# Effect size, sample size and power of forced swim test assays in mice: Guidelines for investigators to optimize reproducibility — Source link ↗

Neil R. Smalheiser, Elena E. Graetz, Zhou Yu, Jing Wang

**Institutions:** University of Illinois at Chicago

Related papers:

- Effect Size Measures and Their Relationships in Stroke Studies

- Detecting Qualitative Interactions in Clinical Trials: An Extension of Range Test

- Power to detect clinically relevant carry-over in a series of cross-over studies.

- Is Noninvasive Testing for Coronary Artery Disease Accurate

- Statistical sampling and hypothesis testing in orthopaedic research.

# Effect size, sample size and power of forced swim test assays in mice: Guidelines for investigators to optimize reproducibility

Neil R. Smalheiser[1][*], Elena E. Graetz[2], Zhou Yu[2], Jing Wang[2]

[1] Department of Psychiatry, University of Illinois School of Medicine,
Chicago, Illinois, United States of America
[2] Department of Mathematics, Statistics, and Computer Science,
University of Illinois at Chicago, Chicago, Illinois, United States of America

*Corresponding author

E-mail: neils@uic.edu

## Abstract

A recent flood of publications has documented serious problems in scientific reproducibility, power, and reporting of biomedical articles, yet scientists persist in their usual practices. Why? We examined a popular and important preclinical assay, the Forced Swim Test (FST) in mice used to test putative antidepressants. Whether the mice were assayed in a naïve state vs. in a model of depression or stress, and whether the mice were given test agents vs. known antidepressants regarded as positive controls, the mean effect sizes seen in the experiments were indeed extremely large (1.5 – 2.5 in Cohen's d units); most of the experiments utilized 7-10 animals per group which did have adequate power to reliably detect effects of this magnitude. We propose that this may at least partially explain why investigators using the FST do not perceive intuitively that their experimental designs fall short -- even though proper prospective design would require ~21-26 animals per group to detect, at a minimum, large effects (0.8 in Cohen's d units) when the true effect of a test agent is unknown. Our data provide explicit

1

30    parameters and guidance for investigators seeking to carry out prospective power estimation for

31    the FST. More generally, altering the real-life behavior of scientists in planning their

32    experiments may require developing educational tools that allow them to actively visualize the

33    inter-relationships among effect size, sample size, statistical power, and replicability in a direct

34    and  intuitive manner.

35

## Keywords

37    reproducibility, power, experimental design, preclinical assays, forced swim test, neuroscience,

38    psychiatry, antidepressants, meta-science, science education

39

## Introduction

41    A recent flood of publications has documented serious problems in scientific reproducibility,

42    power, and reporting of biomedical articles, including psychology, neuroscience, and preclinical

43    animal models of disease [1-16]. The power of published articles in many subfields of

44    neuroscience and psychology hovers around 0.3-0.4, whereas the accepted standard is 0.8 [3, 4,

45    7, 9, 15]. Only a tiny percentage of biomedical articles specify prospective power estimations

46    [e.g., 17]. This is important since under-powered studies have a tendency to over-estimate true

47    effect sizes, and to show a very high false-positive rate [1, 18]. Even when the nominal statistical

48    significance of a finding achieves p= 0.05 or better, the possibility of reporting a false positive

49    finding may approach 50% [1, 3, 19]. In several fields, when attempts have been made to repeat

50    experiments as closely as possible, replication is only achieved about 50% of the time,

51    suggesting that the theoretical critiques are actually not far from the real situation [6, 20].

52

2

53   Why might scientists persist in their usual practices, in the face of objective, clear evidence that

54   their work collectively has limited reproducibility? Most critiques have focused on inadequate

55   education or the incentives that scientists have to perpetuate the status quo. Simply put, scientists

56   are instructed in "usual practice" and rewarded, directly and indirectly, for doing so [2, 3, 16].

57   There are more subtle reasons too; for example, PIs may worry that choosing an adequate

58   number of animals per experimental group as specified by power estimation, if more than the 8-

59   10 typically used in the field, will create problems in animal care committees who are concerned

60   about reducing overall use of animals in research [21]. However, one of the major factors that

61   causes resistance to change may be that investigators honestly do not have the perception that

62   their own findings lack reproducibility [22].

