

Extraordinary claims in the literature on high-intensity interval training (HIIT): I. Bonafide scientific revolution or a looming crisis of replication and credibility?

EKKEKAKIS, P., SWINTON, P and TILLER, N.B.

2023

This is the accepted version of the above article, which is distributed under Springer's [AM terms of use](#). The version of record is available from the journal website: <https://doi.org/10.1007/s40279-023-01880-7>

1
2 RUNNING HEAD: High-Intensity Interval Training and Replication

3
4 Extraordinary Claims in the Literature on High-Intensity Interval Training (HIIT):

5 I. Bonafide Scientific Revolution or a Looming Crisis of Replication and Credibility?

6
7 Panteleimon Ekkekakis

8 Department of Kinesiology, Michigan State University

9 East Lansing, Michigan, U.S.A.

10
11 Paul Swinton

12 School of Health Sciences, Robert Gordon University

13 Aberdeen, Scotland, United Kingdom

14
15 Nicholas B. Tiller

16 The Lundquist Institute for Biomedical Innovation at Harbor-UCLA Medical Center

17 Torrance, California, U.S.A.

18 Address correspondence to:

19 Panteleimon Ekkekakis

20 Department of Kinesiology

21 Michigan State University

22 308 W Circle Dr #134

23 East Lansing, MI 48824

24 United States of America

25 Tel. +1 (517) 355-4730

26 E-mail: ekkekaki@msu.edu

27
28 1st draft Jun 10, 2022; 2nd draft Feb 28, 2023; 3rd draft June 7, 2023

50 **Competing Interests:** The authors declare no competing financial or other interests.

51

52 **Funding:** No funding was received for the preparation and/or publication of this manuscript.

53

54 **Author Contributions:** PE conceived and drafted the manuscript and responses to peer review

55 and editorial comments. PS and NBT revised and edited the original and revised manuscript,

56 each contributing additional, original intellectual content.

57

58

59 **1. Introduction**

60 In the mid-1990s, exercise science underwent what can be characterized as the most
61 consequential paradigmatic shift in its history, expanding its focus from exercise training for
62 fitness enhancement to lifestyle physical activity for the promotion of public health [1,2]. This
63 new perspective resulted in a series of physical activity recommendations from organizations in
64 the United States, including the Centers for Disease Control and Prevention [3], the Surgeon
65 General [4], and the National Institutes of Health [5,6], followed by similar initiatives in other
66 countries. These recommendations converged on a common, easy-to-remember message: adults
67 should accumulate (in short bouts, dispersed throughout the day) at least 30 min of physical
68 activity, performed at least at a moderate intensity, on most, but preferably all, days of the week.

69 At the time, several aspects of these recommendations were criticized for their lack of
70 specificity (e.g., what is "moderate" intensity?) or for relying on a weak empirical basis (e.g.,
71 scant evidence on "accumulated" physical activity). Furthermore, while the recommendations
72 implied that additional health benefits could be obtained with activities of higher-than-moderate
73 intensity, the emphasis was clearly placed on activity options that involve moderate intensity,
74 such as brisk walking, based on the assumption that such options are realistic and non-
75 intimidating for a largely hypoactive adult population [7]. This rationale was supported by a
76 meta-analysis showing that interventions attempting to implement activity of higher intensity
77 were associated with lower participation [8].

78 Despite good intentions, the guidelines had no measurable effect on public participation
79 in physical activity. Accelerometry data from the 2003-2004 National Health and Nutritional
80 Examination Survey (NHANES), a nationally representative study in the United States (with
81 6,329 individuals providing at least one day of data), showed that only 3.5% of individuals 20 to

82 59 years of age and 2.4% of those aged 60 years or older registered at least 30 min of moderate-
83 intensity physical activity per day on at least five days per week [9]. Less than 1% of adults
84 registered 20 min of vigorous-intensity activity on at least three days per week [10]. In the 2005-
85 2006 NHANES, the situation was unchanged, with only 3.2% of adults achieving the
86 recommended dose of moderate-intensity activity [11]. The absence of positive results from
87 population surveys encouraged calls for renewed emphasis on higher intensity [12-14]. Indeed,
88 reformulated physical activity guidelines explicitly offered a choice between moderate intensity
89 (for at least 30 min on five days per week or 150 min per week), vigorous intensity (for at least
90 20-25 min on three days per week or 75 min per week), or an equivalent combination [15,16].

91 In 2005, in the midst of the debate preceding the reformulation of the guidelines and the
92 renewed emphasis on vigorous-intensity activities, researchers published results from a doctoral
93 dissertation [17] in the *Journal of Applied Physiology*. The article reported a remarkable finding,
94 namely that a group of two women and six men doubled their cycling endurance performance
95 (time to fatigue while pedaling at 80% VO_2peak) after a total of only about 15 min of high-
96 intensity interval training (HIIT) over two weeks, without changing their maximal aerobic
97 capacity. An accompanying editorial [18] underscored the "effectiveness and remarkable time
98 efficiency" of high-intensity training but noted that the "price" participants have to pay is a need
99 for "a high level of motivation" and "a feeling of severe fatigue lasting for at least 10–20 min" (p.
100 1983) [18]. Over the next several years, fueled by extensive media coverage in which HIIT was
101 portrayed as a solution for individuals with limited available discretionary time, HIIT became a
102 top trend in the fitness industry worldwide [19]. Moreover, since 2005, HIIT has been the subject
103 of approximately 4,000 articles, with more than 700 new articles being added to the literature
104 each year, 10% of them being meta-analyses (see Figure 1).

105 The data on the fitness and health benefits of HIIT have been characterized as "clear and
106 convincing" (p. 1231) [20]. Nevertheless, as claims about HIIT are now influencing policy on a
107 national and global scale (e.g., through exercise prescription guidelines and physical activity
108 recommendations), it would be prudent to assess whether these claims can withstand statistical
109 scrutiny. Steen [21] has argued that "error and fraud are the main sources of scientific
110 misinformation" but "error is more prevalent than fraud" (p. 501). He insisted that "bias can also
111 result from earnest error, statistical naiveté, or other innocent causes; not all bias is fraud" (p.
112 502). However, it has already been established that some of the extraordinary claims surrounding
113 HIIT cannot be attributed solely to earnest human error. For example, on 14 February 2019, the
114 *British Journal of Sports Medicine* issued a press release, promoting the publication of a meta-
115 analysis entitled "Is interval training the magic bullet for fat loss?" [22], which purportedly
116 showed that, indeed, HIIT results in significantly larger reduction in total absolute fat mass than
117 moderate-intensity continuous exercise (-2.28 kg, 95% CI -4.00 to -0.56, $p = 0.0094$). The press
118 release issued by the journal appeared under the title "Interval training may shed more pounds
119 than continuous moderate intensity workout," and attracted the attention of major news outlets,
120 including the global news agency *Reuters* and influential magazines like *Runner's World*.¹
121 However, the meta-analysis was later retracted because the authors could not explain how they
122 obtained their data (e.g., a larger reduction of body fat by -13.44 kg in HIIT than moderate-
123 intensity continuous exercise, associated with a 12-week study that reported no relevant data).

124 Drawing lists of studies from two recently published systematic reviews, the present
125 critical analysis focuses on statistical concerns emanating from the rapidly expanding literature

¹ See: (1) <https://bjsm.bmj.com/content/bjsports/suppl/2019/02/19/bjsports-2018-099928.DC1/bjsports-2018-099928.pdf>; (2) <https://www.reuters.com/article/us-health-exercise-training/interval-training-burns-off-more-pounds-than-jogging-or-cycling-idUSKCN1Q71TT>; (3) <https://www.runnersworld.com/news/a26339798/interval-training-for-weight-loss-study/>

126 on HIIT. This analysis highlights alarming parallels between prevalent practices in the HIIT
127 literature and the emergence of a replication crisis in other scientific fields. The narrative
128 culminates in a call for a return to fundamental principles of statistics. Unlike some of the more
129 complicated scenarios outlined by Sainani et al. [23], the points raised in the following sections
130 refer to elementary statistical principles, such as the mechanisms that raise the risk of Type I and
131 Type II errors of statistical inference. The analysis culminates in a call not for the
132 implementation of novel, obscure, or advanced statistical methods but rather for a *return* to
133 fundamental statistical principles, along with the *readoption* of the critical outlook that should, in
134 principle, characterize all manner of scientific inquiry.

135 **2. Statistical Preliminaries: (Mis-) Understanding Null-Hypothesis Significance Testing**

136 Studies evaluating the effectiveness of HIIT reach their conclusions following the
137 statistical methodology known as null-hypothesis significance testing (NHST). Despite strong
138 concerns [24,25] and the presence of alternatives (i.e., Bayesian inference and fiducial inference)
139 [26], NHST has been established as the standard method for evaluating statistical tests in most
140 domains of human-science research, including the exercise sciences. Despite its popularity,
141 however, the NHST is frequently misunderstood, misapplied, and misinterpreted [24,25,27].

142 NHST represents the amalgamation of the testing methodologies proposed during the
143 period 1915-1933 by Ronald Aylmer Fisher (1890-1962) and the duo of Jerzy Neyman (1894-
144 1981) and Egon Sharpe Pearson (1895-1980). Fisher on the one hand, and Neyman and Pearson
145 on the other, contributed different pieces of what evolved into the NHST methodology, but it is
146 important to emphasize that, as applied today, the NHST is "essentially an anonymous hybrid"
147 and "a marriage of convenience that neither party would have condoned" (p. 171) [28].

148 Fisher, who emphasized the importance of inductive reasoning (i.e., analyzing samples to

149 draw inferences about the population), is credited with the concept of the null hypothesis (i.e.,
150 data demonstrating random variance) and the use of exact p values as a quantitative measure of
151 the "extremeness" of the data given the null hypothesis. By extension, he considered p values as
152 an indication of the plausibility or implausibility of the null hypothesis. However, although he
153 famously wrote that "we shall not often be astray if we draw a conventional line at .05" (p. 82)
154 [29], for Fisher, a low p value, such as $p < .05$, represented merely a sign that a finding may be
155 worthy of further study, starting with an attempt at replication.

156 In the central point of contention with Fisher, Neyman and Pearson espoused a deductive
157 approach, in which the null hypothesis is either rejected in favor of an alternative or retained for
158 further study (which is not the same as accepting that the null hypothesis is true). Unlike Fisher,
159 who believed that a specific hypothesis can be tested using data from a single study, Neyman and
160 Pearson were not interested in developing a method for drawing inductive inferences about a
161 single hypothesis based on the "statistical significance" of data from a single study. Instead, their
162 goal was to use a deductive approach and probability theory to develop "rules of behavior" (i.e.,
163 rejection vs. non-rejection of a hypothesis) to ensure that the frequency of errors (i.e., the
164 erroneous rejection or non-rejection) would be kept below an acceptably low limit over a series
165 of many studies:

166 But we may look at the purpose of tests from another view-point. Without hoping to
167 know whether each separate hypothesis is true or false, we may search for rules to govern
168 our behaviour with regard to them, in following which we insure that, in the long run of
169 experience, we shall not be too often wrong. Here, for example, would be such a "rule of
170 behaviour": to decide whether a hypothesis, H , of a given type be rejected or not,
171 calculate a specified character, x , of the observed facts; if $x > x_0$ reject H , if $x \leq x_0$ accept

172 H. Such a rule tells us nothing as to whether in a particular case H is true when $x \leq x_0$ or
173 false when $x > x_0$. But it may often be proved that if we behave according to such a rule,
174 then in the long run we shall reject H when it is true not more, say, than once in a
175 hundred times, and in addition we may have evidence that we shall reject H sufficiently
176 often when it is false (p. 291) [30].

177 The Neyman-Pearson approach, therefore, implied two types of errors, called Type I and
178 Type II, with the rate of those errors symbolized by the Greek letters α and β , respectively, as
179 well as the concept of statistical power, symbolized as $1-\beta$ [31]. A Type I error (α) occurs when
180 "if we reject H_0 , we may reject it when it is true," whereas a Type II error (β) occurs when "if we
181 accept H_0 , we may be accepting it when it is false, that is to say, when really some alternative is
182 true" (p. 296) [30]. Statistical power ($1-\beta$) is defined as "the probability of rejecting the
183 hypothesis tested, H_0 , when the true hypothesis is H_i " (p. 498) [32].

184 Fisher [33] concurred with the notion of Type I errors and was keenly aware of the risk of
185 raising the rate of such errors as a result of performing a multitude of tests. For example, he
186 argued that a comparison between two extreme values "picked out from the results, will often
187 appear to be significant, even from undifferentiated material" (p. 66). His proposed remedy was
188 analogous to alpha-splitting, namely making the criterion for evaluating the p value more
189 stringent: "We might, therefore, require the probability of the observed difference to be as small
190 as 1 in 900, instead of 1 in 20, before attaching statistical significance to the contrast" (p. 66). On
191 the other hand, arguing from an inductive standpoint, Fisher rejected the notion of Type II errors
192 because he believed that scientific research is a process of "learning by experience" and, in such
193 a process, a priori knowledge is "almost always absent or negligible" (p. 73) [34]. Thus, although
194 he considered the rate of Type I error "calculable, and therefore controllable," he insisted that

195 Type II error is "incalculable both in frequency and in magnitude" (p. 73).

196 Interestingly, while Fisher rejected the notion of Type II error, he was aware of the
197 importance of statistical power (although he used the term "sensitivity" or "sensitiveness") and
198 the role of sample size and a higher number of repetitions in increasing statistical power: "By
199 increasing the size of the experiment, we can render it more sensitive, meaning by this that it will
200 allow of the detection of a lower degree of sensory discrimination, or, in other words, of a
201 quantitatively smaller departure from the null hypothesis" (p. 25) [33]. Commentators have noted
202 that "Fisher's 'sensitivity' and Neyman-Pearson's 'power' refer to the same concept" (p. 173) [28],
203 but Fisher "denied the possibility of assessing it quantitatively" (p. 1245) [35].

204 The main misinterpretations surrounding the NHST emerged following the merger of the
205 Fisher and Neyman-Pearson approaches by anonymous researchers [35,36], a merger "that
206 neither party would have condoned," to repeat the phrase of Hubbard and Bayarri (p. 171) [28].
207 This anonymous and unsanctioned merger has resulted in several persistent misuses and
208 misinterpretations that have plagued research for decades [24,37,38]. Of these, the following
209 problems are arguably most relevant to research on HIIT.

210 **2.1. The p Value as an Indication of the Plausibility of the Null Hypothesis**

211 First, there is a widespread but mistaken belief that a p value of .05 means that there is
212 only 5% probability of the null hypothesis being true (or, conversely, for $1-p$, that there is 95%
213 probability that the null hypothesis is false). This belief is mistaken because p values are
214 calculated from the data under the assumption that the null hypothesis is true [39]. A p value
215 merely indicates the probability (assuming that the null hypothesis is true) of observing a test
216 statistic (e.g., a t value) as extreme or more extreme than the value observed in the present
217 sample. This can be expressed as $\Pr(\text{data}|\text{H}_0)$ in probability notation. This statement is not

218 equivalent to the interpretation that a p value of .05 means that there is only 5% probability of
 219 the null hypothesis being true, namely $\Pr(H_0|\text{data})$. While the p value does provide some
 220 indication of the plausibility or implausibility of the null hypothesis, a p near .05 "greatly
 221 overstates the evidence against the null hypothesis" (p. 139) [37]. Berger and Sellke [40]
 222 calculated that the lower bound of $\Pr(H_0|\text{data})$ can be estimated as:

$$\Pr(H_0|\text{data}) = (1 + (1 + n)^{-1/2} \exp\{t^2 / [2 (1 + 1/n)]\})^{-1}$$

223 Using a t value that yields $p = .05$ ($t = 1.96$) and a sample size of $n = 50$ per group results in
 224 $\Pr(H_0|\text{data}) = .52$, which surpasses $p = .05$ by more than an order of magnitude [40,41].

225 **2.2. The p Value as an Index of the Risk of Type I Errors**

226 Second, related to the previous point, there is pervasive confusion between a p value,
 227 namely the probability of obtaining a test statistic at least as extreme as that obtained from a
 228 given study under the assumption that the null hypothesis is true, and α , namely the rate of Type
 229 I errors [28]. In actuality, a single number (i.e., a p value) cannot simultaneously serve the dual
 230 function of providing an indication of the "extremeness" of the data from any given study and, at
 231 the same time, an indication of the "long-run" frequency of improperly rejecting the null
 232 hypothesis when it is true [39]. Nevertheless, statisticians [40-42] have estimated that, at least for
 233 the range $p < 1/e$, where e is Euler's constant (2.71828), namely $p < .36787$, the lower bound of α
 234 (i.e., the minimum risk of a Type I error when rejecting the null hypothesis) can be estimated by:

$$\alpha(p) = (1 + [-e p \log(p)]^{-1})^{-1}$$

235 where $\log(p)$ is the natural logarithm of the p value. Substituting $p = .05$ yields $\alpha = .289$. This
 236 means that there is at least 28.9% probability of a Type I error when rejecting the null hypothesis
 237 on the basis of a p value close to .05. In other words, at least 28.9% of p values near .05 can be
 238 expected to come from studies in which the null hypothesis is true.

239 **2.3. The p Value as an Index of Replicability**

240 Third, researchers often mistakenly assume that a low p value (e.g., $p < .05$) entails that,
241 if the same test were performed on a different sample randomly drawn from the same population
242 (e.g., same sample sizes, same treatments), there would be high probability (e.g., $> 95\%$) that the
243 new p value would be similarly low (e.g., $p < .05$) [43]. In fact, except in studies with levels of
244 statistical power over 90%, p values are characterized by extraordinary uncertainty [44,45].
245 Thus, for a comparison between two means resulting in $p < .05$, the probability of finding $p < .05$
246 in a (theoretical) "identical" replication (with the difference between the means being in the same
247 direction) has been estimated as only 50% [46-49].

248 **2.4. A Non-Significant p Value As a Basis for Accepting the Null Hypothesis**

249 Fourth, a widely prevalent and persistent misunderstanding is that obtaining a
250 nonsignificant test result (e.g., $p > .05$) can be interpreted as an indication that the null hypothesis
251 (e.g., $\mu_1 - \mu_2 = 0$) is true or as indication of the absence of an effect [24,37,38,50-52]. Fisher [33]
252 famously asserted that "the null hypothesis is never proved or established, but is possibly
253 disproved, in the course of experimentation" (p. 19). Accordingly, one of the oft-quoted
254 admonitions of statisticians is that "the absence of evidence is not the same as evidence of
255 absence" [53-54]. A non-significant p value cannot provide a basis for accepting the null
256 hypothesis as true or for the rejection of alternatives. It only suggests that a null effect is
257 statistically consistent (or not inconsistent) with the data, along with the range of other effects
258 encompassed within the confidence interval. However, $p > .05$ provides no indication that the
259 null effect, specifically, is the most likely among these. Moreover, using non-significant p values
260 as an indication in support of the null hypothesis is especially precarious in scientific fields, such
261 as the exercise sciences [55], that are characterized by a preponderance of underpowered studies.

262 Authors have warned that "null results are surprisingly easy to obtain by mere statistical
263 artefacts; simply using a small sample or a noisy measure can suffice to produce a false
264 negative" (p. 97) [56].

265 Collectively, the aforementioned misinterpretations suggest that NHST is a potentially
266 useful, but delicate, test methodology. As such, it should be approached cautiously, recognizing
267 and respecting its considerable limitations. The wide prevalence of the misinterpretations and
268 misuses of the NHST across many domains of scientific research cannot be deemed a valid
269 excuse for their ubiquity within the field of exercise science in general and research on HIIT in
270 particular. Likewise, the fact that prestigious journals within the field of exercise science have
271 permitted such practices does not render them any less egregious or harmful.

272 While there is ongoing debate about the causes and potential remedies of these
273 misinterpretations and misuses of the NHST [57], many statistical experts see these
274 misinterpretations and misuses as contributors to the phenomenon of non-replicable research
275 [58-61]. Whether implemented deliberately or inadvertently, questionable statistical practices
276 can result in intriguing, albeit fanciful, findings, with a high probability of attracting the attention
277 of other researchers and the public. Serra-Garcia and Gneezy [62] speculated that, while
278 evaluating manuscripts, journal editors and peer reviewers probably weigh two considerations
279 against each other, namely the likely robustness or reliability of the result on one hand and its
280 interest or curiosity on the other: "when the paper is more interesting, the review team may apply
281 lower standards regarding its reproducibility" (p. 4).

282 **3. Misuses of Null-Hypothesis Significance Testing in Research on HIIT**

283 The following two sections present critical commentaries on two major variants of claims
284 pertaining to HIIT, namely (i) that HIIT is effective in improving a variety of fitness and health

285 outcomes, and (ii) that HIIT is as effective as more time-consuming moderate-intensity
286 continuous exercise. We examine studies contained in two recent systematic reviews to
287 demonstrate that deviating from elementary statistical principles can result in data that can be
288 portrayed as supporting both of these conclusions, but with a high probability that such
289 conclusions reflect errors of statistical inference. It is important to reiterate that the problems to
290 be discussed are certainly not unique to the HIIT literature but have long plagued the broader
291 exercise-science literature [63].

