

Meta-research article

Effect size and statistical power in the rodent fear conditioning
literature – a systematic review

Short title: Effect sizes and statistical power in fear conditioning

Authors: Clarissa F. D. Carneiro^{1*}, Thiago C. Moulin^{1*}, Malcolm R. Macleod², Olavo B. Amaral¹

¹Institute of Medical Biochemistry Leopoldo de Meis, Federal University of Rio de Janeiro, Rio de Janeiro, Brazil and ²Division of Clinical Neurosciences, University of Edinburgh, Edinburgh, UK

* Both authors contributed equally to the work.

Corresponding author:

Olavo B. Amaral, M.D., PhD.

Instituto de Bioquímica Médica Leopoldo de Meis
Av. Carlos Chagas Filho 373, E-38
Cidade Universitária
Rio de Janeiro, RJ, Brazil
CEP 21941-902
Phone: +55-21-39386762
E-mail: olavo@bioqmed.ufrj.br

Abstract

Proposals to increase research reproducibility frequently call for focusing on effect sizes instead of p values, as well as for increasing the statistical power of experiments. However, it is unclear to what extent these two concepts are indeed taken into account in basic biomedical science. To study this in a real-case scenario, we performed a systematic review of effect sizes and statistical power in studies using rodent fear conditioning, a widely used behavioral task to evaluate learning and memory. Our search criteria yielded 410 experiments comparing control and treated groups in 122 articles. Interventions with statistically significant differences had a mean effect size of 45.6%, and amnesia caused by memory-impairing interventions was nearly always partial. Mean statistical power to detect the average effect size observed in well-powered experiments (37.2%) was 65%, and was lower in studies with non-significant results. Only one article reported a sample size calculation, and our estimated sample size to achieve 80% power considering typical effect sizes and variances (15 animals per group) was reached in only 12.2% of experiments. Effect size correlated with textual descriptions of results only when findings were non-significant, and neither effect size nor power correlated with study quality indicators, number of citations or impact factor of articles. In summary, effect sizes and statistical power have a wide distribution in the rodent fear conditioning literature, but do not seem to have a large influence on how results are described or cited. Failure to take these concepts into consideration might limit attempts to improve reproducibility in this field of science.

Introduction

Biomedical research over the last decades has relied heavily on the concept of statistical significance – i.e. the probability that a difference would occur by chance under the null hypothesis – and classifying results as “significant” or “non-significant” on the basis of an arbitrary threshold (usually set at $p < 0.05$) has become standard practice in most fields. This approach, however, has well-described limitations that can lead to erroneous conclusions when researchers rely on p values alone to judge results (1–5). First of all, p values do not measure the magnitude of an effect, and thus are not indicators of biological significance (6). Moreover, the predictive value of a significance test is heavily influenced by factors such as the prior probability of the tested hypothesis, the number of tests performed and their statistical power (7); thus, similar p values can lead to very different conclusions in different scenarios (1).

Recent calls for improving research reproducibility have focused on reporting effect sizes and confidence intervals alongside or in place of p values (5,6,8) and for the use of both informal Bayesian inference (9) and formal data synthesis methods (10) when aggregating data from multiple studies. The concepts of effect size and statistical power are central for such approaches, as how much a given experiment will change a conclusion or an effect estimate will depend on both. However, it is unclear whether they receive much attention from authors in basic science publications. Discussion of effect sizes seems to be scarce, and recent data has shown that sample size and power calculations are very rare in the preclinical literature (11,12). The potential impact of these omissions is large, as reliance on the results of significance tests without consideration of statistical power can decrease the reliability of study conclusions when studies are underpowered (13).

Another issue is that, if effect size is not taken into account, it is difficult to adequately assess the biological significance of a given finding. As p values will be low even for small effect sizes if sample size is sufficiently large, biologically unimportant effects can be found to be statistically significant. Thus, if effect sizes are not considered, it is difficult to dissect essential mechanisms in a given physiological pathway from more peripheral modulatory influences (14). Thus, the wealth of findings in the literature translates poorly into better comprehension of phenomena, and the abundance of statistically significant findings with small effect sizes can eventually do more harm than good. This problem is made much worse when many of these studies have low positive predictive values due to insufficient power, leading a large fraction of them to be false positives (7,15,16).

To analyze how effect sizes and statistical power are taken into account in the description and publication of findings in a real-case scenario of basic biomedical science, we chose to perform a systematic review of articles using rodent fear conditioning, probably the most widely used behavioral task to study learning and memory in animals (17). Focusing on this task provides a major advantage in the fact that the vast majority of articles use the same measure to describe results (i.e. percentage of time spent in freezing behavior during a test session). Thus, effect sizes are comparable across studies, and studying their distribution allows one to estimate the statistical power of individual experiments to detect typical differences. After analyzing the distribution of effect sizes and statistical power in a large sample of experiments, we will evaluate how they correlate with (a) the outcome of significance tests, (b) aspects of experimental design, (c) descriptions of experimental results, (d) risk of bias indicators and (e) citation metrics from the corresponding articles.

Results

Article search and inclusion

As previously described in a protocol published in advance of full data collection (18), we performed a PubMed search for “fear conditioning” AND (“learning” OR “consolidation” OR “acquisition”) AND (“mouse” OR “mice” OR “rat” OR “rats”)” among articles published online in 2013. The search process (**Fig. 1**) yielded 400 search hits, of which 386 were original articles that were full-text screened and included if they (a) described the effects of interventions in fear conditioning acquisition or consolidation in rodents, (b) had a proper control group for comparison, (c) used freezing behavior in a test session as a measure of memory and (d) had available data on mean freezing, standard deviation (SD) or standard error of mean (SEM), sample size and statistical significance. Two investigators examined all included articles, and agreement for exclusions measured on a double-screened sample of 40 articles was 95%. This led to a final sample of 122 articles and 410 experiments, from which various features were extracted to build the database provided as **S1 Data**.

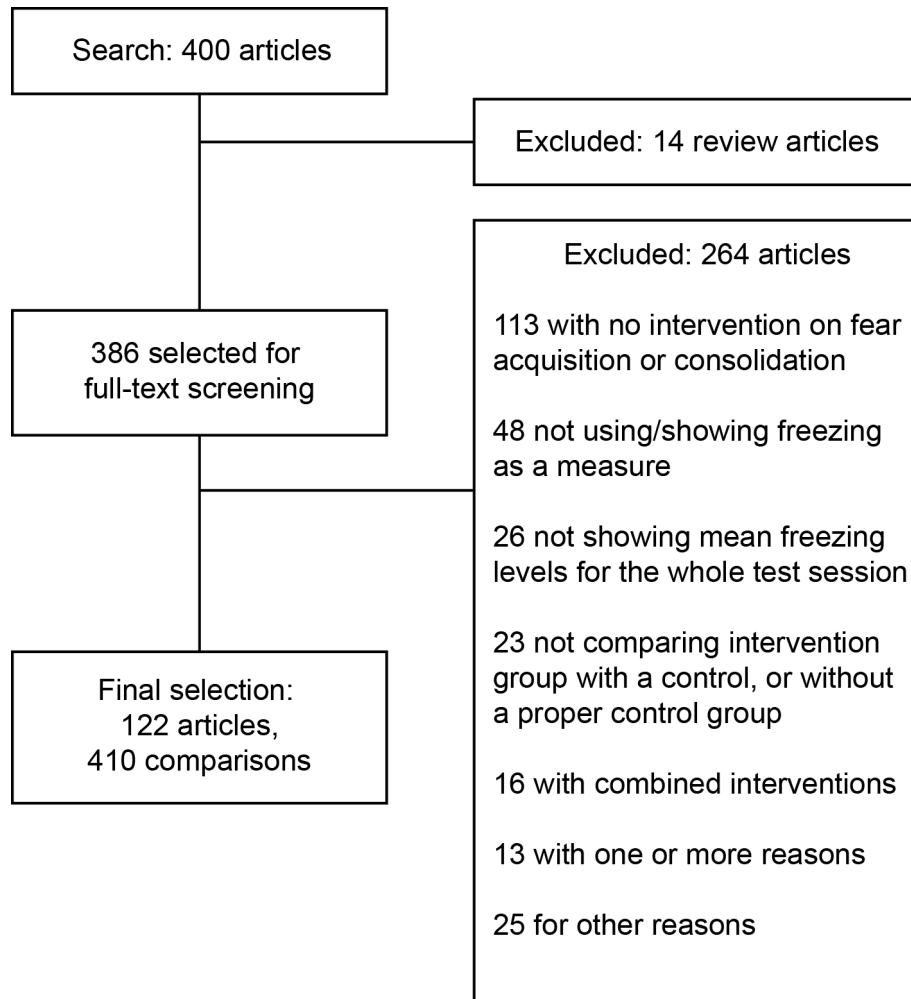


Figure 1. Study flow diagram. Our PubMed search yielded 400 results, of which 14 were excluded based on initial screening of titles and abstracts and 386 were selected for full-text analysis. This final analysis led to the inclusion of 122 articles, containing a total of 410 comparisons (i.e. individual experiments). The main reasons for exclusion are listed in the figure.