63

64   In order to get a more detailed understanding of the current situation of biomedical experiments,

65   particularly in behavioral neuroscience, we decided to focus on a single, popular and important

66   preclinical assay, the Forced Swim Test (FST), which has been widely used to screen

67   antidepressants developed as treatments in humans. Proper design of preclinical assays is

68   important because they are used as the basis for translating new treatments to humans [eg., 21,

69   23]. Recently, Kara et al. presented a systematic review and meta-analysis of known

70   antidepressants injected acutely in adult male mice, and reported extremely large mean effect

71   sizes (Cohen's d ranging from 1.6 to 3.0 units) [24]. However, such antidepressants may have

72   been originally chosen for clinical development (at least in part) because of their impressive

73   results in the FST. Thus, in the present study, we have repeated and extended their analysis:

74   making an unbiased random sampling of the FST literature, considering as separate cases

75   whether the mice were assayed in a naïve state vs. in a model of depression or stress, and

3

76    whether the mice were given test agents vs. known clinically prescribed antidepressants regarded

77    as positive controls.

78

79    Our findings demonstrate that the mean effect sizes seen in the experiments were indeed

80    extremely large; most of the experiments analyzed did have adequate sample sizes and did have

81    the power to detect effects of this magnitude. Our data go further to provide explicit guidelines

82    for investigators planning new experiments using the Forced Swim Test, who wish to ensure that

83    they will have adequate power and reproducibility when new, unknown agents are tested. We

84    also suggest the need to develop tools that may help educate scientists to perceive more directly

85    the relationships among effect size, sample size, statistical power, and replicability.

86
87

# Materials and Methods

89    In this study, searching PubMed using the query ["mice" AND "forced swim test" AND

90    "2014/08/03"[PDat] : "2019/08/01"[PDat]] resulted in 737 articles, of which 40 articles were

91    chosen at random using a random number generator. We only scored articles describing assays in

92    which some test agent(s), e.g. drugs or natural products, postulated to have antidepressant

93    properties, were given to mice relative to some control or baseline. Treatments might either be

94    acute or repeated, for up to 28 days prior to testing. Assays involving both male and female mice

95    were included. Articles were excluded if they did not utilize the most common definition of

96    forced swim test measures (i.e., the mice is in a tank for six minutes and during the last four

97    minutes, the duration of immobility is recorded in seconds). We further excluded assays in rats

98    or other species; assays that did not examine test agents (e.g. FST assays seeking to directly

99   compare genetically modified vs. wild-type mice, or comparing males vs. females); interactional

100  assays (i.e., assays to see if agent X blocks the effects of agent Y); and a few studies with

101  extremely complex designs. When more than one FST assay satisfying the criteria was reported

102  in a paper, all assays included were recorded and analyzed. We thus scored a total of 77 assays

103  across 16 articles (S1 File).

104

105  Mean values and standard error were extracted from online versions of the articles by examining

106  graphs, figures legends, and data in text if available. In addition, sample size, p-values and

107  significance level were recorded. When sample size was not provided directly, it was inferred

108  from t-test or ANOVA parameters and divided equally among treatment and groups, rounding up

109  to the nearest whole number if necessary. If only a range for sample size was provided, the

110  average of the range was assigned to all treatments, and rounded up if needed.

111

112  Control baseline immobility times were documented, indicating whether naïve mice were used or

113  mice subjected to a model of depression or stress.  To normalize effect size across experiments,

114  Cohen's d was used since it is the most widely used measure [25, 26].

115
116
117
## Results
119
120  As shown in Table 1, across all assays, the FST effect sizes of both test agents and known

121  clinically prescribed antidepressants regarded as positive controls had mean values in Cohen's d

122  units of -1.67 (95% Confidence Interval: -2.12 to -1.23) and -2.45 (95% CI: -3.34 to -1.55),

123  respectively. (Although Cohen's d units are defined as positive values, we add negative signs

124  here to indicate that immobility times decreased relative to control values.) These are extremely

125 large effects -- twice as large as the standard definition of a "large" effect, i.e. a Cohen's d value

126 of -0.8 [25, 26]!

127

128 The effect sizes of test agents vs. clinically prescribed antidepressants across all assays were not

129 significantly different (two-tailed t-test for difference of means: t = 1.5859, p-value = 0.1202;

130 Wilcoxon rank sum test for difference of medians: W = 839, p-value = 0.1347). We found no

131 evidence for either ceiling or floor effects in these assays, that is, in no case did immobility times

132 approach the theoretical minimum or maximum. The sample sizes (i.e., number of animals per

133 treatment group) averaged 8-9 (Table 2).