292 **3.1. The "Is Effective" Problem**

293 As evidenced in meta-analyses, [64,65] a striking feature of the research literature on
294 HIIT is an abundance of implausibly large effect sizes (e.g., standardized mean differences over
295 2.0 or 2.5 standard deviations) reportedly demonstrating the extraordinary effectiveness of HIIT
296 compared to control conditions or even compared to active interventions consisting of moderate-
297 intensity continuous exercise training. Some of these can be dismissed as mistakes, such as
298 standardized mean differences (Hedges' g) of 11, 16, or 29 standard deviations [64], which can
299 be readily attributed to computational errors (e.g., mistaking standard errors of the mean as
300 standard deviations). Other cases, however, may be more complicated. For example, a
301 remarkable standardized mean difference in maximal oxygen consumption of 4.59 standard
302 deviations [65] from a 12-week comparison between HIIT and moderate-intensity continuous
303 exercise [66] could be due to a host of well-established but frequently overlooked sources of
304 methodological bias. These include, but are not limited to, the inadequate concealment of the
305 randomization sequence, the absence of intention-to-treat analyses, and the use of unblinded
306 outcome assessors. In addition, exercise researchers are aware of the biasing effect of several
307 exercise-specific factors, such as the lack of control for verbal encouragement during tests of

308 maximal performance [67-69]. When exercise testing is conducted by researchers who are ardent
309 proponents of HIIT (e.g., "HIIT should play a central role in health activity guidelines" because
310 it can "maximize the benefits of physical activity globally," p. 5216) [70], and are unblinded to
311 treatment allocation, finding a standardized mean difference of 4.59 standard deviations in favor
312 of HIIT becomes a plausible occurrence.

313 Such methodological sources of bias are beyond the scope of the present analysis. Here,
314 we focus on statistical mechanisms that can produce similarly extraordinary (and likely non-
315 replicable) results. For example, meta-analyses have reported that HIIT interventions have
316 produced standardized mean differences that exceeded 2.5 standard deviations [71,72]. Closer
317 inspection of the characteristics of the studies that produced these large effect sizes [73-75]
318 reveals certain notable commonalities: (i) small sample sizes (e.g., 10-20 participants per group),
319 resulting in wide confidence intervals and low statistical power to detect even large effects, (ii)
320 long lists of dependent variables, covering several multidimensional domains (e.g.,
321 anthropometric characteristics, inflammatory or immune markers, indices of cardiac, vascular,
322 cardiorespiratory, or metabolic function), (iii) absence of pre-registration that could have allayed
323 concerns about selective reporting, (iv) absence of designation of dependent variables as primary
324 vs. secondary, and (v) numerous statistical tests, each evaluated with the criterion of $p < .05$.
325 Because of sampling variability and the lack of precision associated with small samples,
326 estimates of population values (means, standard deviations) and, therefore, the associated p
327 values "dance around" (p. 1720), as Gandevia [76] put it. Given a long enough list of dependent
328 variables, it becomes almost inevitable that some means will happen to show exaggerated
329 differences, thus resulting in extraordinarily large effect sizes. With a lax criterion such as $p <$
330 $.05$, one or more comparisons will cross the threshold of "statistical significance," increasing the

331 likelihood of publication. A cynic might argue that this approach could be used, deliberately or
332 unwittingly, as a recipe for producing seemingly "significant" and possibly novel or intriguing
333 results, albeit results that are probably non-reproducible.

334 These basic statistical mechanics are taught in undergraduate and postgraduate university
335 courses on research methodology. It is, therefore, surprising and disheartening that studies with
336 the aforementioned characteristics, and attendant risk of producing untenable results, continue to
337 be commonplace in large sections of exercise-science research [77], including research on HIIT.
338 Nosek et al. [57] criticized the "disciplinary incentives" that tend to "inflate the rate of false
339 effects in published science" and "favor novelty over replication" (p. 615). In the following
340 sections, we elaborate on several aspects of this problem.

341 **3.1.1. Multiplicity**

342 Methodologically strong studies, including most well-designed randomized controlled
343 trials, have one outcome variable designated as "primary" and, accordingly, test one main
344 hypothesis, typically using the criterion of $p < .05$. Moreover, methodologically strong studies
345 are pre-registered, which eliminates concerns about "outcome switching" (i.e., replacing the
346 primary outcome of interest if it did not reach statistical significance with a different one that
347 did) or selective reporting (i.e., only reporting the outcome that happened to reach the threshold
348 of statistical significance out of a larger set of tested outcomes). However, in several domains of
349 research, including studies investigating the effects of HIIT, pre-registration remains rare, and
350 researchers report results pertaining to numerous dependent variables, each tested using the
351 criterion of $p < .05$. This scenario is problematic insofar as it can raise the risk of Type I errors
352 (or "false positives"), namely rejecting the null hypothesis when it is true.

353 Besides pre-registration, it is important for the tested hypotheses to be precise (e.g., "it is

354 hypothesized that HIIT will improve outcome X as measured by test Y because of reason Z").
355 Instead, in the HIIT literature, studies often claim to have demonstrated the "effectiveness" of
356 HIIT relative to control treatments or relative to moderate-intensity continuous exercise (despite
357 a smaller time commitment) by testing imprecise hypotheses that refer to broad concepts (e.g.,
358 cardiorespiratory fitness, endurance performance, muscle enzymes, blood pressure, glucose
359 metabolism, inflammatory parameters, cardiometabolic health). In turn, each of these broad
360 concepts is assessed by several variables (e.g., long lists of different indicators of
361 cardiorespiratory fitness, endurance performance, muscle enzymes, and so on). If researchers
362 explicitly follow a "conjunction" approach [78], they need to reject all the constituent null
363 hypotheses (e.g., one for each of the multiple inflammatory parameters) in order to claim that
364 they rejected the joint null hypothesis (i.e., that HIIT has a stronger anti-inflammatory effect, in
365 general, than moderate-intensity continuous exercise). The "conjunction" approach, because of
366 the nature of the joint null hypothesis (i.e., all constituent tests must be significant), gives
367 researchers only a single opportunity to reject the joint null hypothesis at the prespecified level
368 of α (i.e., 5%) and, therefore, despite entailing multiple tests, it does not raise the overall risk of a
369 Type I error. On the other hand, the "conjunction" approach is characterized by low statistical
370 power because researchers would fail to reject the joint null hypothesis if even one of the
371 constituent tests yields a non-significant result. The low statistical power is the likely reason why
372 the "conjunction" approach is rarely encountered in the research literature.

373 In contrast, in the "disjunction" approach, it is only necessary to reject one of multiple
374 constituent null hypotheses in order for researchers to be able to claim that they have rejected the
375 joint null hypothesis [78]. For example, researchers may conclude that HIIT benefits "muscle
376 enzymes" (or "cardiometabolic health" or "arterial stiffness" or "cytokines") if only one or two of

377 the variables that make up this broad category, out of a larger set of tested variables, showed
378 significant results in the expected direction. Consequently, the "disjunction" approach increases
379 the risk of Type I error because researchers have multiple opportunities to *incorrectly* reject the
380 joint null hypothesis (i.e., each test of a constituent null hypothesis is also an opportunity to
381 reject the joint null hypothesis).

382 For two independent events, the probability of observing both of these events together is
383 given by the product of their (separate) probabilities. Therefore, if the probability of making a
384 Type I error is $\alpha = .05$, the probability of *not making* a Type I error (i.e., erroneously rejecting
385 the null hypothesis when it is true) on two independent simultaneous tests would be given by $(1-$
386 $\alpha) \times (1-\alpha) = (1-\alpha)^2 = (1-.05)^2 = .9025$. Conversely, the probability of *making* a Type I error
387 would be given by $1-(1-\alpha)^2 = 1 - .9025 = .0975$. Therefore, more broadly, the formula for the
388 inflation of the Type I error rate due to conducting multiple independent probability tests, often
389 referred to as the Šidák equation, is $\alpha^* = 1-(1-\alpha)^M$, where α^* is the inflated value of α as a result
390 of conducting multiple independent tests, α is the conventionally defined probability of
391 committing a Type I error (typically, $\alpha = .05$), and M is the number of independent probability
392 tests conducted at the level of α [79-81].

393 Applying this formula, one finds, for example, that conducting 14 independent tests
394 following the "disjunction" approach results in $\alpha = .51$, namely more than 10 times the nominal
395 rate of .05. This means that, if 14 independent tests were to be conducted, one should expect the
396 probability of making at least one Type I error to be greater than .50. According to a statistical
397 textbook: "It is especially important to realize that failing to control for multiple testing may play
398 a major role in contributing to a disappointing failure rate in attempts to replicate published
399 studies" (p. 216) [82].

400 As noted, the aforementioned formula relies on the simplifying assumption that the
401 multiple probability tests are independent of each other. This assumption, however, is usually
402 false in practice since, in a common example, several variables within the same data set may
403 examine various facets of the same phenomenon (e.g., different parameters of glucose
404 metabolism, immune function, or health-related quality of life), and will, therefore, probably be
405 intercorrelated. To account for this dependence, researchers have proposed variations of the
406 Šidák equation [83-86]. For example, an approach that originated in the field of genetics [87,88]
407 suggests that, when conducting 14 tests, instead of α rising to .51 when the tests are independent,
408 α would rise to .48, .42, and .32 when the variables are intercorrelated $r = .30$, $r = .50$, and $r =$
409 $.70$, respectively. Thus, while the formula $\alpha^* = 1-(1-\alpha)^M$ represents only the "worst-case
410 scenario," it is nevertheless a useful reminder of the possible deleterious consequences of
411 conducting multiple tests without consideration of the inflation of the Type I error rate.

412 With pre-registration still being a rarity in exercise science [63], there is no guarantee that
413 the dependent variables listed in an article represent a complete accounting of all the variables
414 measured or analyzed. Even with this caveat in mind, it is common in the HIIT literature to
415 encounter studies that follow the "disjunction" approach, hypothesizing joint null hypotheses,
416 each consisting of numerous constituent tests, each tested at $p < .05$ [89-92]. This practice can
417 increase the risk of Type I error to high levels (see Figure 2), even compared to other research
418 within exercise science [55], thus raising serious concerns about the validity and reproducibility
419 of any reported effects.

420 **3.1.2. Sampling variability and the instability of p values**

421 To compound the problem of multiplicity described in the previous section, the samples
422 used in the HIIT literature tend to be small (e.g., with as few as 5 individuals per group). The

423 combination of long lists of dependent variables and small samples creates a statistical "perfect
424 storm," a recipe for non-replicable science [43,44,46,93]. Due to sampling variability, small
425 samples produce highly volatile and imprecise estimates of the "true" population values (e.g.,
426 means, standard deviations, intermean differences, and p values). The combination of instability
427 and imprecision with an extremely lax criterion for determining "statistical significance," given a
428 large enough number of tests, essentially guarantees two outcomes: (i) at least some of the tests
429 will cross the liberal threshold of "statistical significance" and (ii) these findings will have a high
430 likelihood of being non-replicable in different samples.

431 The small samples have occasionally been justified on the basis of the argument that the
432 studies are "pilot" trials that were "not designed to be powered to detect statistically significant
433 differences in small or moderate effects" (p. 2072) [94]. Instead, their purpose is portrayed as
434 estimating "the magnitude of effect to lay the foundation for a fully powered efficacy trial" (p.
435 2072). It should be emphasized, however, that this rationale, although commonly encountered, is
436 flawed, due to the inability of small-sample studies to accurately estimate population parameters
437 [95,96]. This lack of precision can lead to considerable over- or underestimations of the true
438 effect size, with potentially devastating consequences for the design of subsequent larger trials.

439 As noted earlier (Section 2), although some researchers operate under the assumption that
440 a finding of $p < .05$ entails 95% confidence that the same result would re-occur in a subsequent
441 replication study, this is not the case. This misconception has been termed the "replication
442 fallacy" or "replication delusion" [61]. In actuality, following an initial finding of $p < .05$, a
443 subsequent (hypothetical) "perfect" replication study drawing an equal number of participants
444 from the same population has only about 50% chance of resulting in a finding of $p < .05$ with the
445 intergroup difference in the same direction [43]. Based on an empirical analysis of 45,955

446 observed effects derived from the Cochrane Database of Systematic Reviews, van Zwet and
447 Goodman [97] put the estimate considerably lower, at 29%. Many researchers may find these
448 figures surprising, despite numerous relevant warnings having been issued in applied literatures,
449 including in psychology [46], physiology [76,98,99], medicine [93,100], and pharmacology
450 [101].

451 In an effort to understand the implications of p values for replication, statisticians have
452 been analyzing the behavior of p values under various conditions, including different
453 hypothetical population effect sizes, the level of α , and sample size [43,102-105]. These efforts
454 have resulted in formulas that enable researchers to calculate the probability of obtaining
455 statistically significant results (e.g., $p < .05$) in subsequent replication studies [46]. One
456 realization that has emerged from these investigations is that sampling variability renders p
457 values extremely unstable and, therefore, an unreliable basis for drawing inferences about
458 experimental effects in most applied-research contexts (given typical effect sizes and sample
459 sizes), especially inferences regarding the replicability of findings [44,46,106].

460 To illustrate the implications for the HIIT literature, we examined the 48-study database
461 used in a meta-analysis by Mattioni Maturana et al. [65], which concluded that HIIT "was
462 superior to [moderate-intensity continuous training] in improving $VO_2\text{max}$ " (p. 559). In this
463 meta-analysis, the median sample size was $N = 10$ per group, and the pooled effect size for
464 $VO_2\text{max}$ (i.e., the most extensively studied outcome) in comparison to moderate-intensity
465 continuous training was $d = 0.40$. Assuming that the pooled effect size approximates the "true"
466 population effect size δ , the combination of these two numbers results in a noncentrality
467 parameter $z = \delta\sqrt{(N/2)} = .40\sqrt{(10/2)} = .894$, which corresponds to an expected p value of .371
468 (the observed mean p value was slightly lower, at .323, for reasons that will be explained in

469 Section 3.1.4).

470 Under these conditions ($N = 10$ per group, $\alpha = .05$, $\delta = 0.40$), statistical power ($1-\beta$) is
471 only .14 (i.e., 14% of p values are expected to be below .05), much lower than the .80
472 conventionally considered adequate. As shown in Figure 3, while 80% of the studies with $1-\beta =$
473 .81 will yield p values of .047 or less, 80% of the studies with $1-\beta = .14$ will yield p values of
474 .707 or less (which also means that 20% of studies will yield p values higher than .707). Indeed,
475 39 of the 48 p values (81.25%) associated with the studies in the meta-analysis by Mattioni
476 Maturana et al. [65] were lower than .707, whereas 9 of 48 (18.75%) were larger than .707.

477 As a demonstration of the volatility of p values one can expect from this combination of
478 effect sizes and sample sizes, Figure 4 shows that the 48 p values related to $VO_2\max$ [65]
479 covered the range from $p = .00000004$ to $p = 1.000$, and effect sizes exhibited an astounding
480 range of 5.33 standard deviations, from -0.74 to +4.59. In other words, assuming that the effect
481 size of the phenomenon under investigation is in the range between small and medium,
482 attempting to study it with approximately 10 participants per group can lead to any outcome [46].

483 Moreover, as noted earlier and illustrated in Figure 5 and Table 1, if an initial study
484 yields $p < .05$, there is a 50% chance that a subsequent replication will also yield $p < .05$,
485 regardless of whether the population effect size is considered "known" or "unknown." However,
486 if the initial study yields a p value of .371 (i.e., the p value expected from studies with the
487 characteristics of those in the meta-analysis by Mattioni Maturana et al. [65]), the probability that
488 a subsequent replication would yield $p < .05$ is only 14.6%. In other words, 85.4% of direct and
489 exact replications (i.e., without any changes to research protocols, including sample size) would
490 likely yield $p > .05$. Moreover, as noted by Cumming [46] and shown in Table 1, to have 90%
491 confidence that a replication would yield $p < .05$, the initial study would have to produce $p <$

492 .00054.

493 As shown in Table 2 and Figure 6, the p intervals are extremely wide. The two-sided p
 494 interval, from the 10th to the 90th percentile, extends from .006 to .828, whereas the one-sided p
 495 interval from zero to the 80th percentile extends to .662. This means that 80% of replication two-
 496 tail p values would fall between .006 and .828 or between .000 and .662. Indeed, 85.42% of the
 497 two-tail p values associated with the studies in the meta-analysis by Mattioni Maturana et al. [65]
 498 were between .006 and .828, and 79.17% were between .000 and .662. For comparison (see
 499 Table 2), in a hypothetical literature in which one can expect a study to yield $p = .001$, the two-
 500 sided p interval for a replication study, from the 10th to the 90th percentile, extends from
 501 .0000005 to .139, whereas the one-sided p interval from zero to the 80th percentile extends to
 502 .036 (or to .018 in the case of a one-tail test).

503 3.1.3. Positive predictive value and false positive risk

504 Positive predictive value (PPV) is defined as the probability that a "positive" research
 505 finding (e.g., $p < .05$) represents a true effect (i.e., that the finding is a true positive). PPV can be
 506 estimated by the formula [107,108]:

$$PPV = \frac{(1 - \beta)R}{(1 - \beta)R + \alpha}$$

507 where $1 - \beta$ is statistical power, R indicates the prestudy odds (i.e., the odds that an effect is indeed
 508 non-null prior to the study being conducted, based on prior evidence), and α is the probability of
 509 a Type I error. Although R is difficult to estimate, the highest value one can reasonably assume
 510 when there are no prior studies on a given topic is 50% (i.e., a 50-50 chance). Even in the
 511 unrealistic scenario of $R = .50$, using the above formula shows, for example, that conducting 19,
 512 23, 32, or 41 independent tests in underpowered studies (e.g., $1 - \beta = .14$) will result in only 7-

513 10% probability of a true positive (see Figure 7). Under the more realistic scenarios of 1-in-4 or
 514 1-in-5 odds (i.e., $R = .25$ or $.20$), the probability of a true positive drops to 3-5%.

515 As noted in the previous section, in the meta-analysis by Mattioni Maturana et al. [65],
 516 the median sample size was 10 per group (the mean was 13.2) and the pooled effect was $d =$
 517 0.40. As shown in Figure 8, assuming that this effect size approximates the "true" population
 518 effect (although this is likely an overestimate for reasons explained in Section 3.1.5), the median
 519 study exhibited only 14% statistical power (the mean of 16% was slightly higher due to one
 520 study with 75% power). This level of power is even lower than the median power of 21%
 521 highlighted as undermining the reliability of neuroscience [107]. Researchers have found that
 522 between 43% and 57% of studies in different domains of biomedicine have statistical power in
 523 the 0–20% range [109]. Of the 48 studies on VO_{2max} included in the Mattioni Maturana et al.
 524 [65] meta-analysis, considering the pooled effect size of $d = 0.40$ as the effect size of interest, 42
 525 (88%) had statistical power in the 0–20% range and all but one (47 of 48, or 98%) were in the 0–
 526 33% range. The combination of the Type I error rate (α) being allowed to escalate and the
 527 extraordinarily small (i.e., severely underpowered) studies can easily (i.e., in common, entirely
 528 realistic scenarios) lead to false discovery rates that approach 100%.

529 A complementary way to think of this problem is in terms of the False Positive Risk
 530 (FPR), namely the probability that a "significant" result (e.g., $p < .05$) represents a false positive.
 531 The FPR can be estimated by the formula [60]:

$$FPR = \frac{p(1 - R)}{p(1 - R) + (1 - \beta)R}$$

532 where p is the p value of a study, R indicates the prestudy odds (i.e., the odds that an effect is
 533 indeed non-null prior to the study being conducted, based on prior evidence), and $1 - \beta$ is the

534 statistical power of the study. The FPR is associated with efforts [40-42], reviewed in Section 2,
 535 to associate the p value from a single study to the lower bound of the long-run risk of Type I
 536 error (α). Applying the formula to the studies on VO_{2max} that were included in the Mattioni
 537 Maturana et al. [65] meta-analysis, and assuming that $R = .50$, shows that only 3 of the 48 studies
 538 produced FPR lower than .05 (see Figure 9). Given their low level of statistical power (median
 539 .155, mean .169), even under the unrealistic assumption of $R = .50$, the FPR of the 13 studies that
 540 produced $p < .05$ was as high as .245, with a mean of .130 and a median of .123 (recall that the
 541 risk of Type I error associated with $p = .05$ has been estimated as at least .289).

542 **3.1.4. Excess of "significant" results**

543 Assuming that the null hypothesis is false (e.g., that there is a difference between HIIT
 544 and moderate-intensity continuous training in terms of improving VO_{2max}), and the effect size is
 545 $\delta = 0.40$, samples of 10 per group are expected to reject the false null hypothesis in only 14% of
 546 the cases (i.e., statistical power of 14%). Instead, as shown in Figure 10, 13 of the 48 studies
 547 (27.1%) included in the meta-analysis by Mattioni Maturana et al. [65], nearly double the
 548 expected rate, produced results with $p < .05$.

549 This rate indicates an "excess of significant findings" according to the test proposed by
 550 Ioannidis and Trikalinos [110]. This is a χ^2 statistic calculated as:

$$A = [(O-E)^2/E + (O-E)^2/(n-E)]$$

551 where O is the number of studies reporting "statistically significant" results ($p < .05$), E is the
 552 sum of the levels of statistical power in all the studies in the sample to detect the population
 553 effect size (assumed here to equal the pooled effect size from the meta-analysis, namely $d =$
 554 0.40), and n is the number of studies in the sample. For the studies in the meta-analysis by
 555 Mattioni Maturana et al. [65], E is 7.851, $O = 13$, and $n = 48$. Therefore, $\chi^2(1) = 4.038$, $p = .044$,

556 indicating the presence of an excessive proportion of "statistically significant" results.