Distribution of effect sizes and statistical power

For each experiment, we initially calculated effect size as the percentage of increase or decrease in the freezing levels of treated groups when compared to controls. As shown in **Fig. 2A**, this leads interventions that enhance memory acquisition (i.e. those in which freezing is significantly higher in the treated group) to have larger effect sizes than those that impair it (i.e. those in which freezing is significantly lower in the treated group) due to an asymmetry that is inherent to ratios. To account for this and

make effect sizes comparable between both types of interventions, we used a normalized effect size, in which the difference is expressed as a percentage of the highest freezing value, whether in the treatment or control group (**Fig. 2B**) (10). Freezing levels in the reference group correlated negatively with normalized effect size (in %) and pooled coefficient of variation (defined as the ratio between the sample size-weighted pooled SD and pooled mean of each experiment), leading it to correlate positively with statistical power; however, this form of normalization reduced this effect when compared to others (**S1 Fig.**). We also calculated effect sizes in Cohen's *d* (**S2 Fig.**), but chose to use percentages throughout the study, as these are more closely related to the way results are expressed in articles.

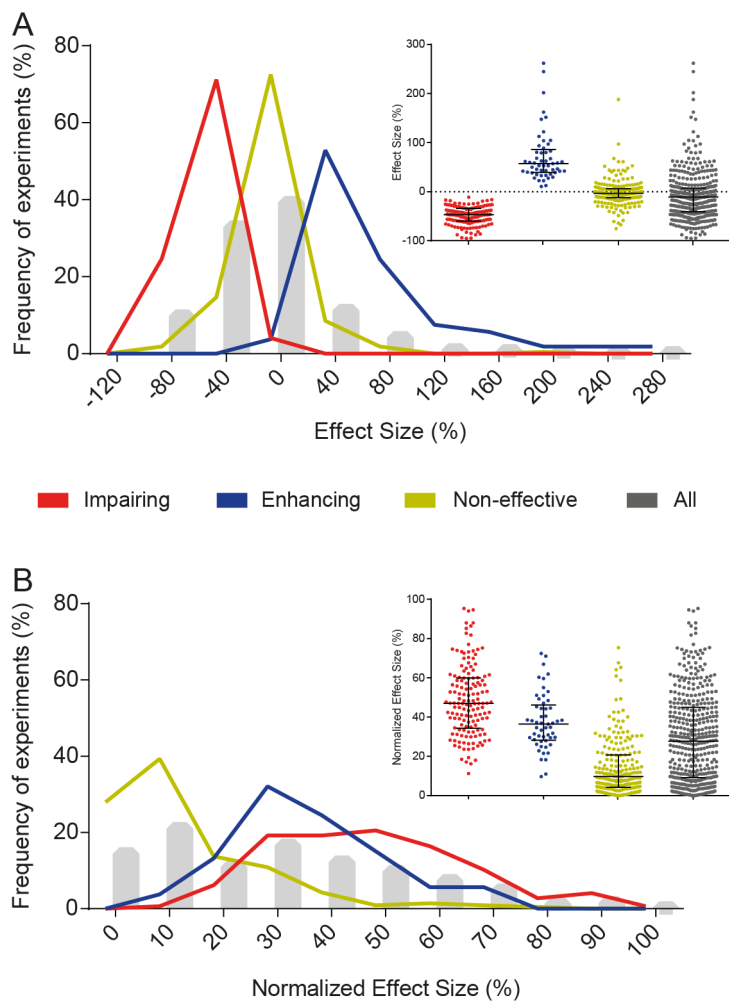


Figure 2. Distribution of effect sizes. (A) Distribution of effect sizes, calculated as percentage of control group freezing. Interventions were divided into memory-impairing ($-48.6 \pm 18.1\%$ [-

51.5 to -45.6], n=146), memory-enhancing ($71.6 \pm 53.2\%$ [56.9 to 86.2], n=53) or non-effective ($-1.8 \pm 26.2\%$ [-5.3 to 1.8], n=211) according to the statistical significance of the comparison. Additionally, the whole sample of experiments is shown in grey ($-9.0 \pm 47.5\%$ [-13.6 to -4.4], n=410). Values are expressed as mean \pm SD [95% confidence interval]. Lines and whiskers in the inset express median and interquartile interval. (B) Distribution of normalized effect sizes, calculated as percentage of the group with the highest mean (i.e. the control group for memory-impairing interventions, or the treated group for memory-enhancing interventions).

Mean normalized effect size was $48.6 \pm 18.1\%$ [45.6 to 51.5] for memory-impairing interventions, $37.6 \pm 14.2\%$ [33.7 to 41.6] for memory-enhancing interventions and $14.4 \pm 14.2\%$ [12.4 to 16.3] for non-effective interventions – i.e. those in which a significant difference between groups was not found (all measures are expressed as mean \pm SD [95% confidence interval]). All 410 experiments combined had a mean normalized effect size of $29.5 \pm 22.4\%$ [27.4 to 31.7]. The distribution of effect sizes shows that the vast majority of memory-impairing interventions cause partial reductions in learning, leaving the treated group with residual freezing levels that are significantly higher than those of a non-conditioned animal. In fact, in all 44 memory-impairing experiments in which pre-conditioning freezing levels were shown for the treated group, these were lower than those observed in the test session – with p values below 0.05 in 32 (82%) out of the 39 cases in which there was enough information for us to perform an unpaired *t* test between sessions.

Among non-significant results, it is worth noting that 26.5% of them had an effect size greater than 20%, suggesting that these experiments might have been underpowered. With this in mind, we sought to evaluate if the distribution of mean sample size and coefficients of variation (both of which are determinants of statistical power) differed between memory-impairing, memory-enhancing and non-effective interventions. As shown in Fig. 3A, most experiments had sample sizes between 8 and 12 animals/group, and this distribution did not vary between enhancing, impairing and non-effective interventions (one-way ANOVA, $p=0.30$). On the other hand, higher

coefficients of variation were more frequent among non-effective interventions (Fig. **3B**, one-way ANOVA, $p=0.001$). This difference in variability was partly explained by freezing levels, which were lower on average for the reference group in non-significant experiments (49.3% vs. 52.9% in memory-impairing and 61.3% in memory-enhancing experiments, one way ANOVA, $p=0.0006$), and was no longer significant after adjustment for this variable (ANCOVA, $p=0.33$).

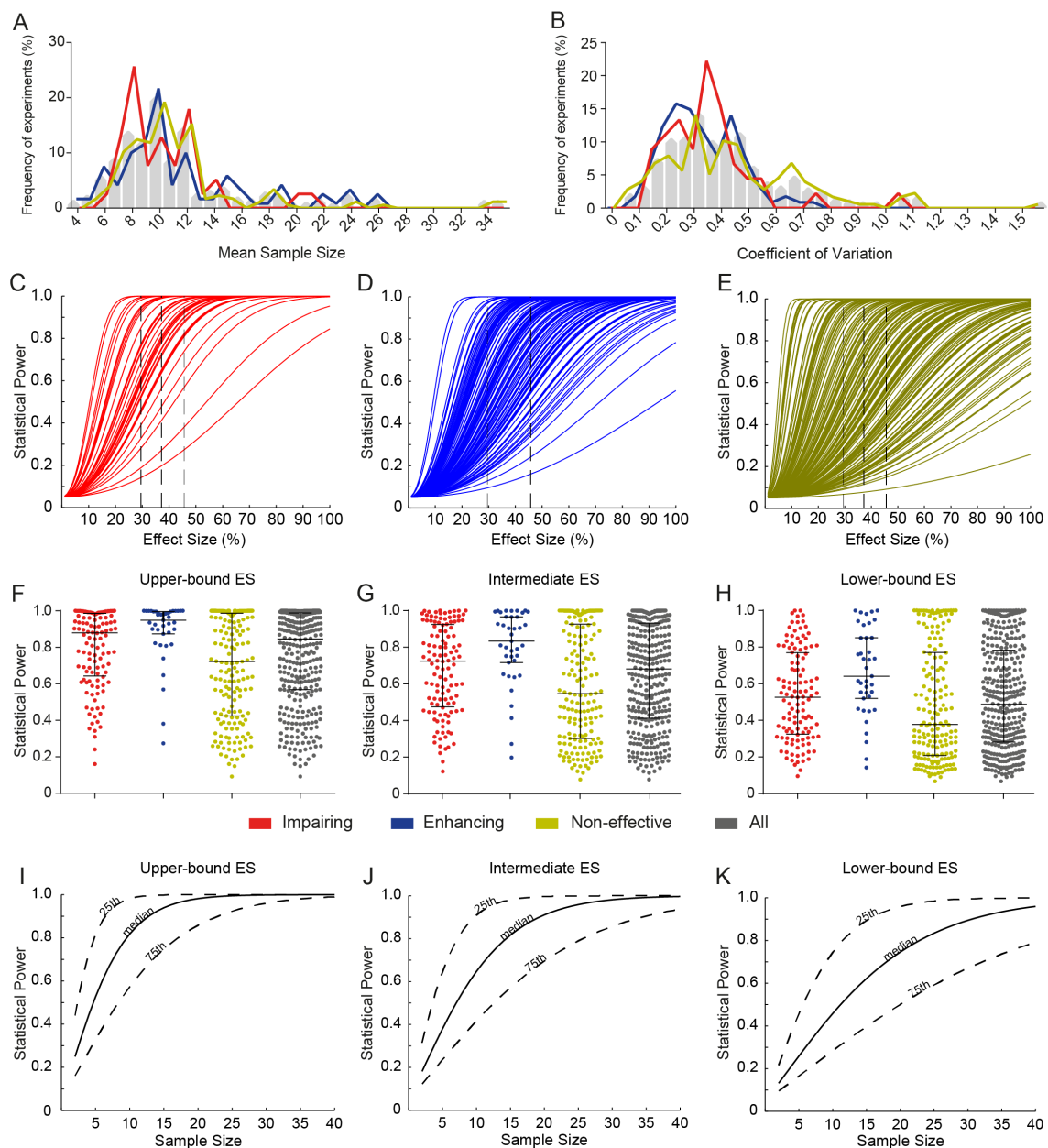


Figure 3. Distribution of sample size, variation and statistical power. (A) Distribution of mean sample size between groups for impairing ($n=120$), enhancing ($n=39$) and non-significant ($n=177$) experiments ($p=0.30$, one-way ANOVA). (B) Distribution of coefficients of variation