134
135
136 **Table 1. Test agents vs. known antidepressants: effect sizes**
137

| | | MEAN | MEDIAN | SD | RANGE | CV |
|---|---|---|---|---|---|---|
| **TEST AGENTS** | **N = 48** | -1.671 | -1.571 | 1.534 | -8.471, 0.759 | 0.918 |
| **ANTIDEPRESSANTS** | **N = 29** | -2.448 | -2.144 | 2.354 | -9.428, 1.702 | 0.961 |

138
139 Shown are effect sizes (in Cohen's d units) for all FST assays that examined test agents and
140 those that examined known clinically prescribed antidepressants regarded as positive controls
141 (regardless of whether the effects achieved statistical significance). The mean effect size,
142 median, range, and coefficient of variation (CV) are shown. The negative signs serve as a
143 reminder that immobility times decreased relative to control values. N refers to the number of
144 assays measured for each category.
145
146 **Table 2. Test agents vs. known antidepressants: sample sizes**
147

| | | MEAN | MEDIAN | SD | RANGE |
|---|---|---|---|---|---|
| **TEST AGENTS** | **N = 48** | 8.31 | 8 | 2.183 | 6, 15 |
| **ANTIDEPRESSANTS** | **N = 29** | 9.12 | 8 | 3.821 | 6, 24 |

148
149 Shown are sample sizes (number of animals per treatment group) for FST assays that examined
150 test agents and those that examined known clinically prescribed antidepressants regarded as
151 positive controls.
152
153

154 **Assays in naïve mice vs. in models of depression or stress**

155

6

156     Agents were tested for antidepressant effects in both naïve mice and mice subjected to various

157     models of depression or stress. To our surprise, although one might expect longer baseline

158     immobility times in "depressed" mice, our data indicate that the mean baseline immobility times

159     of naïve and "depressed" mice (Table 3) did not differ significantly (one tailed t-test: p-value =

160     0.3375).

161
162     **Table 3. Control baseline immobility times in seconds**
163

| | | MEAN | MEDIAN | SD | RANGE |
|---|---|---|---|---|---|
| **NAÏVE** | **N = 63** | 143.817 | 159 | 38.985 | 56, 208 |
| **DEPRESSED** | **N = 14** | 148.643 | 175 | 36.923 | 93, 184 |

164
165
166
167     We then examined the effect sizes of test agents in naïve vs. depressive models (Table 4). There

168     were no significant differences in mean effect size for test agents in naïve vs. depressed mice

169     (two-tailed t-test t = -0.61513, p-value = 0.5423). Interestingly, the assays in depressed models

170     showed a smaller coefficient of variation (i.e., standard deviation divided by the mean) than in

171     naïve mice. A smaller coefficient of variation in depressed models means that they show less

172     intrinsic variability, which in turn means that it is easier for a given effect size to achieve

173     statistical significance.

174
175
176
177     **Table 4. Test agents and known antidepressants in naïve vs. depressed models: Effect sizes**
178

| | | | MEAN | MEDIAN | SD | RANGE | CV |
|---|---|---|---|---|---|---|---|
| **TEST AGENTS** | Naïve | N = 37 | -1.729 | -1.731 | 1.717 | -8.471, 0.759 | 0.993 |
| | Depressed | N = 11 | -1.496 | -1.231 | 0.826 | -3.406, -0.557 | 0.552 |
| **ANTIDEPRESSANTS** | Naïve | N = 26 | -2.554 | -2.389 | 2.492 | -9.428, 1.702 | 0.975 |
| | Depressed | N = 3 | -2.115 | -0.856 | 2.255 | -4.718, -0.771 | 1.066 |

179
180     Shown are effect sizes (in Cohen's d units) for FST assays that examined test agents and those
181     that examined known clinically prescribed antidepressants, in naïve or depressed models,
182     respectively.
183
184

7

## Reporting parameters

None of the 16 randomly chosen articles in our dataset mentioned whether the FST assay was blinded to the group identity of the mouse being tested (although some did use automated systems to score the mice). None presented the raw data (immobility times) for individual mice. None discussed data issues such as removal of outliers, or whether the observed distribution of immobility times across animals in the same group was approximately normal or skewed. Only one mentioned power estimation at all (though no details or parameters were given). All studies utilized parametric statistical tests (t-test or ANOVA), which were either two-tailed or unspecified -- none specified explicitly that they were using a one-tailed test.

