

Boise State University

ScholarWorks

---

Kinesiology Faculty Publications and  
Presentations

Department of Kinesiology

---

4-2023

## On the Reproducibility of Power Analyses in Motor Behavior Research

Brad McKay  
*McMaster University*

Mariane F. B. Bacelar  
*Boise State University*

Michael J. Carter  
*McMaster University*

1      On the reproducibility of power analyses in motor behavior research

Abstract

2

3 Recent metascience suggests that motor behavior research may be underpowered, on average.  
4 Researchers can perform *a priori* power analyses to ensure adequately powered studies.  
5 However, there are common pitfalls that can result in underestimating the required sample  
6 size for a given design and effect size of interest. Critical evaluation of power analyses  
7 requires successful analysis reproduction, which is conditional on the reporting of sufficient  
8 information. Here we attempted to reproduce every power analysis reported in articles ( $k =$   
9 84/635) in three motor behavior journals between January 2019 and June 2021. We  
10 reproduced 7% of analyses using the reported information, which increased to 43% when we  
11 assumed plausible values for missing parameters. Among studies that reported sufficient  
12 information to evaluate, 63% reported using the same statistical test in the power analysis as  
13 in the study itself, and in 77% the test addressed at least one of the identified hypotheses.  
14 Overall, power analyses were not commonly reported with sufficient information to ensure  
15 reproducibility. A non-trivial number of power analyses were also affected by common  
16 pitfalls. There is substantial opportunity to address the issue of underpowered research in  
17 motor behavior by increasing adoption of power analyses and ensuring reproducible  
18 reporting practices.

19 *Keywords:* Motor learning, Motor control, Sample size planning, Metascience,  
20 Reproducibility

21 In statistics, power is the probability of observing a significant effect given the  
22 statistical analysis, sample size, and the true effect size in the population. Recent evidence  
23 suggests that many studies in sports science and motor behavior have been underpowered to  
24 reliably detect the effects researchers are investigating. For example, Mesquida et al. (2022)  
25 estimated the average power of studies sampled from the *Journal of Sports Sciences* to be  
26 48%, albeit with substantial uncertainty. Similarly, Lohse et al. (2016) reported evidence  
27 that motor learning experiments sampled from seven motor behavior journals between 2012  
28 and 2014 were likely underpowered; estimating an average power between 21% and 57%.  
29 Meta-analyses of specific motor learning phenomena have also found evidence of low power  
30 among studies. For example, the average power of experiments ( $k = 75$ ) that compared a  
31 reduced frequency of feedback to a 100% frequency was estimated to be 27%, again with  
32 substantial uncertainty (McKay, Hussien, et al., 2022). Even lower average power estimates  
33 of 6% were reported in meta-analyses of enhanced expectancies (Bacelar et al., 2022; McKay,  
34 Bacelar, et al., 2022) and self-controlled practice (McKay, Yantha, et al., 2022), with an  
35 upper bound estimate of 13%. Despite having a low probability of observing a significant  
36 result *a priori*, positive results in these literatures have been much more frequent than  
37 expected. In fact, the overall positivity rates in exercise and sport science publications in  
38 general, and motor behavior publications specifically, have been estimated between 81%  
39 (Twomey et al., 2021) and 84% (McKay et al., in-press). When individual studies are  
40 unlikely to observe positive results and the published literature is unlikely to contain  
41 negative results, the estimates contained in the published literature are likely to be biased  
42 (Carter et al., 2015; Gelman & Carlin, 2014; Maier et al., 2022). This bias can result in  
43 exaggerated estimates, the appearance of an effect when there is none, or even results in the  
44 wrong direction. Therefore, the combination of low power and selective reporting of positive  
45 results will severely undermine the credibility of a scientific literature.

46 Researchers can design studies with a high probability of observing informative  
47 results (Cohen, 1988; Lakens, 2022). If a study is designed to have 95% power to detect the

48 smallest effect a researcher is interested in, then 95% of the time the researcher will detect  
49 the effect if it is real. If the researcher fails to observe a significant result, they can rule out  
50 effects as large or larger than their smallest effect of interest with an error rate of  $1 - power$ ,  
51 or 5% in this example. Power analysis is therefore a critical tool for designing informative  
52 studies and numerous open-source software packages are available to researchers, including  
53 but not limited to G\*Power (Faul et al., 2009), **Superpower** (Lakens & Caldwell, 2021), and  
54 PANGEA (Westfall, 2015). Despite the widespread availability of power analysis software,  
55 power analyses are not typically reported in sports science research (Abt et al., 2020; Borg et  
56 al., 2022; McCrum et al., 2022; McKay et al., in-press; Robinson et al., 2021; Twomey et al.,  
57 2021). In motor behavior specifically, only 13% of all studies published in *Human Movement*  
58 *Science*, the *Journal of Motor Learning and Development*, and the *Journal of Motor*  
59 *Behavior* between 2019 and June 2021 included a power analysis (McKay et al., in-press). It  
60 is perhaps not surprising that power analyses are uncommon given the low average power  
61 estimates in the literature. However, we argue that this presents an opportunity to the field;  
62 by increasing the use of appropriate and reproducible power analyses, we can improve the  
63 overall reliability of our literature.

64         Conducting a power analysis can be a straightforward task, but new power analysts  
65 may fall victim to some common traps. Each power analysis requires specifying the primary  
66 hypothesis, the effect of interest, the statistical test, and choosing acceptable Type 1 (false  
67 positive) and Type 2 (false negative) error rates. For power calculations to be accurate and  
68 appropriate, it is crucial that the design included in the power analysis addresses the effect  
69 predicted by the primary hypothesis. For example, a study might include both within and  
70 between subject components, but the primary hypothesis may pertain to between subject  
71 differences. In this case, a power analysis based on the within-subjects analysis will  
72 dramatically overestimate the power of the study with respect to the primary hypothesis. It  
73 is also important that the statistical analysis used in the power analysis match that used on  
74 the raw data, otherwise the power calculations can be inaccurate. For example, parametric

75 and non-parametric approaches tend to differ in their power, so it is important that the same  
76 method that will be applied to the data is included in the power analysis. Choice of software  
77 to conduct a power analysis is also important, as different designs may require different  
78 software. For instance, G\*Power cannot, accurately calculate power for mixed factorial  
79 designs that include three or more levels of the within-subjects factor. While other packages,  
80 such as Superpower, can handle this more complex design, there are many possible designs  
81 that will require simulation-based approaches and likely consultation with a statistician. For  
82 example, power analysis for mixed-effects models can be conducted via simulation with the R  
83 package faux (DeBruine et al., 2021), and power analysis for mediation analyses can be  
84 conducted with the package powerMediation (Qiu, 2021). Each of these common pitfalls  
85 can result in conducting an underpowered study, or (less likely) an inefficient study.

