

1 Decline effects, statistical artifacts, and a meta-analytic paradox

2 Richard D. Morey<sup>1</sup>

3 <sup>1</sup> Cardiff University

4 Author Note

5 All code to reproduce this paper can be found at

6 [https://github.com/richarddmorey/decline\\_bias](https://github.com/richarddmorey/decline_bias). This draft was compiled Monday 09 Feb  
7 2026 at 11:24:40 GMT.

8 The authors made the following contributions. Richard D. Morey: Conceptualization,  
9 Writing - Original Draft Preparation, Writing - Review & Editing.

10 Correspondence concerning this article should be addressed to Richard D. Morey, 70  
11 Park Place. E-mail: moreyr@cardiff.ac.uk

12

## Abstract

13 The decline effect (Protzko & Schooler, 2017) is an observed phenomenon where effect sizes  
14 in experiments apparently diminish in size from the first paper demonstrating the effect to  
15 later replications. This has been taken as a symptom of an unhealthy scientific ecosystem,  
16 possibly caused by the “winner’s curse” (selection on significance and regression to the  
17 mean), publication bias or opportunistic analyses. I show that decline effects can arise as  
18 an artifact from a much simpler source: the original article determining the sign of the  
19 effect in a meta-analysis. Moreover, such artifactual decline effects will show correlations  
20 with some of the same experimental properties that one would expect from biases from  
21 poor behavior, such as the sample size of the original study.

22 *Keywords:* decline effect

23 Word count: X

24 Decline effects, statistical artifacts, and a meta-analytic paradox

25 **Introduction**

26 Because scientists often cannot directly observe a whole system, they often make  
27 inferences from phenomena that are deemed signatures of some underlying cause:  
28 morphological similarities for descent from a common ancestor (Darwin, 1859), cosmic  
29 microwave background radiation for the Big Bang (Dicke, Peebles, Roll, & Wilkinson,  
30 1965), or additivity in response times for serial cognitive processing (Sternberg, 1998).  
31 Signatures are consequences of particular ways that a system might work.

32 A weakness of using signatures, however, is that purported signatures may be  
33 observed for reasons other than the putative cause. The worst possible way this could go  
34 wrong is if there is an *artifactual* reason for the signature: one due merely to the manner in  
35 which a phenomenon is studied.

36 Statistical metascience is built on signatures that are used to infer issues in scientific  
37 behaviour. For instance, the observation that effect sizes appear to diminish over time after  
38 an initial discovery—called the “decline effect” (Protzko & Schooler, 2017)—is taken as a  
39 signature of poor scientific behavior (either individual or systemic). Pietschnig, Siegel,  
40 Eder, and Gittler (2019) say that

41 “We show in the present meta-meta-analysis evidence for  
42 overproportional... and stronger effect declines than increases in the published  
43 intelligence literature. Effect misestimations are most likely due to low initial  
44 study power and strategic research and submission behaviors of exploratory  
45 researchers and replicators alike.” (p. 12)

46 They suggest reforms meant to ameliorate the problems they infer.

47 I will show that such signatures can arise as an artifact of the meta-analytic study of  
48 such effects. Seemingly paradoxically, I will show that biased estimates of meta-analytic

49 quantities can arise from unbiased estimates of effect sizes, without any poor behaviour  
50 (assuming the most basic statistical model for the outcomes, every study published, and no  
51 opportunistic reporting). I will then suggest a resolution of the meta-analytical paradox,  
52 proposing that it arises from the stripping of away of the research context in a way that is  
53 central to the meta-scientific perspective.

#### 54 Meta-analytic “coining”

55 “Coining” is the practice of aligning effect sizes by their observed for interpretability.  
56 For instance, Ioannidis (2008) (Figure 2, p. 65) demonstrated how selection for significance  
57 can lead to a relationship between effect size and sample size. He first aligned all observed  
58 effect sizes so they were positive.

59 Coining is also used in the study of decline effects. In this context, the sign of the  
60 effect of an initial study (the earliest one in the set, often the first to report the effect in  
61 question) is flipped to be positive, and then all subsequent studies are re-aligned to be  
62 consistent with that direction (e.g. Fanelli, Costas, & Ioannidis, 2017; Gong & Jiao, 2019;  
63 Nuijten, van Assen, Augusteijn, Crompvoets, & Wicherts, 2020; Pietschnig et al., 2019).  
64 For instance, Fanelli et al. (2017) report in their methods that:

65 “Each dataset was standardized following previously established protocols (14).  
66 Moreover we inverted the sign of (i.e., multiplied by -1) all primary studies  
67 within meta-analyses whose pooled summary estimates were smaller than  
68 zero—a process known as ‘coining.’” (p. 3719)

69 Coining appears to be standard practice in metascience, and may not be explicitly  
70 called “coining”. For instance, Open Science Collaboration (2015) reported that

71 “To be able to compare and analyze correlations across study-pairs, the original  
72 study’s effect size was coded as positive; the replication study’s effect size was

73 coded as negative if the replication study's effect was opposite to that of the  
74 original study." (p. aac4716-3)

75 and reported a "substantial decline" (p. 943) in effect sizes from original to replication (see  
76 also e.g. Camerer et al., 2018; Errington et al., 2021).

77 In the study of decline effects, coining is used to yield a common definition of  
78 decline across studies: a "decline" is when the coined effect size becomes less positive.  
79 Otherwise, a decline of a positive effect (e.g. moving from 1 to .5) might be cancelled by  
80 decline of a negative effect (e.g. moving from -1 to -.5). Coining aligns the signs across all  
81 studies so that many research areas can be studied simultaneously, drawing far-reaching  
82 conclusions across science.

### 83 Data-dependent reporting practices

84 Consider any research context in which there two possible outcomes that can be  
85 supported, and these outcomes are symmetric in the sense that neither outcome is  
86 inherently favored from a reporting perspective. Gender or sex differences are a common  
87 example; one could report a result as favoring men or women.

88 There are different ways to choose how to report a difference in such contexts. If the  
89 research is confirmatory or following up on previous studies that showed a similar effect,  
90 one might choose the to align a pre-registration or a previously reported effect. For  
91 instance, if we predicted that women would score higher on a particular inventory than  
92 men, then when the opposite is observed we might report a negative effect size. Likewise, if  
93 we have a theoretical reason to expect one direction or another, we might choose to report  
94 from the perspective of the theory.

95 But what about novel effects? Another choice is simply to report the effect from the  
96 perspective of the group that scored the highest, i.e. that makes the difference positive.  
97 Johfre and Freese (2021) explicitly recommend "relying on the values of the [differences] to

98 make a decision [on the reference category]. Given that positive numbers are cognitively  
 99 simpler than negative values, the reference category can be chosen such that the presented  
 100 coefficients are positive.” Given that this will have no effect on inferential statistics, a  
 101 data-driven approach is a reasonable choice.

102 I call reporting practices based on outcomes “data-driven reporting practices”  
 103 (DRPs). Note that these are not the same as hypothesizing after the results are known  
 104 (indeed, they may be agnostic to any hypothesis) or cherry picking. Data-driven reports are  
 105 perfectly transparent with respect to the data; they simply choose a *reporting perspective*  
 106 based on the outcome. There is nothing questionable about such reporting *per se*.

107 The practice of coining in meta-analysis is equivalent to re-aligning effect sizes as if  
 108 the original authors had used DRP, and then subsequent authors had reported their effect  
 109 sizes in line with the original study. As a shorthand for both this particular DRP and the  
 110 meta-analytic practice of coining, I will use the term “sign alignment”. For the purpose of  
 111 this paper, it does not matter *who* does it; I will explore its effects on meta-analysis.

112 **Statistical assumptions.** For shorthand, call the experimental context  $A$  and the  
 113 two outcomes/perspectives  $a_1$  and  $a_2$  (e.g.,  $A$  might be “response time in a visual search  
 114 task” and  $a_1$  might be “women/women were faster on average” and  $a_2$  “men/men were  
 115 faster on average”). Of course, there is a difference between the perspective and the  
 116 outcome; which I am referring to will be clear from context.

117 Any experimental outcome  $X_i$  ( $i = 1$  an initial study, and  $i = 2, \dots$  for replications) is  
 118 an unbiased estimate of estimate of some underlying true effect size  $\mu$  with normal error:

$$X_i \sim Normal(\mu, \sigma_i)$$

119 where  $\sigma_i$  is the standard error or estimate  $X_i$ . To begin, we assume that there will be one  
 120 replication, that the initial and replication outcomes are statistically independent of one  
 121 another, that both the initial study and the replication will be reported, and have the same  
 122 true mean. Arbitrarily, we assume that  $\mu$  is the effect size when  $a_1$  is taken as the research

<sup>123</sup> perspective (correspondingly,  $-\mu$  is the effect size when  $a_2$  is taken as the research  
<sup>124</sup> perspective).

<sup>125</sup> We assume that all results are “coined” after being observed. Equivalently, 1) all  
<sup>126</sup> initial studies report their findings by stating the effect size from the perspective that was  
<sup>127</sup> favored by the data, and 2) replications will adopt the perspective of the initial study. The  
<sup>128</sup> initial effect size reported will thus always be positive.