## Discussion

Our literature analysis of the Forced Swim Test in mice agrees with, and extends, the previous meta-analysis of Kara et al [24], which found that known antidepressants exhibit extremely large effect sizes across a variety of individual drugs and mouse strains. The first question that might be asked is whether the effects might be tainted by publication bias, i.e., if negative or unimpressive results were less likely to be published [10]. Ramos-Hryb et al. failed to find evidence for publication bias in FST studies of imipramine [27]. We cannot rule out bias against publishing negative results in the case of FST studies of test agents (i.e. agents not already clinically prescribed as antidepressants in humans), since nearly all articles concerning test agents reported positive statistically significant results (though not every assay in every article was significant). On the other hand, most if not all of the agents tested were not chosen at random, but had preliminary or indirect (e.g., receptor binding) findings in favor of their hypothesis.

209

210    The immobility time measured by the FST may reflect a discontinuous yes/no behavioral

211    decision by mice, rather than a continuous variable like running speed or spontaneous activity.

212    Kara et al [24] observed that the FST test does not exhibit clear dose-response curves in most of

213    the published experiments that looked for them, which further suggests a switch-like rather than

214    graded response of the mice. This phenomenon may partially explain why effects in the FST

215    appear to be very large and robust, and it complicates efforts to assess whether the effect sizes

216    reported in the literature are inflated due to positive publication bias or low statistical power.

217

218    Surprisingly, we found that the baseline immobility time of naïve mice was not significantly

219    different than the baseline immobility time of mice subjected to various models of depression or

220    chronic stress (Table 2). This might potentially be explained by high variability of baseline

221    values across heterogeneous experiments and laboratories. Alternatively, naïve mice housed and

222    handled under routine conditions may be somewhat "depressed" insofar as they have longer

223    immobility times relative to those housed in more naturalistic environments [28].

224

225    ## Guidelines for investigators using FST assays

226    One of the reasons that investigators rarely calculate prospective power estimations is the

227    difficulty in ascertaining the necessary parameters accurately. Our results provide explicit values

228    for these parameters for the FST, at least for the simple designs that are represented in our

229    dataset. For example, for two independent groups of mice treated with an unknown test agent vs.

230    control, one needs to enter a) the baseline immobility time expected in the control group (Table

231    3), b) the expected immobility time for the treated group (at the minimum biologically

232    meaningful effect size that the investigator wishes to detect), c) the standard deviations of each

233    group (Table 1), and d) the relative number of animals in each group (generally 1:1).

234    Alternatively, one can enter the minimum biologically relevant effect size in Cohen's d units that

235    the investigator wants to be able to detect (this encompasses both the difference in immobility

236    times in the two groups as well as their standard deviations) (Table 5). This is sufficient to

237    estimate the required number of animals per group (Table 5), assuming two groups (treated vs.

238    control), standard criteria of power = 0.8, false-positive rate = 0.05, and a parametric statistical

239    test (t-test or ANOVA).

240

## 241    But the power of current FST assays is adequate, isn't it?

242    From Tables 1 and 4, one can see that the observed mean effect sizes across the literature fall

243    into the range of 1.5 to 2.5 Cohen's d units and for the sake of this discussion, we will assume

244    that these values are not inflated. Indeed, if an investigator merely wants to be able to detect

245    effects of this size, only 7-8 animals per group are required, which is in line with the number

246    actually used in these experiments (Table 5). This is likely to explain why scientists in this field

247    have the intuition that the empirical standard sample size of 8-9 (Table 2) is enough to ensure

248    adequate power.

249

250    **Table 5. Prospective Power Estimation for test agents in the FST assay.**

|  | EFFECT SIZE | #ANIMALS REQUIRED PER GROUP |
|---|---|---|
| **MODERATE ES** | -0.5 | 64 |
| **LARGE ES** | -0.8 | 26 |
| **MEAN ES (THIS STUDY)** | -1.671 | 7 |
| **MEDIAN ES (THIS STUDY)** | -1.572 | 7 |

251    These sample size calculations are based on the observed mean and median effect sizes (ES) in
252    Cohen's d units for novel test agents (Table 1), two groups (treated vs. controls), for desired
253    power=0.8, alpha=0.05, and two-sided t-test or ANOVA [25].

10

254

255  However, setting the **minimum** effect size at the observed **mean (or median)** value is clearly

256  not satisfactory since half of the assays fall below that value. When an investigator is examining

257  an unknown test agent, the general guidance is to set the minimum effect size at "moderate" (0.5)

258  if not "large" (0.8) [29], which would require 64 or 26 animals per group, respectively, in order

259  to ensure adequate power (Table 5). Setting the minimum effect size is not something to be

260  fixed, and depends not only on the assay but also on the investigator's hypothesis to be tested

261  [30]. Nevertheless, the appropriate minimum should always be set smaller than the mean

262  observed effect size of the assay as a whole, especially when the agent to be tested lacks

263  preliminary evidence showing efficacy.  From this perspective, a new FST experiment planned

264  using 7-10 animals will be greatly under-powered. Nevertheless, this does shed light on why

265  scientists performing the FST assay may not intuitively perceive that their experiments are

266  under-powered.

