

# Reporting and interpreting non-significant results in animal cognition research

Benjamin G Farrar <sup>Corresp., 1, 2</sup>, Alizée Vernouillet <sup>3</sup>, Elias Garcia-Pelegrin <sup>1</sup>, Edward W Legg <sup>4, 5, 6</sup>, Katharina F Brecht <sup>7</sup>, Poppy Lambert <sup>8</sup>, Mahmoud Elsherif <sup>9</sup>, Shannon Francis <sup>10</sup>, Laurie O'Neill <sup>10</sup>, Nicola S Clayton <sup>1</sup>, Ljerka Ostojic <sup>Corresp. 4, 5, 6</sup>

Corresponding Authors: Benjamin G Farrar, Ljerka Ostojic Email address: farrarbg@gmail.com, lj.ostojic@uniri.hr

How negative results are reported and interpreted following null hypothesis significance testing is often criticized. With small sample sizes and often low number of test trials, studies in animal cognition are prone to producing non-significant p-values, irrespective of whether this is a false negative or true negative result. Thus, we assessed how negative results are reported and interpreted across published articles in animal cognition and related fields. In this study, we manually extracted and classified how researchers report and interpret non-significant p-values and examined the p-value distribution of these non-significant results. We found a large amount of heterogeneity in how researchers report non-significant p-values in the result sections of articles, and how they interpret them in the titles and abstracts. "No Effect" interpretations were common in the titles (84%), abstracts (64%), and results sections (41%) of papers, whereas "Non-Significant" interpretations were less common in the titles (0%) and abstracts (26%), but were present in the results (52%). Discussions of effect sizes were rare (<5% of articles). A p-value distribution analysis was consistent with research being performed with low power research to detect effect sizes of interest.

<sup>1</sup> Department of Psychology, University of Cambridge, Cambridge, United Kingdom

<sup>&</sup>lt;sup>2</sup> Institute for Globally Distributed Open Research and Education (IGDORE), Cambridge, United Kingdom

<sup>3</sup> Department of Experimental Psychology, Universiteit Gent, Gent, Belgium

<sup>4</sup> Department of Psychology, Faculty of Humanities and Social Sciences, University of Rijeka, RIjeka, Croatia

<sup>5</sup> Division of Cognitive Sciences, University of Rijeka, Rijeka, Croatia

<sup>&</sup>lt;sup>6</sup> Centre for Mind and Behaviour, University of Rijeka, Rijeka, Croatia

<sup>7</sup> Institute for Neurobiology, University of Tuebingen, Tuebingen, Germany

<sup>8</sup> Messerli Research Insititute, University of Vienna, Vienna, Austria

Department of Psychology, University of Birmingham, Birmingham, United Kingdom

<sup>10</sup> Comparative Cognition Research Group, Max Planck Institute for Ornithology, Seewiesen, Germany





Gı Ye

1	Reporting and interpreting non-significant results in animal cognition research
2	
3	Benjamin G. Farrar <sup>1, 2</sup> , Alizée Vernouillet <sup>3</sup> , Elias Garcia-Pelegrin <sup>1</sup> , Edward W. Legg <sup>4, 5, 6</sup> ,
4	Katharina F. Brecht <sup>7</sup> , Poppy Lambert <sup>8</sup> , Mahmoud Elsherif <sup>9</sup> , Shannon Francis <sup>10</sup> , Laurie O'Neill <sup>10</sup>
5	Nicola S. Clayton <sup>1</sup> , Ljerka Ostojić <sup>4, 5, 6</sup>
6	
7	<sup>1</sup> Department of Psychology, University of Cambridge, UK
8	<sup>2</sup> Institute for Globally Distributed Open Research and Education (IGDORE), UK
9	<sup>3</sup> Department of Experimental Psychology, Universiteit Gent, Belgium <sup>4</sup> Department of Psychology, Faculty of Hymanities and Social Sciences, University of Piicks
10 11	<sup>4</sup> Department of Psychology, Faculty of Humanities and Social Sciences, University of Rijeka, Croatia
12	<sup>5</sup> Division of Cognitive Sciences, Faculty of Humanities and Social Sciences, University of
13	Rijeka
14	<sup>6</sup> Centre for Mind and Behaviour, University of Rijeka, Croatia
15	<sup>7</sup> Institute for Neurobiology, University of Tübingen, Germany
16	<sup>8</sup> Messerli Research Institute, University of Veterinary Medicine Vienna, Austria
17 18	<sup>9</sup> Department of Psychology, University of Birmingham, UK <sup>10</sup> Comparative Cognition Research Group, Max Planck Institute for Ornithology, Germany
19	Comparative Cognition Research Group, wax Flanck institute for Ornithology, Germany
20	Corresponding Authors:
21	Benjamin G. Farrar
22	20 Bosworth Road
23	Cambridge
24	CB1 8RG
25	United Kingdom
26	Email address: farrarbg@gmail.com
27	
28	
29	Ljerka Ostojić
30	Odsjek za psihologiju
31	Filozofski fakultet u Rijeci
32	Sveučilišna avenija 4
33	51000 Rijeka
34	Croatia
35	Email address: <u>lj.ostojic@uniri.hr</u>
36	
37	
38	Abstract
39	How negative results are reported and interpreted following null hypothesis significance testing
40	is often criticized. With small sample sizes and often low number of test trials, studies in animal
41	cognition are prone to producing non-significant <i>p</i> -values, irrespective of whether this is a false





