi(\alpha_\ell)}.$$

1177

Integration over the study population

1178

1179

1180

Since Δ_i and N_i vary across studies, we integrate over $\delta \sim \mathcal{N}(\mu, \tau^2)$ and average over sample-size strata $\{(n_j, w_j)\}$ (midpoints: $n = 25, 50, 85$ with weights 0.25, 0.50, 0.25):

$$\mathbb{E}[Y \mid \text{sig}] = \frac{\int \sum_j w_j \mathbb{E}[Y \mid Y > Y_{\text{crit}}, \delta, n_j] \pi(\delta, n_j) f(\delta) d\delta}{\int \sum_j w_j \pi(\delta, n_j) f(\delta) d\delta},$$

1181

1182

1183

where $f(\delta)$ is the $\mathcal{N}(\mu, \tau^2)$ density. The non-significant expectation $\mathbb{E}[Y \mid \text{non-sig}]$ is analogous. The marginal significance probability is $\bar{\pi} = \int \sum_j w_j \pi(\delta, n_j) f(\delta) d\delta$.

1184

P-hacking components

1185

1186

1187

1188

1189

1190

1191

When $K\bar{\pi} < s^*$, the shortfall $h = s^* - K\bar{\pi}$ studies require p-hacking (S2). The first $r = \lceil h/2 \rceil$ undergo mild hacking (effect Y_i unchanged), the remaining $h - r$ undergo aggressive hacking. Their expected effects are:

Mild: these studies have $p \in (0.05, 0.30]$; their original effects sit between the $p = 0.30$ and $p = 0.05$ boundaries, so the expected original effect is approximated as

$$\mathbb{E}[Y \mid \text{mild}] \approx \frac{1}{2}(Y_{\text{low}} + Y_{\text{crit}}),$$

1192

1193

1194

where $Y_{\text{low}} = t_{0.85, N-2} \bar{\sigma}$ is the $p = 0.30$ boundary.

Aggressive: following S2, $Y_i^* = T_i^* \sqrt{V_i}$ with $T_i^* = F_{t, \text{df}_i}^{-1}(1 - U_i/2)$, $U_i \sim \text{Unif}(0.01, 0.049)$, giving

$$\mathbb{E}[Y^* \mid \text{aggressive}] = \bar{\sigma} \cdot m(\bar{d}), \quad m(d) = \frac{1}{0.039} \int_{0.01}^{0.049} F_{t, d}^{-1}(1 - \frac{u}{2}) du \quad (\text{S2}).$$

1195

Combined approximation

1196

With $s = \min(K\bar{\pi}, s^*)$ originally significant published studies:

$$\mathbb{E}[\hat{\mu}_{\text{CMA}}] \approx \frac{1}{k} \left[s \mathbb{E}[Y \mid \text{sig}] + r \mathbb{E}[Y \mid \text{mild}] + (h - r) \mathbb{E}[Y^* \mid \text{aggressive}] + \ell^* \mathbb{E}[Y \mid \text{non-sig}] \right],$$

1197

$$\text{Bias}_{\text{CMA}}^{\text{approx}} = \mathbb{E}[\hat{\mu}_{\text{CMA}}] - \mu.$$

1198

1199

1200

1201

1202

1203

1204

When $h = 0$ (no p-hacking needed), this simplifies to $\mathbb{E}[\hat{\mu}_{\text{CMA}}] \approx (s^*/k) \mathbb{E}[Y \mid \text{sig}] + (\ell^*/k) \mathbb{E}[Y \mid \text{non-sig}]$.

These expressions allow researchers to obtain a quick estimate of expected CMA bias for their own field-specific parameters—effect size magnitude, heterogeneity, sample sizes, and publication bias severity—without running the full simulation. An accompanying R script (S3.R) is provided in the repository.

1205

Interpretation

1206

1207

1208

1209

1210

1211

The bias is driven by the ratio of the significance threshold Y_{crit} to the true effect. When $Y_{\text{crit}} \gg \mu$ (small effects, small samples), only studies with large positive sampling errors pass the filter, and the uplift $\phi(\alpha_u)/(1 - \Phi(\alpha_u))$ is large. When μ is large enough that most studies are naturally significant, the truncation point falls into the left tail and the uplift vanishes.

1212

1213

1214

1215

Increasing τ^2 widens the distribution of Δ_i , so more studies achieve significance through genuinely large effects rather than sampling luck, attenuating the selection bias. This matches the simulation finding that RDMA's advantage shrinks with heterogeneity.

1216

CI coverage and error bounds

1217

1218

1219

Under selective publication, the CMA estimate is centred on $\mu + \text{Bias}$ while the CI is constructed as if unbiased. With half-width $w = z_{0.975} \cdot \text{SE}$:

$$\text{Coverage} \approx \Phi\left(\frac{w - \text{Bias}}{\text{SE}}\right) - \Phi\left(\frac{-w - \text{Bias}}{\text{SE}}\right).$$

1220

1221

1222

1223

1224

When $\text{Bias} \gg w$, coverage collapses, consistent with 20% coverage for small effects under severe bias. The MAE decomposes as $|\text{Bias}| + \mathbb{E}[|\varepsilon|]$ where $\varepsilon \sim O(\sqrt{(\bar{V} + \tau^2)/k})$. The bias is a structural property of the publication filter and does not diminish with k ; only the random component shrinks. RDMA eliminates the bias term entirely.

1225

Validation

1226

1227

1228

1229

1230

Table 2 compares the approximation with simulation (Table 1, 10,000 iterations). Mean absolute difference: 0.018; maximum: 0.06. The largest gaps occur for small effects under severe bias with low heterogeneity, where p-hacking is heaviest and the midpoint approximation for borderline studies is coarsest.

Table 2: Approximation vs. simulation for $\mathbb{E}[\hat{\mu}_{\text{CMA}}]$

True ES	Het.	Bias	Approx.	Sim.	Diff.
0.2	Low	65%	0.41	0.37	+0.04
0.2	Low	95%	0.55	0.49	+0.06
0.2	Medium	65%	0.33	0.33	+0.00
0.2	Medium	95%	0.50	0.44	+0.06
0.2	High	65%	0.30	0.30	+0.00
0.2	High	95%	0.45	0.42	+0.03
0.5	Low	65%	0.63	0.64	-0.01
0.5	Low	95%	0.83	0.81	+0.02
0.5	Medium	65%	0.58	0.59	-0.01
0.5	Medium	95%	0.81	0.81	+0.00
0.5	High	65%	0.55	0.56	-0.01
0.5	High	95%	0.77	0.78	-0.01
0.8	Low	95%	1.07	1.07	+0.00
0.8	Medium	95%	1.10	1.10	+0.00
0.8	High	95%	1.07	1.08	-0.01