267

268  ## Possible experimental design strategies for improved power

269  **One tail or two?** Investigators in our dataset never stated that they used one-tailed statistical

270  tests, even though they generally had preliminary or suggestive prior evidence suggesting that

271  the agent being tested may have antidepressant effects in the FST. Using a one-tailed hypothesis

272  in prospective power estimation reduces the number of animals needed per group, for the same

273  power and false-positive rate. For a minimum effect size of 0.8, a two-tailed hypothesis that

274  requires 26 animals per group reduces to 21 animals per group for a one-tailed hypothesis [31].

275

276    In summary, for testing an unknown agent (e.g., chosen without prior experimental evidence or

277    as part of a high-throughput screen), with minimum effect size = 0.8, power = 0.8 and false-

278    positive rate = 0.05, the results suggest that an investigator should use a two-tailed hypothesis

279    and will need ~26 animals per group. (High throughput assays will need additional post hoc

280    corrections for multiple testing.) For a test agent which has preliminary or prior evidence in favor

281    of being an antidepressant, a one-tailed hypothesis is appropriate and ~21 animals per group can

282    be used.  Note that this discussion applies to simple experimental designs only. Interactional

283    assays (e.g., does agent X block the effects of agent Y?) are expected to have larger standard

284    deviations than direct assays and would require somewhat larger sample sizes, as would complex

285    experimental designs of any type.

286

287    **Parametric or nonparametric testing?** All experiments in our dataset employed parametric

288    statistical tests, either ANOVA or t-test. This is probably acceptable when sample sizes of 20 or

289    more are employed, as recommended in the present paper, but not for the usual 7-10 animals per

290    group, as performed by most of the investigators in our dataset. This is for two reasons: First,

291    investigators in our dataset have not presented the raw data for individual animals in each group

292    to verify that the underlying data distribution across individuals resembles a normal distribution.

293    Second, when sample sizes are so small, parametric tests have a tendency to ascribe too much

294    significance to a finding [14], and together with the issue of inflated effect sizes, this results in

295    over-optimistic prospective power estimation.  Nonparametric tests such as the Wilcoxon signed

296    rank test (with either one-tailed or two-tailed hypothesis) are appropriate regardless of normality,

297    and will be more conservative than parametric tests, i.e. will have less tendency to ascribe too

298    much significance to a finding [14]. Popular software including G*Power are able to handle

299    nonparametric testing [31]. A warning though: Using a nonparametric test will result in estimates

300    of required sample sizes larger than those obtained using parametric tests.

301

302    **Within-animal design?** None of the assays in our dataset involved a before/after design in the

303    same animals. This means giving a control vs. an agent to a mouse, observing the immobility

304    time in the FST assay, then repeating the assay in the same mouse with the other treatment.

305    Using an individual mouse as its own control has the advantage of less variability (i.e. no inter-

306    animal variability needs to be considered) and allows the investigator to use paired statistics

307    instead of unpaired tests. Both of these advantages should tend to increase power for the same

308    number of animals, plus, one can divide the number of total animals needed in half since each

309    one is its own control. Unfortunately, control baseline immobility times are not stable on

310    retesting, and investigators have found that the test-retest scheme results in similar effect sizes as

311    the standard assay in some but not all cases [25, 32-34]. Thus, one would need to employ test-

312    retest FST paradigms with some caution and with extra controls.

313

## Limitations of our study

315    Our literature analysis did not examine how effect sizes may vary across mouse strain, or across

316    individual drugs [24].  We also did not undertake a Bayesian analysis to estimate the prior

317    probability that any given test agent chosen at random will have antidepressant effects in the FST

318    assay. We did not consider how power might be affected if animals are not truly independent

319    (e.g. they may be littermates) and if they are not randomly allocated to groups [35]. Our

320    guidelines do not encompass designs in which the sample size is not pre-set at the outset [36].

321    Finally, we did not directly assess the replicability of published FST experiments, i.e., if one

322     publication reports a statistically significant finding, what is the probability that another group

323     examining the same question will also report that the finding is statistically significant?