43

44

45

46

47

48

49

50

51

52

negative or true negative result. Thus, we assessed how negative results are reported and interpreted across published articles in animal cognition and related fields. In this study, we manually extracted and classified how researchers report and interpret non-significant *p*-values and examined the *p*-value distribution of these non-significant results. We found a large amount of heterogeneity in how researchers report non-significant *p*-values in the result sections of articles, and how they interpret them in the titles and abstracts. "No Effect" interpretations were common in the titles (84%), abstracts (64%), and results sections (41%) of papers, whereas "Non-Significant" interpretations were less common in the titles (0%) and abstracts (26%), but were present in the results (52%). Discussions of effect sizes were rare (<5% of articles). A *p*-value distribution analysis was consistent with research being performed with low power research to detect effect sizes of interest.

5354

#### Introduction

5556

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

Null hypothesis significance testing (NHST) is a primary method of statistical analysis in animal cognition research. However, when NHST produces results that are not statistically significant, these are often difficult to interpret. If researchers test null hypotheses of zero effect (i.e., there are no differences between groups or conditions), a non-significant result could result from a lack of any effect in the population (a true negative), or a failure to detect some true difference (a false negative). While current guidance encourages researchers to design studies with high statistical power to detect theoretically interesting effect sizes (Lakens, 2017, 2021) – which can provide context for negative results – power analyses appear infrequent (Fritz et al., 2013). Hence, how negative results are reported and interpreted following non-significant results from NHST has been criticised on several grounds (Gigerenzer et al., 2004; Lambdin, 2012), with some researchers noting that false negative claims can inhibit scientific progress just as much as false positive claims (Fiedler et al., 2012; Vadillo et al., 2016). The most prominent criticism is that researchers often misreport or misinterpret non-significant results as meaning that, i) within the **sample** that was tested there was no effect (a specific report about what happened in the study) and/or that ii) the nonsignificant result means that in the target **population** in general there is no effect (Aczel et al., 2018; Fidler et al., 2006; Hoekstra et al., 2006). This misreporting or misinterpretation even



occurs even when the null hypothesis being considered is very likely to be incorrect (Cohen,

75 1994; Gelman & Carlin, 2014). Given these concerns, this study explored how animal cognition

researchers report and interpret non-significant results using a manually extracted dataset of

77 negative claims following NHST from over 200 articles.

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

76

#### 1. NHST and p-values

When using NHST, researchers attempt to reject a statistical model (the null hypothesis) with their data while controlling the rate at which they will make false-positive decisions in the long-term (Neyman & Pearson, 1933). Most often, this statistical null is that there is absolutely no difference between two groups or conditions (for example a mean difference of 0 for a t-test; 'nil' hypothesis; Cohen, 1994), or, in the case of a one-tailed test, that the difference will not be zero or that it will be not in a certain direction, i.e., researchers make a directional prediction for their alternative hypothesis. A statistical test then produces a p-value, i.e., the probability of observing the researchers' data or more extreme data if the null hypothesis and all its assumptions were true,  $Pr(d(X) \ge d(x_0); H_0)$ . If the p-value is lower than a pre-specified threshold (the  $\alpha$  level), the statistical null hypothesis (H<sub>0</sub>) is rejected in favor of an alternative hypothesis (Neyman & Pearson, 1933), whereas if the p-value is larger than the pre-specified threshold, H<sub>0</sub> should not be rejected. However, how researchers should behave towards their null and alternative hypotheses following a non-significant result has been a continued locus of criticism of NHST (Lambdin, 2012). Formally, researchers *can* make statements about the long-run error probabilities of their test procedures. For example, with an  $\alpha$  level of .05 and if no  $\alpha$ -inflating research practices were used (Simmons et al. 2011), they can say that in the long run they would not reject H<sub>0</sub> more than 5% of the time, if H<sub>0</sub> were true. Similarly, if the design of the study is such that the statistical test had 90% power to detect the smallest effect size of interest, in the long run the researchers would only fail to reject  $H_0$  10% of the time, if the smallest effect size of interest did exist in the population.