<sup>129</sup> Let us call the coined result  $Y_i$ , which can be thought of as a pair of numbers.<sup>1</sup> For  
<sup>130</sup> reports from initial studies  $Y_1$ ,

$$Y_1 = \begin{cases} (|X_1|, a_1) & \text{if } X_1 > 0 \\ (|X_1|, a_2) & \text{if } X_1 < 0 \end{cases}$$

<sup>131</sup> and for replications (where  $i > 1$ ),

$$Y_i = \begin{cases} (X_2 sgn(X_1), a_1) & \text{if } X_1 > 0, i > 1 \\ (X_2 sgn(X_1), a_2) & \text{if } X_1 < 0, i > 1 \end{cases}$$

<sup>132</sup> where  $sgn x$  is the sign function that returns -1 when  $x < 0$  and 1 when  $x > 0$  (ignoring for  
<sup>133</sup> simplicity the zero probability event that an experimental result is exactly 0).

<sup>134</sup> Whenever  $X_1$  and  $X_2$  disagree in sign, the effect sizes in  $Y_1$  and  $Y_2$  will also disagree  
<sup>135</sup> in sign. The difference between the results  $X_1, X_2$  and reported results  $Y_1, Y_2$  are that  
<sup>136</sup>  $Y_1, Y_2$  have been “aligned” to the outcome of the initial study.

<sup>137</sup> It is obvious that we can obtain an unbiased estimate of  $\mu$  by reversing the sign of the  
<sup>138</sup> effect size report in  $Y_i$  when  $a_2$  is reported, then taking the average with all the reports  
<sup>139</sup> where  $a_1$  is reported. An estimate of  $-\mu$  (the effect from the opposite perspective) could  
<sup>140</sup> found in the same way. There is nothing defective about the DRP from the perspective of

---

<sup>1</sup> In a slight abuse of notation, when I refer to  $Y_i$  by itself in mathematical formulae, assume that it is shorthand for the first number: the reported, sign-aligned effect size.

<sup>141</sup> estimating the underlying effect size, as long as one takes care to understand the  
<sup>142</sup> perspective that each study takes.

<sup>143</sup> **Bias in decline effects**

<sup>144</sup> We now switch to a meta-scientific perspective. Suppose we ask the question: *do*  
<sup>145</sup> *effect sizes tend to decrease from the initial to the replication?* Assume that for every initial  
<sup>146</sup> report  $Y_1$ , we also have a report from a replication  $Y_2$ . In order to assess whether reported  
<sup>147</sup> effects decline, we compare the numerical effect size in the report of  $Y_1$  ( $|X_1|$ ) with the  
<sup>148</sup> numerical effect size in the report of  $Y_2$  ( $X_2 \text{sgn}(X_1)$ ) to produce a “decline effect”  $d_s$  (where  
<sup>149</sup> the  $s$  is for “sign” to indicate the effect has been aligned to the sign of the initial study):

$$\begin{aligned} d_s &= Y_1 - Y_2 \\ &= |X_1| - X_2 \text{sgn } X_1 \end{aligned}$$

<sup>150</sup> If  $d_s > 0$ , this is taken to mean that the effect size has “declined” to some extent from  
<sup>151</sup> the initial study. Although I frame this section as comparing an initial study to a  
<sup>152</sup> replication, my critique here applies to any comparison between two kinds of studies  
<sup>153</sup> (e.g. observational vs. randomized control trials, Franklin et al., 2017).

<sup>154</sup> When both  $X_1$  and  $X_2$  have the same true mean  $\mu$ —that is, there is no true decline  
<sup>155</sup> effect—will the expected decline effect  $d_s$  be equal to 0? Surprisingly, no: The expected  
<sup>156</sup> value of  $d_s$  will *always* be larger than 0 under these circumstances. There is always an  
<sup>157</sup> artifactual bias in the decline effect in precisely the direction that metascientists use as a  
<sup>158</sup> signature.

<sup>159</sup> To see why, consider the situation where the true effect size  $\mu$  is 0. The signs of the  
<sup>160</sup> initial and replication will differ in 1/2 of cases. Anytime this occurs,  $Y_1$  is positive and  $Y_2$   
<sup>161</sup> is negative; hence, a decline is observed. When the signs agree (probability 1/2), the

162 probability that  $Y_1 > Y_2$  is  $1/2$  (because they have the same mean). Thus, the probability  
163 of observing a decline in this scenario is  $1/2 + 1/4 = 3/4$ , despite there being no decline  
164 effect.

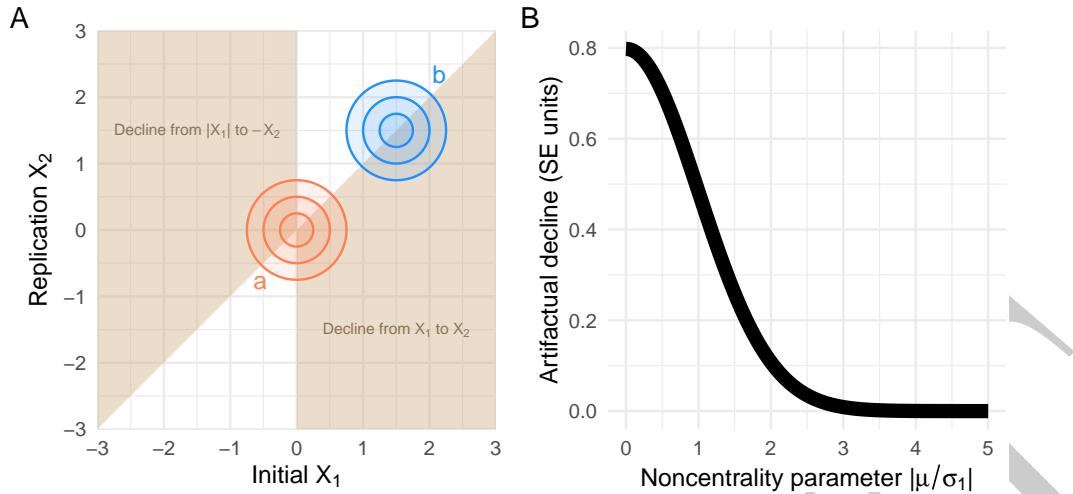
165 Figure 1A depicts the situation graphically for two true effect sizes. The figure shows  
166 the patterns of results that would lead to the identification of a decline effect. As can be  
167 seen, most of the space is occupied by “decline”. When the true effect size is 0 and there is  
168 no decline (“a”),  $3/4$  of the bivariate distribution would count as a “decline”. When the  
169 true effect size is larger (“b”), only about half the bivariate distribution would be in the  
170 “decline” region.

171 The key insight is that the sign alignment happens in one direction or the other with  
172 some probability: whether the result needs to be “coined” is itself random. Unless we have  
173 a specific *perspective* pinned down before the meta-analysis (e.g.,  $a_1$  or  $a_2$ ), the perspective  
174 will thus be random.

175 When there is no decline effect, the expected value of  $d_s$  is  $COV(X_1, sgn X_1)$  (proof  
176 in appendix), which will be a function of how far, in standard error units,  $\mu$  is from 0. The  
177 amount of bias is shown in Figure 1A, as a function of the true effect size in standard  
178 errors of  $X_1$ . The larger the true effect (in either direction), the less probability of the  
179 observed sign differing from the true sign; hence the source of bias diminishes.

## 180 Bias in funnel plots, Egger’s tests, and PET

181 Unlike in decline effect analysis, meta-analysts do not generally use coining. Coining  
182 is used across research context to align all the effect sizes; funnel plots and related tests  
183 (e.g. Egger’s regression test and the precision-effect test, PET) are about looking at a  
184 single research context, so explicit sign alignment is not necessary. However, as I point out,  
185 sign-alignment can also be done by the original authors as a data-dependent reporting  
186 practice. It turns out that sign alignment will bias also funnel plots and the related tests.



*Figure 1.* A: Result combinations that would lead to a conclusion of “decline”, if the error-prone sign of the first experiment determines the interpretation of the effect (shaded regions). Labels “a” and “b” denote two hypothetical bivariate effect size distributions with true effect sizes 0 and 1.5, respectively. See text for explanation. B: Amount of bias in the decline effect as a function of the noncentrality parameter  $\mu/\sigma_1$  of the initial study.

187 To show why sign alignment introduces bias into funnel plots, consider a simple

188 situation in which we perform an initial experiment and a replication. Both experiments  
189 have the same true effect size ( $\mu = 0.20$ ). The initial experiment has a standard error of  
190 0.20 for estimating  $\mu$ , while the replication has 4 times the sample size, yielding a smaller  
191 standard error of 0.10. The initial experiment has a 16% of being observed in the wrong  
192 direction.

193 Consider what happens without sign alignment, but *conditional* on the sign of the

194 initial experiment. Figure 2A depicts the situation graphically. Each experiment/sign  
195 possibility is shown as a distribution. The points show the expectation of the observed  
196 effect size conditional on the sign of the initial experiment. The sizes of the distributions  
197 and points are proportional to the probability of that outcome.