86 Despite the challenges, power analyses can be reproduced quickly and independent of  
87 the final data. This provides collaborators (and even peer reviewers in a registered report)  
88 the opportunity to easily verify and, if necessary, correct a power calculation to ensure an  
89 adequately powered and informative study. Peer reviewers of standard reports can at least  
90 ensure that an accurate power calculation is reported in the final manuscript. While power  
91 analyses can include errors that result in underpowered designs, if reported in a reproducible  
92 fashion, these errors can be caught in time to ensure a better outcome. As a means of  
93 improving the reliability and transparency of the literature, requiring a power analysis for  
94 publication is as easy to implement as simply enforcing the guidelines at most journals.  
95 McKay and colleagues (in-press) reported that 13% of studies in three motor behavior  
96 journals included a power analysis; yet, all three of the journals required a power analysis in  
97 their author guidelines. If power analyses are reported with sufficient information to  
98 reproduce the results, we believe that increasing the adoption of power analyses has the  
99 potential to improve the state of the literature in the long term. However, the largest  
100 benefits to increased usage of reproducible power analyses would likely be seen in  
101 preregistered studies or Registered Reports. Otherwise, power analyses may be conducted

102 post-hoc, limiting (but not eliminating) their usefulness.

103         The goal of this study was to evaluate the reproducibility of power analyses reported  
104 in the motor behavior literature between 2019 and 2021. We attempted to reproduce each  
105 power analysis identified by McKay et al. (in-press) to determine potential areas for  
106 improvement and identify common pitfalls in power analysis reporting. For power analyses  
107 to improve study design, researchers need to conduct them. We have already described  
108 research showing this has not commonly been the case. Power analyses also need to be  
109 conducted properly, but to understand if that is the case, they need to be reported in a  
110 reproducible fashion. Here we sought to answer five preregistered research questions. First,  
111 what proportion of power analyses reported in motor behavior research can be reproduced  
112 using only the information reported in the article or shared as supplementary information?  
113 Second, what proportion of power analyses can be reproduced conditional on making  
114 assumptions for missing parameters in the study article? Third, in what proportion of  
115 studies does the statistical test used in the power analysis match the design used in the data  
116 analysis? Fourth, in what proportion of studies does the design used in the power analysis  
117 address the prediction made by the primary hypothesis? And fifth, what proportion of  
118 studies that used partial eta-squared as the effect size parameter in a power analysis  
119 conducted in G\*Power used the default partial eta-squared settings?

120

## Methods

121         The preregistration, data, and code for this study can be found using either of these  
122 links: [https://github.com/cartermaclab/proj\\_power-reproducibility-motor-behaviour](https://github.com/cartermaclab/proj_power-reproducibility-motor-behaviour) or  
123 <https://osf.io/9a6m8/>

## 124 Piloting

125         We piloted our reproduction and extraction procedures on six papers, two from each  
126 publication year in the sample (2019-2021). During piloting we developed our methods to  
127 account for the diversity of study types and reporting practices we anticipated encountering.

128 The most influential adjustment made during piloting was the removal of a planned code for  
129 the number of primary hypotheses. There was often enough ambiguity about hypothesis  
130 priority that consensus felt arbitrary, so we opted to treat all hypotheses as primary.

### 131 **Sample**

132 The 84 power analyses examined were from studies identified by McKay et al.  
133 (in-press). Inclusion in that project required: a) publication in *Human Movement Science*,  
134 the *Journal of Motor Learning and Development* or the *Journal of Motor Behavior*, b)  
135 published between January 2019 and June 2021, and c) a hypothesis test, including the null.  
136 A total of 635 studies met those inclusion criteria, of which 84 reported a power analysis.

### 137 **Power Analysis Reproduction and Data Extraction**

138 The first and second authors attempted to conduct the power analysis reported in  
139 each study using G\*Power 3.1 (Faul et al., 2009). Although other means of calculating power  
140 are available, all studies in the sample either reported using G\*Power or did not report the  
141 software they used. The authors began by attempting to calculate the power using the  
142 parameters that were reported in the paper. A power analysis was fully reproducible if the  
143 sample size calculation could be confirmed using the reported parameters. If insufficient  
144 parameters were explicitly reported, which was typical, the authors recorded that the power  
145 analysis was not reproducible from the description of the analysis alone. When a study was  
146 not immediately reproducible, we attempted making assumptions for missing parameters.  
147 For example, if the statistical analysis was not reported, we tried assuming the actual  
148 analyses reported in the results section of the study. *All plausible analyses were attempted,*  
149 *but effect size, power, and alphas other than .05 were not guessed.* Studies that could not be  
150 reproduced by assuming parameters were recorded as not reproducible, otherwise they were  
151 considered conditionally reproducible.

152 If the statistical analysis used in the power analysis was reported in a study, it was  
153 assessed whether the analysis tested any of the study's hypotheses. For example, it might be

154 hypothesized that two groups will differ on a measure that is taken twice. If the power  
155 analysis was conducted for the within-subject effect of time, the analysis did not match the  
156 hypothesis. We recorded quotes of the hypotheses from each study and if multiple  
157 hypotheses were made all were considered. We also evaluated whether the analysis used in  
158 the power analysis was consistent with the analysis used in the study. If a *t*-test was used in  
159 the power analysis but an ANOVA was used in the study, the analyses did not match. All  
160 the main analyses reported in a study were considered.