100

101

102

103

104

#### 2. Accepting the null: How much of an error?

Without performing further analyses, it can be an error to conclude that there is evidence in favor of the null hypothesis following a non-significant result. The arbitrary nature of the  $\alpha$  level highlights this: as an example, let us assume that we calculate a p-value of 0.08 with an  $\alpha$ 



105	level of .05. By not rejecting $H_0$ in this instance, we can say that in the long run we would not
106	reject $H_0$ more than 5% of the time, if it were true, when performing this procedure. However, if
107	we had chosen an $\alpha$ level of .10 instead, we would have rejected $H_0$ . Clearly, then, the p-value
108	when using NHST is not a direct indication of the strength of evidence for or against H <sub>0</sub> , but
109	must be interpreted relative to error rates and alternative hypotheses (Lakens et al., 2018).
110	However, despite the <i>p</i> -value not being the probability of the null hypothesis being true, survey
111	studies suggest researchers do interpret $p$ -values in such a way (e.g. Goodman, 2008). Moreover,
112	scientists often misreport non-significant results as evidence of absence of a difference between
113	groups of conditions or evidence of no effect when this inference is not necessarily warranted.
114	For example, Hoekstra et al., (2006) reported that 41% of articles containing non-significant
115	results in 1994 and 1995 in Psychonomic Bulletin & Review included the interpreted of non-
116	significant results as "evidence of no effect", a figure which rose to 60% in 2002 to 2004.
117	Similarly, Fidler et al. (2006) found that 63% of articles in 2000 and 2001 in Conservation
118	Biology and Biological Conservation reported non-significant findings as "evidence of no
119	effect". More recently, Aczel et al. (2018) found that 72% of non-significant results were
120	reported as "no effect" in the abstracts of 2015 articles in Psychonomic Bulletin & Review,
121	Journal of Experimental Psychology: General, and Psychological Science. Such an error might
122	be especially important in animal cognition research, in which a combination of small sample
123	sizes and low trial number may limit the ability of researchers to design studies and statistical
124	test combinations with high power of statistical tests to detect the minimum effect size of
125	theoretical interest (Farrar et al., 2020).
126	While 'accepting the null' may be an error, just how severe an error it is requires
127	discussing. In their survey of 86 Psychonomic Bulletin & Review Hoekstra et al. (2006, p. 1036)
128	reported that: "We found the serious mistake of accepting the null hypothesis and claiming no
129	effect in 60% (CI: 53%, 66%) of the articles that reported statistically nonsignificant results"
130	(emphasis added). And interpreting a non-significant result as if there were no differences
131	between conditions ranks at Number 2 of Goodman's (2008) "Dirty Dozen" p-value
132	misconceptions. However, just because a researcher might report the results of significance tests
133	incorrectly, this does not mean that they themselves, or their readers, necessarily interpreted the
134	significance test incorrectly. In their 1933 paper, Neyman and Pearson often discussed 'accepting
135	H <sub>0</sub> ' following a result that was not statistically significant (Neyman & Pearson, 1933). In fact, as

Mayo (2018, p. 135) writes, Neyman used the term 'acceptance' as shorthand, and even preferred the phrase "No evidence against [the null hypothesis] is found" to "Do not reject [the null hypothesis]" (Neyman, 1976, postscript, p. 749). If scientists equate phrases such as "there were no differences between conditions (p > 0.05)" or "therefore we accept  $H_0$ " with "there was no statistically significant difference between the conditions" or "therefore we do no reject  $H_0$ ", then the "serious mistake" of accepting the null becomes an issue of precision in language, rather than an egregious error. This is exemplified in cases where the observed experimental data are clearly more in-line with the null hypothesis than the alternative hypothesis,  $H_1$ .