557 Various mechanisms may account for this phenomenon [111]. One category includes
558 "researcher degrees of freedom" [112], some of which may be questionable (e.g., "*p*-hacking,"
559 selective outcome reporting, selective removal of data points, failing to account for multiplicity)
560 and some of which may reflect publication bias (e.g., the "file drawer" problem, namely the low
561 probability of studies reporting non-significant results being accepted for publication) [113].

562 **3.1.5. "Winner's curse"**

563 An additional problem, named "winner's curse" [114,115], emerges from underpowered
564 studies. The "winner's curse" refers to the fact that, when an underpowered study happens to
565 correctly reject a null hypothesis, the estimate of the magnitude of the effect derived from such a
566 study will likely be exaggerated. This is because, for a result to satisfy the criterion of statistical
567 significance (even the uncorrected $p < .05$) in an underpowered study, the effect will have to be
568 unusually large. Young et al. [115] described the problem as follows:

569 The average result from multiple studies yields a reasonable estimate of a "true"
570 relationship. However, the more extreme, spectacular results (the largest treatment
571 effects, the strongest associations, or the most unusually novel and exciting biological
572 stories) may be preferentially published. Journals serve as intermediaries and may suffer
573 minimal immediate consequences for errors of over- or mis-estimation, but it is the
574 consumers of these laboratory and clinical results (other expert scientists; trainees
575 choosing fields of endeavour; physicians and their patients; funding agencies; the media)
576 who are "cursed" if these results are severely exaggerated—overvalued and
577 unrepresentative of the true outcomes of many similar experiments (p. 1418).

578 The "winner's curse" can be shown by simulation, following the procedure proposed by

579 Colquhoun [116]. If we consider the pooled effect size reported by Mattioni Maturana et al. [65],
580 namely $d = 0.40$, and run 100,000 simulated "experiments" by drawing random samples of 100
581 per group from populations designed to differ by $d = 0.40$ (i.e., experiments with 80% statistical
582 power), we find that (i) consistent with the theoretical power level of 80.36%, 80.38% of the
583 comparisons satisfy the $p < .05$ criterion of statistical significance, and (ii) importantly, the
584 average observed effect size is $d = 0.45$, which approximates the given effect size of $d = 0.40$.
585 On the other hand, if one runs 100,000 simulated experiments with the same effect size but
586 sample sizes of 10 per group, namely the median sample size of the 48 studies on VO₂max
587 included in the meta-analysis by Mattioni Maturana et al. [65], (i) the statistical power of 13.66%
588 approximates the theoretical value of 13.55% but (ii) the average observed effect size is highly
589 exaggerated, namely $d = 1.04$ instead of the given $d = 0.40$ (see Figure 11). Indeed, after
590 excluding an apparent outlier with a nearly fivefold effect size [66], the average effect size of the
591 remaining 12 studies on VO₂max in the meta-analysis by Mattioni Maturana et al. [65] that
592 produced $p < .05$ was 1.01. In general, larger sample sizes enable the estimation of the
593 population effects with greater precision, whereas small samples increase the risk of greatly
594 exaggerated estimates of effects.

595 **3.1.6. Accuracy of population estimates**

596 Davis-Stober and Dana [117] have proposed an index of the accuracy of population
597 estimates produced by the conventional method of ordinary least squares (used in most of the
598 commonly employed statistical tests, including tests of comparisons between sample means)
599 compared against a "benchmark" method of estimation that uses random estimates for both the
600 direction and the magnitude of treatment effects (called "random least squares"). The index,
601 called the v-statistic, can range from zero to one, with a value of one indicating that the

602 conventional method of estimation (ordinary least squares) is consistently more accurate than the
603 random method, and a value of zero indicating that the random method of estimation is
604 consistently more accurate than ordinary least squares. The values of the v -statistic are
605 influenced by (i) the sample sizes, (ii) the magnitude of the effect being investigated, and (iii) the
606 number of parameters that need to be estimated (i.e., two means in the case of a t -test).

607 Preempting the criticism that comparing the accuracy of statistical tests against a "benchmark" of
608 random guessing sets a meaninglessly "low bar," Davis-Stober and Dana [117] wrote:

609 If one's estimates are less accurate than our guessing benchmark more than half of the
610 time, there is little point in using them to establish treatment effects. As low as this hurdle
611 may seem, we show that $v < .5$, or even $v = 0$, can happen surprisingly often, particularly
612 when researching effect sizes conventionally categorized as small and medium (p. 6)

613 This is precisely the scenario encountered in the HIIT literature: small- to medium-size
614 effects are being studied with small samples. Therefore, to gauge the accuracy of estimates
615 derived from the studies included in the meta-analysis by Mattioni Maturana et al. [65],
616 comparing the effects of HIIT and moderate-intensity continuous exercise on $VO_2\max$, the v -
617 statistic for each study was calculated following the computational method outlined by Lakens
618 and Evers [118]. The average v -statistic was .124 and the median was .000. Nearly all studies (46
619 of 48, or 96%) had values of the v -statistic below .500, and more than half (28 of 48, or 58%)
620 had a v -statistic of zero (see Figure 12). In the words of Lakens and Evers [118], "obviously, if a
621 random estimator is more accurate than the estimator based on the observed data (indicated by a
622 v -statistic smaller than .5), a study does not really reduce the uncertainty about whether the
623 hypothesis is true" (p. 283).

624 **3.1.7. Summary**

625 In summary, when judged by conventional statistical standards, most studies
626 investigating the effects of HIIT on fitness or health have limited informational yield. This is
627 because they are examining small-to-medium effects with small samples, and commonly test a
628 plethora of dependent variables. Estimates of small-to-medium effects derived from small,
629 underpowered studies are characterized by such imprecision and volatility that, given a large
630 enough number of tests, some will probably cross the conventional threshold of statistical
631 significance. Such "statistically significant" results will likely reflect chance and, therefore, entail
632 a low probability of replication. In addition, even if they represent true effects, such results likely
633 overestimate the magnitude of the underlying effects.

634 **3.2. The "Is As Effective As" Problem**

635 As noted in Section 2, statisticians commonly emphasize that "absence of evidence is not
636 evidence of absence" [53,54]. The principle behind this motto is that $p > .05$ (i.e., "absence of
637 evidence") provides no indication that the null effect, namely $\mu_1 - \mu_2 = 0$, is the most likely result
638 (i.e., "evidence of absence"). In other words, finding $p > .05$ for a comparison between two
639 sample means (such as the mean of a group participating in HIIT and a group participating in
640 moderate-intensity continuous exercise training) only permits a researcher to decide not to reject
641 the null hypothesis. Such a result cannot be taken as a basis for *accepting* the null hypothesis
642 (i.e., to conclude that there is "no difference" or that the two treatments being compared have
643 effects that are "same," "equal," "similar," "equivalent," or "comparable").

644 Establishing the "equivalence" of two interventions requires a different hypothesis,
645 different design, different power calculations, and a different statistical approach [50-52]. An
646 equivalence study begins with the difficult decision of determining a difference between the
647 treatments that represents the smallest effect size of interest (e.g., smaller than any effect that can

648 be considered clinically relevant, meaningful, or worthwhile). Then, the null hypothesis is
649 formulated, stating that the difference between the two treatment means, or part of its
650 surrounding confidence interval, falls outside the prespecified margin (i.e., suggesting that the
651 treatments may not be equivalent, or one may be meaningfully more effective than the other).
652 The alternative hypothesis would be that the difference between the treatments, and its
653 surrounding confidence interval, are within the prespecified margin (i.e., that the treatments are
654 equivalent, or one is as effective as the other). Power calculations for an equivalence study are
655 based on the largest treatment difference considered to be practically irrelevant or
656 inconsequential. The hypothesis of equivalence can be tested by specialized procedures, such as
657 the "two one-sided tests" (TOST) method [119-121].

658 Most researchers carefully avoid the use of the adjectives "similar" or "comparable" (let
659 alone "equal" or "same") to describe treatment means following a finding of $p > .05$. This is
660 because a very common scenario is that tests fail to reject the null hypothesis, even though it is
661 false, because of low statistical power (e.g., having too few participants to detect an effect given
662 the magnitude of that effect). Yet, the HIIT literature contains numerous claims that various HIIT
663 protocols have "similar" or "comparable" effects to more time-consuming moderate-intensity
664 continuous exercise. Invariably, these claims are made on the basis of findings of $p > .05$ from
665 studies that are underpowered to detect small ($d = 0.20$, requiring $N = 394$ per group), medium (d
666 $= 0.50$, requiring $N = 64$ per group), or even large effects ($d = 0.80$, requiring $N = 26$ per group).
667 As noted earlier, of the 48 studies included in the Mattioni Maturana et al. meta-analysis [65]
668 comparing HIIT to moderate-intensity continuous exercise on VO_{2max} , all but one (47 of 48, or
669 98%) had statistical power in the 0–33% range. Examples of claims made on the basis of
670 underpowered studies include claims of "equal" changes across a wide range of physiological

671 parameters (samples of 8 and 8) [92], "similar" changes in aerobic capacity (samples of 7 and 7)
672 [122], "similar" metabolic adaptations (samples of 10 and 10) [89], "similar" changes in arterial
673 stiffness (samples of 10 and 10) [123], "similar" cardiometabolic changes (samples of 9, 10, and
674 6) [90], "similar" cardiorespiratory adaptations in patients with heart failure (samples of 8 and 8)
675 [124], "similar" changes in body composition and fitness (samples of 16, 16, and 14) [125],
676 "similar" muscular and performance changes (samples of 8 and 8) [126], and "similar"
677 enjoyment and adherence (samples of 9 and 8) [127]. Likewise, such claims are made on the
678 basis of findings of $p > .05$ from studies using within-subject designs that are also underpowered
679 to detect small ($d = 0.20$, requiring $N = 199$), medium ($d = 0.50$, requiring $N = 34$) or even large
680 effects ($d = 0.80$, requiring $N = 15$). Examples include claims of "similar" adaptations in
681 signaling molecules associated with mitochondrial biogenesis ($N = 10$) [128], "similar"
682 mitochondrial function ($N = 8$) [129], "similar" 24-hour oxygen consumption ($N = 8$) [130],
683 "similar" energy expenditure ($N = 9$) [131], "similar" increases in serum brain-derived
684 neurotrophic factor ($N = 8$) [132], and "similar" enjoyment levels ($N = 7$ [133]; $N = 11$ [134]). To
685 reiterate the essential point, claims of "similar" or "comparable" effects are unjustified on the
686 basis of "non-significant" comparisons between means ($p > .05$). Claims of "similar" or
687 "comparable" effects can only be justified if appropriate hypotheses and associated tests (i.e., of
688 equivalence or non-inferiority) are used [119-121].

689 **3.2.1. Poor reporting of power calculations**

690 By using $p > .05$ as a criterion for establishing equivalence, there is no end to the
691 extraordinary discoveries that researchers can claim. One common approach has been using
692 severely underpowered comparative studies in conjunction with the $p > .05$ criterion in a race to
693 discover the smallest duration or amount of exercise that can still be claimed to be "as effective

694 as" (or "similar" or "comparable" to) either "traditional" HIIT or moderate-intensity continuous
695 exercise. These minimalist forms have been termed "low-volume HIIT," "very low volume
696 HIIT," or "reduced exertion HIIT," among other labels.

697 To illustrate the problems associated with this approach, we examined the studies
698 included in a recent systematic review of "low-volume HIIT," which concluded that it "can
699 induce similar, and at times greater, improvements in cardiorespiratory fitness, glucose control,
700 blood pressure, and cardiac function when compared to more traditional forms of aerobic
701 exercise training including high-volume HIIT and moderate intensity continuous training, despite
702 requiring less time commitment and lower energy expenditure" (p. 1013) [135]. This is a
703 remarkable claim because "low-volume HIIT" was said to differ from regular HIIT solely by
704 entailing a lower total duration of high-intensity intervals (< 15 min). Otherwise, the two
705 modalities of training were said to share common features (e.g., intensity of 80–100% $\text{VO}_{2\text{max}}$
706 or HR_{max} , duration of each high-intensity interval of 1–4 min, work-to-rest ratio of 1:1 to 1:2).
707 In other words, the review concluded that, contrary to conventional wisdom, doing less exercise
708 is "as effective as" (or, remarkably, even "more effective than") doing more exercise while
709 holding other important aspects of the exercise "dose" constant.

710 The review was based on 11 studies (see Table 3) and used the adjective "comparable" to
711 describe the results of the comparisons between the minimalist versions of HIIT to the
712 comparator groups in 9 of the 11 cases [135]. Predictably, the studies had the common
713 denominator of being underpowered (sample size range: 5 to 22 per group, mean: 13.5, mode:
714 12). Using a two-tail test, a two-group comparative study with $N = 12$ per group has 7.6%,
715 21.6%, and 46.6% statistical power to reject a small ($d = 0.20$), medium ($d = 0.50$), and large (d
716 $= 0.80$) false null hypothesis, respectively.

717 Researchers might wonder how this is possible since item 7a of the CONSORT checklist
718 explicitly states that authors must explain "how sample size was determined" [147]. Given the
719 sample size range of 5–22 per group, it is unsurprising that the claimed adequacy of the sample
720 size could not be verified in any of the 11 studies. In four, no information was provided for how
721 the sample size was determined. In the remaining studies, the irregularities ranged from not
722 providing complete information (e.g., not stating the anticipated effect size), citing nonverifiable
723 or incorrect information (e.g., citing effect sizes for within-group changes from previous studies
724 but aiming to conduct between-group comparisons), citing the effect size from an early study
725 [66] that has been identified as an outlier [148], to reporting the required information but
726 claiming that the sample size needed to be only a fraction of what the calculations indicated in
727 order to reach the desired level of statistical power. As one example:

728 Based on a meta-analysis that compared HIIT with continuous endurance training on
729 maximal oxygen uptake ($VO_2\text{max}$) improvements in adults, the estimated standardized
730 mean difference (Cohen's d) between HIIT and [moderate-intensity continuous training]
731 was approximately 0.4. Therefore, it was anticipated that a sample size of 12 participants
732 per group was adequate to detect this difference between groups on our primary outcome
733 (i.e., $VO_2\text{max}$), with a power of 0.8 at an alpha level of 0.05 (pp. 1998–1999) [141].

734 To reach 80% statistical power given an effect size $d = 0.4$ requires 100 participants per
735 group rather than 12. Bonafiglia et al. [149] similarly found that 21 of 27 studies included in a
736 meta-analysis comparing the effects of sprint interval training and continuous training either did
737 not report sample-size calculations or did not provide full information. The reporting of power
738 calculations is suboptimal both in the medical literature [150] and within exercise and sport
739 science [151]. According to Charles et al. [150], only 34% of trials published in medical journals

740 reported all data required to calculate the sample size, had accurate calculations, and were based
741 on accurate assumptions. Of the remaining, 43% did not report all the required parameters to
742 allow readers to verify the calculation, and 5% did not report sample size calculations. Within
743 exercise and sport science, the situation appears worse. An analysis of 120 manuscripts
744 submitted to a prominent disciplinary journal [151] shows that the median sample size was 19.
745 Only 12 of the manuscripts (10%) included any sample-size calculations and, of them, four did
746 not provide a justification for the cited effect size. Similar to the situation in the HIIT literature
747 discussed in this section [135], none of the 12 manuscripts provided all the information required
748 to enable the correct reproduction of the cited sample-size goal (i.e., the statistical test to be
749 conducted, the targeted effect size, the level of α , and the desired level of statistical power). This
750 situation is of grave concern and necessitating urgent change [77].

751 **4. A Crisis of Confidence, a Looming Trainwreck, or an Opportunity for Reform?**

752 Over the past 15 years, the research literature on HIIT has produced some extraordinary
753 claims, which, upon closer inspection, are backed by surprisingly fragile evidence. This
754 phenomenon can be analyzed from several angles. Perhaps the striking discrepancy between the
755 boldness of the claims and the limitations of the experimental evidence is a reflection of a field
756 eager for a scientific breakthrough. As noted in Section 2, journal editors and peer reviewers
757 may, consciously or subconsciously, "apply lower standards" (p. 4) [62] when evaluating
758 manuscripts that purport to report findings that seem highly intriguing or novel. Likewise, the
759 willingness of the press to disseminate, and occasionally amplify, the extraordinary claims
760 surrounding HIIT also suggests that the public at large may be eager for a breakthrough from
761 exercise science, some miraculous discovery that would magnify and accelerate the benefits of
762 exercise while requiring less effort [152].

763 An equally fascinating question pertains to the apparent willingness of exercise science as
764 a research field to enter a state of "suspension of disbelief," accepting and propagating claims
765 that defy conventional wisdom and research choices that directly contradict established
766 methodological and statistical best practices. Like other scientific fields, exercise science will
767 inevitably, sooner or later, have to confront its own crisis of replication and confidence [63].
768 Postponing this conversation will not help avert it. Therefore, it seems ironic that, while a push
769 for more stringent methodologies [112,153] and more responsible reporting [154] is sweeping
770 the scientific landscape, one of the most prominent research lines within exercise science is
771 characterized by a preponderance of studies with questionable statistical standards.

772 In the previous sections, it was shown that most samples in the HIIT literature are small,
773 and thus the studies are underpowered to detect small, medium, or even large effects. This is
774 important because the effect sizes, in most cases (especially when HIIT is compared against
775 moderate-intensity continuous exercise rather than a no-exercise control), are likely to be small.
776 It was also shown that most studies do not have one outcome designated as primary but rather
777 tend to include long lists of dependent variables, all of which are tested at $p < .05$, without
778 consideration for the inflation of α . There is also great flexibility in designs, definitions,
779 outcomes, and analytic approaches, from the definition of HIIT to the selection of variables to
780 represent various domains of physiological function (e.g., metabolism). Moreover, extraordinary
781 claims related to the effectiveness of HIIT, along with claims that HIIT addresses "the most
782 commonly cited reason for not exercising" (p. 212) [155] or "the primary reason for [the] failure
783 to exercise on a regular basis" (p. 61) [156], namely "lack of time," stimulate the interest or
784 curiosity of the public (e.g., the narrative that, contrary to current recommendations, one only
785 needs to exercise for a few seconds per day). The intense interest from the media may encourage

786 or incentivize researchers to produce research results that support compelling narratives but may
787 have low replicability. In particular, claims that smaller and smaller amounts of exercise were
788 found to be "effective" for improving fitness and health are bound to capture the interest of the
789 general public. For example, recent media reports have highlighted that repeated 4-sec spurts of
790 exercise, totaling no more than 2 min per day [157], or a single 3-sec muscular contraction per
791 day [158] have been found to result in "significant" gains in aerobic capacity (by 13%) and
792 muscular strength (by 12%), respectively (based on samples of 11 and 13, respectively).

793 Arguably, there is a striking similarity between the patterns seen in the HIIT literature
794 and what was unfolding in the research field investigating phenomena of behavioral priming
795 within psychology in the 2000s. The literature was being inundated with findings that have been
796 described as "implausible" (p. 13) [159], "spectacular" (p. 19) [160], "fascinating" (p. 20) [161],
797 and "eye-catching and counter-intuitive... the kind of sexy research that popular science writers
798 love to describe" (p. 6) [161]. Failed attempts to replicate several of these widely publicized
799 results led to an ongoing "replication crisis" [162] or "crisis of confidence" [163] in psychology.
800 In response, Nobel laureate Daniel Kahneman wrote an open letter to researchers involved in
801 research on priming, in which he encouraged them to try to remove the question mark that had
802 been attached to their field [164]. He emphasized: "Your problem is not with the few people who
803 have actively challenged the validity of some priming results. It is with the much larger
804 population of colleagues who in the past accepted your surprising results as facts when they were
805 published." Reminding readers that "a posture of defiant denial is self-defeating," Kahneman
806 pointed out what was at stake: "I see a train wreck looming. I expect the first victims to be young
807 people on the job market. Being associated with a controversial and suspicious field will put
808 them at a severe disadvantage in the competition for positions. Because of the high visibility of

809 the issue, you may already expect the coming crop of graduates to encounter problems."

810 Although undertaking the kind of radical reforms advocated by Kahneman is unlikely to
811 be universally appreciated or endorsed, psychology has, to some extent, entered a period of
812 critical self-reflection. Many authors have argued that the replication crisis can be seen as an
813 opportunity for positive change [165-167]. This perspective has grown into a movement [168]
814 that has even been characterized, perhaps optimistically or prematurely, as a "renaissance" [169].
815 The winds of change are reaching other fields, even beyond the social sciences, such as cancer
816 biology and drug development, which are coming to terms with the fact that they, too, are facing
817 a replication crisis [170,171].

818 The replication crisis in psychology offers a potential blueprint for how exercise science
819 could proceed. Arguing that there is no problem is certainly a comforting option but, to echo
820 Kahneman, "a posture of defiant denial is self-defeating." Continuing to overlook the
821 fundamental principles of statistics in pursuit of implausible results that will capture the next
822 headline will predictably lead to poor long-term outcomes. The exorbitant claims in the HIIT
823 literature could serve as a clarion call that should inspire a period of critical self-reflection and
824 positive reform. Recognizing the pitfalls, returning to, and respecting the fundamentals could
825 have a lasting positive influence on the integrity, societal value, and reputation of exercise
826 science.