198 Conditional on a positive initial outcome, the initial effect size is slightly

199 overestimated on average. This is exactly counteracted by the underestimation when the

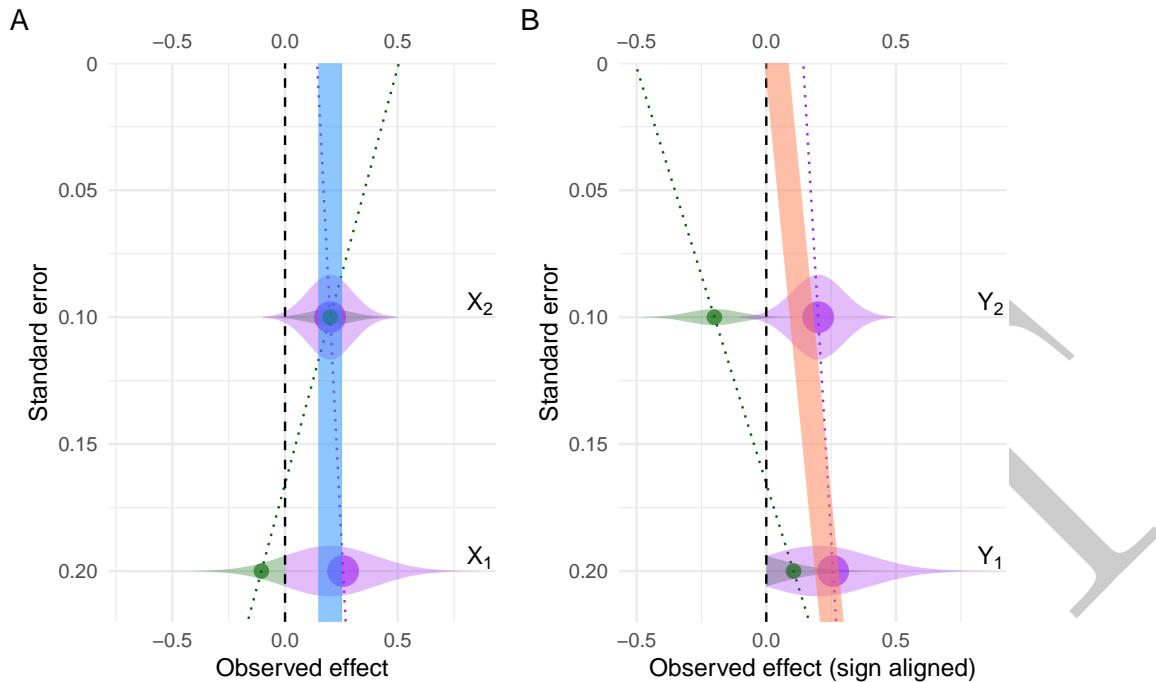
200 initial outcome is negative, weighted by the lower probability of the negative outcome.  
201 With no sign alignment, the sign of the initial experiment has no effect on how the  
202 replication is reported, hence both replications have the same conditional mean: exactly  $\mu$ .

203 Drawing the conditional meta-regression lines (dotted) up to the  $y = 0$ , we can see  
204 that the two lines are biased in opposite directions. When the initial effect is positive, the  
205 slope is slightly negative; when the initial effect is negative, the slope is strongly positive.  
206 The *average* meta-regression line, however (solid, thick blue line), is exactly vertical and  
207 intersects with the  $x$ -axis at exactly  $x = \mu$ . This is the expected behaviour for a funnel plot.

208 Figure 2B shows the situation with sign alignment. Sign alignment does two things:  
209 when the initial experiment has a negative sign, it is flipped to be positive; second,  
210 replication outcome signs are also flipped to be consistent this interpretation. When the  
211 observed effect in the initial experiment is in the correct direction—positive—nothing  
212 changes: the conditional meta-regression line is the same. But when the initial experiment  
213 has a negative outcome, the sign alignment flips the conditional meta-regression line across  
214  $y = 0$ .

215 Both conditional regression lines point in the same direction and they cannot balance  
216 one another out. The average meta-regression line (solid, thick red line) has an  $x$ -intercept  
217 that does not correspond to the true effect size. Moreover, on average it has a negative  
218 slope, which is supposed to be a signature of poor scientific behavior.

219 As the number of studies in the set increases, the bias in the intercept does not  
220 disappear, though the bias in the slope tends to 0 if there is no publication bias. The  
221 source of the bias is the sign alignment of the initial study, and hence all subsequent  
222 studies will be affected in the same way. As with the bias in decline effects, the bias in  
223 funnel plots due to sign alignment decreases for larger effect sizes because the probability  
224 of an initial sign flip gets smaller; see proof in Appendix.



*Figure 2.* Demonstration of the source of sign-alignment bias in funnel plots. See text for explanation.

225

226 In this section we continue using our simple two-experiment (initial/replication)  
 227 setup, but we will vary the sample sizes and true effect sizes. I give details about the  
 228 simulation setup in the Appendix; here, I give an abbreviated version that is enough to  
 229 understand the simulations without burdening the reader with too much formal  
 230 information.

231 Each simulated experiment is assumed to have two independent groups whose mean  
 232 difference will be the effect size of interest. For simplicity, the true variance of each group  
 233 is known to be 1; we can thus treat the mean difference as a standardized effect size. In the  
 234 initial study, the sample size in the two groups are equal and drawn from the distribution  
 235 shown in Figure 3A. For the replication, the sample sizes in the groups were assumed to be  
 236 4 times larger than in the initial study. The true effect sizes for each pair of initial study  
 237 and replication were assumed equal, drawn from the distribution shown in Figure 3B.

238 For each of the  $2$  (initial/replication)  $\times 5000$  simulated studies, an observed difference

239 ( $X$ ) was drawn assuming normal error. These differences were then coined to produce

240 sign-aligned effect sizes ( $Y$ ). The simulation's original/replication simulation setup is

241 similar to the setup of Open Science Collaboration (2015).

242 In the simulations that follow, for each statistic I show that there is no bias if  $X$  is

243 considered; however, using  $Y$  in place of  $X$  will artifactually yield meta-scientific signatures

244 of poor scientific behaviour.

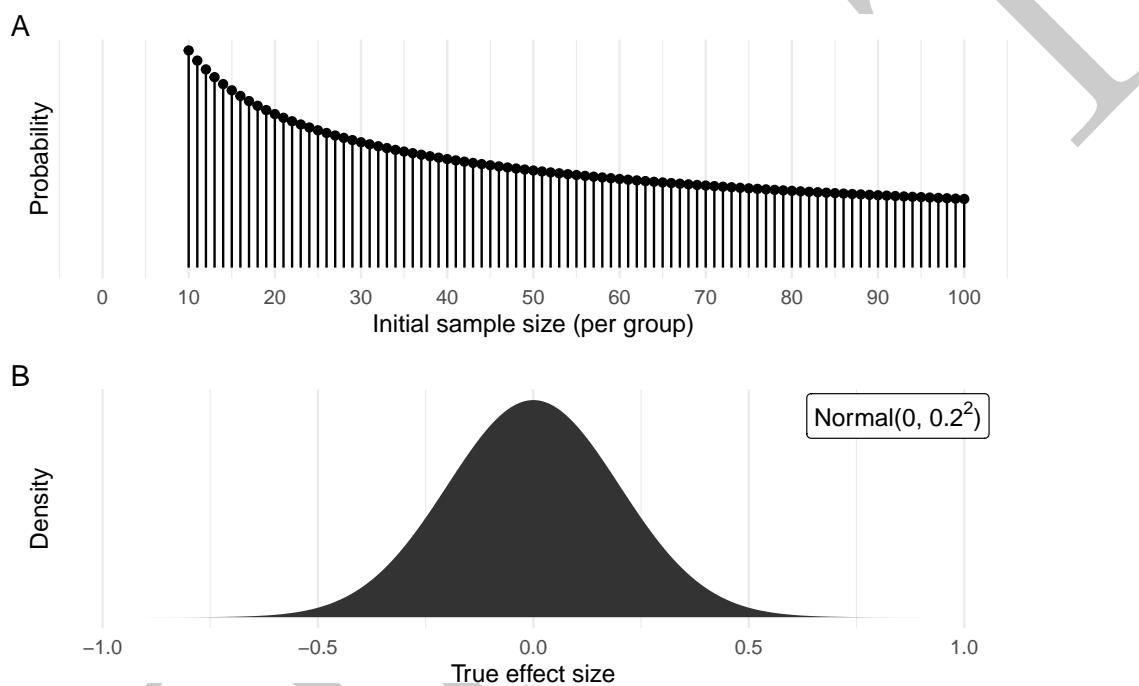


Figure 3. Distributions from which initial per-group sample sizes (A) and true effect sizes (B) were drawn for the reported simulations. All samples were independent.