(pooled standard deviation/pooled mean) for each type of experiment ($p=0.001^*$, one-way ANOVA, $p=0.33$, ANCOVA with freezing levels as a covariate). (C) Distribution of statistical power for memory-impairing interventions: based on each experiment's variance and sample size, power varies according to the difference to be detected for $\alpha=0.05$. Dashed lines show the three effect sizes used for point estimates of power in F, G and H. (D) Distribution of statistical power for memory-enhancing interventions. (E) Distribution of statistical power for non-effective interventions. (F) Distribution of statistical power to detect the upper-bound effect size of 45.6% (right dashed line on C, D and E). Lines and whiskers express median and interquartile interval. Mean power is 0.79 ± 0.21 [0.75 to 0.83] ($n=120$) for memory-impairing, 0.89 ± 0.17 [0.84 to 0.94] ($n=39$) for memory-enhancing and 0.69 ± 0.28 [0.65 to 0.73] ($n=177$) for non-effective interventions. (G) Distribution of statistical power to detect the intermediate effect size of 37.2% (middle dashed line on C, D and E). Mean power is 0.68 ± 0.24 [0.63 to 0.72] ($n=120$) for memory-impairing, 0.8 ± 0.2 [0.73 to 0.87] ($n=39$) for memory-enhancing and 0.59 ± 0.3 [0.55 to 0.64] ($n=177$) for non-effective interventions. (H) Distribution of statistical power to detect the lower-bound effect size of 29.5% (left dashed line on C, D and E). Mean power is 0.54 ± 0.24 [0.49 to 0.58] ($n=120$) for memory-impairing, 0.67 ± 0.23 [0.59 to 0.74] ($n=39$) for memory-enhancing and 0.49 ± 0.31 [0.44 to 0.53] ($n=177$) for non-effective interventions. (I) Sample size vs. statistical power to detect the upper-bound effect size of 45.6%. Continuous lines use the 50th percentile of coefficients of variation for calculations, while dotted lines use the 25th and 75th percentiles. (J) Sample size vs. statistical power to detect the intermediate effect size of 37.2%. (K) Sample size vs. statistical power to detect the lower-bound effect size of 29.5%. Asterisks indicate significant results according to Holm-Sidak correction for 14 experiment-level comparisons.

Based on each experiment's variance and sample size, we built power curves to show how power varies according to the difference to be detected at $\alpha=0.05$ for each individual experiment within each class (**Fig. 3C-E**). To detect the mean effect size of 45.6% found for nominally effective interventions (i.e. those leading to statistically significant differences between groups), mean statistical power in our sample was 0.75 ± 0.26 [0.72 - 0.78] (**Fig. 3F**). This estimate, however, is an optimistic, upper-bound calculation of the mean effect size of biologically effective interventions (from here on referred to as "upper-bound ES"): as only large effects will be detected by underpowered studies, basing calculations on significant results alone leads to effect size inflation (13). A more realistic estimate of effect size was obtained based on experiments achieving statistical power above 0.95 ($n=60$) in the first analysis, leading to a mean effect size of 37.2%. Predictably, mean statistical power to detect this

difference (“intermediate ES”, **Fig. 3G**) fell to 0.65 ± 0.28 [0.62 - 0.68]. Using the mean effect size of all experiments (“lower-bound ES”, 29.5%) led to an even lower estimate of 0.52 ± 0.29 [0.49 - 0.56] (**Fig. 3H**), although this estimate of a typical effect size is likely pessimistic, as it probably includes many true negative effects. Interestingly, using Cohen’s traditionally accepted definitions of small ($d=0.2$), medium ($d=0.5$) and large ($d=0.8$) effect sizes (19), mean power was 0.07 ± 0.01 , 0.21 ± 0.07 and 0.44 ± 0.13 , respectively. These much lower estimates reflect the fact that effect sizes are typically much larger in rodent fear conditioning than in psychology experiments, for which this arbitrary classification was originally devised.

A practical application of these power curves is that we were able to calculate the necessary sample size to achieve desired power for each of these effect size estimates, considering the median coefficient of variation (as well as the 25th and 75th quartiles) of experiments in our sample (**Fig. 3I, J, K**). Thus, for an experiment with typical variation, around 15 animals per group are needed to achieve 80% power to detect our ‘intermediate effect size’ of 37.2%, which we consider to be our more realistic estimate for a typical effect size in the field. Nevertheless, only 12.2% of comparisons in our sample had a sample size of 15 or above in each experimental group, suggesting that this kind of calculation is seldom performed.

We also analyzed the distributions of both effect sizes and statistical power at the level of articles instead of individual experiments, using the mean effect size or power for each class of experiment (memory-impairing, memory-enhancing or non-effective) in each article as the observational unit. Results for these analyses are shown in **S3 Fig.** and **S4 Fig.**, and are generally similar to those obtained for the experiment-level analysis, except that the long tail of non-significant experiments with large coefficients of variation is not observed. This suggests that experiments with large

variation and low power are frequently found alongside others with adequate power within the same articles. It is unclear, however, whether this means that the low power of some experiments is a consequence of random fluctuations of experimental variance, or if these experiments use protocols that lead to larger coefficients of variation – for example, by generating lower mean levels of freezing (see **S1 Fig.**).

Correlation between effect sizes and statistical power/sample size

We next sought to correlate normalized effect size and statistical power. When analyzing effective and non-effective experiments separately, a correlation is expected mathematically from the definition of power; this is not the case, however, when analyzing the whole set of experiments together, where the presence of a negative correlation has been considered an indirect measure of publication bias (20). In our analysis, no correlation was found between effect size and sample size (**Fig. 4A**, $r=0.0007$, $p=0.99$); thus, relationships between effect size and power were driven mostly by a positive correlation between effect size and coefficients of variation (**Fig. 4B**, $r=0.37$, $p<0.0001$). Part of this correlation was mediated by the association of both variables with freezing levels (**S1 Fig.**), but the correlation remained significant after adjustment for this variable ($r=0.32$, $p<0.001$). Weak negative correlations between effect size and power were observed for the three effect size estimates used (**Figs. 4C-E**), although they were larger for the lower-bound estimate (**Fig. 4E**, $r=-0.21$, $p<0.0001$) than for the intermediate (**Fig. 4D**, $r=-0.16$, $p=0.003$) and upper-bound (**Fig. 4C**, $r=-0.12$, $p=0.03$) ones due to a ceiling effect on power. These, however, were no longer significant after adjustment for freezing levels ($r=-0.10$, $p=0.06$ for lower-bound, $r=-0.04$, $p=0.48$ for intermediate, $r=0.11$, $p=0.84$ for upper-bound effect size estimate).