827 It is, therefore, encouraging that the first signs of reform within exercise science have
828 started to appear. Statistical experts [23,77] and journal editors [76,99,151,172] are making
829 strong cases about the need to improve the quality of research designs and statistical analyses.
830 Newly created organizations, such as the Consortium for Transparency in Exercise Science [63]
831 and the Society for Transparency, Openness, and Replication in Kinesiology, are spearheading

832 educational initiatives aimed at promoting stronger research practices. In psychology, arguably
833 one of the most consequential reform efforts has been the push to expand the practice of study
834 preregistration [173-176]. Therefore, the growing number of journals within exercise science that
835 encourage preregistration and welcome registered reports represents a particularly promising
836 development [177]. Beyond these efforts, curricular reforms will be necessary, with the goal of
837 significantly improving statistical literacy at both the undergraduate and postgraduate levels. At
838 the undergraduate level, courses intended to promote critical appraisal skills, specifically
839 designed for consumers of research information (i.e., future exercise professionals), should be
840 considered a necessity for a field aspiring to fully transition to a model of evidence-based
841 practice. At the postgraduate level, where most students are prospective producers of research
842 information, the teaching of statistical skills should be combined with efforts to cultivate a
843 mindset that welcomes openness and transparency while resisting the "disciplinary incentives" to
844 "favor novelty over replication" (p. 615) [57]. Finally, an important issue that the extraordinary
845 claims surrounding HIIT have brought to the surface is that the field of exercise science must
846 critically reexamine its relationship with the mass media. Researchers, university press offices,
847 and journal editors should also resist the temptation to construct and disseminate media-friendly
848 narratives that are based on statistically questionable or fragile evidence.

849 **References**

- 850 1. Haskell WL. Health consequences of physical activity: understanding and challenges
851 regarding dose-response. *Med Sci Sports Exerc.* 1994;26(6):649-660. [https://doi.org/
852 10.1249/00005768-199406000-00001](https://doi.org/10.1249/00005768-199406000-00001)
- 853 2. Pate RR. Physical activity and health: dose-response issues. *Res Q Exerc Sport.*
854 1995;66(4):313-317. [https://doi.org/ 10.1080/02701367.1995.10607917](https://doi.org/10.1080/02701367.1995.10607917)

- 855 3. Pate RR, Pratt M, Blair SN, Haskell WL, Macera CA, Bouchard C, Buchner D, Ettinger
 856 W, Heath GW, King AC, et al. Physical activity and public health: a recommendation
 857 from the Centers for Disease Control and Prevention and the American College of Sports
 858 Medicine. *JAMA*. 1995;273(5):402-407. [https://doi.org/ 10.1001/jama.273.5.402](https://doi.org/10.1001/jama.273.5.402)
- 859 4. US Department of Health and Human Services. Physical activity and health: a report of
 860 the Surgeon General. Atlanta, Georgia: US Department of Health and Human Services,
 861 Public Health Service, Centers for Disease Control and Prevention, National Center for
 862 Chronic Disease Prevention and Health Promotion; 1996.
- 863 5. Leon AS. Physical activity and cardiovascular health: a national consensus. Champaign,
 864 Illinois: Human Kinetics; 1997.
- 865 6. NIH Consensus Development Panel on Physical Activity and Cardiovascular Health.
 866 Physical activity and cardiovascular health. *JAMA*. 1996;276(3):241-246.
 867 <https://doi.org/10.1001/jama.1996.03540030075036>
- 868 7. Blair SN, LaMonte MJ, Nichaman MZ. The evolution of physical activity
 869 recommendations: how much is enough? *Am J Clin Nutr*. 2004;79(5):913S-920S.
 870 <https://doi.org/10.1093/ajcn/79.5.913S>
- 871 8. Dishman RK, Buckworth J. Increasing physical activity: a quantitative synthesis. *Med*
 872 *Sci Sports Exerc*. 1996;28(6):706-719. [https://doi.org/10.1097/00005768-199606000-](https://doi.org/10.1097/00005768-199606000-00010)
 873 [00010](https://doi.org/10.1097/00005768-199606000-00010)
- 874 9. Troiano RP, Berrigan D, Dodd KW, Mâsse LC, Tilert T, McDowell M. Physical activity
 875 in the United States measured by accelerometer. *Med Sci Sports Exerc*. 2008;40(1):181-
 876 188. <https://doi.org/10.1249/mss.0b013e31815a51b3>
- 877 10. Metzger JS, Catellier DJ, Evenson KR, Treuth MS, Rosamond WD, Siega-Riz AM.

- 878 Patterns of objectively measured physical activity in the United States. *Med Sci Sports*
 879 *Exerc.* 2008;40(4):630-638. <https://doi.org/10.1249/MSS.0b013e3181620ebc>
- 880 11. Tudor-Locke C, Brashear MM, Johnson WD, Katzmarzyk PT. Accelerometer profiles of
 881 physical activity and inactivity in normal weight, overweight, and obese U.S. men and
 882 women. *Int J Behav Nutr Phys Act.* 2010;7:60. <https://doi.org/10.1186/1479-5868-7-60>
- 883 12. Winett RA. Developing more effective health-behavior programs: analyzing the
 884 epidemiological and biological bases for activity and exercise programs. *Appl Prev*
 885 *Psychol.* 1998;7(4):209-224. [https://doi.org/10.1016/S0962-1849\(98\)80025-5](https://doi.org/10.1016/S0962-1849(98)80025-5)
- 886 13. Swain DP, Franklin BA. Comparison of cardioprotective benefits of vigorous versus
 887 moderate intensity aerobic exercise. *Am J Cardiol.* 2006;97(1):141-147.
 888 <https://doi.org/10.1016/j.amjcard.2005.07.130>
- 889 14. O'Donovan G, Shave R. British adults' views on the health benefits of moderate and
 890 vigorous activity. *Prev Med.* 2007;45(6):432-435.
 891 <https://doi.org/10.1016/j.ypmed.2007.07.026>
- 892 15. Haskell WL, Lee IM, Pate RR, Powell KE, Blair SN, Franklin BA, Macera CA, Heath
 893 GW, Thompson PD, Bauman A; American College of Sports Medicine; American Heart
 894 Association. Physical activity and public health: updated recommendation for adults from
 895 the American College of Sports Medicine and the American Heart Association.
 896 *Circulation.* 2007;116(9):1081-1093.
 897 <https://doi.org/10.1161/CIRCULATIONAHA.107.185649>
- 898 16. O'Donovan G, Blazeovich AJ, Boreham C, Cooper AR, Crank H, Ekelund U, Fox KR,
 899 Gately P, Giles-Corti B, Gill JM, Hamer M, McDermott I, Murphy M, Mutrie N, Reilly
 900 JJ, Saxton JM, Stamatakis E. The ABC of Physical Activity for Health: a consensus

- 901 statement from the British Association of Sport and Exercise Sciences. *J Sports Sci.*
 902 2010;28(6):573-591. <https://doi.org/10.1080/02640411003671212>
- 903 17. Burgomaster KA, Hughes SC, Heigenhauser GJ, Bradwell SN, Gibala MJ. Six sessions
 904 of sprint interval training increases muscle oxidative potential and cycle endurance
 905 capacity in humans. *J Appl Physiol.* 2005;98(6):1985-1990.
 906 <https://doi.org/10.1152/jappphysiol.01095.2004>
- 907 18. Coyle EF. Very intense exercise-training is extremely potent and time efficient: a
 908 reminder. *J Appl Physiol.* 2005;98(6):1983-1984.
 909 <https://doi.org/10.1152/jappphysiol.00215.2005>
- 910 19. Thompson WR. Worldwide survey of fitness trends for 2014. *ACSM Health Fitness J.*
 911 2013;17(6):10-20.
- 912 20. Gray SR, Ferguson C, Birch K, Forrest LJ, Gill JM. High-intensity interval training: key
 913 data needed to bridge the gap from laboratory to public health policy. *Br J Sports Med.*
 914 2016;50(20):1231-1232. <https://doi.org/10.1136/bjsports-2015-095705>
- 915 21. Steen RG. Misinformation in the medical literature: what role do error and fraud play? *J*
 916 *Med Ethics.* 2011;37(8):498-503. <https://doi.org/10.1136/jme.2010.041830>.
- 917 22. Viana RB, Naves JPA, Coswig VS, de Lira CAB, Steele J, Fisher JP, Gentil P. Is interval
 918 training the magic bullet for fat loss? A systematic review and meta-analysis comparing
 919 moderate-intensity continuous training with high-intensity interval training (HIIT). *Br J*
 920 *Sports Med.* 2019;53(10):655-664. <https://doi.org/10.1136/bjsports-2018-099928>
- 921 23. Sainani KL, Borg DN, Caldwell AR, Butson ML, Tenan MS, Vickers AJ, Vigotsky AD,
 922 Warmenhoven J, Nguyen R, Lohse KR, Knight EJ, Bargary N. Call to increase statistical
 923 collaboration in sports science, sport and exercise medicine and sports physiotherapy. *Br*

- 924 J Sports Med. 2021;55(2):118-122. <https://doi.org/10.1136/bjsports-2020-102607>
- 925 24. Nickerson RS. Null hypothesis significance testing: a review of an old and continuing
926 controversy. Psychol Methods. 2000;5(2):241-301. [https://doi.org/10.1037/1082-](https://doi.org/10.1037/1082-989x.5.2.241)
927 989x.5.2.241
- 928 25. Sterne JA, Davey Smith G. Sifting the evidence: what's wrong with significance tests?
929 BMJ. 2001;322(7280):226-231. <https://doi.org/10.1136/bmj.322.7280.226>
- 930 26. Christensen R. Testing Fisher, Neyman, Pearson, and Bayes. Am Stat. 2005;59(2):121-
931 126. <https://doi.org/10.1198/000313005X20871>
- 932 27. Lakens D. The practical alternative to the p value is the correctly used p value. Perspect
933 Psychol Sci. 2021;16(3):639-648. <https://doi.org/10.1177/1745691620958012>
- 934 28. Hubbard R, Bayarri MJ. Confusion over measures of evidence (p's) versus errors (α 's) in
935 classical statistical testing. Am Stat. 2003;57(3):171-178.
936 <https://doi.org/10.1198/0003130031856>
- 937 29. Fisher RA. Statistical methods for research workers (5th ed.). Edinburgh: Oliver and
938 Boyd; 1934.
- 939 30. Neyman J, Pearson ES. IX. On the problem of the most efficient tests of statistical
940 hypotheses. Philos Trans R Soc Lond A. 1933;231:289-337.
941 <https://doi.org/10.1098/rsta.1933.0009>
- 942 31. Szucs D, Ioannidis JPA. When null hypothesis significance testing is unsuitable for
943 research: a reassessment. Front Hum Neurosci. 2017;11:390.
944 <https://doi.org/10.3389/fnhum.2017.00390>
- 945 32. Neyman J, Pearson ES. The testing of statistical hypotheses in relation to probabilities a
946 priori. Math Proc Camb Philos Soc. 1933;29(4):492-510.

- 947 <https://doi.org/10.1017/S030500410001152X>
- 948 33. Fisher RA. The design of experiments. Edinburgh: Oliver and Boyd; 1935.
- 949 34. Fisher R. Statistical methods and scientific induction. *J R Stat Soc Series B Methodol.*
950 1955;17:69-78. <https://doi.org/10.1111/j.2517-6161.1955.tb00180.x>
- 951 35. Lehmann EL. The Fisher, Neyman-Pearson theories of testing hypotheses: one theory or
952 two? *J Am Stat Assoc.* 1993;88(424):1242-1249.
953 <https://doi.org/10.1080/01621459.1993.10476404>
- 954 36. Lehmann EL. Fisher, Neyman, and the creation of classical statistics. New York:
955 Springer; 2011. <https://doi.org/10.1007/978-1-4419-9500-1>
- 956 37. Goodman S. A dirty dozen: twelve p-value misconceptions. *Semin Hematol.*
957 2008;45(3):135-140. <https://doi.org/10.1053/j.seminhematol.2008.04.003>
- 958 38. Greenland S, Senn SJ, Rothman KJ, et al. Statistical tests, p values, confidence intervals,
959 and power: a guide to misinterpretations. *Eur J Epidemiol.* 2016;31(4):337-350.
960 <https://doi.org/10.1007/s10654-016-0149-3>
- 961 39. Goodman SN. Toward evidence-based medical statistics. 1: The p value fallacy. *Ann*
962 *Intern Med.* 1999;130(12):995-1004. [https://doi.org/10.7326/0003-4819-130-12-](https://doi.org/10.7326/0003-4819-130-12-199906150-00008)
963 [199906150-00008](https://doi.org/10.7326/0003-4819-130-12-199906150-00008)
- 964 40. Berger JO, Sellke T. Testing a point null hypothesis: the irreconcilability of p values and
965 evidence. *J Am Stat Assoc.* 1987;82(397):112-122.
966 <https://doi.org/10.1080/01621459.1987.10478397>
- 967 41. Sellke T, Bayarri MJ, Berger JO. Calibration of p values for testing precise null
968 hypotheses. *Am Stat.* 2001;55(1):62-71. <https://doi.org/10.1198/000313001300339950>
- 969 42. Berger JO. Could Fisher, Jeffreys and Neyman have agreed on testing? *Stat Sci.*

- 970 2003;18(1):1-32. <https://doi.org/10.1214/ss/1056397485>
- 971 43. Goodman SN. A comment on replication, p-values and evidence. *Stat Med*.
972 1992;11(7):875-879. <https://doi.org/10.1002/sim.4780110705>
- 973 44. Halsey LG, Curran-Everett D, Vowler SL, Drummond GB. The fickle P value generates
974 irreproducible results. *Nat Methods*. 2015;12(3):179-185.
975 <https://doi.org/10.1038/nmeth.3288>
- 976 45. Lazzeroni LC, Lu Y, Belitskaya-Lévy I. Solutions for quantifying p-value uncertainty
977 and replication power. *Nat Methods*. 2016;13(2):107-108.
978 <https://doi.org/10.1038/nmeth.3741>
- 979 46. Cumming G. Replication and p intervals: p values predict the future only vaguely, but
980 confidence intervals do much better. *Perspect Psychol Sci*. 2008;3(4):286-300.
981 <https://doi.org/10.1111/j.1745-6924.2008.00079.x>
- 982 47. Killeen PR. An alternative to null-hypothesis significance tests. *Psychol Sci*.
983 2005;16(5):345-353. <https://doi.org/10.1111/j.0956-7976.2005.01538.x>
- 984 48. Lecoutre B, Lecoutre MP, Poitevineau J. Killeen's probability of replication and
985 predictive probabilities: how to compute, use, and interpret them. *Psychol Methods*.
986 2010;15(2):158-171. <https://doi.org/10.1037/a0015915>
- 987 49. Sanabria F, Killeen PR. Better statistics for better decisions: rejecting null hypotheses
988 statistical tests in favor of replication statistics. *Psychol Sch*. 2007;44(5):471-481.
989 <https://doi.org/10.1002/pits.20239>
- 990 50. Amrhein V, Greenland S, McShane B. Scientists rise up against statistical significance.
991 *Nature*. 2019;567(7748):305-307. <https://doi.org/10.1038/d41586-019-00857-9>
- 992 51. Hoekstra R, Finch S, Kiers HA, Johnson A. Probability as certainty: dichotomous

- 993 thinking and the misuse of p values. *Psychon Bull Rev.* 2006;13(6):1033-1037.
 994 <https://doi.org/10.3758/bf03213921>
- 995 52. Smith RJ. $P > .05$: The incorrect interpretation of "not significant" results is a significant
 996 problem. *Am J Phys Anthropol.* 2020;172(4):521-527. <https://doi.org/10.1002/ajpa.24092>
- 997 53. Alderson P. Absence of evidence is not evidence of absence. *BMJ.* 2004;328(7438):476-
 998 477. <https://doi.org/10.1136/bmj.328.7438.476>
- 999 54. Altman DG, Bland JM. Absence of evidence is not evidence of absence. *BMJ.*
 1000 1995;311(7003):485. <https://doi.org/10.1136/bmj.311.7003.485>
- 1001 55. Speed HD, Andersen MB. What exercise and sport scientists don't understand. *J Sci Med*
 1002 *Sport.* 2000;3(1):84-92. [https://doi.org/10.1016/s1440-2440\(00\)80051-1](https://doi.org/10.1016/s1440-2440(00)80051-1)
- 1003 56. Vadillo MA, Konstantinidis E, Shanks DR. Underpowered samples, false negatives, and
 1004 unconscious learning. *Psychon Bull Rev.* 2016;23(1):87-102.
 1005 <https://doi.org/10.3758/s13423-015-0892-6>
- 1006 57. Nosek BA, Spies JR, Motyl M. Scientific utopia: II. Restructuring incentives and
 1007 practices to promote truth over publishability. *Perspect Psychol Sci.* 2012;7(6):615-631.
 1008 <https://doi.org/10.1177/1745691612459058>
- 1009 58. Anderson SF. Misinterpreting p: the discrepancy between p values and the probability the
 1010 null hypothesis is true, the influence of multiple testing, and implications for the
 1011 replication crisis. *Psychol Methods.* 2020;25(5):596-609.
 1012 <https://doi.org/10.1037/met0000248>
- 1013 59. Colling LJ, Szucs D. Statistical inference and the replication crisis. *Rev Phil Psychol.*
 1014 2021;12(1):121-147. <https://doi.org/10.1007/s13164-018-0421-4>
- 1015 60. Colquhoun D. The reproducibility of research and the misinterpretation of p-values. *R*

- 1016 Soc Open Sci. 2017;4(12):171085. <https://doi.org/10.1098/rsos.171085>
- 1017 61. Gigerenzer G. Statistical rituals: the replication delusion and how we got there. *Adv*
1018 *Methods Pract Psychol Sci.* 2018;1(2):198-218.
1019 <https://doi.org/10.1177/2515245918771329>
- 1020 62. Serra-Garcia M, Gneezy U. Nonreplicable publications are cited more than replicable
1021 ones. *Sci Adv.* 2021;7(21):eabd1705. <https://doi.org/10.1126/sciadv.abd1705>
- 1022 63. Caldwell AR, Vigotsky AD, Tenan MS, Radel R, Mellor DT, Kreutzer A, Lahart IM,
1023 Mills JP, Boisgontier MP; Consortium for Transparency in Exercise Science (COTES)
1024 Collaborators. Moving sport and exercise science forward: a call for the adoption of more
1025 transparent research practices. *Sports Med.* 2020;50(3):449-459.
1026 <https://doi.org/10.1007/s40279-019-01227-1>
- 1027 64. Bauer N, Sperlich B, Holmberg HC, Engel FA. Effects of high-intensity interval training
1028 in school on the physical performance and health of children and adolescents: a
1029 systematic review with meta-analysis. *Sports Med Open.* 2022;8(1):50.
1030 <https://doi.org/10.1186/s40798-022-00437-8>
- 1031 65. Mattioni Maturana F, Martus P, Zipfel S, Nieß AM. Effectiveness of HIIE versus MICT
1032 in improving cardiometabolic risk factors in health and disease: a meta-analysis. *Med Sci*
1033 *Sports Exerc.* 2021;53(3):559-573. <https://doi.org/10.1249/MSS.0000000000002506>
- 1034 66. Wisløff U, Støylen A, Loennechen JP, Bruvold M, Rognum Ø, Haram PM, Tjønnå AE,
1035 Helgerud J, Slørdahl SA, Lee SJ, Videm V, Bye A, Smith GL, Najjar SM, Ellingsen Ø,
1036 Skjaerpe T. Superior cardiovascular effect of aerobic interval training versus moderate
1037 continuous training in heart failure patients: a randomized study. *Circulation.*
1038 2007;115(24):3086-3094. <https://doi.org/10.1161/CIRCULATIONAHA.106.675041>

- 1039 67. Andreacci JL, LeMura LM, Cohen SL, Urbansky EA, Chelland SA, Von Duvillard SP.
1040 The effects of frequency of encouragement on performance during maximal exercise
1041 testing. *J Sports Sci.* 2002;20(4):345-352. <https://doi.org/10.1080/026404102753576125>
- 1042 68. Halperin I, Pyne DB, Martin DT. Threats to internal validity in exercise science: a review
1043 of overlooked confounding variables. *Int J Sports Physiol Perform.* 2015;10(7):823-829.
1044 <https://doi.org/10.1123/ijsp.2014-0566>
- 1045 69. Midgley AW, Marchant DC, Levy AR. A call to action towards an evidence-based
1046 approach to using verbal encouragement during maximal exercise testing. *Clin Physiol
1047 Funct Imaging.* 2018;38(4):547-553. <https://doi.org/10.1111/cpf.12454>
- 1048 70. Wisløff U, Coombes JS, Rognum Ø. CrossTalk proposal: high intensity interval training
1049 does have a role in risk reduction or treatment of disease. *J Physiol.* 2015;593(24):5215-
1050 5217. <https://doi.org/10.1113/JP271041>
- 1051 71. Khalafi M, Symonds ME. The impact of high-intensity interval training on inflammatory
1052 markers in metabolic disorders: A meta-analysis. *Scand J Med Sci Sports.*
1053 2020;30(11):2020-2036. <https://doi.org/10.1111/sms.13754>
- 1054 72. Solera-Martínez M, Herraiz-Adillo Á, Manzanares-Domínguez I, De La Cruz LL,
1055 Martínez-Vizcaíno V, Pozuelo-Carrascosa DP. High-intensity interval training and
1056 cardiometabolic risk factors in children: a meta-analysis. *Pediatrics.*
1057 2021;148(4):e2021050810. <https://doi.org/10.1542/peds.2021-050810>
- 1058 73. Gerosa-Neto J, Antunes BM, Campos EZ, et al. Impact of long-term high-intensity
1059 interval and moderate-intensity continuous training on subclinical inflammation in
1060 overweight/obese adults. *J Exerc Rehabil.* 2016;12(6):575-580.
1061 <https://doi.org/10.12965/jer.1632770.385>