#### 245 Relationship between effect size and sample size

246 A relationship between effect size and sample size is often cited as evidence for

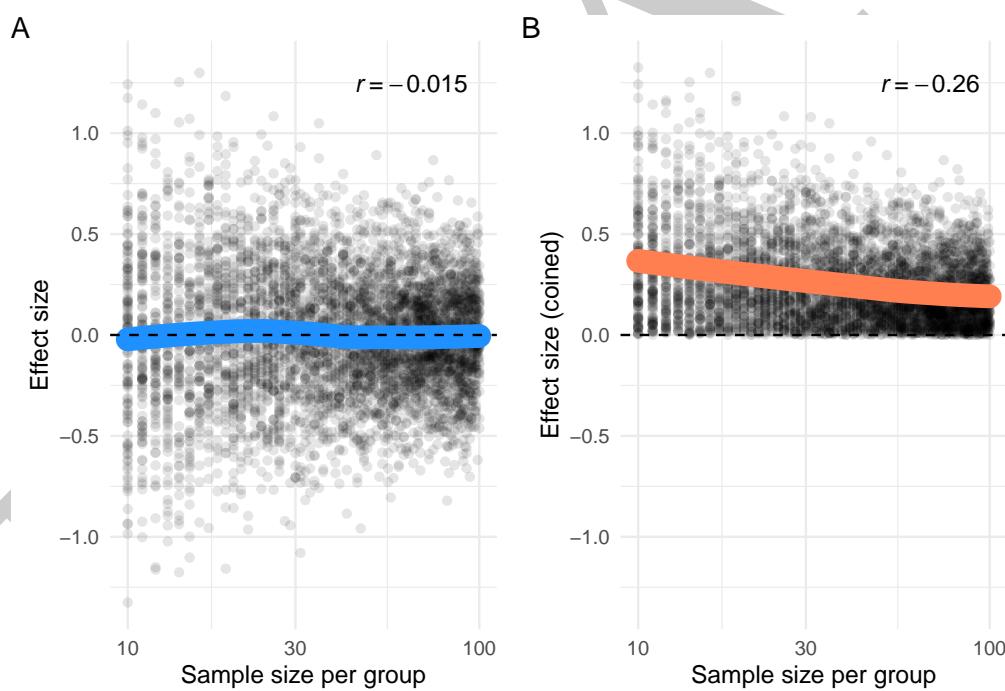
247 publication bias; for example, Stanley, Carter, and Doucouliagos (2018) cited such a

248 correlation as part of the “quite clear” evidence for publication bias, because “inverse

249 correlation between the magnitude of the effect size and sample size would be expected

when there is selective reporting for statistical significance” (p. 1328). Ioannidis (2008) purported to demonstrate the effects of publication bias by showing a relationship between effect size and total sample size in a set of meta-analyses. Likewise, the logic of funnel plots and meta-regression techniques such as Egger’s regression and PET-PEESE depend on finding such relationships.

Figure 4A shows the relationship between the initial effect size  $X_{1j}$  and the sample size in the 5000 simulations. By design, these were sampled to be independent, so there is no relationship between them. Figure 4B, however, shows that once the results are sign-aligned, the relationship appears. Although Ioannidis (2008) intended to show that selection on statistical significance would show this effect, he would have seen the same effect even had he not selected on significance simply due to his explicit coining of the effect sizes.



*Figure 4.* Relationship between the initial observed effect sizes and initial sample sizes for (A) non-sign-aligned results and (B) sign-aligned results. Each point is a single simulated initial result. The thick line in each plot shows the LOESS nonparametric regression curve.

262 **Decline effects**

263 To demonstrate the effect of sign-alignment on estimates of the decline effect, I  
 264 computed the difference between the initial and replication estimates  $X_{1j} - X_{2j}$   
 265 (unaligned) or  $Y_{1j} - Y_{2j}$  (sign-aligned).

266 Figure 5A shows the unaligned differences. Because these simulated results have the  
 267 same mean, as expected, the average difference is 0. However, as Figure 5B shows, when  
 268 the results are sign aligned, the results show artifactual decline effects. Although there is  
 269 no true decline effect in any simulation, the apparent decline effect is large: on average, the  
 270 effect sizes “decline” by over 50%.

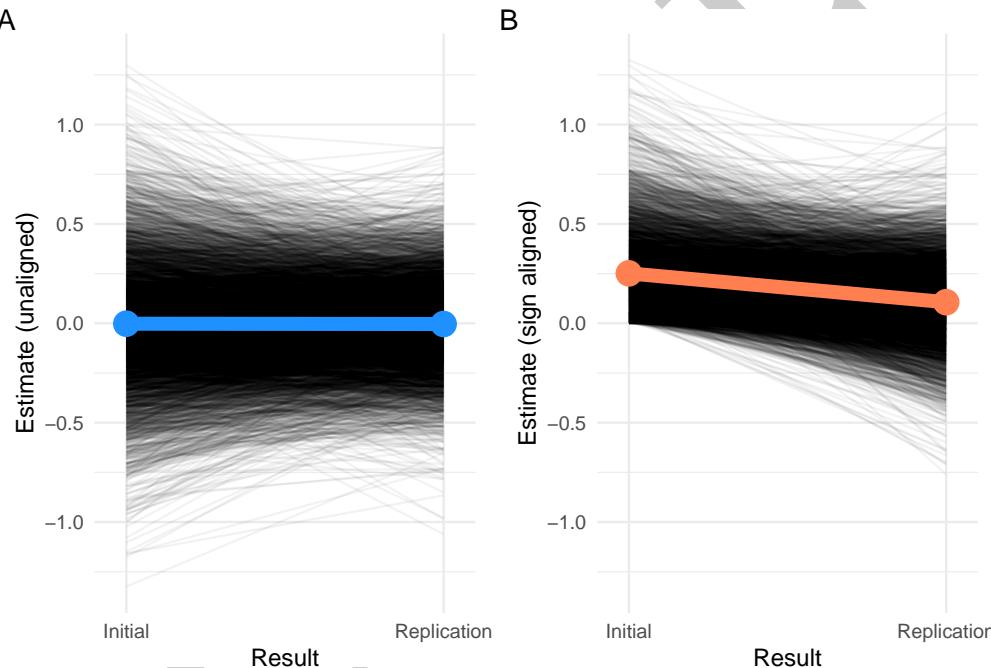
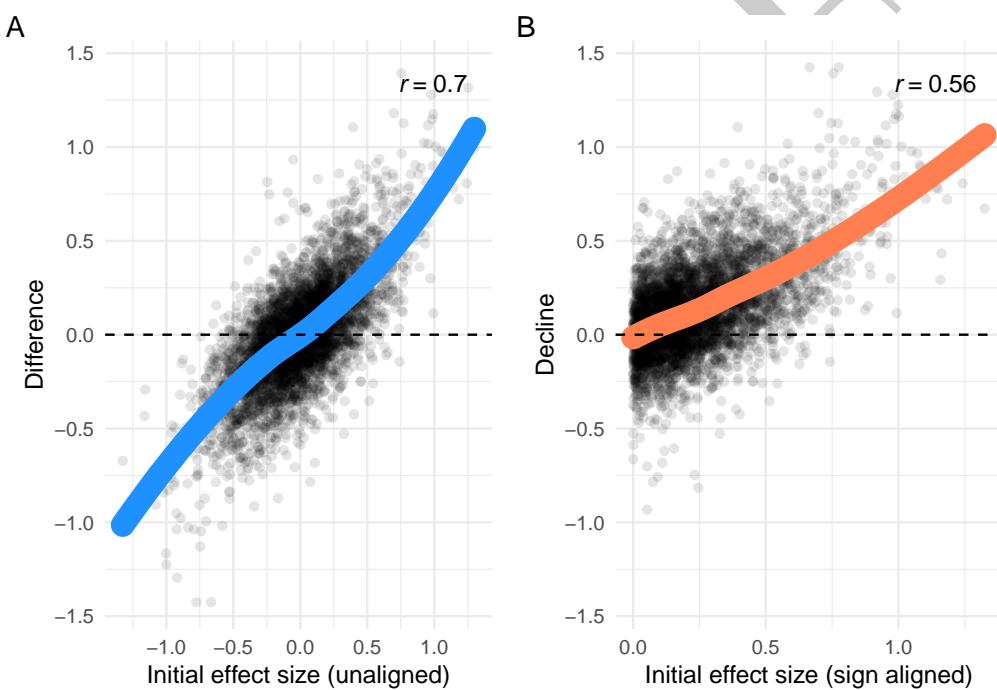


Figure 5. Initial and replication estimates without sign-alignment (A;  $X_{1j}$  and  $X_{2j}$ ) and with sign-alignment (B;  $Y_{1j}$  and  $Y_{2j}$ ). Each thin black line shows a single simulation. The thick line in each graph shows the average.

271 **Correlation between decline effect and sample size.** Meta-analysts look for  
 272 relationships between decline effects and other properties of initial papers to understand  
 273 “risk factors” for bias in literatures. For instance, Pietschnig et al. (2019) report that

<sup>274</sup> “[e]ffect misestimations were more substantial when initial studies had smaller sample sizes  
<sup>275</sup> and reported larger effects, thus indicating suboptimal initial study power as the main  
<sup>276</sup> driver of effect misestimations in initial studies” (p. 1).