#### 3. Exploring non-significant result reporting and interpretation in animal cognition

Understanding how animal cognition researchers have reported and interpreted nonsignificant findings in their published articles is an important step to, i) identify how often
negative conclusions in animal cognition might be a result of NHST misreporting or
misinterpretation, and ii) highlight areas in which animal cognition researchers can improve their
statistical inferences and statistical reporting. In this study we explored how authors in fields
related to animal cognition report and interpret non-significant results by building on the
methods used in similar studies in psychology and conservation biology (Aczel et al., 2018;
Fidler et al., 2006; Hoekstra et al., 2006). Furthermore, we also extracted the *p*-values associated
with the negative results in our sample. We compared the distribution of these to four simulated *p*-value distributions in which research was performed with no publication bias and either 10%,
33%, 50%, or 80% power. Comparing the observed literature p-value distribution to the four
simulated distribution can provide cues to the average statistical power of the research we
extracted data from, as well as the presence of any biasing effects like publication bias (but note
there is likely a large degree of heterogeneity here – see Nord et al. (2017) for a discussion of
this in neuroscience research).

#### 4. Reporting and interpreting non-significant results in animal cognition

F

 $<sup>^1</sup>$  For example, consider a study in which birds' latencies to approach a novel object are compared between 10,000 wild and 10,000 hand reared birds (with 99% power to detect a pre-specified effect size of interest of 2 seconds), and in which a difference of 0.02 seconds was observed. This difference may even be statistically significant, but the minimum effect size of interest. Here, saying "there was no difference between the latency of wild and hand reared birds to approach the novel object (p > .05)", although literally incorrect, does not seem to be an error of great consequence.

F



163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

In order to investigate how animal cognition researchers report negative results, we manually extracted text and data of non-significant results from 18 journals in animal cognition, behavior, and welfare, one pre-print server, and from articles recommended by PCI: Animal Science. We extracted data from articles reporting non-significant findings in their titles, abstracts and results sections and classified how the authors interpreted them. Our classification was descriptive and aimed to characterize the different ways in which researchers reported the results of nonsignificant findings, and how these were translated into claims about populations and thus into substantive claims. Specifically, for reporting of negative results in the abstract or results section about the specific sample tested in the study, we classified the negative result text into three categories: 1) "Non-Significant" statements that either reported that there was no significant difference between two conditions, or words to that effect, or reported a correct directional statement; 2) "No Effect" statements that here was no difference within the sample, when in fact there was — it was just not significant in the analysis; 3) "Ambiguous" statements that neither suggest that samples were the same, nor that there was no significant difference between them. Similarly, for substantive claims about the population in the abstracts or titles, we had three related categories: 1) "Justified": An interpretation that commented on statistical power, uses equivalence tests or otherwise justifies why a non-significant result suggests that there is no theoretically important difference in the population, or that the study provides no strong evidence of a difference, 2) "Caveated, Ambiguous or Similar": An interpretation of the non-significant results as suggesting/indicating etc. that X and Y do not differ in the population, or showing that they are similar, or 3) "No Effect": An interpretation of the non-significant result as showing that X and Y do not differ in the population.

184 185 186

187 188

189 190

191

192

193

194

#### **Materials & Methods**

#### **Sample**

We extracted data from a total of 20 sources, comprising 18 peer-reviewed journals, one pre-print server, and articles recommended through Peer Communities In. The 20 sources are detailed in *Table 1*.



#### **Data extraction and Classification**

BGF, AV, KB, EGP, LoN, PL, SF, EL, and ME performed the coding and were each
assigned two journals, except BGF who conducted the coding for four journals. Each coder
screened the abstracts of each article of their assigned journals and identified any negative
statements about either, i) the specific sample tested in that study or, ii) a wider population. If a
negative statement was present, the coder then recorded the paper's information (title, first
author, journal, and year) and the negative statement. For articles with multiple negative
statements for either the sample or the population, the coder recorded the negative statement that
they thought was most clearly related to the paper's main claim, such that for each article, we
had a maximum of one negative sample statement and one negative population statement. Next,
the coder verified that the statements were based on results from NHST. If verified, the coder
then extracted the text of the NHST that corresponded to the abstract claim from the results
section of the manuscript, including the associated <i>p</i> -value. If there was more than one
corresponding statistical test within an experiment, the coder extracted the test result that they
thought was most relevant to the claim. If the abstract claim was equally supported by multiple
studies or experiments, the coder extracted the information from the first study or experiment
presented.
After the title, abstract claims (sample and population), result text and p-value had been
extracted, the coder categorized how each negative statement was reported. Through piloting,
discussion, from looking at previous studies (Aczel et al., 2018; Fidler et al., 2006; Hoekstra et
al., 2006), we developed three categories. For the sample claims and result text, these were: 1)
"Non-Significant" statements that either there was no significant difference between two
conditions, or words to that effect, or a correct directional statement; 2) "No Effect" statements
that there was not a difference within the sample, when in fact there was — it was just not
significant in the analysis; 3) "Ambiguous, Similar or Small Effect Size" statements about the
results that neither suggest that samples were the same, nor that there was no significant
difference between them (which were later split into "Ambiguous" and "Similar or Small Effect
Size" categories). In addition to these descriptions, we developed a table of hypothetical
statements that are detailed in Table 2, which were available to the coders during the project.
Similarly, the title, if it contained a negative statement, and population claims from the abstracts
were categorized into three categories: 1) "Justified": An interpretation that commented on