- 1062 74. Oh S, So R, Shida T, et al. High-intensity aerobic exercise improves both hepatic fat
1063 content and stiffness in sedentary obese men with nonalcoholic fatty liver disease. *Sci*
1064 *Rep.* 2017;7:43029. <https://doi.org/10.1038/srep43029>
- 1065 75. Paahoo A, Tadibi V, Behpoor N. Effectiveness of continuous aerobic versus high-
1066 intensity interval training on atherosclerotic and inflammatory markers in boys with
1067 overweight/obesity. *Pediatr Exerc Sci.* 2021;33(3):132-138.
1068 <https://doi.org/10.1123/pes.2020-0138>
- 1069 76. Gandevia S. Publications, replication and statistics in physiology plus two neglected
1070 curves. *J Physiol.* 2021;599(6):1719-1721. <https://doi.org/10.1113/JP281360>
- 1071 77. Sainani K, Chamari K. Wish list for improving the quality of statistics in sport science.
1072 *Int J Sports Physiol Perform.* 2022;17(5):673-674. <https://doi.org/10.1123/ijsp.2022->
1073 0023
- 1074 78. Rubin M. When to adjust alpha during multiple testing: a consideration of disjunction,
1075 conjunction, and individual testing. *Synthese* 2021;199(3-4):10969–11000.
1076 <https://doi.org/10.1007/s11229-021-03276-4>
- 1077 79. Albers C. The problem with unadjusted multiple and sequential statistical testing. *Nat*
1078 *Commun.* 2019;10(1):1921. <https://doi.org/10.1038/s41467-019-09941-0>
- 1079 80. Forstmeier W, Wagenmakers EJ, Parker TH. Detecting and avoiding likely false-positive
1080 findings: a practical guide. *Biol Rev Camb Philos Soc.* 2017;92(4):1941-1968.
1081 <https://doi.org/10.1111/brv.12315>
- 1082 81. Streiner DL. Best (but oft-forgotten) practices: the multiple problems of multiplicity -
1083 whether and how to correct for many statistical tests. *Am J Clin Nutr.* 2015;102(4):721-
1084 728. <https://doi.org/10.3945/ajcn.115.113548>

- 1085 82. Maxwell SE, Delaney HD, Kelley K. Designing experiments and analyzing data: a model
 1086 comparison perspective. 3rd ed. Oxfordshire, England: Routledge; 2018.
- 1087 83. Blakesley RE, Mazumdar S, Dew MA, Houck PR, Tang G, Reynolds CF 3rd, Butters
 1088 MA. Comparisons of methods for multiple hypothesis testing in neuropsychological
 1089 research. *Neuropsychology*. 2009;23(2):255-264. <https://doi.org/10.1037/a0012850>
- 1090 84. Sankoh AJ, Huque MF, Dubey SD. Some comments on frequently used multiple
 1091 endpoint adjustment methods in clinical trials. *Stat Med*. 1997;16(22):2529-2542.
 1092 [https://doi.org/10.1002/\(sici\)1097-0258\(19971130\)16:22<2529::aid-sim692>3.0.co;2-j](https://doi.org/10.1002/(sici)1097-0258(19971130)16:22<2529::aid-sim692>3.0.co;2-j)
- 1093 85. Vickerstaff V, Omar RZ, Ambler G. Methods to adjust for multiple comparisons in the
 1094 analysis and sample size calculation of randomised controlled trials with multiple
 1095 primary outcomes. *BMC Med Res Methodol*. 2019;19(1):129.
 1096 <https://doi.org/10.1186/s12874-019-0754-4>
- 1097 86. Sankoh AJ, D'Agostino RB Sr, Huque MF. Efficacy endpoint selection and multiplicity
 1098 adjustment methods in clinical trials with inherent multiple endpoint issues. *Stat Med*.
 1099 2003;22(20):3133-3150. <https://doi.org/10.1002/sim.1557>
- 1100 87. Cheverud JM. A simple correction for multiple comparisons in interval mapping genome
 1101 scans. *Heredity*. 2001;87(Pt 1):52-58. <https://doi.org/10.1046/j.1365-2540.2001.00901.x>
- 1102 88. Nyholt DR. A simple correction for multiple testing for single-nucleotide polymorphisms
 1103 in linkage disequilibrium with each other. *Am J Hum Genet*. 2004;74(4):765-769.
 1104 <https://doi.org/10.1086/383251>
- 1105 89. Burgomaster KA, Howarth KR, Phillips SM, Rakobowchuk M, Macdonald MJ, McGee
 1106 SL, Gibala MJ. Similar metabolic adaptations during exercise after low volume sprint
 1107 interval and traditional endurance training in humans. *J Physiol*. 2008;586(1):151-160.

- 1108 <https://doi.org/10.1113/jphysiol.2007.142109>
- 1109 90. Gillen JB, Martin BJ, MacInnis MJ, Skelly LE, Tarnopolsky MA, Gibala MJ. Twelve
1110 weeks of sprint interval training improves indices of cardiometabolic health similar to
1111 traditional endurance training despite a five-fold lower exercise volume and time
1112 commitment. *PLoS ONE*. 2016;11(4):e0154075.
1113 <https://doi.org/10.1371/journal.pone.0154075>
- 1114 91. Robinson E, Durrer C, Simtchouk S, Jung ME, Bourne JE, Voth E, Little JP. Short-term
1115 high-intensity interval and moderate-intensity continuous training reduce leukocyte TLR4
1116 in inactive adults at elevated risk of type 2 diabetes. *J Appl Physiol*. 2015;119(5):508-
1117 516. <https://doi.org/10.1152/jappphysiol.00334.2015>
- 1118 92. Cocks M, Shaw CS, Shepherd SO, Fisher JP, Ranasinghe A, Barker TA, Wagenmakers
1119 AJ. Sprint interval and moderate-intensity continuous training have equal benefits on
1120 aerobic capacity, insulin sensitivity, muscle capillarisation and endothelial
1121 eNOS/NAD(P)H oxidase protein ratio in obese men. *J Physiol*. 2016;594(8):2307-2321.
1122 <https://doi.org/10.1113/jphysiol.2014.285254>
- 1123 93. McGiffin DC, Cumming G, Myles PS. The frequent insignificance of a "significant" p-
1124 value. *J Card Surg*. 2021;36(11):4322-4331. <https://doi.org/10.1111/jocs.15960>
- 1125 94. Locke SR, Bourne JE, Beauchamp MR, Little JP, Barry J, Singer J, Jung ME. High-
1126 intensity interval or continuous moderate exercise: a 24-week pilot trial. *Med Sci Sports
1127 Exerc*. 2018;50(10):2067-2075. <https://doi.org/10.1249/MSS.0000000000001668>
- 1128 95. Albers C, Lakens D. When power analyses based on pilot data are biased: inaccurate
1129 effect size estimators and follow-up bias. *J Exp Soc Psychol*. 2018;74:187-195.
1130 <https://doi.org/10.1016/j.jesp.2017.09.004>

- 1131 96. Kraemer HC, Mintz J, Noda A, Tinklenberg J, Yesavage JA. Caution regarding the use of
1132 pilot studies to guide power calculations for study proposals. *Arch Gen Psychiatry*.
1133 2006;63(5):484-489. <https://doi.org/10.1001/archpsyc.63.5.484>
- 1134 97. van Zwet EW, Goodman SN. How large should the next study be? Predictive power and
1135 sample size requirements for replication studies. *Stat Med*. 2022;41(16):3090-3101.
1136 <https://doi.org/10.1002/sim.9406>
- 1137 98. Curran-Everett D. Explorations in statistics: statistical facets of reproducibility. *Adv*
1138 *Physiol Educ*. 2016;40(2):248-252. <https://doi.org/10.1152/advan.00042.2016>
- 1139 99. Gandevia S, Cumming G, Amrhein V, Butler A. Replication: do not trust your p-value,
1140 be it small or large. *J Physiol*. 2021;599(11):2989-2990.
1141 <https://doi.org/10.1113/JP281614>
- 1142 100. Gorroochurn P, Hodge SE, Heiman GA, Durner M, Greenberg DA. Non-replication of
1143 association studies: "pseudo-failures" to replicate? *Genet Med*. 2007;9(6):325-331.
1144 <https://doi.org/10.1097/gim.0b013e3180676d79>
- 1145 101. Gibson EW. The role of p-values in judging the strength of evidence and realistic
1146 replication expectations. *Stat Biopharm Res*. 2021;13(1):6-18.
1147 <https://doi.org/10.1080/19466315.2020.1724560>
- 1148 102. Boos DD, Stefanski LA. P-value precision and reproducibility. *Am Stat*. 2011;65(4):213-
1149 221. <https://doi.org/10.1198/tas.2011.10129>
- 1150 103. Hung HM, O'Neill RT, Bauer P, Köhne K. The behavior of the P-value when the
1151 alternative hypothesis is true. *Biometrics*. 1997;53(1):11-22.
1152 <https://doi.org/10.2307/2533093>
- 1153 104. Sackrowitz H, Samuel-Cahn E. P values as random variables: expected p values. *Am*

- 1154 Stat. 1999;53(4):326-331. <https://doi.org/10.1080/00031305.1999.10474484>
- 1155 105. Shao J, Chow SC. Reproducibility probability in clinical trials. *Stat Med*.
1156 2002;21(12):1727-1742. <https://doi.org/10.1002/sim.1177>
- 1157 106. Amrhein V, Korner-Nievergelt F, Roth T. The earth is flat ($p > 0.05$): significance
1158 thresholds and the crisis of unreplicable research. *PeerJ*. 2017;5:e3544.
1159 <https://doi.org/10.7717/peerj.3544>
- 1160 107. Button KS, Ioannidis JP, Mokrysz C, Nosek BA, Flint J, Robinson ES, Munafò MR.
1161 Power failure: why small sample size undermines the reliability of neuroscience. *Nat Rev*
1162 *Neurosci*. 2013;14(5):365-376. <https://doi.org/10.1038/nrn3475>
- 1163 108. Ioannidis JP. Why most published research findings are false. *PLoS Med*.
1164 2005;2(8):e124. <https://doi.org/10.1371/journal.pmed.0020124>
- 1165 109. Dumas-Mallet E, Button KS, Boraud T, Gonon F, Munafò MR. Low statistical power in
1166 biomedical science: a review of three human research domains. *R Soc Open Sci*.
1167 2017;4(2):160254. <https://doi.org/10.1098/rsos.160254>
- 1168 110. Ioannidis JP, Trikalinos TA. An exploratory test for an excess of significant findings.
1169 *Clin Trials*. 2007;4(3):245-253. <https://doi.org/10.1177/1740774507079441>
- 1170 111. Masicampo EJ, Lalande DR. A peculiar prevalence of p values just below .05. *Q J Exp*
1171 *Psychol*. 2012;65(11):2271-2279. <https://doi.org/10.1080/17470218.2012.711335>
- 1172 112. Simmons JP, Nelson LD, Simonsohn U. False-positive psychology: undisclosed
1173 flexibility in data collection and analysis allows presenting anything as significant.
1174 *Psychol Sci*. 2011;22(11):1359-1366. <https://doi.org/10.1177/0956797611417632>
- 1175 113. Rosenthal R. The file drawer problem and tolerance for null results. *Psychol Bull*.
1176 1979;86(3):638-641. <https://doi.org/10.1037/0033-2909.86.3.638>

- 1177 114. Ioannidis JP. Why most discovered true associations are inflated. *Epidemiology*.
1178 2008;19(5):640-648. <https://doi.org/10.1097/EDE.0b013e31818131e7>
- 1179 115. Young NS, Ioannidis JP, Al-Ubaydli O. Why current publication practices may distort
1180 science. *PLoS Med*. 2008;5(10):e201. <https://doi.org/10.1371/journal.pmed.0050201>
- 1181 116. Colquhoun D. An investigation of the false discovery rate and the misinterpretation of p-
1182 values. *R Soc Open Sci*. 2014;1(3):140216. <https://doi.org/10.1098/rsos.140216>
- 1183 117. Davis-Stober CP, Dana J. Comparing the accuracy of experimental estimates to guessing:
1184 a new perspective on replication and the "crisis of confidence" in psychology. *Behav Res*
1185 *Methods*. 2014;46(1):1-14. <https://doi.org/10.3758/s13428-013-0342-1>
- 1186 118. Lakens D, Evers ER. Sailing from the seas of chaos into the corridor of stability: practical
1187 recommendations to increase the informational value of studies. *Perspect Psychol Sci*.
1188 2014;9(3):278-292. <https://doi.org/10.1177/1745691614528520>
- 1189 119. Lakens D. Equivalence tests: a practical primer for t tests, correlations, and meta-
1190 analyses. *Soc Psychol Personal Sci*. 2017;8(4):355-362.
1191 <https://doi.org/10.1177/1948550617697177>
- 1192 120. Parkhurst DF. Statistical significance tests: equivalence and reverse tests should reduce
1193 misinterpretation. *BioScience*. 2001;51(12):1051-1057. [https://doi.org/10.1641/0006-3568\(2001\)051\[1051:SSTEAR\]2.0.CO;2](https://doi.org/10.1641/0006-3568(2001)051[1051:SSTEAR]2.0.CO;2)
- 1194
- 1195 121. Mazzolari R, Porcelli S, Bishop DJ, Lakens D. Myths and methodologies: the use of
1196 equivalence and non-inferiority tests for interventional studies in exercise physiology and
1197 sport science. *Exp Physiol*. 2022;107(3):201-212. <https://doi.org/10.1113/EP090171>
- 1198 122. McRae G, Payne A, Zelt JG, Scribbans TD, Jung ME, Little JP, Gurd BJ. Extremely low
1199 volume, whole-body aerobic-resistance training improves aerobic fitness and muscular

- 1200 endurance in females. *Appl Physiol Nutr Metab.* 2012;37(6):1124-1131.
1201 <https://doi.org/10.1139/h2012-093>
- 1202 123. Rakobowchuk M, Tanguay S, Burgomaster KA, Howarth KR, Gibala MJ, MacDonald
1203 MJ. Sprint interval and traditional endurance training induce similar improvements in
1204 peripheral arterial stiffness and flow-mediated dilation in healthy humans. *Am J Physiol
1205 Regul Integr Comp Physiol.* 2008;295(1):R236-R242.
1206 <https://doi.org/10.1152/ajpregu.00069.2008>
- 1207 124. Iellamo F, Manzi V, Caminiti G, Vitale C, Castagna C, Massaro M, Franchini A, Rosano
1208 G, Volterrani M. Matched dose interval and continuous exercise training induce similar
1209 cardiorespiratory and metabolic adaptations in patients with heart failure. *Int J Cardiol.*
1210 2013;167(6):2561-2565. <https://doi.org/10.1016/j.ijcard.2012.06.057>
- 1211 125. Martins C, Kazakova I, Ludviksen M, Mehus I, Wisloff U, Kulseng B, Morgan L, King
1212 N. High-intensity interval training and isocaloric moderate-intensity continuous training
1213 result in similar improvements in body composition and fitness in obese individuals. *Int J
1214 Sport Nutr Exerc Metab.* 2016;26(3):197-204. <https://doi.org/10.1123/ijsnem.2015-0078>
- 1215 126. Gibala MJ, Little JP, van Essen M, Wilkin GP, Burgomaster KA, Safdar A, Raha S,
1216 Tarnopolsky MA. Short-term sprint interval versus traditional endurance training: similar
1217 initial adaptations in human skeletal muscle and exercise performance. *J Physiol.*
1218 2006;575(Pt 3):901-911. <https://doi.org/10.1113/jphysiol.2006.112094>
- 1219 127. Vella CA, Taylor K, Drummer D. High-intensity interval and moderate-intensity
1220 continuous training elicit similar enjoyment and adherence levels in overweight and
1221 obese adults. *Eur J Sport Sci.* 2017;17(9):1203-1211.
1222 <https://doi.org/10.1080/17461391.2017.1359679>

- 1223 128. Bartlett JD, Hwa Joo C, Jeong TS, Louhelainen J, Cochran AJ, Gibala MJ, Gregson W,
1224 Close GL, Drust B, Morton JP. Matched work high-intensity interval and continuous
1225 running induce similar increases in PGC-1 α mRNA, AMPK, p38, and p53
1226 phosphorylation in human skeletal muscle. *J Appl Physiol.* 2012;112(7):1135-1143.
1227 <https://doi.org/10.1152/jappphysiol.01040.2011>
- 1228 129. Trewin AJ, Parker L, Shaw CS, Hiam DS, Garnham A, Levinger I, McConell GK, Stepto
1229 NK. Acute HIIE elicits similar changes in human skeletal muscle mitochondrial H₂O₂
1230 release, respiration, and cell signaling as endurance exercise even with less work. *Am J*
1231 *Physiol Regul Integr Comp Physiol.* 2018;315(5):R1003-R1016.
1232 <https://doi.org/10.1152/ajpregu.00096.2018>
- 1233 130. Hazell TJ, Olver TD, Hamilton CD, Lemon P WR. Two minutes of sprint-interval
1234 exercise elicits 24-hr oxygen consumption similar to that of 30 min of continuous
1235 endurance exercise. *Int J Sport Nutr Exerc Metab.* 2012;22(4):276-283.
1236 <https://doi.org/10.1123/ijsnem.22.4.276>
- 1237 131. Skelly LE, Andrews PC, Gillen JB, Martin BJ, Percival ME, Gibala MJ. High-intensity
1238 interval exercise induces 24-h energy expenditure similar to traditional endurance
1239 exercise despite reduced time commitment. *Appl Physiol Nutr Metab.* 2014;39(7):845-
1240 848. <https://doi.org/10.1139/apnm-2013-0562>
- 1241 132. Saucedo Marquez CM, Vanaudenaerde B, Troosters T, Wenderoth N. High-intensity
1242 interval training evokes larger serum BDNF levels compared with intense continuous
1243 exercise. *J Appl Physiol.* 2015;119(12):1363-1373.
1244 <https://doi.org/10.1152/jappphysiol.00126.2015>
- 1245 133. Sagelv EH, Hammer T, Hamsund T, Rognmo K, Pettersen SA, Pedersen S. High

- 1246 intensity long interval sets provides similar enjoyment as continuous moderate intensity
1247 exercise: the Tromsø Exercise Enjoyment Study. *Front Psychol.* 2019;10:1788.
1248 <https://doi.org/10.3389/fpsyg.2019.01788>
- 1249 134. Crisp NA, Fournier PA, Licari MK, Braham R, Guelfi KJ. Optimising sprint interval
1250 exercise to maximise energy expenditure and enjoyment in overweight boys. *Appl*
1251 *Physiol Nutr Metab.* 2012;37(6):1222-1231. <https://doi.org/10.1139/h2012-111>
- 1252 135. Sabag A, Little JP, Johnson NA. Low-volume high-intensity interval training for
1253 cardiometabolic health. *J Physiol.* 2022;600(5):1013-1026.
1254 <https://doi.org/10.1113/JP281210>
- 1255 136. Tjønnå AE, Leinan IM, Bartnes AT, Jenssen BM, Gibala MJ, Winnett RA, Wisløff U.
1256 Low- and high-volume of intensive endurance training significantly improves maximal
1257 oxygen uptake after 10-weeks of training in healthy men. *PLoS ONE.* 2013;8(5):e65382.
1258 <https://doi.org/10.1371/journal.pone.0065382>
- 1259 137. Ramos JS, Dalleck LC, Borrani F, Beetham KS, Wallen MP, Mallard AR, Clark B,
1260 Gomersall S, Keating SE, Fassett RG, Coombes JS. Low-volume high-intensity interval
1261 training is sufficient to ameliorate the severity of metabolic syndrome. *Metab Syndr Relat*
1262 *Disord.* 2017;15(7):319-328. <https://doi.org/10.1089/met.2017.0042>
- 1263 138. Oh S, So R, Shida T, Matsuo T, Kim B, Akiyama K, Isobe T, Okamoto Y, Tanaka K,
1264 Shoda J. High-intensity aerobic exercise improves both hepatic fat content and stiffness
1265 in sedentary obese men with nonalcoholic fatty liver disease. *Sci Rep.* 2017;7:43029.
1266 <https://doi.org/10.1038/srep43029>
- 1267 139. Winding KM, Munch GW, Iepsen UW, Van Hall G, Pedersen BK, Mortensen SP. The
1268 effect on glycaemic control of low-volume high-intensity interval training versus

- 1269 endurance training in individuals with type 2 diabetes. *Diabetes Obes Metab.*
1270 2018;20(5):1131-1139. <https://doi.org/10.1111/dom.13198>
- 1271 140. Abdelbasset WK, Tantawy SA, Kamel DM, Alqahtani BA, Elnegamy TE, Soliman GS,
1272 Ibrahim AA. Effects of high-intensity interval and moderate-intensity continuous aerobic
1273 exercise on diabetic obese patients with nonalcoholic fatty liver disease: a comparative
1274 randomized controlled trial. *Medicine.* 2020;99(10):e19471.
1275 <https://doi.org/10.1097/MD.00000000000019471>
- 1276 141. Poon ET, Little JP, Sit CH, Wong SH. The effect of low-volume high-intensity interval
1277 training on cardiometabolic health and psychological responses in overweight/obese
1278 middle-aged men. *J Sports Sci.* 2020;38(17):1997-2004.
1279 <https://doi.org/10.1080/02640414.2020.1766178>
- 1280 142. Sabag A, Way KL, Sultana RN, Keating SE, Gerofi JA, Chuter VH, Byrne NM, Baker
1281 MK, George J, Caterson ID, Twigg SM, Johnson NA. The effect of a novel low-volume
1282 aerobic exercise intervention on liver fat in type 2 diabetes: a randomized controlled trial.
1283 *Diabetes Care.* 2020;43(10):2371-2378. <https://doi.org/10.2337/dc19-2523>
- 1284 143. Ryan BJ, Schleh MW, Ahn C, Ludzki AC, Gillen JB, Varshney P, Van Pelt DW,
1285 Pitchford LM, Chenevert TL, Gioscia-Ryan RA, Howton SM, Rode T, Hummel SL,
1286 Burant CF, Little JP, Horowitz JF. Moderate-intensity exercise and high-intensity interval
1287 training affect insulin sensitivity similarly in obese adults. *J Clin Endocrinol Metab.*
1288 2020;105(8):e2941-e2959. <https://doi.org/10.1210/clinem/dgaa345>
- 1289 144. Matsuo T, Saotome K, Seino S, Shimojo N, Matsushita A, Iemitsu M, Ohshima H,
1290 Tanaka K, Mukai C. Effects of a low-volume aerobic-type interval exercise on VO₂max
1291 and cardiac mass. *Med Sci Sports Exerc.* 2014;46(1):42-50.