161 Two software considerations were probed during data collection. First, we recorded  
162 whether the software used to conduct the power analysis was appropriate for the type of  
163 analysis. Second, if partial eta-squared was used in G\*Power, we recorded the setting  
164 required to reproduce the power analysis if it was reproducible.

165 The first and second authors met frequently throughout data collection to discuss the  
166 extracted studies and resolve coding conflicts. There were a wide range of study designs,  
167 hypotheses, and reporting language in the sample, so meeting frequently ensured consistency  
168 and allowed for quick updating of policies when faced with unexpected scenarios. Power  
169 analyses could be reproduced quickly when reporting was clear (1 to 4 minutes), but it could  
170 take much longer when reporting was unclear (15 to 30 minutes).

## 171 Data Analysis

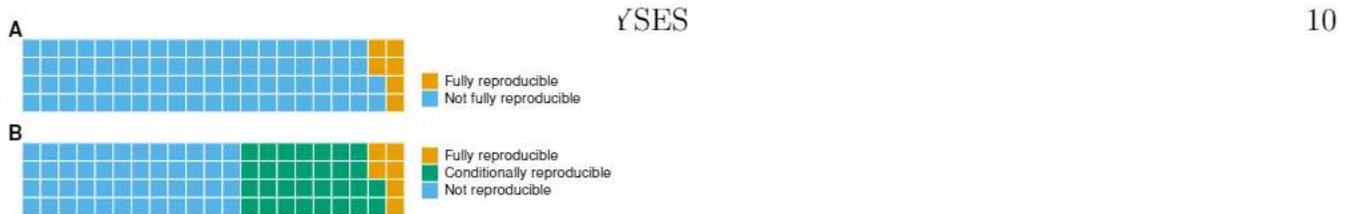
172 Each research question was addressed descriptively by calculating proportions. All  
173 analyses were conducted using R (Version 4.1.2; R Core Team, 2021) and the R-packages  
174 *daff* (Version 0.3.5; Fitzpatrick et al., 2019), *extrafont* (Version 0.18; Chang, 2022), *papaja*  
175 (Version 0.1.1; Aust & Barth, 2020), *renv* (Version 0.15.5; Ushey, 2022), *tidyverse* (Version  
176 1.3.1; Wickham et al., 2019), and *waffle* (Version 1.0.1; Rudis & Gandy, 2019) were used in  
177 this project.

## Results

178

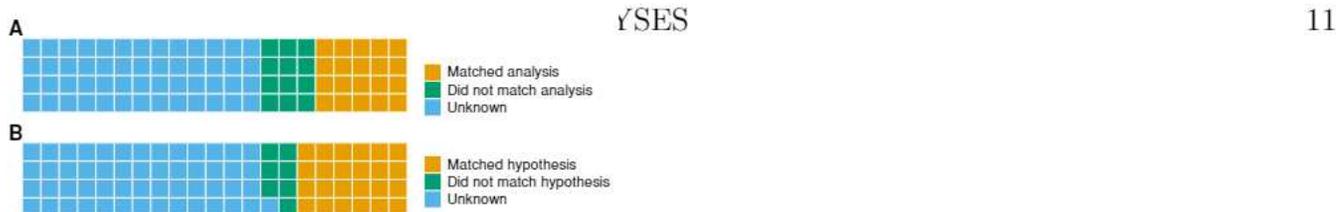
### 179 Preregistered Analyses

180 Of the 84 power analyses reported in 83 articles, 7% ( $n = 6$ ) were fully reproducible  
181 (see Figure 1A) and 36% ( $n = 30$ ) were conditionally reproducible (see Figure 1B). The  
182 statistical test used in the power analysis matched the one used in the data analysis in 24%  
183 of the power analyses ( $n = 20$  experiments), did not match in 14% ( $n = 12$  experiments),  
184 and in the remaining 62% ( $n = 52$  experiments) the statistical test used in the power analysis  
185 could not be accurately identified, precluding an assessment of the congruence between  
186 power analysis design and data analysis design (see Figure 2A). The design used in the  
187 power analysis addressed the experiment's hypothesis in 23% of the experiments ( $n = 19$ ), at  
188 least one of the hypotheses in 6% of the experiments ( $n = 5$ ), none of the hypotheses in 8%  
189 of the experiments ( $n = 7$ ), and in 63% of the experiments ( $n = 53$ ), congruence between  
190 power analysis design and the experiment's hypothesis could not be assessed mainly due to a  
191 lack of information about the design used in the power analysis (see Figure 2B). Finally, of  
192 12 studies that reported using partial eta-squared as the effect size parameter in a power  
193 analysis, 10 reported using G\*Power. Of the studies that used G\*Power, 8 used the default  
194 setting in (80%), one used the *As in SPSS* setting (10%), and one was not reproducible  
195 (10%), precluding an assessment of which setting was used (see Figure 3A). Neither of the  
196 power analyses that did not report using G\*Power could be reproduced with either setting.



### Figure 1

(A) Proportion of power analyses that were fully reproducible (orange) using the information provided in the article or supplemental materials and those that could not be reproduced (blue) based on the provided information. (B) Same data as that shown in (A); however, the power analyses that were conditionally reproducible (green) when certain assumptions were made regarding missing parameters are now highlighted. Each square represents one power analysis in the sample.



**Figure 2**

(A) Proportion of power analyses wherein the statistical test used in the power analysis matched the one used in the data analysis (orange), did not match (green), or was not reported with sufficient information to determine if the analyses matched (blue). (B) Proportion of power analyses that included a statistical test that addressed one of the hypotheses in the study (orange), included a test that did not address any hypotheses in the study (green), or was not reported with sufficient detail to determine if the test addressed a hypothesis (blue). Each square represents one power analysis in the sample.