statistical power, use of equivalence tests or otherwise a justification why a non-significant result
suggests that there is no theoretically important difference in the population, or that the study
provides no strong evidence of a difference, 2) "Justified, Caveated or Ambiguous": An
interpretation of the non-significant results as suggesting/indicating etc. that X and Y do not
differ in the population, or showing that they are similar, and 3) "No Effect": An interpretation of
the non-significant result as showing that X and Y do not differ in the population. In addition to
these descriptions, we developed a table of hypothetical statements that are detailed in <i>Table 3</i> .
Reliability and Quality Control
Twenty-four articles (8.5%) were double-blind coded in order to assess the likely
reliability of our coding scheme, and all articles underwent a quality control procedure involving
a second coder to identify any mistakes or inconsistencies.
Double-Blind Extraction
BGF independently coded 24 articles, namely the first four articles from six randomly
chosen journals, blind to the results of the original coders. From this, we computed inter-rate
agreement for each variable that we extracted (Title Population Claim Level; Title Sample Claim
Level; Abstract Sample Claim Text; Abstract Sample Claim Level; Abstract Population Claim
Text; Abstract Population Claim Level; Result Text; Result Level; p-value).
Quality Control
All articles underwent the quality control procedure. Here, a second coder reviewed the
data extracted from each article. BGF, AV, KB, EGP, LoN, PL, SF, EL, ME, and LO served as
second coders, and each was assigned one other coder's original set of articles to quality control.
The quality controller verified 1) that a negative claim from the title/abstract has been extracted,
2) that any negative claim extracted was really a negative claim, 3) that the result that was
extracted corresponded to the claim that was extracted, and 4) that they agreed with the

classification of each claim. If the quality controller identified a mistake, they classified this as a

major disagreement, whereas if the quality controller disagreed but was uncertain about this

254





judgment, for example in the case of borderline claims, they classified this as a minor
disagreement. BGF reviewed all disagreements and made a final decision on what entered the
final dataset, returning to the original article if necessary.

#### **Analysis**

The primary analysis was descriptive. The percentage of claims in each category across the titles, abstract population claims, abstract sample claims, and result texts are presented. To illustrate the types of claims placed in each category, examples that we felt were particularly representative of each are provided in tables. In addition, every classification can also be viewed in the open dataset. We used a Chi-squared test to test whether, if a "No Effect" interpretation was made in the results, it was more likely that a "No Effect" interpretation would also be made in the abstract than when a correct interpretation was made in the results. All data and code as well as our coding guidelines are openly available at <a href="https://osf.io/84puf/">https://osf.io/84puf/</a>.

#### Results

We extracted data from 302 articles. Of these, 18 were excluded due to their identified claim having no corresponding negative result of NHST (e.g., only descriptive statistics used, or only a Bayesian analysis performed) and one was excluded due to excessive ambiguity in how the results were described. This left a final sample of 283 articles for analysis.

#### **Reliability and Quality Control**

#### 278 Double-Blind Coding

For 24 articles (8.5% of the total sample), two authors (BGF and the author originally assigned to the journal) extracted all the data independently of each other. Only five of the article titles were identified as containing negative statements by either of the two coders, and from this, the two coders agreed on only one out of five (20%) of the articles about whether the title statement was about the tested sample or the population. Following discussion with the whole group, we agreed that it was often ambiguous whether the titles of articles were referring to the



specific sample tested or a wider population, and so we decided to combine these measures and have no sub-group analysis for the title claim, deviating from our original plan. When considering the category (Justified; Caveated, Similar or Ambiguous; No Effect) of the title claim, the two coders agreed on two out of six articles (33%). Three of the four disagreements occurred when one coder did not interpret the title as a negative claim, e.g. as in "Evidence that novel flavors unconditionally suppress weight gain in the absence of flavor-calorie associations" (Seitz et al., 2020), and one where a coder appeared to have made an error. From discussion within the group, it was evident that these ambiguous cases — where the statements were not clearly written as negative statistical results but involved an interpretation that did not directly corresponds to a specific non-significant result from the article — proved the most difficult during the whole coding procedure, and this affected the reliability of the title claims and population claims from the abstract.