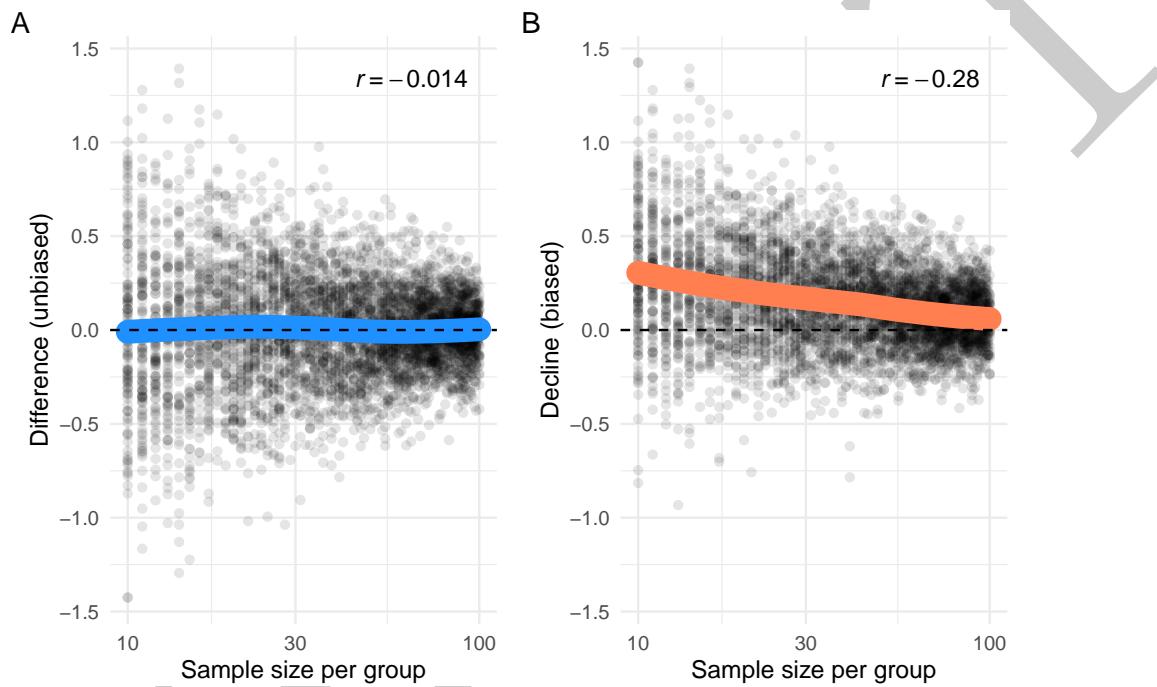
<sup>277</sup> Relationships between decline effects either initial effect sizes or sample sizes can be  
<sup>278</sup> explained artifactually. A relationship between initial effect size and both unaligned or  
<sup>279</sup> aligned decline is explainable by simple regression to the mean. Figure 6A and B show this  
<sup>280</sup> relationship in the simulated data. We replicate the strong relationship Pietschnig et al.  
<sup>281</sup> (2019) with initial effect size without any difference between initial and replication true  
<sup>282</sup> effect sizes. This has relationship has nothing to do with “suboptimal initial study power”;  
<sup>283</sup> it is a mere statistical artifact.



*Figure 6.* Correlation between initial effect size and the difference/decline from initial to replication study in unaligned (A) and sign-aligned (B) results. Each point is the result of a single simulation. The thick line in each plot is a LOESS nonparametric regression curve.

<sup>284</sup> Relationships between initial sample size and the decline effect can be equally  
<sup>285</sup> attributed to a statistical artifact, but the cause is somewhat subtler. As Figure 7A shows,

no such relationship is apparent in the simulated data before sign-alignment; indeed, the data were simulated in such a way that these quantities were independent. However, Figure 7B shows that after sign alignment, the relationship hypothesized and reported by Pietschnig et al. (2019) appears. This artifactual relationship can be attributed to the fact that the bias in computing the decline effect is a decreasing function of the noncentrality parameter (the effect size in standard error units), and the noncentrality parameter is larger when  $N$  is larger.



*Figure 7.* Correlation between sample size and the difference/decline from initial to replication study in unaligned (A) and sign-aligned (B) results. Each point is the result of a single simulation. The thick line in each plot is a LOESS nonparametric regression curve.

### 293 Asymmetric funnel plots

As previously demonstrated, sign-alignment is expected to have the effect of biasing funnel plots to 1) be asymmetric, and 2) have  $x$ -intercepts away from the true effect size even in the absence of publication bias.

297 For the purposes of simulating funnel plots, I added a second simulated replication  
 298 ( $X_{3j}$ ) of the same size as the first replication to prevent the meta-regressions from being  
 299 trivial (only consisting of two points). For visualization purposes, instead of showing the  
 300 individual points, I show the meta-regression line for each simulation, along with the  
 301 average meta-regression line.

302 Meta-regression lines were built from least squares estimates for the model:

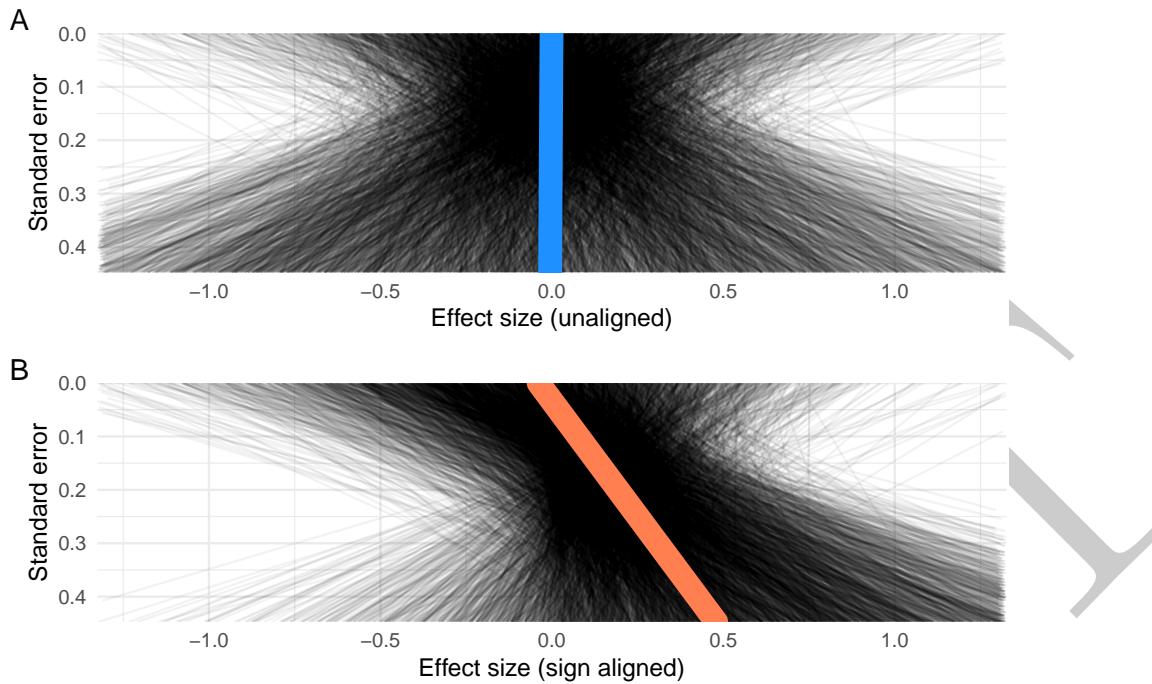
$$X_{ij} = b_{0j} + b_{1j}s_{ij}$$

303 where  $s_{ij}$  is the standard error for the  $i$ th study ( $i = 1, 2, 3$ ) in the  $j$ th simulation.  
 304 Sign-aligned meta-regressions were obtained through corresponding meta-regressions for  
 305  $Y_{ij}$ . For the purposes of this simple example, weighting the points by precision is  
 306 unnecessary.

307 Figure 8A shows the meta-regression lines of standard error onto observed effect size  
 308 for the unaligned data points ( $X_{ij}, j = 1, 2, 3$ ). On average, the funnels are symmetric and  
 309 have an  $x$ -intercept of 0, which is the average effect size in the data.

310 Figure 8B shows the corresponding meta-regression lines for the sign-aligned data.  
 311 The lines are, on average, skewed to have negative slopes. The average  $x$ -intercept is not 0,  
 312 but is instead less than 0. As is clear from the previous discussion (particularly Figure 2),  
 313 this  $x$ -intercept is a function of the true average effect size and the standard errors of the  
 314 replications; it could be shifted right or left by changing the relationship between the initial  
 315 and replication standard errors.

316 We can also assess the bias over all simulations in the  $x$ -intercept as an estimate of  
 317 the true effect size. Figure 9A shows the relationship between the true effect size and error  
 318 in the meta-regression estimate of the effect size for the unaligned results. For the range of  
 319 true effect sizes, the meta-regression estimate appears to be unbiased (though quite  
 320 variable, because we only used three studies).