197 **Exploratory Analyses**

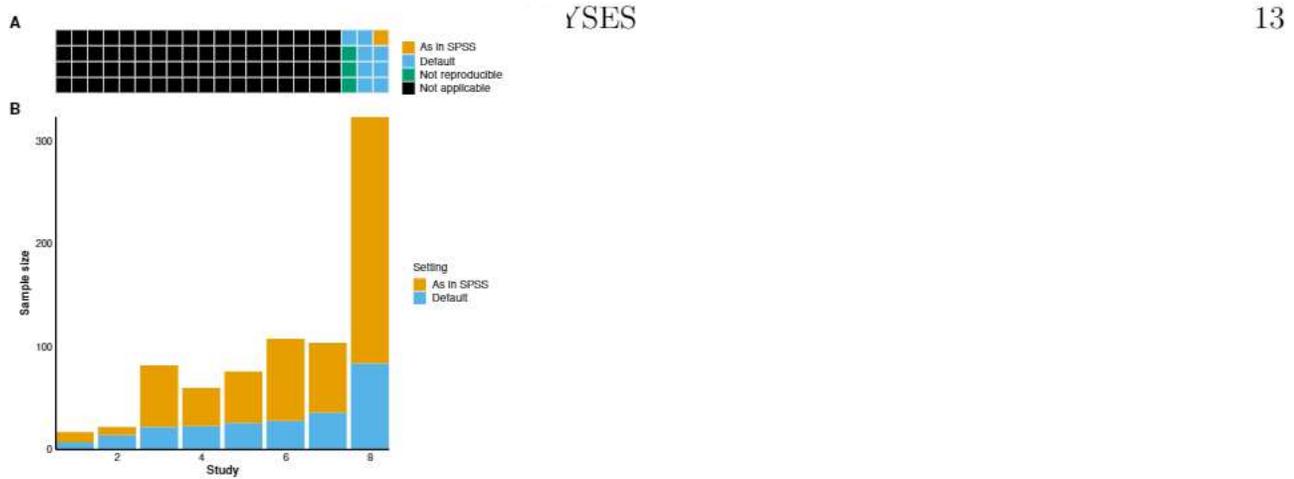
198         Several exploratory analyses were conducted to gather more information about the  
199 current state of the reproducibility of power analyses in motor behavior research.

200 ***Trouble Spots***

201         We noted that critical information required to reproduce power analyses was  
202 frequently missing: The statistical test and information about the effect size. We observed  
203 that 62% ( $n = 52$ ) of the power analyses did not include the statistical test, 48% ( $n = 40$ )  
204 did not include the type of effect size (e.g.,  $d$ ,  $f^2$ ,  $r$ ), and 17% ( $n = 14$ ) did not include the  
205 value of the effect size.

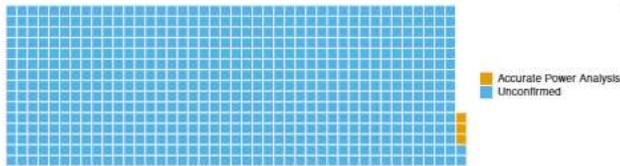
206 ***G\*Power Considerations***

207         G\*Power (Faul et al., 2009) was the chosen software in all studies that reported  
208 which software was used (74%;  $n = 62$ ). However, in at least 7% ( $n = 6$ ) of those studies,  
209 G\*Power does not provide an accurate power calculation for the statistical design of the  
210 study. Further, although G\*Power's user-friendly interface facilitates the process of  
211 conducting power analyses, the software's settings require careful use. For example, when  
212 partial eta-squared is used as the effect size in a power analysis in G\*Power, but was  
213 calculated in SPSS, then failing to change the settings from default to *As in SPSS* can result  
214 in considerably smaller sample sizes. We investigated the impact of this setting on sample  
215 size estimation across the 8 experiments that reported using partial eta-squared as the effect  
216 size and used G\*Power with the default setting to conduct the analysis. As seen in Figure  
217 3B, sample size estimation increased across all experiments when the *As in SPSS* setting was  
218 used, with the number of additional subjects needed ranging from 8 (Carnegie et al., 2020)  
219 to 240 (Uiga et al., 2020).



**Figure 3**

(A) Proportion of power analyses that included partial eta-squared ( $\eta_p^2$ ) as the effect size measure and used the *As in SPSS* setting in G\*Power (orange), the default setting (blue), were not reproducible (green), or did not include partial eta-squared as an effect size measure (black). Each square represents one power analysis in the sample. (B) A comparison of the required sample size based on chosen setting in G\*Power when using partial eta-squared as an effect size measure. The sample size calculated by the eight studies that used the default settings and partial eta-squared as an effect size measure is shown in blue. In contrast, if the partial eta-squared was originally calculated in SPSS, then using the appropriate *As in SPSS* setting would have resulted in substantially larger sample sizes for each study, with the difference represented by the orange bars.



#### Figure 4

Proportion of accurate power analyses (orange). An accurate power analysis had to be 1) reproducible, 2) include a statistical test that addressed at least one hypothesis and was used in the data analysis, and 3) were conducted with the appropriate software and settings. All other studies from the full sample of articles surveyed failed to meet these criteria (blue). Each square represents one study.

#### 220 *Rare Air*

221 Ideally, power analyses should be a) fully reproducible, b) the statistical test used in  
222 the power analysis should match the test used in the data analysis and c) at least one of the  
223 hypotheses, and d) the appropriate software with e) the appropriate settings should be used  
224 to obtain an accurate sample size estimation. Only three studies (4%; see Figure 4) met all  
225 five of these criteria (Daou et al., 2019; Harry et al., 2019; Rhoads et al., 2019).

## Discussion

226

227 *A priori* power analyses are a critical tool for designing informative studies and an  
228 important step toward high quality research. Inaccurate power analyses, however, can have  
229 the opposite effect as they may lead to underpowered study designs. Detecting, and even  
230 preventing, power analysis errors depends on the ability to successfully reproduce a given  
231 analysis, which requires reporting of pertinent information. The goal of the present study  
232 was to assess the current state of power analysis reproducibility in the motor behavior  
233 domain by evaluating 84 power analyses reported in 83 research articles published in the  
234 *Journal of Motor Behavior*, *Human Movement Science*, and the *Journal of Motor Learning*  
235 *and Development* between January 2019 and June 2021. Specifically, following a preregistered  
236 analysis plan, we assessed the proportion of power analyses that could be reproduced with  
237 the information reported in the article or supplementary material, the proportion of power  
238 analyses that could be reproduced conditional on making assumptions for missing  
239 parameters in the article, the proportion of studies wherein the statistical test used in the  
240 power analysis matched the test used in the data analysis, the proportion of studies wherein  
241 the statistical test used in the power analysis addressed the study's primary hypothesis, and  
242 finally, the proportion of studies that conducted a power analysis in G\*Power and used the  
243 default settings when computing the effect size parameter from partial eta-squared.