The coders identified 24 sample claims from the abstracts of the papers, from which they coded the same claim on 22 out of 24 occasions (91.6%). Of these 22 claims, the two coders agreed on 19 of their classification (86.3%). In contrast, the coders identified only eight population claims from the abstracts of articles, from which they agreed on three occasions (37.5%), and of these three, agreed on two of their classification (66.7%). From the results, the coders recorded the same text for 16 of the 22 (72.7%) abstract claims that they had coded the same, and of these 16, they agreed on 13 of their classification (81.3%) and extracted the exact same p-value for 10 of these 13 (76.9%).

In sum, the double-blind coding demonstrated good inter-rater consistency for how the abstract sample claims and associated results and *p*-values were extracted, even before our quality control procedures had been implemented. In contrast, inter-rater consistency was low for the title claims and population claims from the abstracts. This matched our subjective experience of the coding procedure, where we experienced many cases of population claims as vague and about a theoretical hypothesis that did not closely correspond to any particular negative result from the article. In contrast, the negative sample claims could often be easily mapped onto a particular negative result in the text.

Quality Control





Each article was checked by a quality controller. The initial coders identified 67 possible
negative statements in the titles of papers, and the quality controller agreed with the
classification of 39 (58%) of these statements, had a minor disagreement with six statements
(9%), and a major disagreement with 22 statements (33%). Of note, 16 of these 22 major
disagreements came from a single repeated error in which one individual coder coded
'ambiguous' for titles containing no negative statement. In the abstract, coders identified 281
negative statements about the specific sample tested in the paper. Of these, the quality controllers
agreed with the classification of 250 (89%), had minor comments about 16 (6%), and major
disagreements with 15 (5%). Coders identified a much smaller number of negative inferences
about populations in the articles and disagreed more frequently: Of the 82 identified statements,
the quality controllers agreed with the classification of 44 (53%), had minor comments about 18
(22%) and major disagreements with 20 (24%). Regarding the result texts from the article
bodies, coders identified 282 results, of which the quality controller agreed with the
classification and extracted p-value for 252 (89%), had minor comments for 13 (5%), and major
disagreements for 17 (6%).

The quality control process allowed us to, i) identify any clear errors in the data extraction process, ii) highlight borderline cases where our coding scheme could not clearly categorize certain statements, and iii) assess the robustness of the coding procedure. In line with the results from the double-blind coding, the quality control process demonstrated a high interrater agreement and consistency with identifying and classifying negative sample statements from abstracts, and the corresponding results and *p*-values from the main text, yet greater inconsistency in deciding, i) whether titles and population statements were truly "negative" in the sense of being the result of a non-significant NHST, and ii) whether the authors were claiming the absence of an effect from these negative results. This inconsistency occurred mainly because many titles and population claims referred not to a certain statistical result but made a vague theoretical statement.

#### **Title Claims**

Forty-four titles contained negative statements resulting from non-significant results of NHSTs. Of these, 37 (84%) were classified as interpreting the non-significant result as evidence