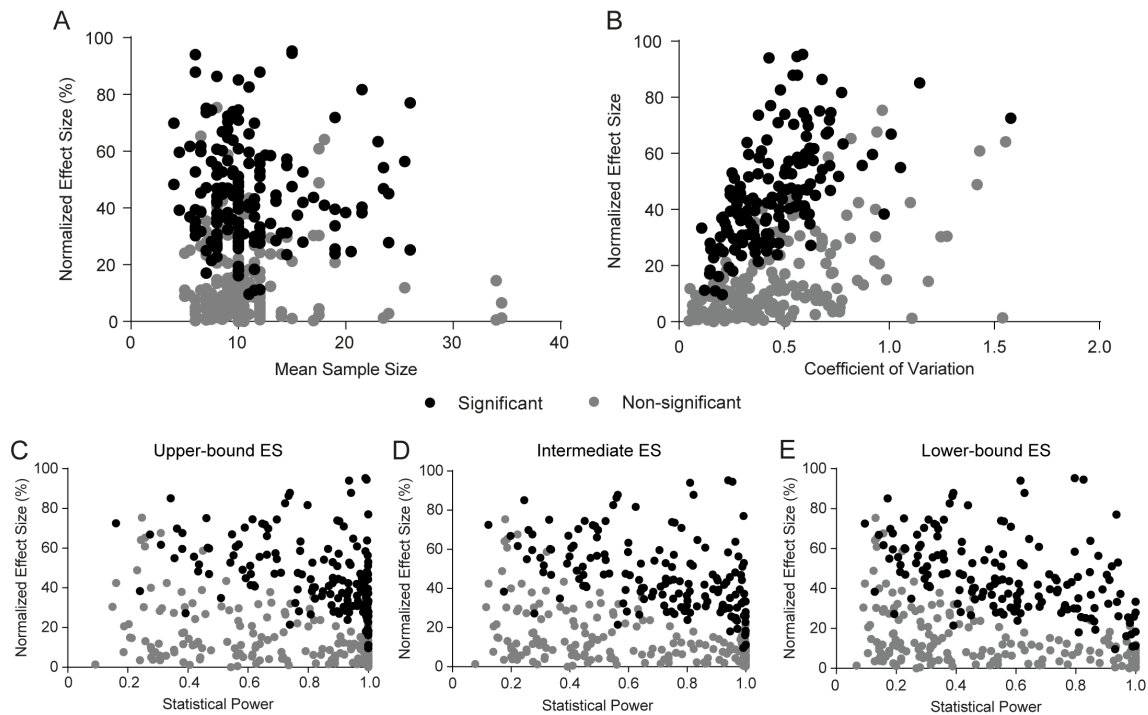


Figure 4. Correlations between effect size, variation and statistical power. (A) Correlation between normalized effect size and mean sample size. No correlation is found ($r=0.0007$, $p=0.99$, $r=-0.26$, $p=0.64$ after adjustment), although sample size variation was limited. (B) Correlation between normalized effect size and coefficient of variation. Correlation of the whole sample of experiments yields $r=0.37$, $p<0.0001^*$ ($n=336$; $r=0.32$, $p<0.001$ after adjustment for freezing levels). (C) Correlation between normalized effect size and statistical power based on upper-bound effect size of 45.6%. Correlation of the whole sample of experiments yields $r=-0.12$, $p=0.03$ ($r=0.11$, $p=0.84$ after adjustment for freezing levels), but distribution is skewed due to a ceiling effect on power. (D) Correlation between normalized effect size and statistical power based on intermediate effect size of 37.2%; $r=-0.16$, $p=0.003^*$ ($r=-0.16$, $p=0.48$ after adjustment). (E) Correlation between normalized effect size and statistical power based on lower-bound effect size of 29.5%; $r=-0.21$, $p<0.0001^*$ ($r=-0.1$, $p=0.06$ after adjustment). Asterisks indicate significant results according to Holm-Sidak correction for 18 experiment-level correlations.

Interestingly, the correlation between effect size and power was driven by a scarcity of experiments with large effect size and high power. This could be due to the fact that both are correlated with freezing levels – thus, relative differences tend to be smaller if freezing levels are high, as occurs in the majority of well-powered experiments. It also raises the possibility, however, that truly large effects are unusual in fear conditioning, and that some of the large effect sizes among low-powered experiments in our sample are inflated. On the other hand, a pattern classically

suggesting publication bias – i.e. a scarcity of low-powered experiments with small effects (20) – is not observed. It should be noted, however, that our analysis focused on individual experiments within an article, meaning that non-significant results were usually presented alongside other experiments with significant differences (either using fear conditioning or other methods); thus, this analysis does not allow us to assess publication bias at the level of articles.

Effects of methodological variables on the distribution of effect sizes and coefficients of variation.

We next examined whether effect sizes and coefficients of variation were influenced by type of conditioning, species or sex of the animals (**Fig. 5**). Mean normalized effect size was slightly larger in contextual than in cued fear conditioning (33.2% vs. 24.4%, Student's t test $p < 0.0001$) and markedly larger in males than in females (30.3% vs. 18.9% vs. 34.2% for experiments using both, one-way ANOVA, $p = 0.004$), but roughly equivalent between mice and rats (29.8% vs. 29.1%, $p = 0.76$). Coefficients of variation were higher in contextual conditioning (0.51 vs 0.41, Student's t test $p = 0.001$), in experiments using animals of both sexes (0.62 vs. 0.44 in males and 0.41 in females, one-way ANOVA, $p < 0.0001$), and in those using mice (0.50 vs. 0.42, Student's t test, $p = 0.008$), although the latter difference was not statistically significant after correction for multiple comparisons. All of these associations should be considered correlational and not causal, as specific types of conditioning or animals of a particular species or sex might be more frequently used for testing interventions with particularly high or low effect sizes. Also of note is the fact that experiments using males were 7.7 times more common than those using females in our sample (277 vs. 36).

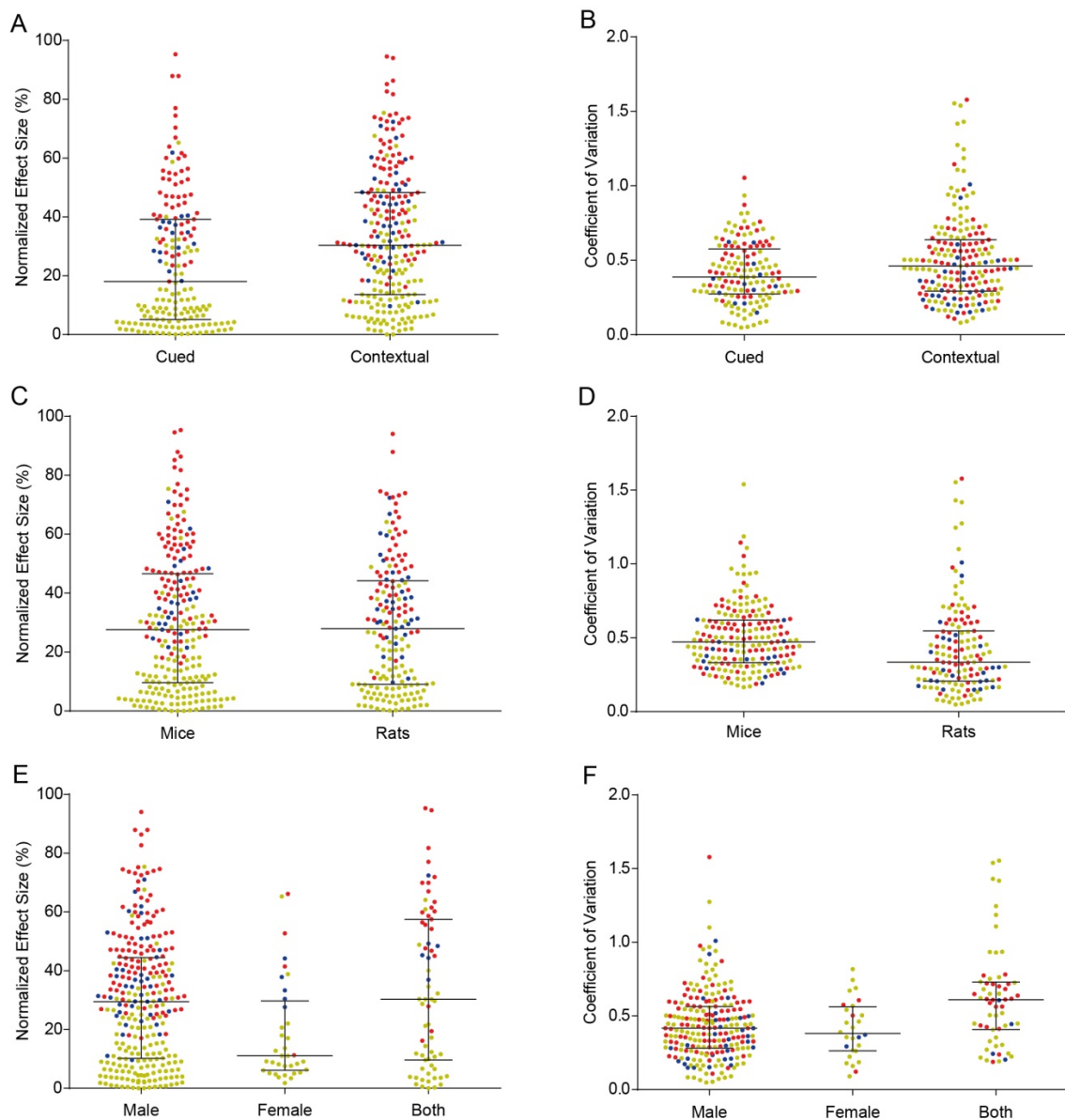


Figure 5. Effect sizes and coefficients of variation across different protocols, species and sexes. Colors indicate memory-enhancing (red), memory-impairing (blue) or non-effective (yellow) experiments. Lines and whiskers express median and interquartile interval. (A) Distribution of effect sizes across cued (n=171) and contextual (n=239) conditioning protocols. Student's t test, $p < 0.0001^*$. (B) Coefficients of variation across cued (n=145) and contextual (n=191) conditioning protocols. Student's t test, $p = 0.001^*$. (C) Distribution of effect sizes across experiments using mice (n=237) or rats (n=173). Student's t test, $p = 0.76$. (D) Coefficients of variation across experiments using mice (n=193) or rats (n=143). Student's t test, $p = 0.008$. (E) Distribution of effect sizes across experiments using male (n=277), female (n=36) or both (n=67) sexes. One-way ANOVA, $p = 0.004^*$; Tukey's post-hoc test, male vs. female $p = 0.01$, male vs. both $p = 0.40$, female vs. both $p = 0.003$. 30 experiments were excluded from this analysis for not stating the sex of animals. (F) Coefficients of variation across experiments using male (n=233), female (n=28) or both (n=60) sexes. One-way ANOVA, $p < 0.0001^*$; Tukey's test, male vs. female $p = 0.85$, male vs. both $p < 0.0001$, female vs. both $p = 0.0006$. For coefficient of variation analyses, 74 experiments were excluded due to lack of information on sample size for individual groups. Asterisks indicate significant results according to Holm-Sidak correction for 14 experiment-level comparisons.

We also examined whether effect sizes and coefficients of variation differed systematically according to the type, timing or anatomical site of intervention (**S5 Fig.**). Effect sizes did not differ significantly between surgical, pharmacological, genetic and behavioral interventions (38.7% vs. 28.1% vs. 30.5% vs. 25.8% one-way ANOVA, $p=0.12$), although there was a trend for greater effects with surgical interventions (which were uncommon in our sample). No differences were found between the effect sizes of systemic and intracerebral interventions (28.7% vs. 30.3%, Student's t test, $p=0.45$) or between those of pre- and post-training interventions (30.5% vs. 25.4%, Student's t test, $p=0.07$), although pre-training interventions had slightly higher coefficients of variation (0.49 vs 0.37, Student's t test $p=0.0015$). Coefficients of variation did not differ significantly between surgical, pharmacological, genetic and behavioral interventions (0.41 vs. 0.43 vs. 0.50 vs. 0.50, one-way ANOVA $p=0.08$) or between systemic and intracerebral interventions (0.49 vs. 0.45, Student's t test $p=0.15$). Once again, these differences can only be considered correlational and not causal.