244 We were unable to reproduce 93% of the power analyses in the sample using only the  
245 information provided in the article or shared as supplementary information. By making  
246 assumptions for missing parameters, we were able to reproduce 43% of the power analyses,  
247 although this of course comes with caveats. Different parameters can yield the same sample  
248 size estimation, so despite our efforts to make plausible assumptions this approach does not  
249 guarantee that the original analyses adopted the same parameters we assumed. Therefore,  
250 43% represents the upper bound on reproducibility with the truth likely being even more  
251 concerning. Common reasons as to why power analysis reproducibility failed include lack of

252 information regarding the design used in the power analysis, the type of effect size, and the  
253 effect size value. A missing effect size value is particularly problematic because one cannot  
254 simply guess what effect size authors are targeting.

255         The process of conducting power analyses is facilitated by an abundance of  
256 user-friendly and openly available programs, including G\*Power (Faul et al., 2009), which is  
257 commonly used in social and behavioral research. In our sample, all studies ( $n = 62$ ) that  
258 reported the software used G\*Power, establishing a preference for this program in the motor  
259 behavior domain. While conducting a power analysis in G\*Power can be straightforward,  
260 easy-to-make mistakes when using the software can lead to inaccurate power calculations.  
261 For instance, G\*Power is not suitable for calculating power for mixed factorial designs with  
262 three or more within-subject factors, which require the use of other packages such as  
263 **Superpower** (Lakens & Caldwell, 2021). In our sample, at least 7% of the power analyses  
264 adopted designs that are too complex for G\*Power. More critically, G\*Power's method to  
265 compute the effect size partial eta-squared differs from the method used in SPSS. If  
266 researchers are basing their effect size target on previous estimates of partial eta-squared,  
267 and those estimates were calculated in SPSS, they need to change the effect size specification  
268 under *Options* from *Default* to *As in SPSS* (G\*Power version: 3.1.9.7). Across the power  
269 analyses assessed in the present study, 10 used partial eta-squared as the effect size  
270 parameter in G\*Power but only one used the *As in SPSS* setting. All 8 experiments that  
271 originally used the default setting would have been underpowered to detect the effect of  
272 interest if it was originally calculated in SPSS.

273         A lack of thoroughly reported and vetted power analyses contributes to the  
274 proliferation of underpowered studies, which combined with selection for significant results  
275 threatens the credibility of our literature. The impact of low power and selection bias is well  
276 illustrated by the growing body of metascience calling into question the reliability of research  
277 paradigms long considered robust (Carter et al., 2015; e.g., Maier et al., 2022; Vohs et al.,

278 2021), such as self-controlled practice in the motor learning domain (McKay, Yantha, et al.,  
279 2022). In a recent meta-analysis, McKay and colleagues estimated the benefit to motor  
280 learning of giving learners control over an aspect of their environment is trivially small, if  
281 existent, after correcting for publication bias. Nevertheless, the average effect size in the  
282 published literature was  $g = .54$ , suggesting apparent benefits. Similarly, another  
283 meta-analysis (McKay, Bacelar, et al., 2022) investigated the second motivational factor in  
284 OPTIMAL theory (Wulf & Lewthwaite, 2016), enhanced expectancies. The analysis found  
285 that despite an average benefit of  $g = .54$  in the published literature, the true effect of  
286 enhanced expectancies is likely much smaller, if it exists at all. The studies examined in  
287 these meta-analyses had median sample sizes of  $n = 14$  (McKay, Yantha, et al., 2022) and  $n$   
288  $= 18$  (McKay, Bacelar, et al., 2022), requiring effects larger than  $g = .8$  to achieve  
289 significance with an independent  $t$ -test. Therefore, selectively publishing significant results in  
290 these literatures meant publishing an abundance of large effects, making it possible for even  
291 null effects to appear moderately beneficial on average.

292 It is not only the extant but the future literature that is affected by underpowered  
293 studies. Small studies with positive results generate inflated effect sizes (Gelman & Carlin,  
294 2014). When these inflated effect sizes are used in power calculations for future studies,  
295 those studies become underpowered as well. This snowball effect can lead to uncertainty,  
296 research waste, and overall issues with replication as additional studies that are unlikely to  
297 be informative continue to be conducted and discarded, or reported when positive (Open  
298 Science Collaboration, 2015).

299 We have reviewed evidence that power analyses have been reported infrequently in  
300 the motor behavior literature (McKay et al., in-press). When power analyses were reported,  
301 they were rarely reproducible without making assumptions, and even then, most power  
302 analyses could not be reproduced. Meanwhile, there is growing evidence that the average  
303 power among motor behavior studies is low, making the literature vulnerable to more severe

304 bias from various selective reporting mechanisms (e.g., Lohse et al., 2016; McKay, Hussien,  
305 et al., 2022; McKay, Yantha, et al., 2022; Mesquida et al., 2022). Here, we argue that power  
306 analyses can easily be reported in a reproducible fashion and doing so is a progressive step  
307 toward improved research quality overall.

### 308 **Power Analysis Reproducibility: Recommendations for Future Studies**

309 Two simple practices can ensure power analysis reproducibility: 1) complete reporting  
310 and 2) sharing of code (see Table 1). The minimum parameters required to reproduce a  
311 power analysis are the type of effect size and its value (e.g.,  $d$ ,  $f^2$ ,  $r$ ), the accepted  
312 false-positive rate (i.e., alpha), the target power value (e.g., 80%), the specific statistical test,  
313 and the required sample size. Several additional parameters may be required to reproduce a  
314 specific analysis. A helpful strategy for G\*Power users is to report every possible input  
315 variable. Although one can technically reproduce a power analysis without knowing the  
316 primary hypothesis, we argue that researchers should also explicitly state their primary  
317 hypothesis(es) so others (e.g., collaborators, peer-reviewers, and readers) can assess whether  
318 a given study was powered to detect the primary effect(s) of interest.