345	of no effect, whereas seven (16%) were classified as caveated claims or claims about two groups
346	or conditions being 'similar'. Table 4 provides examples of these claims.
347	
348	Abstract Claims
349	Abstract Sample Claims
350	We extracted 278 negative claims about a sample result. Of these, 174 (63%) were
351	classified as claiming evidence of no effect, 71 (26%) as making formally correct statements that
352	there were no statistically significant differences between groups or conditions, 17 (6%) as
353	making claims about an effect being 'similar' between groups or conditions, or as describing a
354	small effect size, and 16 (6%) were classified as ambiguous. Table 5 provides examples of these
355	claims.
356	
357	Abstract Population Claims
358	We extracted 63 negative claims about a population that followed on from the negative result
359	within a sample. Of these, 45 (71%) were classified as caveated and 18 as claiming that there
360	was no effect (29%). Table 6 provides examples of these claims.
361 362	Result Text
363	In the results sections, 276 non-significant results of NHST were coded. Of these, 140 (52%)
364	were classified as reporting the results as "Non-Significant", 113 (41%) as reporting that there
365	was "No Effect", 12 (4%) as reporting groups or conditions being "Similar", 10 (4%) were
366	classified as "Ambiguous", and one $(0.4\%)$ as reporting a "trend" in the opposite direction to the
367	prediction. Several of the classifications as ambiguous were due to authors' use of "main effect"
368	when interpreting ANOVA terms, where we thought that saying there was "no main effect of X"
369	was different enough to saying "no effect of X" to not be included in the "No Effect" category,
370	although this highlights the somewhat arbitrariness of our categories. Table 7 provides examples
371	of the different types of result reporting.
372	Notably, if a sentence reporting the results in the results section was classified as "No
373	Effect", it was more likely that a no effect interpretation would also be made in the abstract,





374	compared to when the result was classified as "Non-Significant" ( $\chi 2(1, N = 211) = 21.65, p < 1.65$
375	.0001). Limiting the data to just those with responses in the abstract and results classified as
376	"Non-Significant" or "No Effect", of the 92 statements in the results classified as "No Effect",
377	80 (87%) of the corresponding sample statements were classified as interpreting there being no
378	effect. In contrast, of the 119 statements in the results classified as "non-significant", only 67
379	(56%) of the corresponding sample statements were classified as interpreting there being no
380	effect. Nevertheless, "No Effect" interpretations in the abstracts were absolutely the most likely
381	classifications for both "No Effect" and "Non-Significant" results statements.
382	
383	p-value distributions
384	In total, 202 of the 283 papers reported exact <i>p</i> -values, with the other 81 reporting either
385	inequalities or not reporting the $p$ -values at all. Of these 202 $p$ -values, four were below .05 and
386	non-significant due to a lower $\alpha$ level. The distribution of the 198 non-significant $p$ -values in the
387	interval .05-1 is displayed in Figure 1. This distribution significantly differs from a uniform
388	distribution (two-sided Kolmogorov-Smirnov test, $D = 0.12$ , $p = .0087$ ).
389	Error! Reference source not found. contrasts the distribution of Figure 1 with the four
390	simulated distributions of bodies of research performed where 80% of alternative hypotheses
391	were correct, and studies had either 10, 33, 50 or 80% power to detect the true effect size of H1
392	if it was true. Notably, p-values in the interval from .05 to .10 were underrepresented in the
393	manually extracted data, making up only 5.6% of observations compared to 8.2% (10% power
394	simulation), 15% (33% power simulation), 19% (50% power simulation), and 20% (80% power
395	simulation). Similarly, very high $p$ -values (.95-1.0) were overrepresented in our manual dataset
396	(7.6% of observations, compared to 4.3%, 3.2%, 2.4% and 3.4% for the 10, 33, 50 and 80%
397	power simulations respectively), which likely reflects either the use of multiple correction
398	procedures, or small sample non-parametric statistics that produce non-uniform distributions
399	under the null hypothesis.
400	
401	Discussion
402	We extracted and classified how animal cognition researchers reported the results of non-
403	significant null hypothesis significance tests in 253 articles between 2019 and 2021. Across
404	titles, abstracts, and results, we classified non-significant results as often being reported with the



405	No Effect' phrasing that has often been labelled as erroneous (titles 84%; abstract sample
406	results 63%; result text 41%). Reporting negative results as "Non-Significant" was less common
407	in titles and abstracts, but as prevalent as "No Effect" phrasings in the results (titles 16%;
408	abstract sample results 26%; result text 52%). The other, albeit less frequently classified method
409	of reporting non-significant results was to comment on the similarity between two groups or
410	conditions (abstract sample results 6%; result text 4%).
411	Overall, these results demonstrate considerable heterogeneity in how animal cognition
412	researchers report and potentially interpret non-significant results in published articles (according
413	to our classification of these reports, made by other researchers in the field). However, it was
414	often difficult to confidently categorize results due to the heterogeneity in how negative results
415	were reported. are Nevertheless, our results suggest that negative results are at risk of being
416	misreported and misinterpreted in animal cognition publications. It remains a question, however,
417	what the consequences of such misreporting might be, i.e., how readers of scientific articles
418	interpret "No Effect" statements, and this could be studied through analyzing how these studies
419	are cited, in other publications but also in media reports and student essays. Possibly
420	encouragingly, when researchers extended "No Effect" statements from the sample to the
421	population, they routinely opted for qualifiers to caveat inference to the populations (e.g.,
422	"these results suggest that there is no effect at the population level"2). Again, however, more
423	research is needed to understand how such statements are interpreted and implemented by
424	scientists and the wider community. One way in which researchers might reduce the ambiguity
425	of their negative statements would be to use more formal methods of assessing evidence against
426	informative null hypotheses, such as by testing against theoretically interesting effect sizes using
427	as equivalence tests or comparing plausible null and alternative hypotheses using Bayes factors.
428	Although beyond the scope of the current project, Lakens (2017) provides a detailed tutorial for
429	equivalence testing in psychological research, and Rose et al., (2018) in animal behavior, and
430	Rouder et al. (2009) provide an introduction to Bayes Factors.