Correlation between effect sizes/statistical power and description of results

Given the wide distribution of both effect sizes and statistical power in the fear conditioning literature, we tried to determine whether these were taken into account in the description of results. For each comparison we extracted the terms (words or phrases) describing the results of that experiment in the text or figure legends, and asked 14 researchers with experience in behavioral neuroscience and fluent or native levels of English to classify them. For comparisons with significant differences, we asked them to classify terms as implying strong (i.e. large effect size) or weak (i.e. small effect size) effects, or as neutral terms (i.e. those from which effect size could not be deduced). For comparisons with non-significant differences, terms were to be classified as implying similarity between groups, as suggesting a trend towards difference, or as neutral terms

(i.e. those from which the presence or absence of a trend could not be deduced). From the average of these classifications, we defined a score for each term (**S1 and S2 Tables**) and sought to correlate these scores with the actual effect size and statistical power of experiments.

Agreement between researchers over the classification of terms was quite low, especially in the case of significant interventions: single measures intraclass correlation coefficients (reflecting the reliability of individual researchers when compared to the whole sample) were .234 for significant interventions and .597 for non-significant ones, while average measures coefficients (reflecting the aggregated reliability of the whole sample) were .839 and .962, respectively. This, along with a trend for the use of terms with little effect size information (“increase”, “decrease”, “significantly more”, “significantly less”, etc.), led most terms to receive intermediate scores approaching 1 (i.e. neutral). For significant differences, no correlations were observed between this score and either effect size ($r=-0.05$, $p=0.48$) or statistical power ($r=0.03$, $p=0.73$) (**Fig. 6A and 6B**). For non-effective interventions, a significant correlation between description score and effect size was observed (**Fig 6C**, $r=0.28$, $p=0.0002$), mostly because larger effect sizes were associated with terms indicating a trend for difference. Still, no correlation was observed between textual descriptions of results and power (**Fig 6D**, $r=0.03$, $p=0.74$). Moreover, statistical power was rarely mentioned in the textual description of results – the term “power” was used in this context in only 4 articles in our sample – suggesting that it is largely ignored when discussing findings, as previously shown in other areas of research (21).

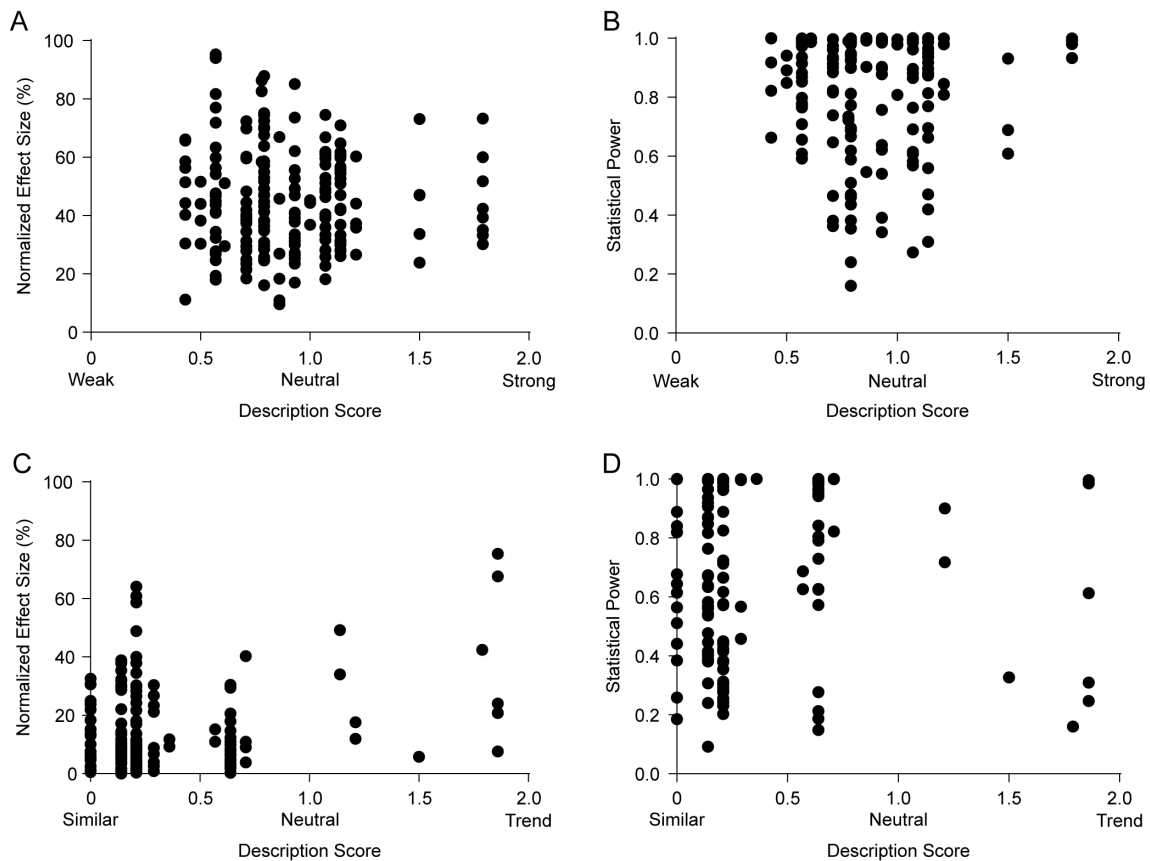


Figure 6. Correlation between description of results and effect size/statistical power. Description scores refer to the mean score given by 14 neuroscience researchers who rated terms as “weak” (0), “neutral” (1) or “strong” (2) in the case of those describing significant differences, or as “similar” (0), “neutral” (1) or “trend” (2) in the case of those describing non-significant ones. (A) Correlation between normalized effect size and description score for significant results. $r=-0.05$, $p=0.48$ ($n=195$). (B) Correlation between statistical power and description score for significant results. $r=0.03$, $p=0.73$ ($n=155$). (C) Correlation between normalized effect size and description score for non-significant results. $r=0.28$, $p=0.0002^*$ ($n=174$). (D) Correlation between upper-bound estimate of statistical power and description score for non-significant results. $r=0.03$, $p=0.74$ ($n=146$). Asterisk indicates significant result according to Holm-Sidak correction for 18 experiment-level correlations.

Risk of bias indicators and their relationship with effect size and power

As previous studies have shown that measures to reduce risk of bias are not widely reported in animal research (11,12), we sought to evaluate the prevalence of these measures in our sample of fear conditioning articles. **Table 1** shows the percentage of articles reporting 7 items thought to reduce risk of bias in animal studies, adapted and expanded from the CAMARADES checklist (22). As shown in previous studies, although some items were reported for fear conditioning experiments in most

articles (statement of compliance with animal regulations, adequate description of sample size, blinding), others were virtually inexistent, such as the presence of a sample size calculation (1 article) and compliance with the ARRIVE guidelines (23) (0 articles). Contrary to previous reports in other areas (24–27), however, no significant association was found between reporting of these indicators and either the percentage of significant experiments, the mean effect size of effective interventions or the mean statistical power of experiments in our sample (**S6 Fig.**). The region of origin of the article also had no correlation with either of these variables (**S7 Fig.**). Nevertheless, it should be noted that mean effect size, percentage of significant experiments and mean power in this case refer only to fear conditioning experiments, which were not necessarily the only results or the main findings presented in these articles. Thus, it is possible that other results in the article might have shown higher correlation with risk of bias indicators.

Quality assessment item	Randomization of allocation	Blinded or automated assessment	Sample size calculation	Exact sample size description	Statement of compliance with regulatory requirements	Statement on conflict of interest	Statement of compliance with ARRIVE
Number of articles (%)	18/77 (23.4%)	92/122 (75.4%)	1/122 (0.8%)	98/122 (80.3%)	118/122 (96.7%)	66/122 (54.1%)	0/122 (0%)

Table 1. Number of articles including quality assessment items. Percentages were calculated using all 122 articles, except in the case of randomization, which was calculated based on 77 articles, as it is not applicable to genetic interventions. In the case of blinding, 72 articles used automated analysis and 20 used blinded observers, totaling 92 articles scored for this item.

Correlations of effect size, power and study quality with article citations

Finally, we investigated whether the percentage of significant experiments, mean effect size for effective interventions, mean statistical power or a composite study quality score (aggregating the 7 risk of bias indicators described in Table 1) correlated

with article impact, as measured by the number of citations (obtained in August 26th, 2016) (**Fig. 7**) and the impact factor of the publication venue (**S8 Fig.**). None of the correlations was significant after adjustment for multiple comparisons, although a weak positive correlation was observed between study quality score and impact factor ($r=0.22$, $p=0.01$), driven by associations of higher impact factors with blinding (Student's *t* test with Welch's correction, $p=0.0001$), conflict of interest reporting (Student's *t* test with Welch's correction, $p=0.03$) and exact sample size description (Student's *t* test, $p=0.03$). It should be noted that the distribution of impact factors and citations is heavily skewed, limiting the use of linear correlations as planned in the original protocol – nevertheless, exploratory non-parametric analysis of the data confirmed the lack of significance of correlations. Once again, our data refers to fear conditioning experiments only – therefore, other data in the articles could feasibly account for the variation in impact factors and citations.