- 1292 <https://doi.org/10.1249/MSS.0b013e3182a38da8>
- 1293 145. Wilson GA, Wilkins GT, Cotter JD, Lamberts RR, Lal S, Baldi JC. HIIT improves left
1294 ventricular exercise response in adults with type 2 diabetes. *Med Sci Sports Exerc.*
1295 2019;51(6):1099-1105. <https://doi.org/10.1249/MSS.0000000000001897>
- 1296 146. Way KL, Sabag A, Sultana RN, Baker MK, Keating SE, Lanting S, Gerofi J, Chuter VH,
1297 Caterson ID, Twigg SM, Johnson NA. The effect of low-volume high-intensity interval
1298 training on cardiovascular health outcomes in type 2 diabetes: a randomised controlled
1299 trial. *Int J Cardiol.* 2020;320:148-154. <https://doi.org/10.1016/j.ijcard.2020.06.019>
- 1300 147. Schulz KF, Altman DG, Moher D; CONSORT Group. CONSORT 2010 Statement:
1301 Updated guidelines for reporting parallel group randomised trials. *J Clin Epidemiol.*
1302 2010;63(8):834-840. <https://doi.org/10.1016/j.jclinepi.2010.02.005>
- 1303 148. Pattyn N, Beulque R, Cornelissen V. Aerobic interval vs. continuous training in patients
1304 with coronary artery disease or heart failure: an updated systematic review and meta-
1305 analysis with a focus on secondary outcomes. *Sports Med.* 2018;48(5):1189-1205.
1306 <https://doi.org/10.1007/s40279-018-0885-5>
- 1307 149. Bonafiglia JT, Islam H, Preobrazenski N, Gurd BJ. Risk of bias and reporting practices in
1308 studies comparing VO₂max responses to sprint interval vs. continuous training: a
1309 systematic review and meta-analysis. *J Sport Health Sci.* 2021.
1310 <https://doi.org/10.1016/j.jshs.2021.03.005>
- 1311 150. Charles P, Giraudeau B, Dechartres A, Baron G, Ravaud P. Reporting of sample size
1312 calculation in randomised controlled trials: review. *BMJ.* 2009;338:b1732.
1313 <https://doi.org/10.1136/bmj.b1732>
- 1314 151. Abt G, Boreham C, Davison G, Jackson R, Nevill A, Wallace E, Williams M. Power,

- 1315 precision, and sample size estimation in sport and exercise science research. *J Sports Sci.*
1316 2020;38(17):1933-1935. <https://doi.org/10.1080/02640414.2020.1776002>
- 1317 152. Cheval B, Boisgontier MP. The theory of effort minimization in physical activity. *Exerc*
1318 *Sport Sci Rev.* 2021;49(3):168-178. <https://doi.org/10.1249/JES.0000000000000252>
- 1319 153. Ioannidis JP. How to make more published research true. *PLoS Med.*
1320 2014;11(10):e1001747. <https://doi.org/10.1371/journal.pmed.1001747>
- 1321 154. Fanelli D. Redefine misconduct as distorted reporting. *Nature.* 2013;494(7436):149.
1322 <https://doi.org/10.1038/494149a>
- 1323 155. Gibala MJ. High-intensity interval training: a time-efficient strategy for health
1324 promotion? *Curr Sports Med Rep.* 2007;6(4):211-213.
- 1325 156. Gibala MJ, McGee SL. Metabolic adaptations to short-term high-intensity interval
1326 training: a little pain for a lot of gain? *Exerc Sport Sci Rev.* 2008;36(2):58-63.
1327 <https://doi.org/10.1097/JES.0b013e318168ec1f>
- 1328 157. Satiroglu R, Lalande S, Hong S, Nagel MJ, Coyle EF. Four-second power cycling
1329 training increases maximal anaerobic power, peak oxygen consumption, and total blood
1330 volume. *Med Sci Sports Exerc.* 2021;53(12):2536-2542.
1331 <https://doi.org/10.1249/MSS.0000000000002748>
- 1332 158. Sato S, Yoshida R, Murakoshi F, Sasaki Y, Yahata K, Nosaka K, Nakamura M. Effect of
1333 daily 3-s maximum voluntary isometric, concentric, or eccentric contraction on elbow
1334 flexor strength. *Scand J Med Sci Sports.* 2022;32(5):833-843.
1335 <https://doi.org/10.1111/sms.14138>
- 1336 159. Wagenmakers EJ. Defiant denial is self-defeating. *Psychol Inq.* 2021;32(1):12-16.
1337 <https://doi.org/10.1080/1047840X.2021.1889314>

- 1338 160. Harris C, Rohrer D, Pashler H. A train wreck by any other name. *Psychol Inq.*
 1339 2021;32(1):17-23. <https://doi.org/10.1080/1047840X.2021.1889317>
- 1340 161. Sherman, JW, Rivers AM. There's nothing social about social priming: derailing the
 1341 "train wreck". *Psychol Inq.* 2021;32(1):1-11.
 1342 <https://doi.org/10.1080/1047840X.2021.1889312>
- 1343 162. Wiggins BJ, Christopherson CD. The replication crisis in psychology: an overview for
 1344 theoretical and philosophical psychology. *J Theoret Philos Psychol.* 2019;39(4):202-217.
 1345 <https://doi.org/10.1037/teo0000137>
- 1346 163. Pashler H, Wagenmakers EJ. Editors' introduction to the special section on replicability
 1347 in psychological science: a crisis of confidence? *Perspect Psychol Sci.* 2012;7(6):528-
 1348 530. <https://doi.org/10.1177/1745691612465253>
- 1349 164. Yong E. Nobel laureate challenges psychologists to clean up their act. *Nature.* 2012.
 1350 <https://doi.org/10.1038/nature.2012.11535>
- 1351 165. Rodgers JL, Shrout PE. Psychology's replication crisis as scientific opportunity: a précis
 1352 for policymakers. *Policy Insights Behav Brain Sci.* 2018;5(1):134-141.
 1353 <https://doi.org/10.1177/2372732217749254>
- 1354 166. Sharpe D, Goghari VM. Building a cumulative psychological science. *Can Psychol.*
 1355 2020;61(4):269-272. <https://doi.org/10.1037/cap0000252>
- 1356 167. Shrout PE, Rodgers JL. Psychology, science, and knowledge construction: broadening
 1357 perspectives from the replication crisis. *Annu Rev Psychol.* 2018;69:487-510.
 1358 <https://doi.org/10.1146/annurev-psych-122216-011845>
- 1359 168. Munafò MR, Nosek BA, Bishop DVM, Button KS, Chambers CD, du Sert NP,
 1360 Simonsohn U, Wagenmakers EJ, Ware JJ, Ioannidis JPA. A manifesto for reproducible

- 1361 science. *Nat Hum Behav.* 2017;1:0021. <https://doi.org/10.1038/s41562-016-0021>
- 1362 169. Nelson LD, Simmons J, Simonsohn U. Psychology's renaissance. *Annu Rev Psychol.*
1363 2018;69:511-534. <https://doi.org/10.1146/annurev-psych-122216-011836>
- 1364 170. Begley CG, Ellis LM. Drug development: raise standards for preclinical cancer research.
1365 *Nature.* 2012;483(7391):531-533. <https://doi.org/10.1038/483531a>
- 1366 171. Begley CG, Ioannidis JP. Reproducibility in science: improving the standard for basic
1367 and preclinical research. *Circ Res.* 2015;116(1):116-126.
1368 <https://doi.org/10.1161/CIRCRESAHA.114.303819>
- 1369 172. Impellizzeri, F. M., McCall, A., & Meyer, T. (2019). Registered reports coming soon: our
1370 contribution to better science in football research. *Science and Medicine in Football*, 3(2),
1371 87-88. <https://doi.org/10.1080/24733938.2019.1603659>
- 1372 173. Nosek BA, Ebersole CR, DeHaven AC, Mellor DT. The preregistration revolution. *Proc*
1373 *Natl Acad Sci U S A.* 2018;115(11):2600-2606. <https://doi.org/10.1073/pnas.1708274114>
- 1374 174. Nosek BA, Hardwicke TE, Moshontz H, et al. Replicability, Robustness, and
1375 Reproducibility in Psychological Science. *Annu Rev Psychol.* 2022;73:719-748.
1376 <https://doi.org/10.1146/annurev-psych-020821-114157>
- 1377 175. Scheel AM, Schijen MRMJ, Lakens D. An excess of positive results: comparing the
1378 standard psychology literature with registered reports. *Adv Meth Pract Psychol Sci.*
1379 2021;4(2). <https://doi.org/10.1177/25152459211007467>
- 1380 176. Schäfer T, Schwarz MA. The meaningfulness of effect sizes in psychological research:
1381 differences between sub-disciplines and the impact of potential biases. *Front Psychol.*
1382 2019;10:813. <https://doi.org/10.3389/fpsyg.2019.00813>
- 1383 177. Abt G, Boreham C, Davison G, Jackson R, Wallace E, Williams AM. Registered Reports

1384 in the Journal of Sports Sciences. J Sports Sci. 2021;39(16):1789-1790.

1385 <https://doi.org/10.1080/02640414.2021.1950974>

1386

1387 **Figure Captions**

1388

1389 Figure 1.

1390 The number of entries per year in PubMed that include the strings "high intensity interval" or

1391 "sprint interval" are shown in the line chart. The number of meta-analyses (subsample) is shown

1392 in bars.

1393

1394 Figure 2.

1395 The inflation of the risk of Type I error as a function of the number of probability tests (at $p <$ 1396 $.05$). The estimates shown include the theoretical case of statistically independent (uncorrelated)

1397 variables (using the Šidák equation), as well as hypothetical cases in which the variables being

1398 analyzed are intercorrelated at levels of $r = .3$, $r = .5$, and $r = .7$ (using the M_{eff} method) [87,88].

1399 Note: DV – dependent variables.

1400

1401 Figure 3.

1402 The probability distribution of two-tailed p for three hypothetical studies: (i) an adequately1403 powered study, with population effect size $\delta = 0.5$ and $N = 64$ per group ($1-\beta = .81$), (ii) the1404 example shown by Cumming [46] (p. 289), with population effect size $\delta = 0.5$ and $N = 32$ per1405 group ($1-\beta = .52$), and (iii) an example consistent with the studies included in the meta-analysis1406 by Mattioni Maturana et al. [65], with population effect size $\delta = 0.4$ and $N = 10$ per group ($1-\beta =$ 1407 $.14$). The 80th percentiles indicate that 80% of the area under each curve (the probability of two-1408 tail p values) lies to the left of the marker and the figure indicated is the upper limit of the 80%1409 percentile p interval (with a lower limit of zero). The probabilities associated with conventional

1410 intervals of p (i.e., .05, .01, .001) are shown as percentages in the histograms.

1411

1412 Figure 4.

1413 The p values associated with the 48 studies comparing VO₂max between HIIT and moderate-

1414 intensity continuous exercise groups that were included in the meta-analysis by Mattioni

1415 Maturana et al. [65], illustrating the range from 0.000 to 1.000.

1416

1417 Figure 5.

1418 Probability (y axis) that a hypothetical "perfect" replication study (i.e., drawing samples of equal

1419 size from the same population as the original, and applying identical treatment and assessment

1420 methods) would obtain $p < .05$, as a function of the p value obtained in the original study (under

1421 two assumptions: that the population effect size is known, and equal to the effect size obtained in

1422 the initial study, or not). It can be seen that if the initial study yielded $p < .05$, there is only a 50%

1423 chance that a replication would also obtain $p < .05$. If the initial study yielded $p = .371$ (i.e., the p

1424 value expected from studies with the characteristics of those included in the meta-analysis by

1425 Mattioni Maturana et al. [65], given $\delta = 0.40$ and $N = 10$ per group), the probability of obtaining

1426 $p < .05$ from a replication would be only .15 and .25, respectively.

1427

1428 Figure 6.

1429 P intervals estimated to indicate the probability of obtaining $p < .05$ in a replication study as a

1430 function of the (two-tail) p value in an initial study. The two-sided p intervals, extending from

1431 the 10th to the 90th percentile, are shown on the left, whereas the one-sided p intervals,

1432 extending from zero to the 80th percentile, are shown on the right. Estimates are shown for both

1433 two-tail and one-tail tests in the replication study. The upper limits of the 90th percentile (left)
1434 and 80th percentile (right) p intervals associated with an initial study yielding $p = .371$ (i.e., the p
1435 value expected from studies with the characteristics of those included in the meta-analysis by
1436 Mattioni Maturana et al. [65], given $\delta = 0.40$ and $N = 10$ per group) are highlighted.

1437

1438 Figure 7.

1439 Positive predictive value (PPV), namely the probability that a "positive" research finding
1440 represents a true effect (i.e., that the finding is a true positive), as a function of the Type I error
1441 rate (α), when statistical power ($1-\beta$) is sufficient (i.e., $1-\beta = .80$) and when it is the median of the
1442 power of studies included in the meta-analysis by Mattioni Maturana et al. [65] comparing HIIT
1443 and moderate-intensity continuous training on VO_{2max} (i.e., $1-\beta = .14$). When α is allowed to
1444 escalate to high levels, even under the unrealistic scenario of $R = .50$, the PPV drops to $< .10$.

1445

1446 Figure 8.

1447 Levels of statistical power ($1-\beta$) for each of the 48 studies included in the Mattioni Maturana et
1448 al. [65] meta-analysis comparing the effects of HIIT and moderate-intensity continuous exercise
1449 on VO_{2max} . Power was calculated from the reported sample sizes, assuming that the pooled
1450 effect ($d = 0.40$) represents the "true" population effect and $\alpha = .05$. The median study exhibited
1451 14% statistical power, 42 of 48 studies (88%) had statistical power in the 0–20% range and all
1452 but one (47 of 48, or 98%) were in the 0–33% range.

1453

1454 Figure 9

1455 The estimated false positive risk (FPR) of the studies on VO_{2max} that were included in the

1456 Mattioni Maturana et al. [65] meta-analysis, assuming $R = .50$. Only 3 of the 48 studies (6.25%)
1457 produced FPR lower than .05. The FPR of the 13 studies that produced $p < .05$ was as high as
1458 .245, with a mean of .130 and a median of .123. Two related figures are highlighted for
1459 reference: (i) the minimum risk of Type I error (α) associated with $p = .05$ has been estimated as
1460 .289; (ii) the relationship between p values and α holds until $p < 1/e$, namely $p < .368$, after
1461 which α reaches a plateau.

1462

1463 Figure 10.

1464 The expected and observed frequencies of p values, in intervals ranging from $p < .05$ to $.95 < p <$
1465 1.00 , resulting from the studies on $VO_2\text{max}$ included in the meta-analysis by Mattioni Maturana
1466 et al. [65], illustrating the presence of an excessive proportion of studies with $p < .05$.

1467

1468 Figure 11

1469 Results of simulated experiments (100,000 simulated tests per data point) illustrating the
1470 phenomenon of "winner's curse," namely the inflation of the apparent effect size (d) compared to
1471 the known population effect size (δ) from studies with various sample sizes resulting in $p < .05$.
1472 For sample sizes of 10 per group, namely the median sample size of the 48 studies on $VO_2\text{max}$
1473 included in the meta-analysis by Mattioni Maturana et al. [65], a small effect ($\delta = 0.20$) can
1474 appear as large ($d = 0.80$), while a population effect size of $\delta = 0.40$ (the pooled effect from the
1475 meta-analysis by Mattioni Maturana et al. [65]) can appear highly exaggerated, namely $d = 1.04$.
1476 Notice that samples of $N = 100$ per group suffice to eliminate the inflation of medium population
1477 effect sizes ($\delta = d = .50$) but samples of $N = 700$ per group are required to eliminate the inflation
1478 for small population effect sizes ($\delta = d = .20$).

1479

1480 Figure 12.

1481 Values of the v-statistic proposed by Davis-Stober and Dana [116] for each of the 48 studies on

1482 VO₂max included in the meta-analysis by Mattioni Maturana et al. [65], comparing the effects of

1483 HIIT and moderate-intensity continuous exercise. The v-statistic is an index of the relative

1484 accuracy of population estimates produced by the traditional method of ordinary least squares

1485 compared to "random least squares" (i.e., random estimates for both the direction and the

1486 magnitude of treatment effects). The average v-statistic was .124 and the median was .000.

1487 Nearly all studies (46 of 48, or 96%) had values of the v-statistic below .500, and more than half

1488 (28 of 48, or 58%) had a v-statistic of zero, suggesting that random estimates were consistently

1489 more accurate than estimates based on the observed data.

1490

1491 Table 1. Probability of obtaining $p < .05$ from a replication as a function of the p value obtained
 1492 in an initial experiment (p obt) under two assumptions (i.e., that the population effect size is
 1493 known, and equal to the effect size obtained in the initial study, or not). The column labeled
 1494 "Goodman" contains the values calculated by Goodman [43] (Table 1, p. 877), presented here as
 1495 evidence of validation. The p value of .371 (i.e., the expected p value from the meta-analysis by
 1496 Mattioni Maturana et al. [65], given $\delta = 0.40$ and $N = 10$ per group) is also included, to highlight
 1497 the low probabilities of obtaining $p < .05$ from a replication study.
 1498

p obt	Assuming δ is known ($\delta = d$)			Assuming δ is unknown		
	2-tail	Goodman	1-tail	2-tail	Goodman	1-tail
.001	.908	.91	.950	.827	.78	.878
.005	.802	.80	.877	.726	.71	.794
.010	.731	.73	.824	.669	.66	.745
.030	.583	.58	.700	.561	.56	.645
.050	.500	.50	.624	.503	.50	.588
.100	.376	.37	.500	.417	.41	.500
.200	.249		.358	.327		.399
.371	.146		.227	.247		.298
.400	.134		.211	.238		.285
.600	.082		.131	.195		.214

1499
 1500
 1501

1502 Table 2. Two-sided (extending from the 10th to the 90th percentile) and one-sided (extending
 1503 from zero to the 80th percentile) p intervals for two- and one-tail single-study replications as a
 1504 function of the p value obtained in an initial (two-tail) study (p obt). P intervals indicate the
 1505 probability of obtaining $p < .05$ in a single, identical replication study. Compare to the values
 1506 calculated by Cumming [46] (Table 1, p. 292) for validation. As noted by Cumming [46], "for
 1507 the 90% p interval [one-tail] to be [0, .05], p obt must equal .00054" (p. 293). The p value of .371
 1508 (i.e., the expected p value from the studies included in the meta-analysis by Mattioni Maturana et
 1509 al. [65], given $\delta = 0.40$ and $N = 10$ per group) is also included, to highlight the extraordinarily
 1510 wide p interval associated with it.
 1511

p obt	10-90th percentile interval, two-tail	10-90th percentile interval, one-tail	0-80th percentile interval, two-tail	0-80th percentile interval, one-tail
.00054	[.0000005, .099]	[.0000001, .050]	[.000, .023]	[.000, .011]
.001	[.0000005, .139]	[.0000005, .070]	[.000, .036]	[.000, .018]
.010	[.000012, .408]	[.000006, .223]	[.000, .162]	[.000, .083]
.020	[.000035, .517]	[.000018, .304]	[.000, .242]	[.000, .128]
.050	[.000162, .648]	[.000081, .441]	[.000, .379]	[.000, .221]
.100	[.000544, .728]	[.000273, .567]	[.000, .491]	[.000, .325]
.200	[.001924, .789]	[.000988, .702]	[.000, .591]	[.000, .464]
.371	[.005998, .828]	[.003397, .821]	[.000, .662]	[.000, .616]
.400	[.006848, .832]	[.003978, .834]	[.000, .669]	[.000, .636]
.600	[.013091, .849]	[.009726, .901]	[.000, .701]	[.000, .747]

1512
 1513

1514 Table 3. Synopsis of the sample-size calculations of the 11 studies included in the review by
 1515 Sabag et al. [135], comparing the effects of low-volume HIIT to traditional HIIT or moderate-
 1516 intensity continuous exercise.
 1517