319 A common trouble spot among studies in our sample was the description of the  
320 statistical test. We suggest making use of standardized language in power analysis software.  
321 This is a straightforward approach that offers researchers a clear way to describe the power  
322 analysis components, which is not only helpful from a practical standpoint, but it also  
323 reduces uncertainty. For instance, if a researcher reports the use of a test from the ANOVA  
324 family in G\*Power, five different options are possible. However, if she reports the use of the  
325 statistical test *ANOVA: Repeated measures, within-between interaction*, only one option is  
326 available. Reporting the exact language used in the software will clarify the statistical test  
327 for readers.

328 The second simple practice that will ensure power analysis reproducibility is sharing  
329 the protocol output or code. It is easy to save the exact protocol used in the power analysis

330 in software such as G\*Power, **Superpower**, and **R**. In G\*Power, the *Protocol of power*  
331 *analyses* tab includes all the details of the power analysis and can be saved as a PDF.  
332 Researchers can make this file available online in a repository such as the Open Science  
333 Framework (<https://osf.io>) or as part of supplementary material. Sharing code is a great  
334 strategy for ensuring the reproducibility of power analyses and primary analyses alike.

335         The benefits of adopting the practices we have presented go beyond power analysis  
336 reproducibility. For one, these practices increase research transparency, a key goal of the  
337 Open Science movement. Clear reporting can also assist other researchers in determining  
338 parameters for their own power analyses, which is especially helpful for researchers  
339 conducting their first power analysis for a given hypothesis. Although power analyses are  
340 best used for study planning, they can be conducted at any time. Therefore, the most  
341 informative power analyses are not just reproducible, but preregistered. Fortunately, another  
342 benefit of completing a reproducible power analysis while planning a study is that it  
343 represents a huge step toward preregistration. The study's primary hypothesis, smallest  
344 effect size of interest, statistical test to answer the research question, desired error rates, and  
345 the intended sample size comprise at least 50% of a preregistration form (e.g.,  
346 <https://aspredicted.org> form). To illustrate the potential symbiotic relationship between  
347 reproducible power analysis reporting and preregistration, in our sample, 50% of the  
348 experiments considered fully reproducible had a preregistered analysis plan, while only 0.47%  
349 of the overall sample was preregistered.

**Table 1**

*Sample checklist for improving the reproducibility of power analyses based on our two main recommendations.*

---

Item
<b>Reporting practices</b>
Report all input parameters for your selected analysis. <sup>a</sup>
<input type="checkbox"/> Specific statistical test e.g., Means: Difference between two independent means (two groups)
<input type="checkbox"/> Type of power analysis e.g., A priori: Compute required sample size - given $\alpha$ , power, and effect size
<input type="checkbox"/> Tails e.g., Two
<input type="checkbox"/> Effect size of interest and type e.g., $d = 0.4$
<input type="checkbox"/> Accepted false-positive rate e.g., $\alpha = 0.05$
<input type="checkbox"/> Target power e.g., 80%
<input type="checkbox"/> Allocation ratio N2/N1 e.g., 1
<input type="checkbox"/> Required (i.e., total) sample size e.g., 200
<input type="checkbox"/> Primary hypothesis e.g., We predict Group A to have lower total error in retention than Group B
<b>Sharing practices</b>
Can use repositories like Open Science Framework, Github, etc or include as supplementary material.
<input type="checkbox"/> G*Power: Click the <i>Protocol of power analyses</i> tab → Right click in window → <i>Save protocol to file</i>
<input type="checkbox"/> R package: Save as a .R file <sup>b</sup>

---

*Note.* <sup>a</sup> We have used an two-sample *t*-test and the input parameters from G\*Power as an example. Other statistical tests from the same or from a different ‘Test family’ may require fewer or several additional parameters. This process can be streamlined by sharing the protocol output from G\*Power or the R code.

<sup>b</sup> Other options for reproducible documents in R include RMarkdown (.Rmd) and Quarto (.qmd).

## 350 **Limitations**

351           Given we were unable to reproduce most of the power analyses, we cannot assess  
352 whether the primary deficit among studies is in power analysis quality or in reporting quality.  
353 When power analyses were reproducible, we made no effort to evaluate the quality of the  
354 evidence produced by those studies. Although we are optimistic that increased adoption of  
355 reproducible power analyses will benefit the quality of research in our field, we recognize that  
356 power analyses are not a panacea for bias in research. Although we recommend powering  
357 studies to detect the smallest effect size of interest (Lakens, 2022), we give no guidance on  
358 how to select this value. This is no small challenge for researchers and future metascience  
359 should focus on developing methods for choosing which effects are likely to be important in  
360 each study. In the meantime, it is important for researchers to think carefully about the  
361 specific effects they are investigating and not rely on effect size benchmarks to inform their  
362 power analyses. In fact, the benchmarks recommended by Cohen (1988) and used in  
363 G\*Power change depending on the type of analysis, rendering them inconsistent and illogical  
364 for use in sample size planning (Correll et al., 2020). Instead, researchers should think about  
365 raw differences they would not want to miss to help arrive at a smallest effect of interest.

## 366 **Conclusion**

367           Eighty-four motor behavior studies out of a sample of 635 included a power analysis,  
368 and of those we found three that were both appropriate and reproducible. There is  
369 converging evidence that motor behavior research is underpowered; perhaps because power  
370 analyses are not being leveraged to ensure a study produces informative results. Researchers  
371 can improve this situation by reporting all details of their power analyses and sharing their  
372 protocol output or code. Journals can improve this situation by asking for reproducible  
373 power analyses as a condition of publication. Peer reviewers can improve this situation by  
374 double-checking that the power analysis reported in a submission can be reproduced and has  
375 been appropriately conducted. Together, the sports science community can improve the  
376 quality of our research with relatively simple adjustments to the research workflow.