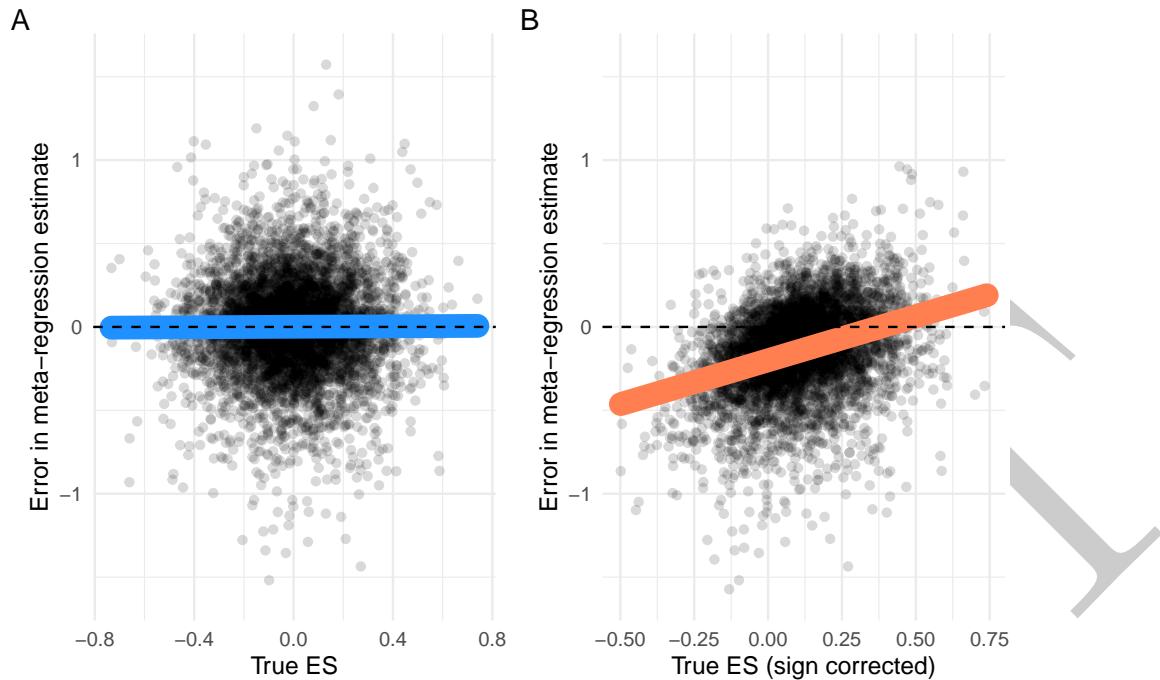


*Figure 8.* Meta-regression lines for each of the simulations using the non-sign-aligned data (A) or sign-aligned data (B).

321        The same is not true of the sign-aligned estimates. After sign-correcting the true  
 322   effect sizes (i.e. when the sign is flipped, we account for the fact that we are estimating  
 323    $-\mu$ ), the sign-corrected estimates are not only biased on average, as shown in Figure 8B;  
 324   Figure 9B shows that bias appears to be an increasing function of the true effect size.

325        Finally, I note that all of the sign-alignment bias noted above is due to the inclusion  
 326   of the initial sign-aligned study in the meta-analysis itself. If we drop the initial study from  
 327   the meta-analysis and consider only the first replication (for instance), the replications are  
 328   unbiased estimates of the true effect size (once the true effect sizes have been suitably  
 329   sign-corrected as well).

330        Figure 10 shows the relationship between the true effect size and observed,  
 331   sign-aligned replication effect size. The bias caused by the initial study and sign-alignment  
 332   has disappeared. This reveals that a data-dependent reporting strategy itself is not a  
 333   problem; the problem is meta-analysts assuming that the effect sizes of the initial studies



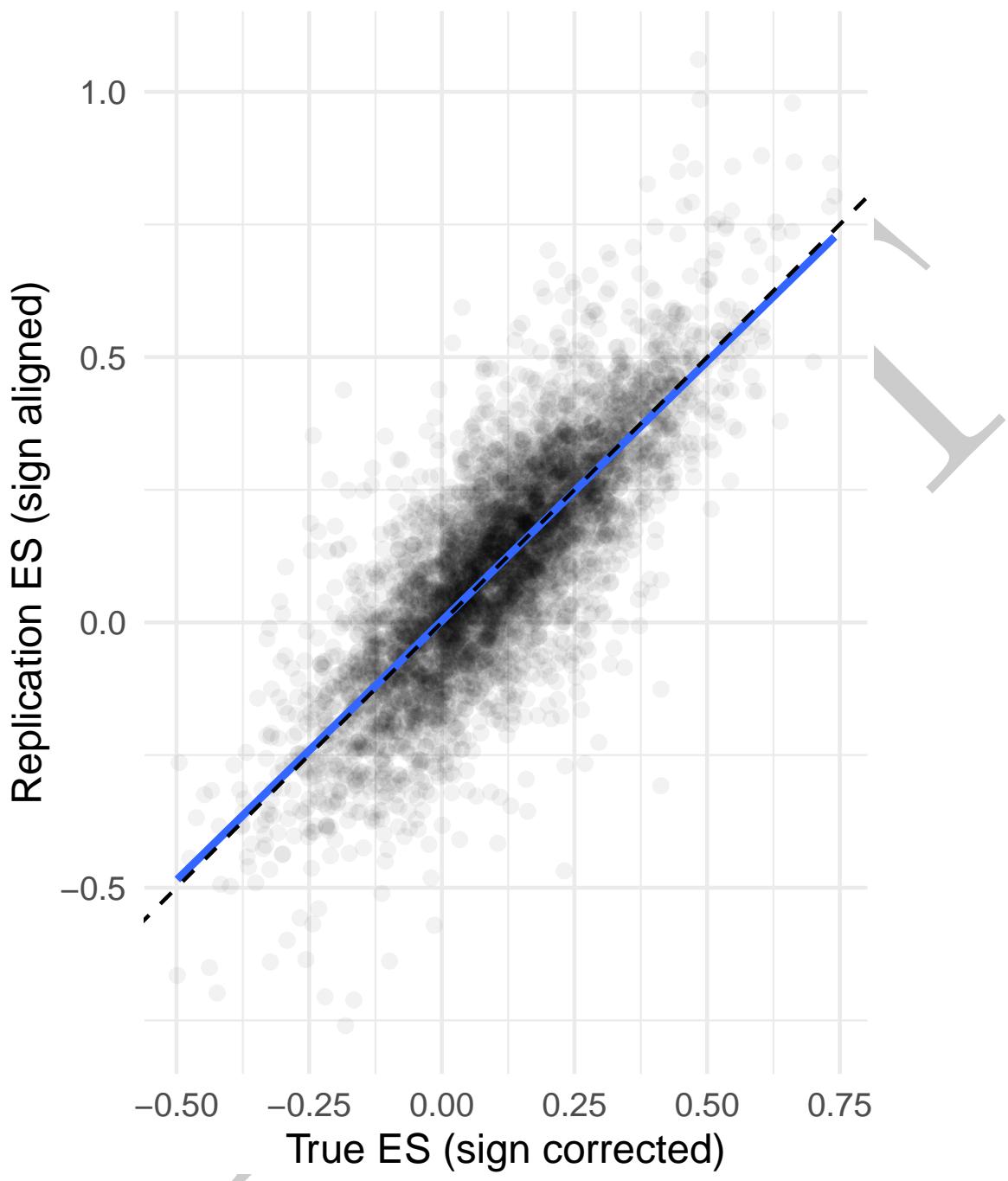
*Figure 9.* Error in the meta-regression estimate of the effect size in the simulations (intercept of the regression of effect size onto standard error, minus the true effect size) as a function of true effect size in non-sign-aligned simulations (A) and sign-aligned simulations (B). In B, the true effect size is corrected when the sign of the initial study was flipped (i.e. the true effect size is  $-\mu$  in these cases).

334 can be modeled in the same way as later replications, when in fact this may not be true.

335 **A meta-analytic paradox**

336 I have shown that sign-alignment leads to analysis artifacts of the sort that one would  
 337 expect under poor scientific behavior. This is clearly undesirable. “Coining” by  
 338 meta-analysts and sign-alignment as a DRP are equivalent, yet I have called the DRP  
 339 “innocuous”. Is this contradictory? Also, how can reports by original authors that I am  
 340 claiming lead to biased differences be themselves unbiased? At first glance this appears to  
 341 be a paradox.

342 The resolution of the paradox must lie in the different goals of those reporting the



*Figure 10.* Relationship between the observed, sign-aligned replication effect size and the true sign-aligned effect size for all simulations. The solid blue line is the least-squares regression line; the dashed black line is the line  $y = x$ . The true effect size is corrected when the sign of the initial study was flipped (i.e. the true effect size is  $-\mu$  in these cases).

343 statistical results and those meta-analysing them. A person reporting an effect size in a  
344 novel experiment has reported the result transparently and accurately even if they  
345 preferentially report it in the positive direction. They have met the goals that someone  
346 reporting results should have.

347 The meta-analyst, on the other hand, has a different goal: to combine results across  
348 studies to come to a well-supported conclusion. By coining—or by assuming that the  
349 results have not been sign-aligned, when in fact they have been—they have undermined  
350 that goal. Their conclusions will be biased.

351 The mistake that the meta-analyst has made is in assuming something that is not  
352 true about the data: that sets of results can be treated as sequences of independent  
353 random variables with symmetric error. Sign-alignment—and perhaps other  
354 DRPs—violates that assumption. The results are not independent because the data from  
355 the first study was used to perform an operation on all studies.

356 But there is a way around the bias induced by sign-alignment: the effect size in an  
357 original report must be interpreted more holistically. In the notation introduced earlier,  
358 when sign alignment happens, the effect size is actually a pair of values  $(|X_1|, a_1)$  - the  
359 first value denoting the effect size, and the second the perspective from which it is  
360 reported. An unbiased estimate of the effect size from either perspective can be obtained  
361 by aligning all effects on  $a_1$  or  $a_2$ . The perspective by which one aligns must be  
362 independent of the data (e.g., from a previous study, a theory, a pre-registration, a  
363 Bayesian prior, or a design asymmetry). Without a fixed perspective, the perspective  
364 becomes random and this is where the trouble begins.