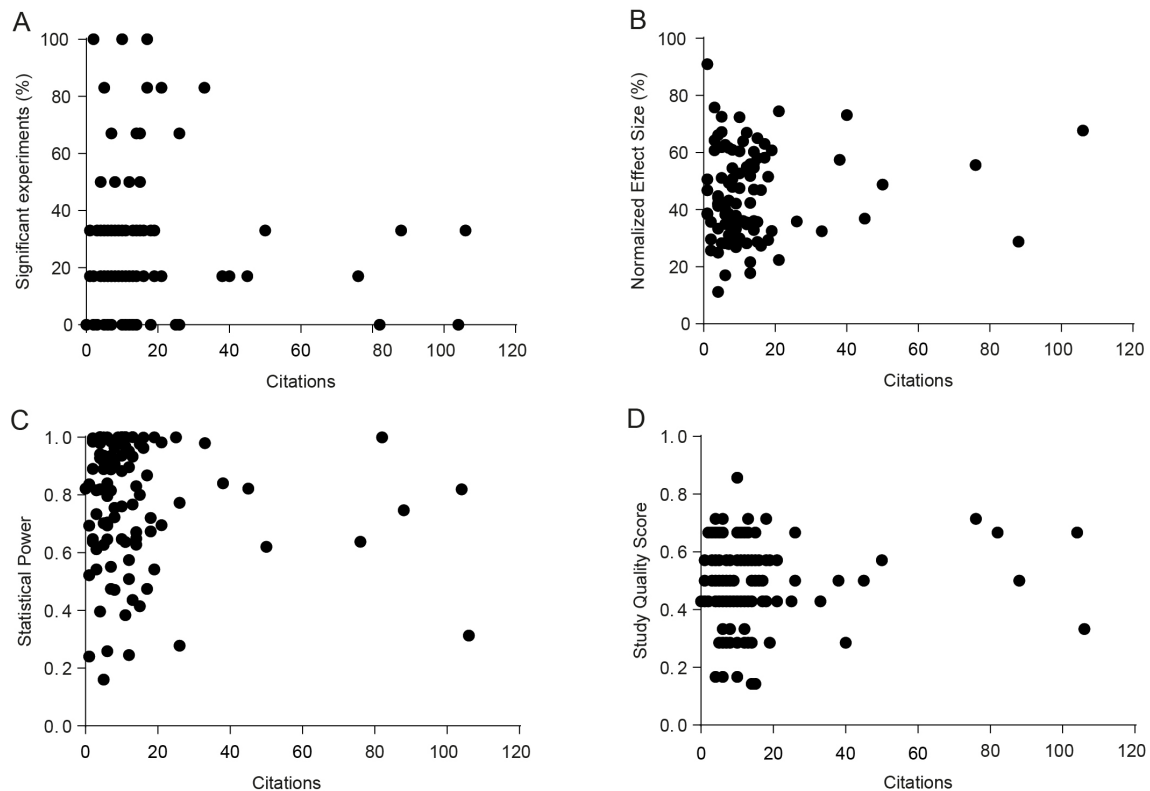


Figure 7. Correlation between citations and percentage of significant experiments, effect size and statistical power. Citations were obtained for all articles on August 26th, 2016. (A) Correlation between % of significant results per article and citations. $r=-0.03$, $p=0.75$ ($n=121$). (B) Correlation between mean normalized effect size of effective interventions and citations. $r=0.097$, $p=0.34$ ($n=98$). (C) Correlation between mean statistical power (upper-bound estimate) and citations. $r=-0.08$, $p=0.40$ ($n=104$). (D) Correlation between study quality score and citations. $r=0.09$, $p=0.31$ ($n=121$). According to Holm-Sidak correction for 8 article-level correlations, none is significant.

Discussion

In light of the low reproducibility of published studies in various fields of biomedical research (28–30), which is thought by many to be a consequence of low statistical power and excessive reliance on significance tests (7,16), calls have been made to report effect sizes and confidence intervals alongside or in place of p values (5,6,8) and to increase power by performing adequate sample size calculations (13,21,31). However, it is unclear whether these proposals have had any impact on most fields of basic science. We have taken one particular memory task in rodents – in this case, fear conditioning, in which outcomes and effect sizes are described in a standardized way and are thus comparable across studies – in order to analyze how these two concepts are dealt with in behavioral neuroscience.

Our first main finding, which should come as no surprise for those following the memory literature, is that most amnesic interventions in fear conditioning cause partial effects, with residual freezing usually remaining way above pre-conditioning levels. Moreover, most of the large effect sizes in our sample were found in underpowered studies, suggesting that they could represent inflated estimates. This in itself is not unexpected: fear memories depend on a well distributed network, both anatomically and molecularly (17), and it seems natural that most interventions directed at a specific site or pharmacological target will not fully block learning. Moreover, fear memory

formation is a complex process that is influenced by many modulatory influences, both within and outside the brain. This becomes a problem, however, when effect sizes are not considered when analyzing results. As even small decreases in freezing can be found to be significant, it is challenging to tear apart essential mechanisms of memory formation from modulatory influences by looking at statistical significance alone. This can lead to a situation in which accumulating evidence, even if correct, can confuse rather than advance understanding, as has been suggested to occur in fields such as long-term potentiation (14) and apoptosis (32).

Matters are complicated further by the possibility that many of these findings are false positives and/or false negatives. The prevalence of both in relation to true positives and negatives depends crucially on the statistical power of experiments, which in turn depends on sample size. Calculating the actual power of published experiments is tricky, as the difference to be used should not be based on the observed results – which leads power analysis to become circular (33). Thus, statistical power depends on expected effect sizes, which are arbitrary by nature – although they can sometimes be estimated from meta-analyses (13). However, by considering the mean effect size for well-powered experiments in our sample, we arrived at an estimate of around 37.2% that might be considered “typical” for a published fear conditioning experiment. Using this number and the median coefficient of variation, our conclusion was that the ideal sample size to reach a statistical power of 80% would be around 15 animals per group. This sample size, however, was reached in only 12.2% of experiments, most of which had samples of around 8 to 12 animals per group, considered to be standard in the field.

As sample size calculations are exceedingly rare, and insufficient power seems to be the norm in other fields of neuroscience as well (13), it is quite possible that classically used sample sizes in behavioral neuroscience (and perhaps in other fields of

basic science) might be insufficient. Considering median variances and our intermediate effect size estimate, a sample size of 8 will lead power to be just above 50% – not much better than a coin toss – and would provide adequate (i.e. 80%) power only for differences of 66.8% or greater, which are very uncommon in fear conditioning (i.e. 7% of experiments in our sample). The high prevalence of these low sample sizes seems therefore to confirm recent predictions from mathematical models suggesting that current incentives in science tend to favor the publication of underpowered studies (16,34), although they could also be due to restrictions on animal use imposed by ethical regulations. That said, average power in our sample for typical effect sizes was higher than those that have been described in other areas of neuroscience by Button et al. (13); however, this could reflect the fact that effect sizes in their study were calculated by meta-analysis, and might be smaller than those derived by our method of estimation.

Another old-established truism in the behavioral neuroscience field – as well as in other fields of basic science – is that experiments in females tend to yield more variable results than in males due to estrous cycle variations (35). This belief is a likely cause for the vast preponderance of experiments on male animals, which were nearly 8 times more common than those in females in our sample – a sex bias even greater than that described in previous reports in other areas of science (36). However, at least in our analysis, coefficients of variation were roughly similar between experiments in males and females (and predictably higher in experiments using both), as has been found in other areas (37,38), suggesting this truism to be a myth. On the other hand, effect sizes were markedly lower in females in our sample. We do not have an explanation for this finding, and it should not be taken to imply a causal relationship between sex and effect sizes. It is possible, for example, that female animals are used to test lower-yield

hypotheses, or that they are used only when results have already been shown to be significant in males, being thus less susceptible to effect inflation due to publication bias. In any case, it would be interesting to test, either through meta-regression or by direct empirical comparison, whether the effects of specific interventions on fear conditioning are indeed more marked in males when other variables are adequately controlled.

Our work also calls attention to the problem of sex bias in the memory literature. Previous work in clinical (39) and preclinical (35,36) data has pointed out the problem of concentrating experiments in male populations, as these might not be automatically translatable to females. However, despite calls for increasing the number of studies on females (40), this problem remains strikingly present in the fear conditioning field, and although we cannot assess its consequences, it seems sensible to call for more memory studies to be performed on female animals.

Concerning risk of bias indicators, the prevalence found in our sample was roughly similar to previous reports on animal studies for randomization and conflict of interest reporting (11), but were distinctly higher for blinded assessment of outcome, largely because 56.6% of articles used automated software to measure freezing, which we considered to be equivalent to blinding in terms of eliminating observer bias. As described previously in many fields (11,12,21), sample size calculations were almost non-existent, a fact that helps to explain why many experiments are underpowered. Interestingly, although we analyzed a sample of papers published 3 years after the ARRIVE guidelines (23), these guidelines were not mentioned in any of the articles, suggesting that their impact, at least in the field of behavioral neuroscience, has been rather limited.