324     Replicability is related to adequate statistical power but also involves multiple aspects of

325     experimental design not considered here [2, 5, 8, 11, 13, 19, 37]. Nevertheless, adequate power is

326     essential for experiments to be replicable, because under-powered studies tend to over-estimate

327     effect sizes and have inflated false-positive rates [4, 38].

328

329

## Conclusions

331     In the case of the Forced Swim Test used to assess antidepressant actions of test agents in mice,

332     we found that the mean effect size is extremely large (i.e., 1.5 - 2.5 in Cohen's d units), so large

333     that only 7-10 animals per group are needed to reliably detect a difference from controls. This

334     may shed light on why scientists in neuroscience, and preclinical biomedical research in general,

335     have the intuition that their usual practice (7-10 animals per group) provides adequate statistical

336     power, when many meta-science studies have shown that the overall field is greatly under-

337     powered. The large mean effect size may at least partially explain why investigators using the

338     FST do not perceive intuitively that their experimental designs fall short. It can be argued that

339     when effects are so large, relatively small sample sizes may be acceptable [39]. The Forced

340     Swim Test is not unique – to name one example, rodent fear conditioning is another popular

341     preclinical assay that exhibits extremely large effect sizes [40]. Nevertheless, we showed that

342     adequate power to detect minimum biologically relevant large effects in this assay actually

343     requires at least ~21-26 animals per group when the true effect of a test agent is unknown.

344

345  We suggest that investigators are not able to perceive intuitively whether or not a given sample

346  size is adequate for a given experiment, and this contributes to a mindset that is skeptical of

347  theoretical or statistical arguments.  Apart from other educational and institutional reforms [2, 3,

348  10, 11, 13, 19, 21, 37, 41], altering the real-life behavior of scientists in planning their

349  experiments may require developing tools that allow them to actively visualize the inter-

350  relationships among effect size, sample size, statistical power, and replicability in a direct and

351  intuitive manner.

352

353  **COMPETING INTERESTS**

354
355  The authors attest that they have no competing interests.

356

357

358
359  **FUNDING**

360  This work was supported by NIH Grants R01LM010817 and P01AG039347. The study sponsor

361  had no role in study design; in the collection, analysis and interpretation data; in the writing of

362  the report; or in the decision to submit the paper for publication.

363
364  **CONTRIBUTORSHIP STATEMENT**

365

366  NS – Conceived of the study, supervised data extraction, and wrote the initial draft of the paper.

367  EG - Assisted with data extraction and carried out data analysis. JW - Supervised data analysis.

368  ZY - Assisted with data extraction. All authors participated in writing the paper.

369

370  # References

15

371    1. Ioannidis JP. Why most published research findings are false. PLoS Med. 2005

372        Aug;2(8):e124. doi: 10.1371/journal.pmed.0020124.

373    2. Ioannidis JP. How to make more published research true. PLoS Med. 2014 Oct

374        21;11(10):e1001747. doi: 10.1371/journal.pmed.1001747.

375    3. Higginson AD, Munafò MR. Current Incentives for Scientists Lead to Underpowered

376        Studies with Erroneous Conclusions. PLoS Biol. 2016 Nov 10;14(11):e2000995. doi:

377        10.1371/journal.pbio.2000995. – half are wrong

378    4. Button KS, Ioannidis JP, Mokrysz C, Nosek BA, Flint J, Robinson ES, Munafò MR.

379        Power failure: why small sample size undermines the reliability of neuroscience. Nat Rev

380        Neurosci. 2013 May;14(5):365-76. doi: 10.1038/nrn3475.

381    5. Curran-Everett D. Explorations in statistics: statistical facets of reproducibility. Adv

382        Physiol Educ. 2016 Jun;40(2):248-52. doi: 10.1152/advan.00042.2016.

383    6. Open Science Collaboration. PSYCHOLOGY. Estimating the reproducibility of

384        psychological science. Science. 2015 Aug 28;349(6251):aac4716. doi:

385        10.1126/science.aac4716.

386    7. Dumas-Mallet E, Button KS, Boraud T, Gonon F, Munafò MR. Low statistical power in

387        biomedical science: a review of three human research domains. R Soc Open Sci. 2017

388        Feb 1;4(2):160254. doi: 10.1098/rsos.160254.

389    8. Tsilidis KK, Panagiotou OA, Sena ES, Aretouli E, Evangelou E, Howells DW, Al-Shahi

390        Salman R, Macleod MR, Ioannidis JP. Evaluation of excess significance bias in animal

391        studies of neurological diseases. PLoS Biol. 2013 Jul;11(7):e1001609. doi:

392        10.1371/journal.pbio.1001609.