365

## Discussion

366 The bias due to coining in meta-analyses was previously pointed out by Franklin et  
367 al. (2017; see also commentary by Sterne, 2018), but does not appear to be widely known.

368 Currently, Franklin et al. (2017)'s excellent commentary has fewer than twenty citations  
369 according to Scopus, and meta-analysts continue to use the technique. Morrissey (2016)  
370 discusses naïve meta-analyses of absolute values and other nonlinear transformations of  
371 data and points out that meta-analytic “findings” may be mere artifacts of poor statistical  
372 models. Here I show that coining introduces these kinds of artifacts, and DRPs may  
373 introduce them as well.

374 Although the problem with coining (and related data transforms) have been noted  
375 before in other literatures, the point that the innocuous DRPs may cause equivalent effects  
376 appears to be novel. This raises the general possibility that the simple assumptions  
377 underlying most meta-analyses—that effect size reports can be treated as independent  
378 observations with symmetric error around a fixed true value—may often be false in  
379 important ways even when there is no problematic behaviour among scientists.

380 How much does this effect any particular meta-analysis? This is unknown, and I am  
381 not arguing here that this artifact accounts for any particular previous finding. Other  
382 issues (e.g. publication bias) can also cause the effects, and not all research contexts will be  
383 as susceptible to sign alignment. I do, however, take it for granted that meta-analysts  
384 should not use methods *known* to be biased. If a meta-analyst wants to draw an inference  
385 in any research context, the burden of evidence is squarely on them to argue that their  
386 inference is not plagued by artifacts such as those from sign-alignment or other DRPs. The  
387 use and interpretation of methods that may be sensitive to the issues I describe either  
388 implicitly or explicitly assumes that they are *not* a problem.

389 It has been argued that many experiments in psychological science have effect sizes  
390 that are small enough to lead to a high probability that a replication may invert the sign  
391 merely by chance (see e.g. Open Science Collaboration (2015), which reported that 17 of  
392 97 replications showed an opposite sign to the initial study). Small effects with reports  
393 depending on the data will lead to artifacts when studying decline effects. Interestingly, the

394 reformers using the biased methods in meta-analyses are the ones that argue that the  
395 conditions exist that would lead to maximum bias.

396 **Recommendations**

397 **Meta-analysis should never coin effect sizes.** I believe it should go without  
398 saying that meta-analysts should not introduce a source of bias into their meta-analyses,  
399 even if it appears to improve interpretability. Meta-analysts must find a way outside the  
400 data itself to align effect sizes. In some cases this will be straightforward (e.g.,  
401 treatment/control designs or based on theory); in others, it will be less straightforward. In  
402 any case, however, it seems that meta-analysts cannot avoid engaging with the research  
403 contexts of the results they study.

404 **Meta-analysts should not include the initial study in most meta-analyses.**

405 If it is uncertain whether an initial report was affected by DRPs, it is simply safer to  
406 exclude it. The simple act of excluding the initial study will entirely mitigate the effect I  
407 have pointed out, because it is caused by the error in the initial study. Exclusion of the  
408 initial study would also mitigate regression to the mean (sometimes called the “winner’s  
409 curse,” Young, Ioannidis, & Al-Ubaydli, 2008), which seems unavoidable otherwise.

410 **Meta-analysts should seek to understand the role of DRPs in**  
411 **meta-analysis.** As I have shown, reporting practices that have no effect on the inferences  
412 in individual studies may nevertheless have an effect on meta-analyses. If the the simple  
413 practice of reporting a novel effect in the direction in which it is observed can bias  
414 meta-analyses—and this has remained unexplored for decades—this raises the possibility  
415 that other DRPs may effect meta-analyses.

416 As I have previously argued (Morey & Davis-Stober, 2025), it is crucial for the health  
417 of science that meta-scientific tools should have good properties (e.g. not find evidence for  
418 poor behaviour) when scientific behaviour is not problematic. We do not have to agree on

<sup>419</sup> every aspect of good or poor scientific behavior to agree that methods yielding artifactual  
<sup>420</sup> decline effects should not be used. If meta-analysts continue to use biased methods, the  
<sup>421</sup> rhetoric used to attack poor-quality science now can (and will) be used against  
<sup>422</sup> good-quality science in the future.

<sup>423</sup>

## Appendix

### <sup>424</sup> Proofs

<sup>425</sup> **Bias in decline effect.** Let  $X_1$  and  $X_2$  be two independent random variables; the  
<sup>426</sup> mean and variance of  $X_1$  will be denoted  $\mu_1$  and  $\sigma_1^2$ , and likewise for  $X_2$ .  $X_1$  and  $X_2$  are  
<sup>427</sup> assumed to be effect size estimates from initial and a replication experiment.

<sup>428</sup> Let  $s(x)$  denote the sign function  $\operatorname{sgn} x$ . Then the decline effect is a random variable  
<sup>429</sup>  $D_s$  defined as:

$$D_s = s(X_1)X_2 - s(X_1)X_2$$

<sup>430</sup> Let  $s_1$  be  $E[s(X_1)] = 2p - 1$  where  $p = \Pr(X_1 > 0)$ . The expectation of  $D_s$  is

$$\begin{aligned} E[D_s] &= E[s(X_1)X_1 - s(X_1)X_2] \\ &= E[s(X_1)X_1] - E[s(X_1)X_2] \\ &= E[s(X_1)X_1] - s_1\mu_2 \\ &= COV[s(X_1), X_1] + s_1\mu_1 - s_1\mu_2 \\ &= COV[s(X_1), X_1] + s_1(\mu_1 - \mu_2) \end{aligned}$$

<sup>431</sup> Under the assumption of no decline effect ( $\mu_1 - \mu_2 = 0$ ),  $E[D_s] = COV[s(X_1), X_1]$ .<sup>2</sup>

<sup>432</sup> Of course,  $s(X_1)$  and  $X_1$  will be positively correlated, so the decline effect estimate will be

<sup>433</sup> biased, in general (for all  $SE < \infty$ ).

<sup>434</sup> When there *is* a difference between  $\mu_1$  and  $\mu_2$ , and assuming  $X_1$  has a normal

<sup>435</sup> distribution (so  $s_1 = 2\Phi(-\mu_1) - 1$  where  $\Phi$  is the CDF of the normal distribution), the bias

<sup>436</sup> in the estimate of the difference  $D_s$  relative to  $\mu_1 - \mu_2$  will be

$$COV[s(X_1), X_1] - 2(\mu_1 - \mu_2)\Phi(\mu_1)$$

<sup>437</sup> However, one may object to this definition of “bias” when  $\mu_1 \neq \mu_2$  because it does not

<sup>438</sup> capture the logic of “attenuation” (that is, if both effect sizes are negative when  $\mu_1 - \mu_2$  is

<sup>439</sup> positive it does not represent an attenuation; it is an increase).

<sup>440</sup> Another option that better captures attenuation might be; e.g.,

$$d_a = |X_1| - |X_2|$$

<sup>441</sup> as a measure of decline. However, this estimator, too, will suffer from bias because

<sup>442</sup>  $E(|X|) > |E(X)|$  when  $X$  can be negative and there is variability in  $X$ . The more

<sup>443</sup> variability, the greater bias; hence, if  $X_2$  is a more precise estimator than  $X_1$ —which would

<sup>444</sup> often be the case if  $X_2$  arises from a replication or a meta-analysis—the decline effect will

<sup>445</sup> again be overestimated, and will not be 0 when the true decline effect is 0. One may also

<sup>446</sup> object to this definition of decline because extreme sign differences may not be picked up as

<sup>447</sup> problematic (i.e. from -1 to 1 is a large change, but the “decline” as the difference in

<sup>448</sup> absolute values is 0).

---

<sup>2</sup> If one is unwilling to assume independence of  $X_1$  and  $X_2$  (e.g. if  $X_2$  is a meta-analytic effect estimate including  $X_1$ ), the bias when there is truly no decline effect will be  $COV[s(X_1), X_1] - COV[s(X_1), X_2]$ . The first term will dominate when  $X_1$  is a small component of  $X_2$ .

449 Any useful definition of “decline” will likely be more complicated than simple  
 450 arithmetic operations, and may depend on the scientific context. Regardless of how decline  
 451 is defined, I propose the following necessary condition: If a pair of published results  
 452 (original/replication) would be defined as a decline, then the reversed pair (treating the  
 453 replication as the original, and vice versa) should *not* indicate a decline. Sign-alignment  
 454 can violate this simple rule.