377 **Author Contributions (CRediT Taxonomy)**

378 Conceptualization: BM, MFBB, MJC

379 Data curation: BM, MJC

380 Formal analysis: BM, MFBB

381 Funding acquisition: MJC

382 Investigation: BM, MFBB

383 Methodology: BM, MFBB

384 Project administration: BM, MJC

385 Software: BM, MJC

386 Supervision: MJC

387 Validation: BM, MJC

388 Visualization: BM, MJC

389 Writing – original draft: BM, MFBB, MJC

390 Writing – review & editing: BM, MFBB, MJC

391 **Preprint**

392 [An unrefereed version of this paper can be found on SportRxiv:](#)

393 <https://doi.org/10.51224/SRXIV.184>.

394 **Open Science Practices**

395 The preregistration, data, and code for this study can be accessed using either of these links:

396 [https://github.com/cartermaclab/proj\\_power-reproducibility-motor-behaviour](https://github.com/cartermaclab/proj_power-reproducibility-motor-behaviour) or

397 <https://osf.io/9a6m8/>.

398 **Conflicts of Interest**

399 All authors declare no conflicts of interest.

400 **Funding**

401 This work was supported by the Natural Sciences and Engineering Research Council

402 (NSERC) of Canada (RGPIN-2018-05589; MJC) and McMaster University (MJC).

References

403

404 Abt, G., Boreham, C., Davison, G., Jackson, R., Nevill, A., Wallace, E., & Williams, M.  
405 (2020). Power, precision, and sample size estimation in sport and exercise science  
406 research. *Journal of Sports Sciences*, *38*(17), 1933–1935.

407 <https://doi.org/10.1080/02640414.2020.1776002>

408 Aust, F., & Barth, M. (2020). *papaja: Prepare reproducible APA journal articles with R*  
409 *Markdown*. <https://github.com/crsh/papaja>

410 Bacelar, M. F. B., Parma, J. O., Murrah, W. M., & Miller, M. W. (2022). Meta-analyzing  
411 enhanced expectancies on motor learning: Positive effects but methodological concerns.  
412 *International Review of Sport and Exercise Psychology*, *0*(0), 1–30.

413 <https://doi.org/10.1080/1750984X.2022.2042839>

414 Borg, D. N., Barnett, A., Caldwell, A. R., White, N., & Stewart, I. (2022). *The bias for*  
415 *statistical significance in sport and exercise medicine*.

416 <https://doi.org/10.31219/osf.io/t7yfc>

417 Carnegie, E., Marchant, D., Towers, S., & Ellison, P. (2020). Beyond visual fixations and  
418 gaze behaviour. Using pupillometry to examine the mechanisms in the planning and  
419 motor performance of a golf putt. *Human Movement Science*, *71*, 102622.

420 <https://doi.org/10.1016/j.humov.2020.102622>

421 Carter, E. C., Kofler, L. M., Forster, D. E., & McCullough, M. E. (2015). A series of  
422 meta-analytic tests of the depletion effect: Self-control does not seem to rely on a limited  
423 resource. *Journal of Experimental Psychology: General*, *144*(4), 796–815.

424 <https://doi.org/10.1037/xge0000083>

425 Chang, W. (2022). *Extrafont: Tools for using fonts*.

426 <https://CRAN.R-project.org/package=extrafont>

427 Cohen, J. (1988). *Statistical power analysis for the behavioral sciences* (2nd ed.). Routledge.

428 Correll, J., Mellinger, C., McClelland, G. H., & Judd, C. M. (2020). Avoid Cohen’s “small,”  
429 “medium,” and “large” for power analysis. *Trends in Cognitive Sciences*, *24*(3), 200–207.

- 430 <https://doi.org/10.1016/j.tics.2019.12.009>
- 431 Daou, M., Rhoads, J. A., Jacobs, T., Lohse, K. R., & Miller, M. W. (2019). Does limiting  
432 pre-movement time during practice eliminate the benefit of practicing while expecting to  
433 teach? *Human Movement Science, 64*, 153–163.  
434 <https://doi.org/10.1016/j.humov.2018.11.017>
- 435 DeBruine, L., Krystalli, A., & Heiss, A. (2021). *faux: Simulation for factorial designs*.  
436 <https://CRAN.R-project.org/package=faux>
- 437 Faul, F., Erdfelder, E., Buchner, A., & Lang, A.-G. (2009). Statistical power analyses using  
438 G\*power 3.1: Tests for correlation and regression analyses. *Behavior Research Methods,*  
439 *41*(4), 1149–1160. <https://doi.org/10.3758/BRM.41.4.1149>
- 440 Fitzpatrick, P., de Jonge, E., & Warnes, G. R. (2019). *Daff: Diff, patch and merge for*  
441 *data.frames*. <https://CRAN.R-project.org/package=daff>
- 442 Gelman, A., & Carlin, J. (2014). Beyond power calculations: Assessing type S (sign) and  
443 type M (magnitude) errors. *Perspectives on Psychological Science, 9*(6), 641–651.
- 444 Harry, J. R., Lanier, R., Nunley, B., & Blinch, J. (2019). Focus of attention effects on lower  
445 extremity biomechanics during vertical jump landings. *Human Movement Science, 68*,  
446 102521. <https://doi.org/10.1016/j.humov.2019.102521>
- 447 Lakens, D. (2022). Sample size justification. *Collabra: Psychology, 8*(1), 33267.
- 448 Lakens, D., & Caldwell, A. R. (2021). Simulation-based power analysis for factorial analysis  
449 of variance designs. *Advances in Methods and Practices in Psychological Science, 4*(1),  
450 2515245920951503. <https://doi.org/10.1177/2515245920951503>
- 451 Lohse, K., Buchanan, T., & Miller, M. (2016). Underpowered and overworked: Problems  
452 with data analysis in motor learning studies. *Journal of Motor Learning and*  
453 *Development, 4*(1), 37–58. <https://doi.org/10.1123/jmld.2015-0010>
- 454 Maier, M., Bartoš, F., Stanley, T. D., Shanks, D. R., Harris, A. J. L., & Wagenmakers, E.-J.  
455 (2022). No evidence for nudging after adjusting for publication bias. *Proceedings of the*  
456 *National Academy of Sciences, 119*(31), e2200300119.