393    9.  Szucs D, Ioannidis JP. Empirical assessment of published effect sizes and power in the

394        recent cognitive neuroscience and psychology literature. PLoS Biol. 2017 Mar

395        2;15(3):e2000797. doi: 10.1371/journal.pbio.2000797.

396    10. Sena ES, van der Worp HB, Bath PM, Howells DW, Macleod MR. Publication bias in

397        reports of animal stroke studies leads to major overstatement of efficacy. PLoS Biol.

398        2010 Mar 30;8(3):e1000344. doi: 10.1371/journal.pbio.1000344.

399    11. Howells DW, Sena ES, Macleod MR. Bringing rigour to translational medicine. Nat Rev

400        Neurol. 2014 Jan;10(1):37-43. doi: 10.1038/nrneurol.2013.232.

401    12. Lazic SE, Clarke-Williams CJ, Munafò MR. What exactly is 'N' in cell culture and

402        animal experiments? PLoS Biol. 2018 Apr 4;16(4):e2005282. doi:

403        10.1371/journal.pbio.2005282.

404    13. Munafò MR, Davey Smith G. Robust research needs many lines of evidence. Nature.

405        2018 Jan 25;553(7689):399-401. doi: 10.1038/d41586-018-01023-3.

406    14. Smalheiser NR. Data literacy: How to make your experiments robust and reproducible.

407        Academic Press; 2017 Sep 5.

408    15. Nord CL, Valton V, Wood J, Roiser JP. Power-up: A Reanalysis of 'Power Failure' in

409        Neuroscience Using Mixture Modeling. J Neurosci. 2017 Aug 23;37(34):8051-8061. doi:

410        10.1523/JNEUROSCI.3592-16.2017.

411    16. Smaldino PE, McElreath R. The natural selection of bad science. R Soc Open Sci. 2016

412        Sep 21;3(9):160384. doi: 10.1098/rsos.160384.

413    17. Vankov I, Bowers J, Munafò MR. On the persistence of low power in psychological

414        science. Q J Exp Psychol (Hove). 2014 May;67(5):1037-40. doi:

415        10.1080/17470218.2014.885986.

416    18. Ioannidis, J. P. A. Why most discovered true associations are inflated. *Epidemiology* **19**,

417         640–648 (2008).

418    19. Simmons JP, Nelson LD, Simonsohn U. False-positive psychology: undisclosed

419         flexibility in data collection and analysis allows presenting anything as significant.

420         Psychol Sci. 2011 Nov;22(11):1359-66. doi: 10.1177/0956797611417632.

421    20. Mullard A. Cancer reproducibility project yields first results. Nat Rev Drug Discov. 2017

422         Feb 2;16(2):77. doi: 10.1038/nrd.2017.19.

423    21. Steward O, Balice-Gordon R. Rigor or mortis: best practices for preclinical research in

424         neuroscience. Neuron. 2014 Nov 5;84(3):572-81. doi: 10.1016/j.neuron.2014.10.042.

425    22. Fitzpatrick BG, Koustova E, Wang Y. Getting personal with the "reproducibility crisis":

426         interviews in the animal research community. Lab Anim (NY). 2018 Jul;47(7):175-177.

427         doi: 10.1038/s41684-018-0088-6.

428    23. Steckler T. Editorial: preclinical data reproducibility for R&D - the challenge for

429         neuroscience. Springerplus. 2015 Jan 13;4(1):1. doi: 10.1186/2193-1801-4-1.

430    24. Kara NZ, Stukalin Y, Einat H. Revisiting the validity of the mouse forced swim test:

431         Systematic review and meta-analysis of the effects of prototypic antidepressants.

432         Neurosci Biobehav Rev. 2018 Jan;84:1-11. doi: 10.1016/j.neubiorev.2017.11.003.

433    25. Lakens D. Calculating and reporting effect sizes to facilitate cumulative science: a

434         practical primer for t-tests and ANOVAs. Front Psychol. 2013 Nov 26;4:863. doi:

435         10.3389/fpsyg.2013.00863.

436    26. Cumming, G. (2012). Understanding the New Statistics: Effect sizes, Confidence

437         Intervals, and Meta-Analysis. New York, NY: Routledge.

438    27. Ramos-Hryb AB, Harris C, Aighewi O, Lino-de-Oliveira C. How would publication bias

439         distort the estimated effect size of prototypic antidepressants in the forced swim test?