455 **Bias in funnel plot meta-regression.** Let  $K \geq 2$  be the total number of points  
 456 in the meta-regression. We assume that the first study ( $X_1$ ) is the one that was used for  
 457 sign alignment, and our vector of observed effect sizes is

$$\mathbf{Y} = [Y_1, Y_2, \dots, Y_K]'$$

458 The expected value of  $\mathbf{Y}$  is

$$E(\mathbf{Y}) = [E(|X_1|), (2p - 1)\mu \mathbf{1}_{K-1}']'$$

459 Let  $e_i$  be the standard error of study  $i$  and  $S = \sum_i (e_i - \bar{e})^2$ . Applying the least squares  
 460 solution the expected slope  $b_1$  and intercept  $b_0$  are:

$$\begin{aligned} E(b_1) &= \frac{e_1 - \bar{e}}{S} (E(|X_1|) - (2p - 1)\mu) \\ E(b_0) &= qE(|X_1|) + (1 - q)(2p - 1)\mu \end{aligned}$$

461 where  $q = 1/K - \bar{e}(e_1 - \bar{e})/S$ . If we assume normality so that

$$E(|X_1|) = \frac{\exp\{-\mu^2/2\}}{\sqrt{\pi/2}} + (2p - 1)\mu,$$

462 we obtain

$$\begin{aligned} E(b_1) &= \frac{e_1 - \bar{e}}{S\sqrt{\pi/2}} \exp\{-\mu^2/2\} \\ E(b_0) &= \frac{q \exp\{-\mu^2/2\}}{\sqrt{\pi/2}} + (2p - 1)\mu \end{aligned}$$

<sup>463</sup> It is obvious that as  $\mu \rightarrow \infty$  and thus  $p \rightarrow 1$ ,  $E(b_1) \rightarrow 0$  and  $E(b_0) - \mu \rightarrow 0$ . Likewise

<sup>464</sup> when  $\mu \rightarrow -\infty$  and thus  $p \rightarrow 0$ ,  $E(b_1) \rightarrow 0$  and  $E(b_0) - (-\mu) \rightarrow 0$  (i.e.,  $E(b_0)$  approaches

<sup>465</sup> the effect size with flipped sign).

<sup>466</sup> Because  $S \rightarrow \infty$  as  $K \rightarrow \infty$ ,  $E(b_1) \rightarrow 0$  and  $E(b_0) \rightarrow (2p - 1)\mu$  also as  $K \rightarrow \infty$ .

<sup>467</sup> Thus the bias in the slope diminishes, but not the bias in the intercept.

## <sup>468</sup> Simulation details

<sup>469</sup> Let  $N_{ij}$  be the sample size per group for the  $i$ th study in the  $j$ th experimental

<sup>470</sup> context ( $i = 1$  for the initial study,  $i > 1$  for replications). We sampled the initial sample

<sup>471</sup> sizes  $N_{1j}$ , then based the replication sample sizes on these:

$$\sqrt{N_{1j}} \stackrel{\text{indep.}}{\sim} \text{Uniform}(\sqrt{10 - 1/2}, \sqrt{100 + 1/2}),$$

$$N_{ij} = 4N_{1j}, i > 1,$$

<sup>472</sup> and  $N_{1j}$  was rounded to the nearest integer. The quantity  $\sqrt{N_{1j}}$  was sampled from a  
<sup>473</sup> uniform distribution so that, once squared, small sample sizes would be more common than  
<sup>474</sup> larger ones. Replication sample sizes were assumed to be 4 times larger than the initial  
<sup>475</sup> studies, though this does not matter much because the bias in the decline effect is only a  
<sup>476</sup> function of the initial study properties.

<sup>477</sup> The standardized effect size  $\mu_j$  for all studies in an experimental context was assumed  
<sup>478</sup> to be the same:

$$\mu_j \stackrel{\text{indep.}}{\sim} \text{Normal}(0, 0.20^2).$$

<sup>479</sup> A mean effect size of 0 was chosen to respect the symmetry in the assumption that either  
<sup>480</sup> group could act as a reference.

<sup>481</sup> The effect size  $X_{ij}$  before sign alignment was then sampled from a Normal with mean

<sup>482</sup>  $\mu_j$  and standard error  $\sqrt{2/N_{ij}}$ :

$$X_{ij} \stackrel{i\text{ndep.}}{\sim} \text{Normal}(\mu_j, 2/N_{ij}).$$

<sup>483</sup> Sign-aligned effect sizes were then computed:

$$Y_{ij} = \begin{cases} |X_{ij}| & i = 1 \\ X_{ij} \text{sgn } X_{1j} & i > 1 \end{cases}$$

484

## References

- 485 Camerer, C. F., Dreber, A., Holzmeister, F., Ho, T.-H., Huber, J., Johannesson, M., ...  
486 Wu, H. (2018). Evaluating the replicability of social science experiments in Nature and  
487 Science between 2010 and 2015. *Nature Human Behaviour*, 2(9), 637–644.  
488 <https://doi.org/10.1038/s41562-018-0399-z>
- 489 Darwin, C. (1859). *On the Origin of Species by Means of Natural Selection*. John Murray.
- 490 Dicke, R. H., Peebles, P. J. E., Roll, P. G., & Wilkinson, D. T. (1965). Cosmic Black-Body  
491 Radiation. *The Astrophysical Journal*, 142, 414–419. <https://doi.org/10.1086/148306>
- 492 Errington, T. M., Mathur, M., Soderberg, C. K., Denis, A., Perfito, N., Iorns, E., & Nosek,  
493 B. A. (2021). Investigating the replicability of preclinical cancer biology. *eLife*, 10,  
494 e71601. <https://doi.org/10.7554/eLife.71601>
- 495 Fanelli, D., Costas, R., & Ioannidis, J. P. A. (2017). Meta-assessment of bias in science.  
496 *Proceedings of the National Academy of Sciences*, 114(14), 3714–3719.  
497 <https://doi.org/10.1073/pnas.1618569114>
- 498 Franklin, J. M., Dejene, S., Huybrechts, K. F., Wang, S. V., Kulldorff, M., & Rothman, K.  
499 J. (2017). A Bias in the Evaluation of Bias Comparing Randomized Trials with  
500 Nonexperimental Studies. *Epidemiologic Methods*, 6(1), 20160018.  
501 <https://doi.org/10.1515/em-2016-0018>
- 502 Gong, Z., & Jiao, X. (2019). Are Effect Sizes in Emotional Intelligence Field Declining? A  
503 Meta-Meta Analysis. *Frontiers in Psychology*, 10.  
504 <https://doi.org/10.3389/fpsyg.2019.01655>
- 505 Ioannidis, J. P. A. (2008). Why Most Discovered True Associations Are Inflated.  
506 *Epidemiology*, 19(5), 640–648. Retrieved from <https://www.jstor.org/stable/25662607>
- 507 Johfre, S. S., & Freese, J. (2021). Reconsidering the Reference Category. *Sociological  
508 Methodology*, 51(2), 253–269. <https://doi.org/10.1177/0081175020982632>
- 509 Morey, R. D., & Davis-Stober, C. P. (2025). On the Poor Statistical Properties of the  
510 P-Curve Meta-Analytic Procedure. *Journal of the American Statistical Association*,

- 511         $\theta(0)$ , 1–13. <https://doi.org/10.1080/01621459.2025.2544397>
- 512        Morrissey, M. B. (2016). Meta-analysis of magnitudes, differences and variation in  
513        evolutionary parameters. *Journal of Evolutionary Biology*, 29(10), 1882–1904.  
514        <https://doi.org/10.1111/jeb.12950>
- 515        Nuijten, M. B., van Assen, M. A. L. M., Augusteijn, H. E. M., Crompvoets, E. A. V., &  
516        Wicherts, J. M. (2020). Effect Sizes, Power, and Biases in Intelligence Research: A  
517        Meta-Meta-Analysis. *Journal of Intelligence*, 8(4), 36.  
518        <https://doi.org/10.3390/jintelligence8040036>
- 519        Open Science Collaboration. (2015). Estimating the reproducibility of psychological  
520        science. *Science*, 349(6521), 943.
- 521        Pietschnig, J., Siegel, M., Eder, J. S. N., & Gittler, G. (2019). Effect Declines Are  
522        Systematic, Strong, and Ubiquitous: A Meta-Meta-Analysis of the Decline Effect in  
523        Intelligence Research. *Frontiers in Psychology*, 10.  
524        <https://doi.org/10.3389/fpsyg.2019.02874>
- 525        Protzko, J., & Schooler, J. W. (2017). Decline effects: Types, mechanisms, and personal  
526        reflections. In *Psychological science under scrutiny: Recent challenges and proposed*  
527        *solutions* (pp. 85–107). Hoboken, NJ, US: Wiley Blackwell.  
528        <https://doi.org/10.1002/9781119095910.ch6>
- 529        Stanley, T. D., Carter, E. C., & Doucouliagos, H. (2018). What meta-analyses reveal about  
530        the replicability of psychological research. *Psychological Bulletin*, 144(12), 1325–1346.  
531        <https://doi.org/10.1037/bul0000169>
- 532        Sternberg, S. (1998). *Discovering Mental Processing Stages: The Method of Additive*  
533        *Factors*. <https://doi.org/10.7551/mitpress/3967.003.0017>
- 534        Sterne, J. (2018). Commentary: Does the selective inversion approach demonstrate bias in  
535        the results of studies using routinely collected data? *BMJ*, 362, k3259.  
536        <https://doi.org/10.1136/bmj.k3259>
- 537        Young, N. S., Ioannidis, J. P. A., & Al-Ubaydli, O. (2008). Why Current Publication

- <sup>538</sup> Practices May Distort Science. *PLOS Medicine*, 5(10), e201.
- <sup>539</sup> <https://doi.org/10.1371/journal.pmed.0050201>

DRAFT