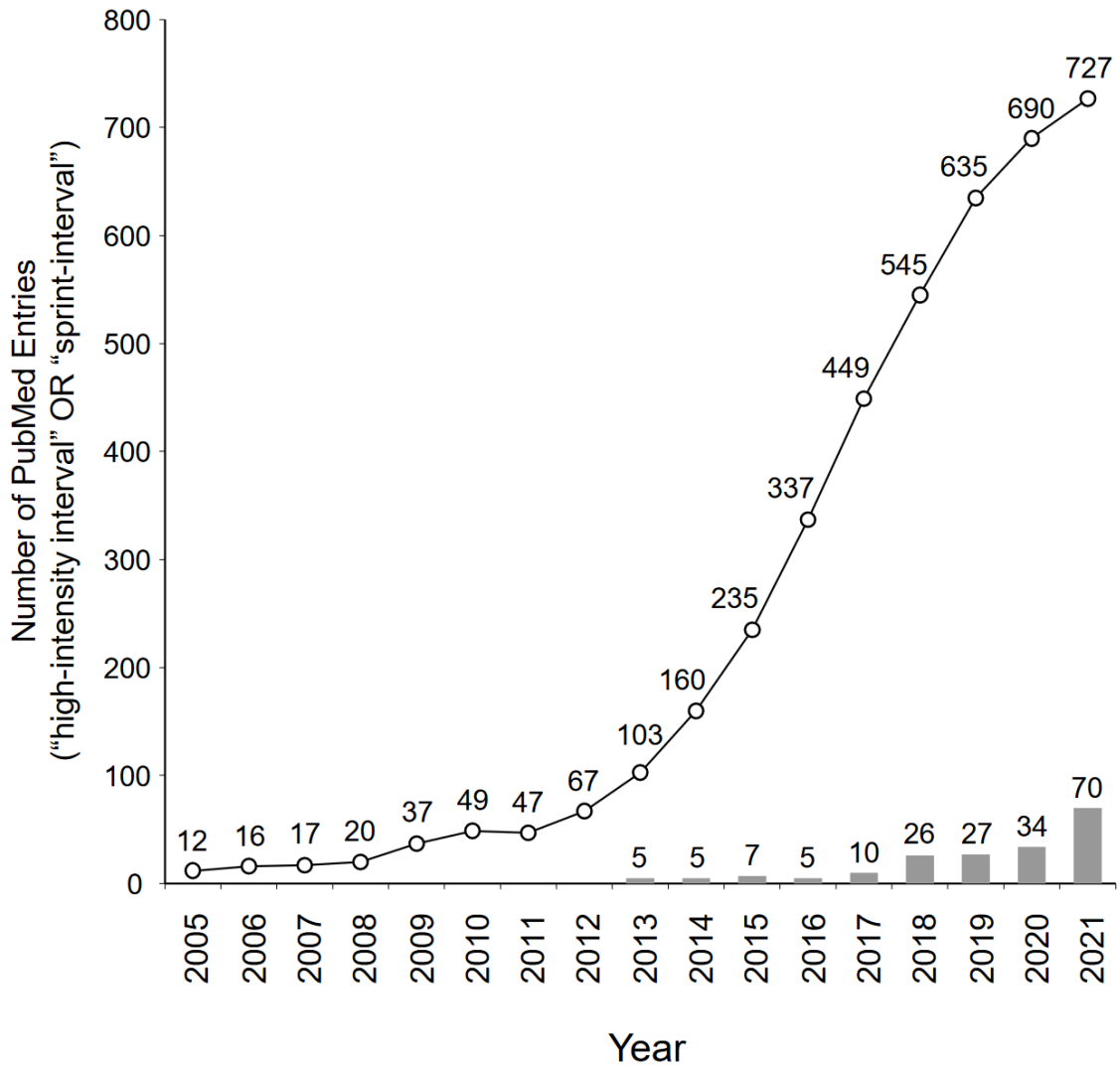
Study	Samples	Verbatim section on power	Comment
Tjønnå et al., 2013 [136]	13 & 13	Prior experience suggests a standard deviation (SD) of about 2.0–3.0 ml/kg/min. According to sample size tables for clinical studies, we needed 10 subjects in each group (we included 13 in case of drop out). With a standardized within-group difference of 1.0 differences may be detected using a paired t-test with 80% power, at a significance level of 5%. Clinically, this corresponds to a detectable difference for VO ₂ max of 3 ml/kg/min (p. 3).	While the calculations for a matched-pair t-test with $d = 1.0$ indeed yields a required sample size of $N = 10$, the cited source did not yield an effect size of $d = 1.0$. Also, the focus of this study was not on "within group" changes but rather inter-group comparisons, and the analysis was not based on a matched-pair t-test but rather on "mixed linear model analyses with group and time interaction."
Ramos et al., 2017 [137]	21 & 22	Sample size for the substudy was calculated using an anticipated mean difference in MetS z-score reduction of 0.60 (power = 0.80, alpha = 0.05 for two-tailed test) between HIIT and MICT groups. This was based on a previous study showing a similar mean difference in reduction of MetS z-score between HIIT and MICT (From: Supplementary material).	The information provided lacks standard deviation. The cited source does not report a mean difference in MetS z-score reduction "similar" to 0.60 but rather 0.46 \forall 1.55, and $d = 0.29$. This entails a total sample size of $N = (188 + 188) = 376$.
Oh et al., 2017 [138]	20 & 13	Our study design did not consider sampling size calculation to estimate the effect of sample size. Therefore, the small sample size might have limited the statistical power of the study (p. 10).	No sample-size calculation.
	13 & 12	A limitation of the present study is the relatively small number of	

<p>Winding et al., 2018 [139]</p>		<p>participants, which may have masked differences between HIIT and END (p. 1138).</p>	<p>No sample-size calculation.</p>
<p>Abdelbasset et al., 2020 [140]</p>	<p>16 & 15</p>	<p>For sample size estimation, an initial power analysis was applied (2-tailed test with statistical power of 0.80, a error=0.05, and effect size = 0.5). Estimates of mean difference and standard deviation for the [intrahepatic triglycerides] value from the previous study assessed 19 patients who received aerobic exercise. According to that study measures, 13 patients were required in each group. Forty-eight patients were included [for three groups] in the study to account for the dropout rate of 20% (p. 3).</p>	<p>Given the cited assumptions ($d = 0.5$, $\alpha = 0.05$, power = 0.80), the required sample is $N = 64$ per group (128 for two groups, 192 for three groups, 230 with 20% oversampling for dropout). However, the cited source (which did not include power calculations) did not yield $d = 0.5$ for the comparison between exercise and placebo for hepatic triglyceride concentration but rather $d = 0.3$. This entails $N = 170$ per group, (340 for two groups, 510 for three, 612 with 20% oversampling for dropout).</p>
<p>Poon et al., 2020 [141]</p>	<p>12 & 12</p>	<p>Based on a meta-analysis that compared HIIT with continuous endurance training on maximal oxygen uptake ($VO_2\max$) improvements in adults, the estimated standardized mean difference (Cohen's d) between HIIT and MICT was approximately 0.4. Therefore, it was anticipated that a sample size of 12 participants per group was adequate to detect this difference between groups on our primary outcome (i.e., $VO_2\max$), with a power of 0.8 at an alpha level of 0.05 (pp. 1998-1999).</p>	<p>The cited source (meta-analysis) reported non-standardized results (i.e., not Cohen's d). When converted to d using the information given (mean difference, 95% confidence limits), d was not 0.4. More importantly, the sample size required for $d = 0.4$, $\alpha = .05$, power = .8 is $N = 100$ per group, not 12.</p>

Sabag et al., 2020 [142]	12 & 12	An a priori, two-tailed power calculation at an α of 0.05 and β of 0.8 gave an actual power of 0.813 for a sample size of 11 in each group. This calculation was determined using the effect size (ES) of 1.28 of a similar exercise intervention from a previous study, which detected significant improvements in liver fat within groups (p. 2373).	Besides confusing β and $1-\beta$ (power), the researchers referred to an effect size "within groups" as the basis for power calculations for a between-groups comparison (also, the reported effect size for high-intensity, low volume exercise was 1.42 for intrahepatic lipids, not 1.28). In the cited source, the effect size for the comparison between high-intensity, low-volume exercise and low-intensity high-volume exercise was $d = 0.19$, requiring $N = (436 + 436) = 872$.
Ryan et al., 2020 [143]	16 & 14		No sample-size calculation.
Matsuo et al., 2014 [144]	14 & 14	A priori power analysis was performed to determine the sample size. The primary outcome variable of this study was the increase of VO ₂ max achieved through three types of exercise intervention. On the basis of data from both a previous study and our preliminary study on changes in VO ₂ max, we assumed a 15% difference in the training effect between the three groups with an SD estimate of 10%. With an alpha error rate of 0.017 (with Bonferroni adjustment for post hoc tests) and statistical power of 80%, the minimal sample size in each group was estimated to be 11 subjects (33 subjects in total).	Assuming a large effect $d = 1.5$ with an adjusted $\alpha = (0.05 / 3) = 0.017$ indeed requires only $N = 11$ per group. However, the cited "preliminary study" only reported within-subjects changes in VO ₂ max in two participants, not intergroup differences or standard deviations. Moreover, the "previous study" was conducted on a patient population (heart failure), with low baseline levels of VO ₂ max and there is no indication of a "15% difference in the training effect" (the cited study reported increases

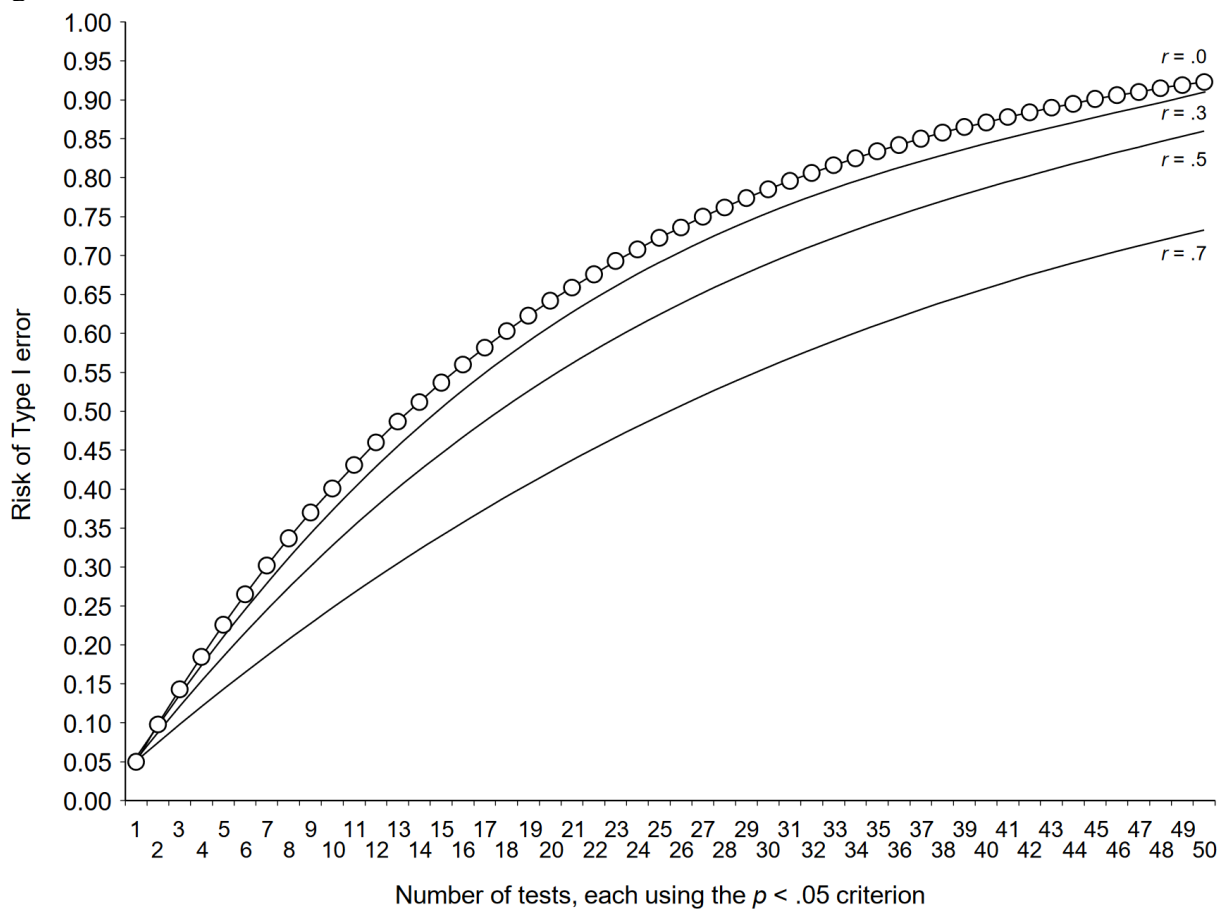
		Assuming subject attrition such as dropout, we recruited 14 subjects for each group (42 subjects in total) in this study (p. 46).	of 46% vs. 14%, for interval and moderate continuous exercise, respectively).
Wilson et al., 2019 [145]	11 & 5		No sample-size calculation.
Way et al., 2020 [146]	12 & 12	Sample size was calculated based on a projected change in peripheral arterial stiffness [pulse wave velocity] with [moderate-intensity continuous training] in adults with [type 2 diabetes] similar to the [moderate-intensity continuous training] protocol in our study. A priori, two-tailed power calculation of $\alpha = 0.05$ and $\beta = 0.20$ gave a power of 0.82 for a total sample size of 45 ($n = 15$ per group) (p. 150).	The researchers did not cite an anticipated effect size, so the calculations cannot be verified. Solving for the missing effect size shows that the study was sufficiently powered only for a large between-group effect ($d = 1.2$). The researchers reported basing their calculations on within-group changes but their analyses were for inter-group comparisons. The cited source reported $d = 0.80$ (radial) and $d = 0.50$ (femoral) for within-group changes and $d = 1.10$ (radial) and $d = 0.84$ (femoral) for inter-group comparisons.

1520 Figure 1



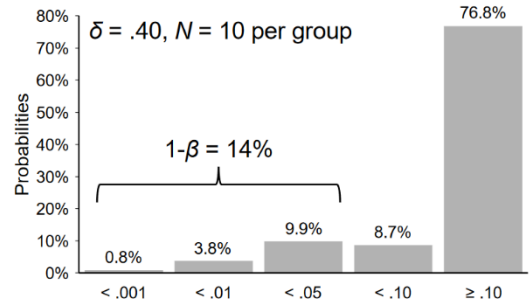
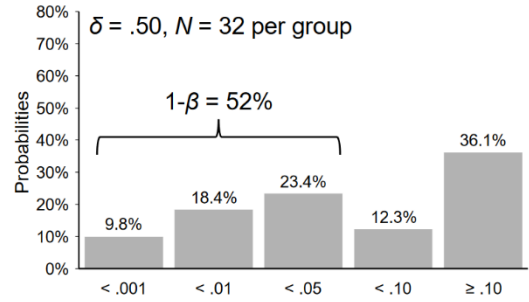
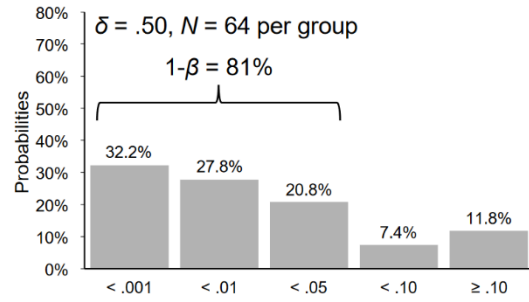
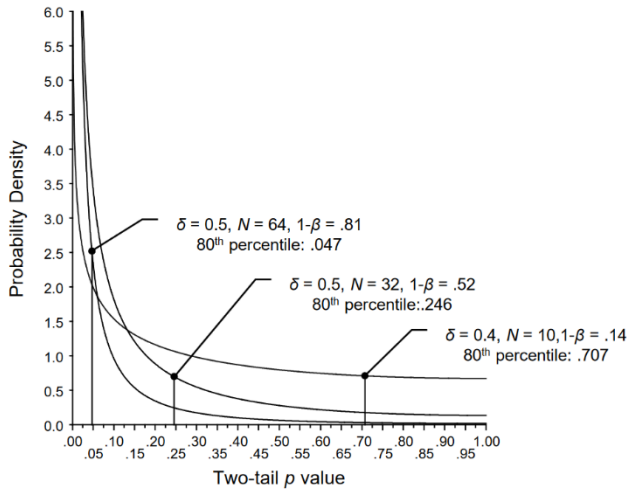
1521
1522

1523 Figure 2



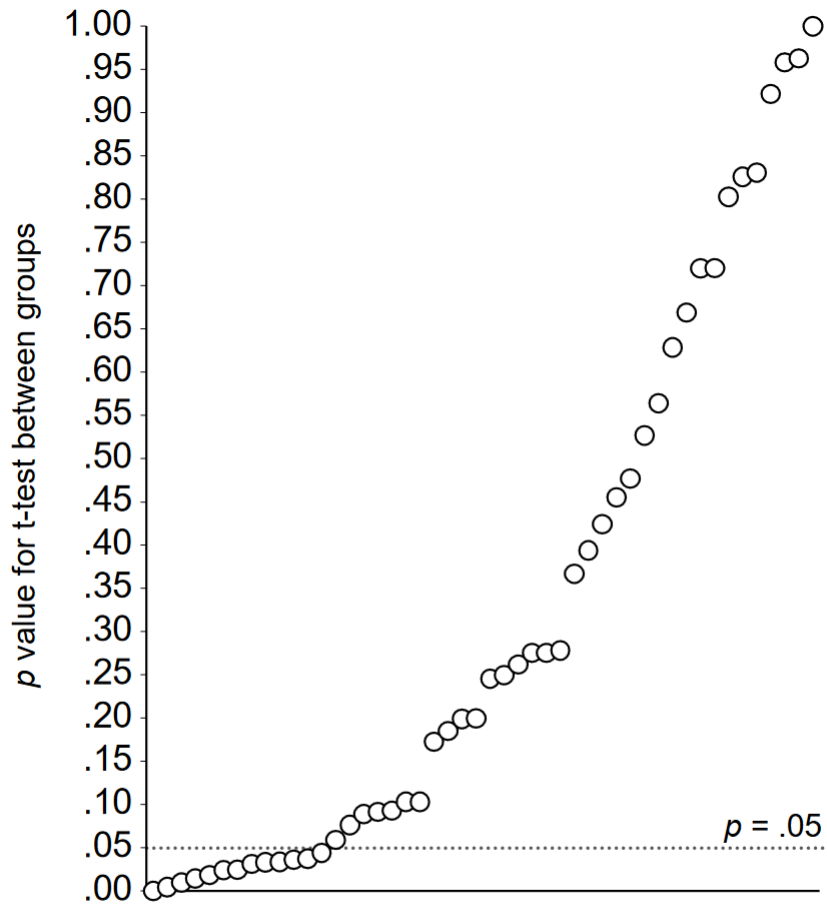
1524
1525

1526 Figure 3



1527
1528

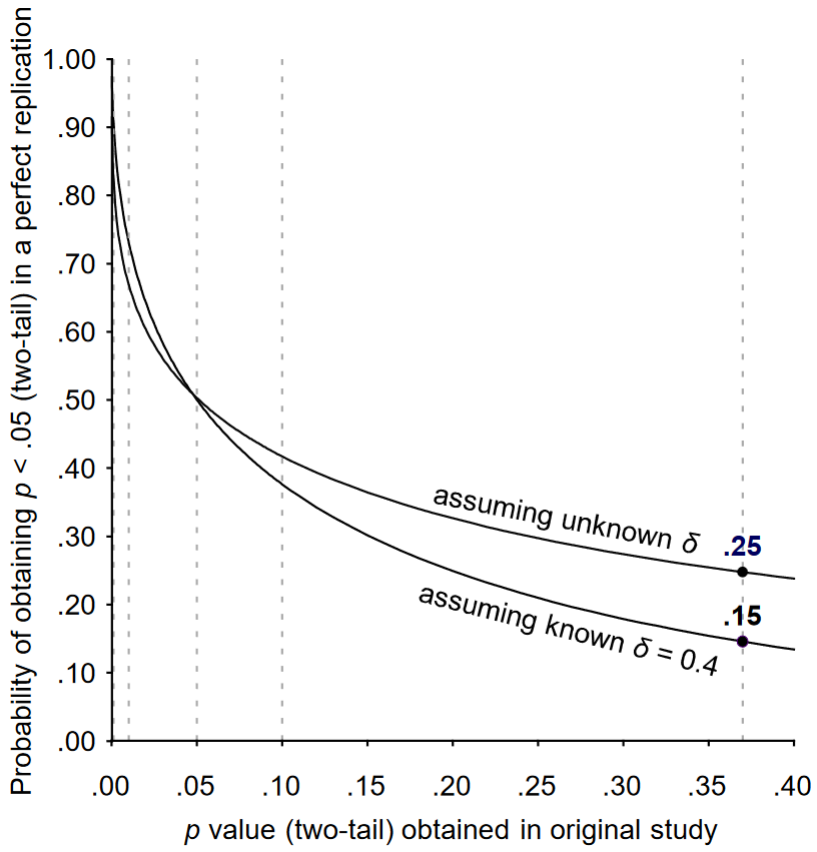
1529 Figure 4



48 studies investigating the effects of HIIT vs. moderate-intensity continuous training on VO₂max in the meta-analysis by Mattioni Maturana et al. [65]