440         Neurosci Biobehav Rev. 2018 Sep;92:192-194. doi: 10.1016/j.neubiorev.2018.05.025.

441    28. Bogdanova OV, Kanekar S, D'Anci KE, Renshaw PF. Factors influencing behavior in the

442         forced swim test. Physiol Behav. 2013 Jun 13;118:227-39. doi:

443         10.1016/j.physbeh.2013.05.012.

444    29. Calin-Jageman RJ. The New Statistics for Neuroscience Majors: Thinking in Effect

445         Sizes. J Undergrad Neurosci Educ. 2018 Jun 15;16(2):E21-E25.

446    30. Ashton JC. Experimental power comes from powerful theories - the real problem in null

447         hypothesis testing. Nat Rev Neurosci. 2013 Aug;14(8):585. doi: 10.1038/nrn3475-c2.

448    31. Faul F, Erdfelder E, Lang AG, Buchner A. G*Power 3: a flexible statistical power

449         analysis program for the social, behavioral, and biomedical sciences. Behav Res

450         Methods. 2007 May;39(2):175-91. doi: 10.3758/bf03193146.

451    32. Su J, Hato-Yamada N, Araki H, Yoshimura H. Test-retest paradigm of the forced

452         swimming test in female mice is not valid for predicting antidepressant-like activity:

453         participation of acetylcholine and sigma-1 receptors. J Pharmacol Sci. 2013;123(3):246-

454         55. doi: 10.1254/jphs.13145fp.

455    33. Mezadri TJ, Batista GM, Portes AC, Marino-Neto J, Lino-de-Oliveira C. Repeated rat-

456         forced swim test: reducing the number of animals to evaluate gradual effects of

457         antidepressants. J Neurosci Methods. 2011 Feb 15;195(2):200-5. doi:

458         10.1016/j.jneumeth.2010.12.015.

459    34. Calil CM, Marcondes FK. The comparison of immobility time in experimental rat

460         swimming models. Life Sci. 2006 Sep 27;79(18):1712-9. doi: 10.1016/j.lfs.2006.06.003.

19

461    35. Lazic SE. Four simple ways to increase power without increasing the sample size. Lab

462        Anim. 2018 Dec;52(6):621-629. doi: 10.1177/0023677218767478.

463    36. Neumann K, Grittner U, Piper SK, Rex A, Florez-Vargas O, Karystianis G, Schneider A,

464        Wellwood I, Siegerink B, Ioannidis JP, Kimmelman J, Dirnagl U. Increasing efficiency

465        of preclinical research by group sequential designs. PLoS Biol. 2017 Mar

466        10;15(3):e2001307. doi: 10.1371/journal.pbio.2001307.

467    37. Snyder HM, Shineman DW, Friedman LG, Hendrix JA, Khachaturian A, Le Guillou I,

468        Pickett J, Refolo L, Sancho RM, Ridley SH. Guidelines to improve animal study design

469        and reproducibility for Alzheimer's disease and related dementias: For funders and

470        researchers. Alzheimers Dement. 2016 Nov;12(11):1177-1185. doi:

471        10.1016/j.jalz.2016.07.001.

472    38. Marino MJ. How often should we expect to be wrong? Statistical power, P values, and

473        the expected prevalence of false discoveries. Biochem Pharmacol. 2018 May;151:226-

474        233. doi: 10.1016/j.bcp.2017.12.011.

475    39. Dumas-Mallet E, Button K, Boraud T, Munafo M, Gonon F. Replication Validity of

476        Initial Association Studies: A Comparison between Psychiatry, Neurology and Four

477        Somatic Diseases. PLoS One. 2016 Jun 23;11(6):e0158064. doi:

478        10.1371/journal.pone.0158064.

479    40. Carneiro CFD, Moulin TC, Macleod MR, Amaral OB. Effect size and statistical power in

480        the rodent fear conditioning literature - A systematic review. PLoS One. 2018 Apr

481        26;13(4):e0196258. doi: 10.1371/journal.pone.0196258.

482   41. Wass MN, Ray L, Michaelis M. Understanding of researcher behavior is required to

483       improve data reliability. Gigascience. 2019 May 1;8(5):giz017. doi:

484       10.1093/gigascience/giz017.

485

## Supporting information

487

488   **S1 File. Excel spreadsheet of the FST assays scored in this paper.** A total of 77 FST assays

489   across 16 randomly chosen articles were scored. See text for details.