Contrary to previous studies, however (24–27), we did not detect an impact of these risk of bias indicators on article-level measures such as percentage of fear conditioning experiments with significant results, mean effect size of significant experiments and mean statistical power. This could mean that, compared to preclinical studies, bias towards positive results is lower in fear conditioning studies, most of which pertain to discovery research. However, it seems more likely that, as we selected particular experiments within papers containing other results, we were not as likely to detect effects of bias on article-level measures. As basic science articles typically contain numerous results using different methods, it is perhaps less likely that all comparisons will be subject to bias towards positive findings. Moreover, the experiments in our sample probably included negative controls for other findings, which might have been expected to yield non-significant results. Thus, although our results do not indicate an impact of bias on article-level results, they should not be taken as evidence that this does not occur.

One of the most interesting findings of our article was the lack of correlation of textual description of results with the actual effect sizes of significant experiments, as well as with statistical power. Although this suggests that these measures are not usually considered in the interpretation of results, there are important caveats to this data. First of all, agreement between what words describe a “strong” or “weak” effect between researchers in the area was strikingly low, suggesting that written language is a poor descriptor for quantitative data in these articles. Moreover, the fact that most terms used to describe differences were neutral to effect sizes (e.g. “significantly higher”, “significantly lower”, etc.) limited our ability to detect a correlation. That said, the high prevalence of neutral terms by itself is evidence that effect sizes are not usually taken

into account when reporting results, as differences tend to be described in the text by their statistical significance only.

This point is especially important to consider in the light of recent calls for basic science to use data synthesis tools such as meta-analysis (10) and formal or informal Bayesian inference (2,7,9,41). In both of these cases, the incremental effect of each new experiment on researchers' beliefs on the veracity of a finding is dependent both on the effect size of the result and on the probability that it might be found by chance under the null hypothesis (i.e. its p value). However, even exact p values were uncommon in our sample, with the majority of articles describing p as being above or below a threshold value for each experiment. This seems to suggest that researchers in the field indeed tend to consider statistical significance as a binary outcome, and might not be quite ready or willing to move towards Bayesian logic, which would require a major paradigm shift in the way results are reported and discussed.

Concerning article impact metrics, our results are in line with previous work showing that journal impact factor does not correlate either with statistical power (13) or with most risk of bias indicators (11). Furthermore, we showed that, in the case of fear conditioning, this lack of correlation also occurs for the percentage of significant experiments and the mean effect size for significant differences, and that it extends to citations measured over 2 subsequent years. That said, our article-level analysis was limited by the fact that, for many articles, the included experiments (e.g. those using fear conditioning and fulfilling inclusion criteria) represented a minority of the findings. Moreover, most articles tend to cluster around intermediate impact factors (i.e. between 3 and 6) and relatively low (< 20) citation numbers. Thus, our methodology might not have been adequate to detect correlations between these metrics with article-wide effect

size and power estimates, which could potentially have been observed if all experiments were considered.

The choice to focus on a particular type of experiment – in this case, rodent fear conditioning – is both one of the main strengths and the major limitation of our findings. On one hand, it allows us to look at effect sizes that are truly on the same scale, as fear conditioning protocols tend to be reasonably similar across laboratories, and all included experiments described their results using the same metric. Thus, the effect sizes we studied are not abstract and have real-life meaning – allowing us to state, for example, that most amnesic effects in the field are partial. On the other hand, however, this decision limits our conclusions to rodent fear conditioning, and also weakens our article-level conclusions, as most articles had only a fraction of their experiments analyzed.

Dealing with multiple experiments using different outcomes presents a major challenge for meta-research in basic science, and all alternatives present limitations. A radically opposite approach of converting all effect sizes in a field to a single metric (e.g. Pearson's r , Cohen's d , etc.) has been used by other researchers investigating similar topics (15,20,21,31). Although normalizing effect sizes allows one to obtain a representative sample of effect sizes from a wider field, it also leads them to be abstract in nature, and not as readily understandable for experimental researchers. Moreover, there is no reason to think that typical effect sizes will be similar in different areas – as an example, a whopping 48.2% of effect sizes in our sample (and 88.7% of statistically significant ones) fell under the “large” category by Cohen's original proposal for psychology studies (19), which has been adopted by other authors as well (15,21). This constitutes a strong argument against using fixed values as arbitrary descriptors for effect sizes across areas, as what constitutes a “small” or “large” effect depends on the

kind of experiment; moreover, ignoring this fact can lead to major distortions in power calculations and interpretation of results.

In our case, on the contrary, studying the concrete scenario of a specific methodology leads to more readily applicable suggestions for experimental researchers, such as the rule-of-thumb recommendation that the average number of animals per group in a fear conditioning experiments to achieve 80% power would be around 15 for typical effect sizes and variances. Our approach also allowed us to detect correlations between results and specific methodological factors (e.g. context vs. cued conditioning, female vs. male animals) that might not have been apparent if multiple types of experiments were pooled together. Still, to provide more solid conclusions on the influence of these factors on experimental results, even our methodology has too wide a focus, as analyzing multiple interventions limits our possibilities to perform meta-analysis and meta-regression, and thus control for confounding variables. Thus, follow-up studies with more specific aims (i.e. meta-analyses of specific interventions in fear conditioning) are warranted to understand the variation between results in the field.

Finally, it is important to note that, while our study has led to some illuminating conclusions, they are inherently limited to the methodology under study. Thus, extrapolating our findings to other types of behavioral studies, not to mention other fields of science, requires data to be collected for each specific subfield. While this area-by-area approach might appear herculean at first glance, it is easily achievable if scientists working within specific domains start to design and perform their own systematic reviews as we have done. Only through this dissemination of meta-research across different areas of science will we be able to develop solutions that, by respecting the particularities of individual subfields, will be accepted enough to have an impact on research reproducibility.

Materials and Methods

The full protocol of data selection, extraction and analysis was initially planned on the basis of a pilot analysis of 30 papers, and was registered, reviewed and published ahead of full data extraction (18). In brief, we searched PubMed for the term “fear conditioning” AND (“learning” OR “consolidation” OR “acquisition”) AND (“mouse” OR “mice” OR “rat” OR “rats”)” to obtain all articles published online in 2013. Titles and abstracts were first scanned for articles presenting original results involving fear conditioning in rodents and that were written in English. Selected articles underwent full-text screening for selection of experiments that (a) described the effects of a single intervention on fear conditioning acquisition or consolidation, (b) had a proper control group to which the experimental group is compared to, (c) used freezing behavior as a measure of conditioned fear in a test session and (d) had available data on mean freezing, SD or SEM, as well as on the significance of the comparison. Articles were screened by one of two investigators (C.F.D.C. or T.C.M.) for relevant data and were analyzed by the other – thus, all included experiments were dual-reviewed.

Only experiments analyzing the effect of interventions performed before or up to 6 hours after the training session (i.e. those affecting fear conditioning acquisition or its immediate consolidation) were included. Data on mean freezing and SD or SEM were obtained for each group from the text when available; otherwise, it was extracted using Gsys 2.4.6 software (Hokkaido University Nuclear Reaction Data Centre). When exact sample size for each group was available, the experiment was used for the analysis of effect size and statistical power – otherwise, only effect size was obtained, and the

experiment was excluded from power analysis. For individual experiments, study design characteristics were also obtained, including species and sex of the animals, type of conditioning protocol, type, timing and site of intervention.

From each comparison, we also obtained the description term used by the authors in the results section of the paper. Classification of the terms used to describe effects (**S1 and S2 Tables**) was based on a blinded assessment of words or phrases by a pool of 14 researchers who were fluent or native speakers of English and had current or past experience in the field of behavioral neuroscience. Categories were given a score from 0 to 2 in order of magnitude (i.e. 0 = weak, 1 = neutral, 2 = strong for significant results; 0 = similar, 1 = neutral, 2 = trend for non-significant results), and the average results for all researchers was used as a continuous variable for analysis.

Apart from experiment-level variables, we also extracted article-level data such as impact factor of the journal in which it was published (based on the 2013 Journal Citations Report), number of citations (obtained for all articles on August 26th 2016), country of origin (defined by the corresponding author's affiliation) and the 7 risk of bias indicators described on **Table 1**. For article-level correlations, we compiled these measures into a normalized score.

After completion of data extraction, all calculations and analyses were performed according to the previously specified protocol. Specific details of calculations (as well as the raw data used) can be found in **S1 Data**. After this, the following additional analyses were performed in an exploratory fashion:

(a) To confirm that residual freezing levels after memory-impairing interventions were indeed above training values, demonstrating that most amnesic

intervention have partial effects, we extracted pre-conditioning freezing levels from training sessions when these were available. These levels were obtained for pre-shock periods only, and separated as baselines for contextual (i.e. values in the absence of tone) or tone conditioning (i.e. values in the presence of a tone, but before shock). These were compared to the corresponding contextual or tone test sessions for treated groups in memory-impairing interventions by an unpaired t test based on the extracted means, SD or SEM and sample size.

(b) In the original protocol, only the mean of all effective interventions (i.e. upper-bound effect size) was planned as a point estimate to be used for power calculations, although we acknowledged this to be optimistic (18). We later decided to perform power calculations based on the mean effect size of the experiments achieving power above 0.95 on the first analysis (i.e. intermediate effect size) to avoid effect size inflation, as we reached the conclusion that this would provide a more realistic estimate. Additionally, we also calculated power based on the mean effect size of the whole sample of experiments as a lower-bound estimate, and presented all three estimates in the results section and figures.