- 457 <https://doi.org/10.1073/pnas.2200300119>
- 458 McCrum, C., Beek, J. van, Schumacher, C., Janssen, S., & Van Hooren, B. (2022). Sample  
459 size justifications in Gait & Posture. *Gait & Posture*, *92*, 333–337.  
460 <https://doi.org/10.1016/j.gaitpost.2021.12.010>
- 461 McKay, B., Bacelar, M., Parma, J. O., Miller, M. W., & Carter, M. J. (2022). *The*  
462 *combination of reporting bias and underpowered study designs have substantially*  
463 *exaggerated the motor learning benefits of self-controlled practice and enhanced*  
464 *expectancies: A meta-analysis*. PsyArXiv. <https://doi.org/10.31234/osf.io/3nhtc>
- 465 McKay, B., Corson, A., Vinh, M.-A., Jeyarajan, G., Tandon, C., Brooks, H., Hubley, J., &  
466 Carter, M. J. (in-press). Low prevalence of a priori power analyses in motor behavior  
467 research. *Journal of Motor Learning & Development*.
- 468 McKay, B., Hussien, J., Vinh, M.-A., Mir-Orefice, A., Brooks, H., & Ste-Marie, D. M. (2022).  
469 Meta-analysis of the reduced relative feedback frequency effect on motor learning and  
470 performance. *Psychology of Sport and Exercise*, 102165.  
471 <https://doi.org/10.1016/j.psychsport.2022.102165>
- 472 McKay, B., Yantha, Z. D., Hussien, J., Carter, M. J., & Ste-Marie, D. M. (2022).  
473 Meta-analytic findings in the self-controlled motor learning literature: Underpowered,  
474 biased, and lacking evidential value. *Meta-Psychology*, *6*, 1–32.  
475 <https://doi.org/10.15626/MP.2021.2803>
- 476 Mesquida, C., Murphy, J., Lakens, D., & Warne, J. (2022). *Replication concerns in sports*  
477 *science: A narrative review of selected methodological issues in the field*. SportRxiv.  
478 <https://sportrxiv.org/index.php/server/preprint/view/127>
- 479 Open Science Collaboration. (2015). Estimating the reproducibility of psychological science.  
480 *Science*, *349*(6251), aac4716. <https://doi.org/10.1126/science.aac4716>
- 481 Qiu, W. (2021). *powerMediation: Power/Sample size calculation for mediation analysis*.  
482 <https://CRAN.R-project.org/package=powerMediation>
- 483 R Core Team. (2021). *R: A language and environment for statistical computing*. R

- 484 Foundation for Statistical Computing. <https://www.R-project.org/>
- 485 Rhoads, J. A., Daou, M., Lohse, K. R., & Miller, M. W. (2019). The effects of expecting to  
486 teach and actually teaching on motor learning. *Journal of Motor Learning &  
487 Development*, 7(1), 84–105.
- 488 Robinson, M. A., Vanrenterghem, J., & Pataky, T. C. (2021). Sample size estimation for  
489 biomechanical waveforms: Current practice, recommendations and a comparison to  
490 discrete power analysis. *Journal of Biomechanics*, 122, 110451.  
491 <https://doi.org/10.1016/j.jbiomech.2021.110451>
- 492 Rudis, B., & Gandy, D. (2019). *Waffle: Create waffle chart visualizations*.  
493 <https://gitlab.com/hrbrmstr/waffle>
- 494 Twomey, R., Yingling, V., Warne, J., Schneider, C., McCrum, C., Atkins, W., Murphy, J.,  
495 Medina, C. R., Harley, S., & Caldwell, A. (2021). The nature of our literature: A  
496 registered report on the positive result rate and reporting practices in kinesiology.  
497 *Communications in Kinesiology*, 1(3). <https://doi.org/10.51224/cik.v1i3.43>
- 498 Uiga, L., Poolton, J. M., Capio, C. M., Wilson, M. R., Ryu, D., & Masters, R. S. W. (2020).  
499 The role of conscious processing of movements during balance by young and older adults.  
500 *Human Movement Science*, 70, 102566. <https://doi.org/10.1016/j.humov.2019.102566>
- 501 Ushey, K. (2022). *renv: Project environments*. <https://CRAN.R-project.org/package=renv>
- 502 Vohs, K. D., Schmeichel, B. J., Lohmann, S., Gronau, Q. F., Finley, A. J., Ainsworth, S. E.,  
503 Alquist, J. L., Baker, M. D., Brizi, A., Bunyi, A., Butschek, G. J., Campbell, C., Capaldi,  
504 J., Cau, C., Chambers, H., Chatzisarantis, N. L. D., Christensen, W. J., Clay, S. L.,  
505 Curtis, J., . . . Albarracín, D. (2021). A multisite preregistered paradigmatic test of the  
506 ego-depletion effect. *Psychological Science*, 32(10), 1566–1581.  
507 <https://doi.org/10.1177/0956797621989733>
- 508 Westfall, J. (2015). PANGEA: Power analysis for general ANOVA designs. *Unpublished  
509 Manuscript*. Available at <Http://Jakewestfall.org/Publications/Pangea.pdf>, 4.
- 510 Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L. D., François, R., Golemund,

- 511 G., Hayes, A., Henry, L., Hester, J., Kuhn, M., Pedersen, T. L., Miller, E., Bache, S. M.,  
512 Müller, K., Ooms, J., Robinson, D., Seidel, D. P., Spinu, V., ... Yutani, H. (2019).  
513 Welcome to the tidyverse. *Journal of Open Source Software*, 4(43), 1686.  
514 <https://doi.org/10.21105/joss.01686>
- 515 Wulf, G., & Lewthwaite, R. (2016). Optimizing performance through intrinsic motivation  
516 and attention for learning: The OPTIMAL theory of motor learning. *Psychonomic*  
517 *Bulletin & Review*, 23(5), 1382–1414.