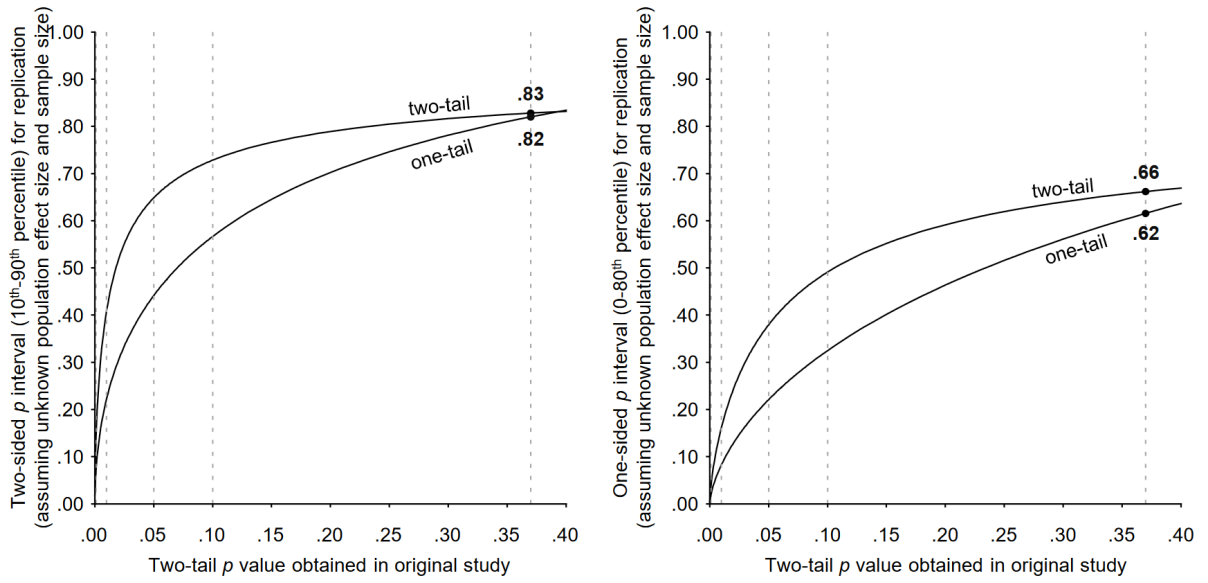
1530
1531

1532 Figure 5



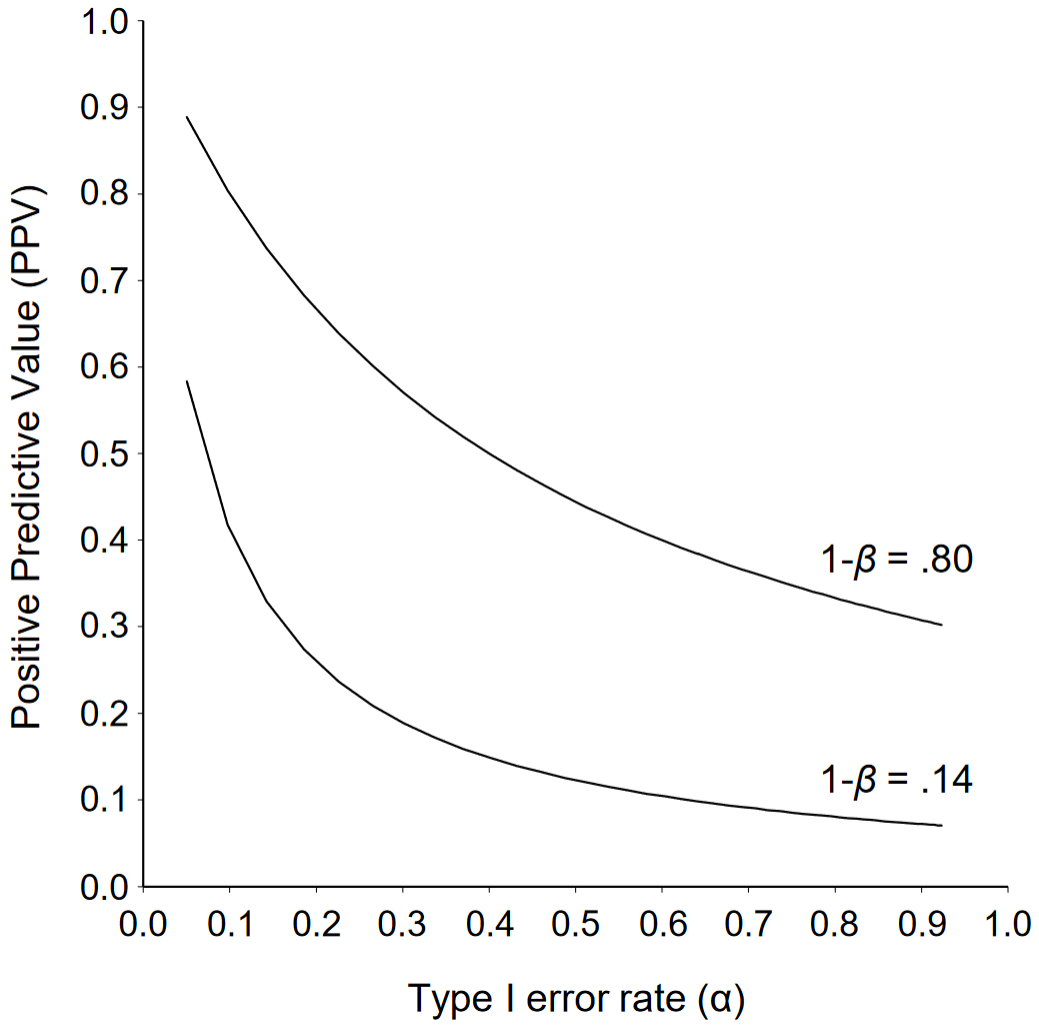
1533
1534

1535 Figure 6



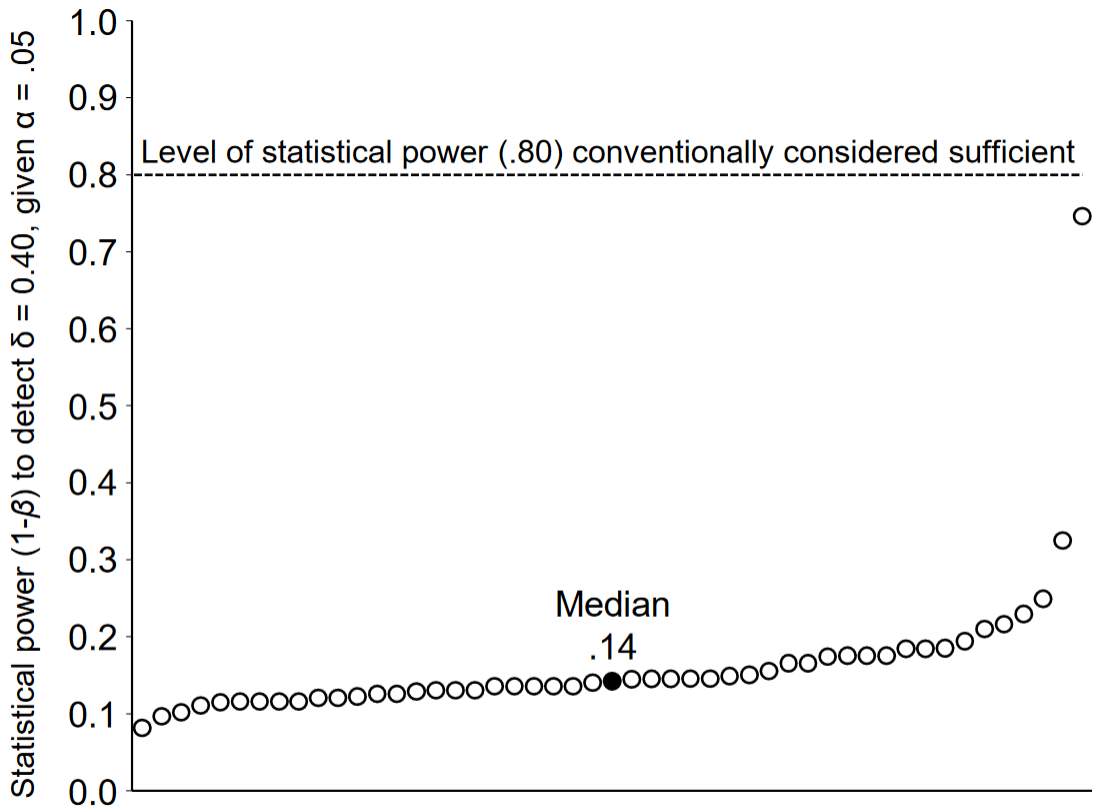
1536
1537

1538 Figure 7



1539
1540

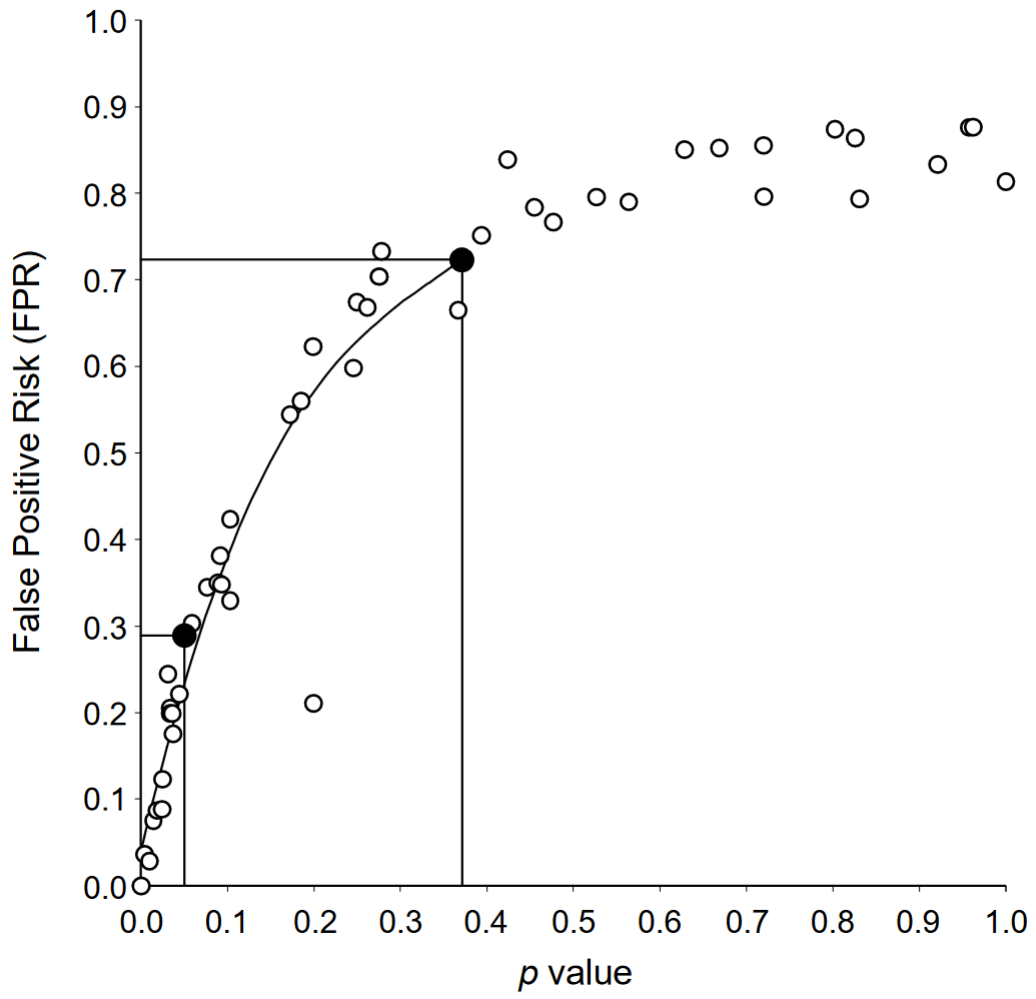
1541 Figure 8



48 studies investigating the effects of HIIT vs. moderate-intensity continuous training on $VO_2\text{max}$ in the meta-analysis by Mattioni Maturana et al. [65]

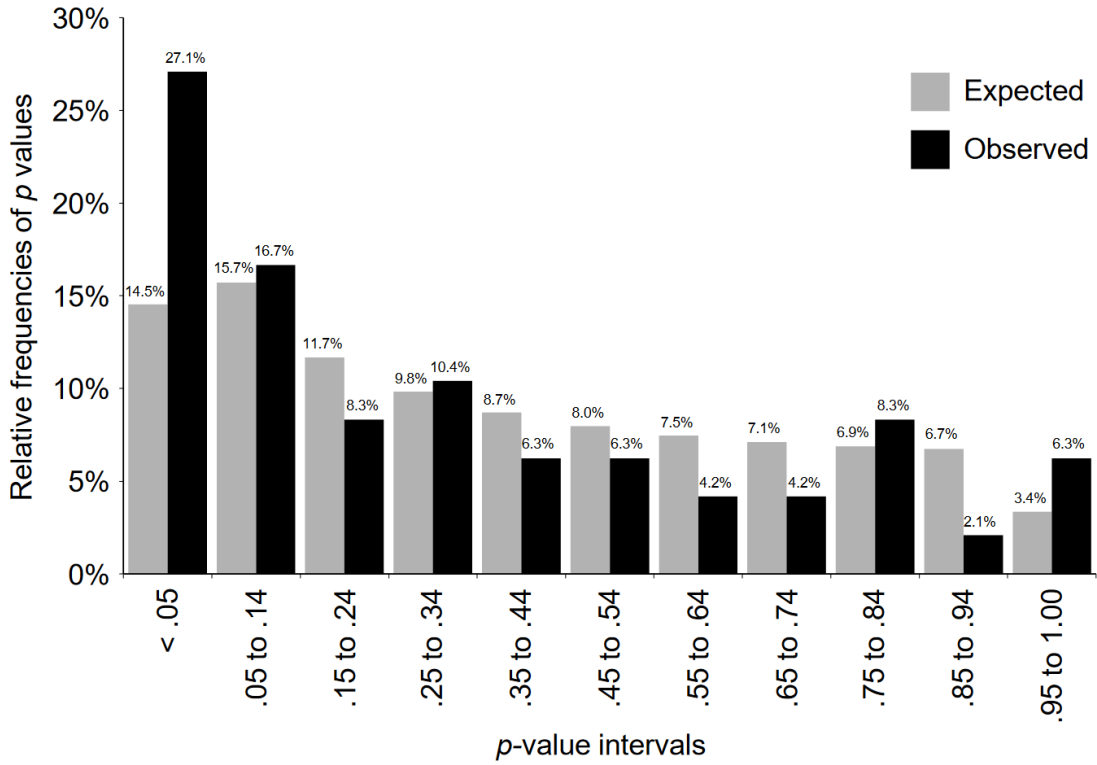
1542
1543

1544 Figure 9



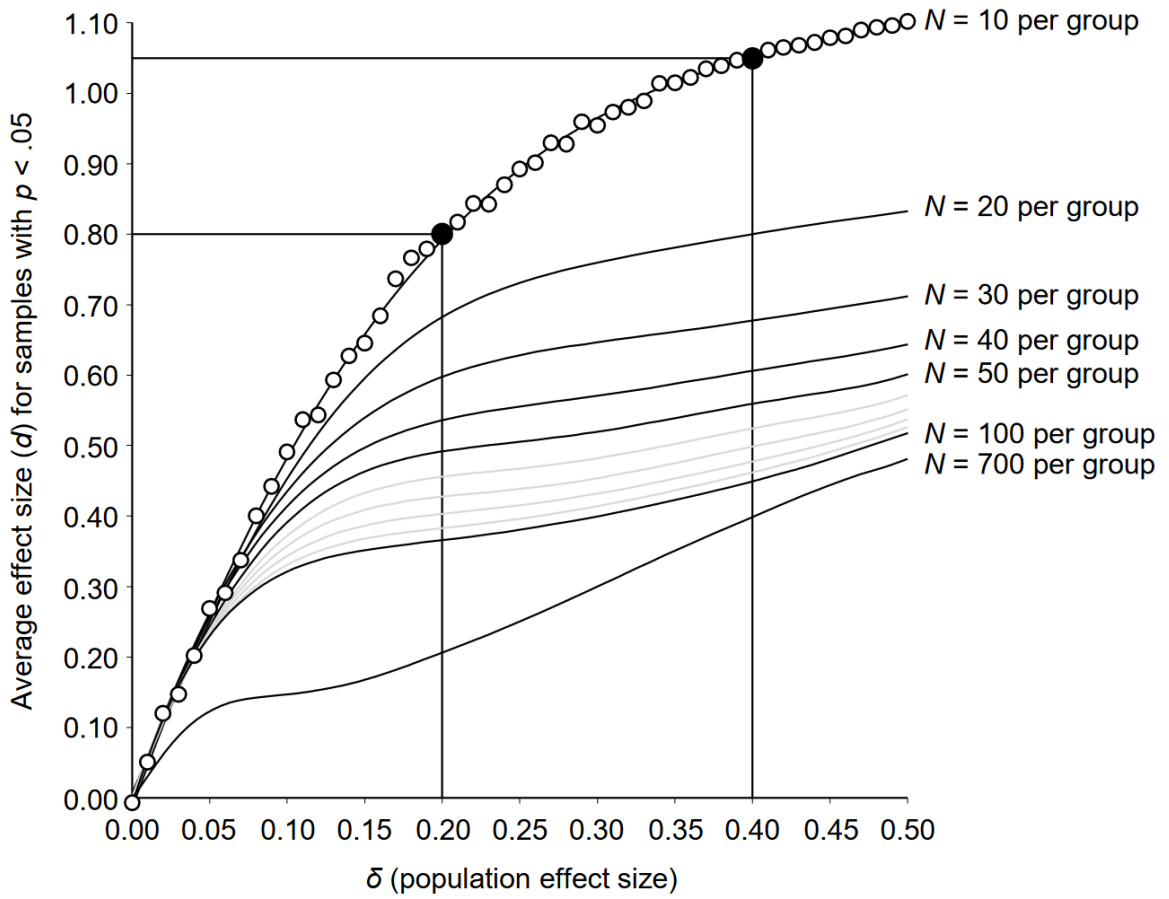
1545
1546

1547 Figure 10



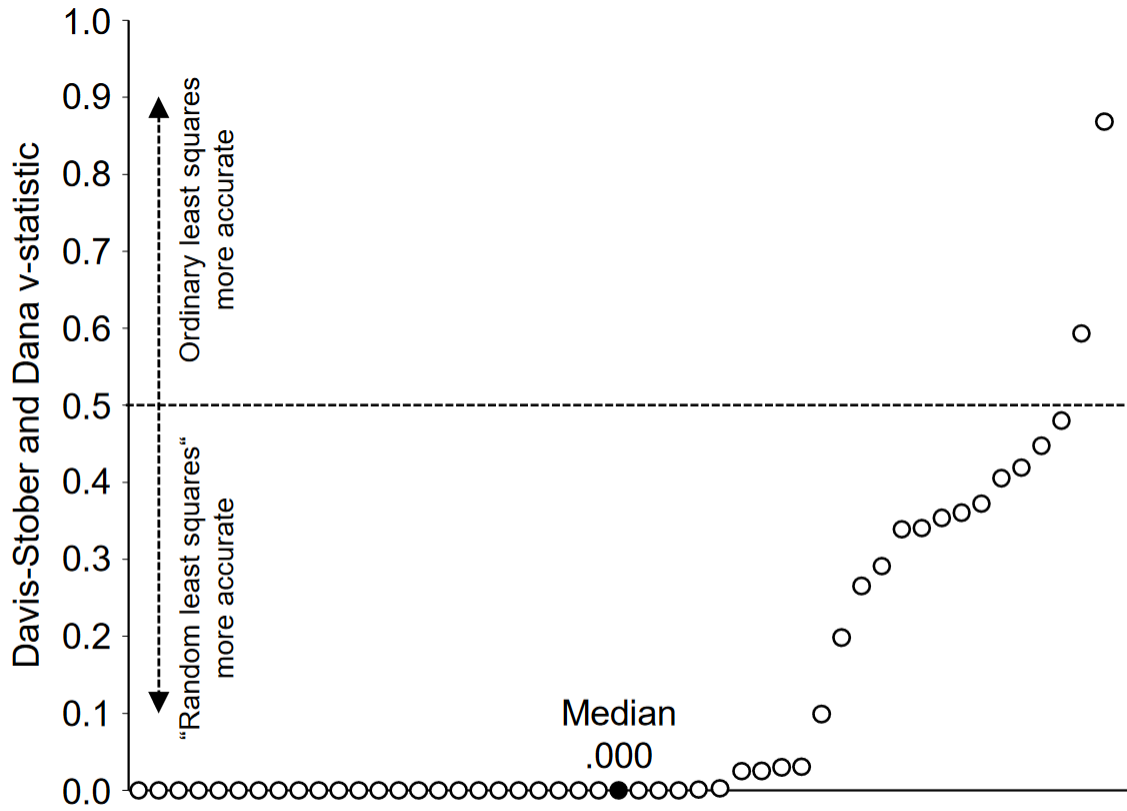
1548
1549

1550 Figure 11



1551
1552

1553 Figure 12



48 studies investigating the effects of HIIT vs. moderate-intensity continuous training on VO₂max in the meta-analysis by Mattioni Maturana et al. [65]

1554