(c) As coefficients of variation showed a strong negative correlation with freezing levels, we decided to perform an ANCOVA using the highest freezing between both groups (i.e. the one used as the reference for effect size normalization) as a covariate for the results on **Fig. 3B**, as well as an ANOVA to compare freezing levels across different classes of experiments. Moreover, as both power and effect size correlated significantly with freezing levels, we decided to perform partial correlations between both variables using freezing as a covariate, in order to analyze how much of the correlations presented in **Fig. 4** were mediated by freezing levels. We also checked

whether adding freezing levels as a covariate influenced the statistical analyses in **Fig. 3A, Fig. 5, Fig. 6** and **S5 Fig.**, but as this did not have a significant impact on the results in these figures, we only reported the originally planned analyses.

(d) All of our planned analyses were parametric; after extraction, however, it was clear that some of the data deviated from a normal distribution (especially in the case of power estimates, citation counts and impact factor). Because of this, we also performed non-parametric analysis for the correlations of citations and impact factor with the percentage of significant results, mean normalized effect sizes, statistical power and study quality score.

(e) In the protocol, we had planned to test correlations between normalized effect sizes and statistical power, mean sample size and absolute freezing levels (using the group with the highest freezing). After analyzing the results, we also decided to correlate normalized effect sizes with coefficients of variation (as this, rather than sample size, seemed to explain the lower power of non-significant results), additional power estimates (as using our original estimate led to a ceiling effect) and different estimates of freezing based on the control group or on the mean freezing of both groups (to compare these forms of normalization with the one we chose).

(f) Due to the correlation of study quality assessment with journal impact factor, we performed an exploratory analysis of the correlation of this metric with each of the individual quality assessment items by performing a Student's t test (corrected for unequal variances by Welch's correction) between the impact factors of studies with and without each item.

(f) Because of the additional analyses above, we adjusted the number of comparisons/correlations used as the basis of the Holm-Sidak correction for multiple comparisons. The total numbers used for each correction were 14 for experiment-level comparisons, 17 for article-level comparisons, 14 for experiment-level comparisons and 8 for article-level correlations, leading to significance thresholds between 0.003 and 0.05.

Author Contributions

O.B.A and T.C.M initially designed the study. C.F.D.C and T.C.M. collected data. O.B.A. solved controversies between investigators. All authors designed the protocol after collection of pilot data. C.F.D.C. and T.C.M. analyzed data. C.F.D.C. prepared figures and tables. O.B.A. wrote the initial version of the manuscript. All authors critically revised the text.

Acknowledgements

This work was supported by FAPERJ grants E-26/111.277/2014 and E-26/201.544/2014 to O.B.A., by RCUK NC3Rs grant NC/L000970/1 to M.R.M., and by CNPq scholarships to T.C.M. and C.F.D.C.

References

1. Nuzzo R. Statistical errors. *Nature*. 2014;506: 150–152.
2. Colquhoun D. An investigation of the false discovery rate and the misinterpretation of p-values. *R Soc Open Sci*. 2014;1: 140216.
3. Altman N, Krzywinski M. Points of significance: P values and the search for significance. *Nat Methods*. 2016;14: 3–4.

4. Wasserstein RL, Lazar NA. The ASA's statement on p-values: context, process, and purpose. *Am Stat.* 2016. <http://dx.doi.org/10.1080/00031305.2016.1154108>
5. Greenland S, Senn SJ, Rothman KJ, Carlin JB, Poole C, Goodman SN, et al. Statistical tests, p-values, confidence intervals, and power: a guide to misinterpretations. *Am Stat.* 2016;15: 1–31.
6. Nakagawa S, Cuthill IC. Effect size, confidence interval and statistical significance: a practical guide for biologists. *Biol Rev.* 2007;82: 591–605.
7. Ioannidis JPA. Why most published research findings are false. *PLoS Med.* 2005;2: e124.
8. Trafimow D, Marks M. Editorial. *Basic Appl Soc Psych.* 2015;37: 1–2.
9. Goodman SN, Fanelli D, Ioannidis JPA. What does research reproducibility mean? *Sci Transl Med.* 2016;8: 1–6.
10. Vesterinen HM, Sena ES, Egan KJ, Hirst TC, Churolov L, Currie GL, et al. Meta-analysis of data from animal studies: A practical guide. *J Neurosci Methods.* 2014;221: 92–102.
11. Macleod MR, Lawson McLean A, Kyriakopoulou A, Serghiou S, de Wilde A, Sherratt N, et al. Risk of bias in reports of in vivo research: a focus for improvement. *PLOS Biol.* 2015;13: e1002273.
12. Kilkenny C, Parsons N, Kadyszewski E, Festing MFW, Cuthill IC, Fry D, et al. Survey of the quality of experimental design, statistical analysis and reporting of research using animals. *PLoS One.* 2009;4: e7824.
13. Button KS, Ioannidis JP a, Mokrysz C, Nosek B a, Flint J, Robinson ESJ, et al. Power failure: why small sample size undermines the reliability of neuroscience. *Nat Rev Neurosci.* 2013;14: 365–376.

14. Sanes JR, Lichtman JW. Can molecules explain long-term potentiation? *Nat Neurosci.* 1999;2: 597–604.
15. Szucs D, Ioannidis JPA. Empirical assessment of published effect sizes and power in the recent cognitive neuroscience and psychology literature. *PLoS Biol.* 2017;15: e2000797.
16. Higginson AD, Munafò MR. Current incentives for scientists lead to underpowered studies with erroneous conclusions. *PLOS Biol.* 2016;14(11): e2000995.
17. Maren S. Neurobiology of Pavlovian fear conditioning. *Annu Rev Neurosci.* 2001;24: 897–931.
18. Moulin TC, Carneiro CFD, Macleod MR, Amaral OB. Protocol for a systematic review of effect sizes and statistical power in the rodent fear conditioning literature. *Evid Based Preclin Med.* 2016;3: 24– 32.
19. Cohen J. *Statistical power analysis for the behavioral sciences* (2nd ed.). New York: Academic Press; 1977.
20. Kühberger A, Fritz A, Scherndl T. Publication bias in psychology: a diagnosis based on the correlation between effect size and sample size. *PLoS One.* 2014;9: e105825.
21. Sedlmeier P, Gigerenzer G. Do Studies of statistical power have an effect on the power of studies? *Psychol Bull.* 1989;105: 309–316.
22. Sena E, van der Worp HB, Howells D, Macleod M. How can we improve the pre-clinical development of drugs for stroke? *Trends Neurosci.* 2007; 30:433–439.
23. Kilkenny C, Browne WJ, Cuthill IC, Emerson M, Altman DG. Improving

- bioscience research reporting: the ARRIVE guidelines for reporting animal research. *PLoS Biol.* 2010;8: e1000412.
24. Macleod MR, van der Worp HB, Sena ES, Howells DW, Dirnagl U, Donnan GA. Evidence for the efficacy of NXY-059 in experimental focal cerebral ischaemia is confounded by study quality. *Stroke.* 2008;39: 2824–2829.
 25. Currie GL, Delaney A, Bennett MI, Dickenson AH, Egan KJ, Vesterinen HM, et al. Animal models of bone cancer pain: Systematic review and meta-analyses. *Pain.* 2013;154: 917-926.
 26. Vesterinen HM, Sena ES, French-Constant C, Williams A, Chandran S, Macleod MR. Improving the translational hit of experimental treatments in multiple sclerosis. *Mult Scler.* 2010;16: 1044-1055.
 27. Rooke EDM, Vesterinen HM, Sena ES, Egan KJ, Macleod MR. Dopamine agonists in animal models of Parkinson's disease: A systematic review and meta-analysis. *Parkinsonism Relat Disord.* 2011;17: 313–320.
 28. Scott S, Kranz JE, Cole J, Lincecum JM, Thompson K, Kelly N, et al. Design, power, and interpretation of studies in the standard murine model of ALS. *Amyotroph Lateral Scler.* 2008;9: 4–15.
 29. Begley CG, Ellis LM. Drug development: Raise standards for preclinical cancer research. *Nature.* 2012;483: 531–533.
 30. Prinz F, Schlange T, Asadullah K. Believe it or not: how much can we rely on published data on potential drug targets? *Nat Rev Drug Discov.* 2011;10: 712.
 31. Cohen J. The statistical power of abnormal-social psychological research: a review. *J Abnorm Soc Psychol.* 1962;65: 145–153.
 32. Lazebnik Y. Can a biologist fix a radio?--Or, what I learned while studying

- apoptosis. *Cancer Cell*. 2002;2: 179–182.
33. Goodman SN, Berlin JA. The use of predicted confidence intervals when planning experiments and the misuse of power when interpreting results. *Ann Intern Med*. 1994;121: 200-206.
 34. Smaldino PE, McElreath R. The natural selection of bad science. *R Soc Open Sci*. 2016;3: 160384.
 35. Wald C, Wu C. Of mice and women: The bias in animal models. *Science*. 2010;327: 1571–1572.
 36. Beery AK, Zucker I. Sex bias in neuroscience and biomedical research. *Neurosci Biobehav Rev*. 2011;35: 565–572.
 37. Mogil JS, Chanda ML. The case for the inclusion of female subjects in basic science studies of pain. *Pain*. 2005;117: 1–5.
 38. Prendergast BJ, Onishi KG, Zucker I. Female mice liberated for inclusion in neuroscience and biomedical research. *Neurosci Biobehav Rev*. 2014;40: 1–5.
 39. Wizemann TM. Sex-specific reporting of scientific research. Washington: National Academies Press; 2012.
 40. Clayton JA, Collins FS. Policy: NIH to balance sex in cell and animal studies. *Nature*. 2014;509: 282–283.
 41. Nuzzo R. How scientists fool themselves – and how they can stop. *Nature*. 2015;526: 182–185.