claims researchers wish to test and the actual statistical hypotheses that are tested, i.e., rarely can

Notably, the coding team found it difficult to identify and classify negative population

statements in the abstracts of articles. This likely reflects the distance between the theoretical

<sup>&</sup>lt;sup>2</sup> Although we did not study this, it is likely this type of caveating is not unique to negative results but used to caveat positive findings, too.





434 a theoretical prediction about an animal's cognition be reduced to a single decision between a 435 null and alternative hypothesis in a null hypothesis significance test. 436 Finally, we classified "No Effect" interpretations more commonly in abstracts and titles than "No Effect" reporting of results in the results section. That is, authors who have written out 437 "Non-Significant" results in the results section nevertheless wrote "No Effect" interpretations in 438 the abstracts and titles. This could be due to two factors, namely word limits and incentives to 439 440 make bolder claims. If this is correct, then the former should be considered by journal editorial 441 boards when setting their policy. 442 The p-value distribution likely differed from a uniform distribution for two reasons: the 443 cumulative frequency was greater in the observed distribution for smaller p-values (p < .3) and 444 was also greater for large p-values (p > .95). The larger density of smaller p-values is consistent 445 with research with low-powered statistical tests in which the null hypothesis was incorrect, but 446 which produces p-values that did not reach statistical significance. The density of very large p-447 values is consistent with researchers applying corrections that might increase p-values, such as 448 Bonferroni corrections, or by using statistical tests with small sample sizes that produce non-449 uniform p-value distributions under the null. An interesting contrast between the observed and 450 simulated p-value distributions is that, unlike in the manual distribution, p-values in the range .05 451 to .10 were much more common than p-values in the range .10 to .15 in the simulated 452 distributions. This is likely because we extracted results that researchers had interpreted as 453 negative for the manual dataset, but p-values in the range .05-0.1 are often interpreted as "trends" 454 or "marginally significant". 455 456 Conclusions This study explored reporting and interpretation of negative result in animal cognition literature 457 458 through classification by other researchers in the field. In line with previous studies in other 459 disciplines (Aczel et al., 2018; Fidler et al., 2006), we found that non-significant results were 460 often reported as if there were no differences observed in the sample, and this was the case in the 461 titles, abstracts and result sections of papers, although it was most frequent in the titles and 462 abstracts. Because of the distance between statistical hypotheses and theoretical claims, and 463 uncertainty around how no difference statements are interpreted, the consequences of this 464 putative error are uncertain. Nevertheless, these results suggest that researchers should pay close

attention to the evidence used to support claims of absence of effects in the animal cognition



466	literature, and prospectively seek to, i) report non-significant results clearly, and ii) use more
467	formal methods of assessing the evidence against theoretical predictions.
468	
469	
470	Acknowledgements
471 472 473	We would like to thank Balazs Aczel for discussions and clarifications about previous research in this area
474	
475	D. formana
476	References
477 478	Aczel, B., Palfi, B., Szollosi, A., Kovacs, M., Szaszi, B., Szecsi, P., Zrubka, M., Gronau, Q. F., van den Bergh, D., & Wagenmakers, EJ. (2018). Quantifying support for the null
479	hypothesis in Psychology: An empirical investigation. Advances in Methods and Practices in
480	Psychological Science, 1(3), 357–366. https://doi.org/10.1177/2515245918773742
481 482	Anselme, P., & Robinson, M. J. F. (2019). Evidence for motivational enhancement of
482 483	sign-tracking behavior under reward uncertainty. <i>Journal of Experimental Psychology: Animal</i>
484	Learning and Cognition, 45(3), 350–355. https://doi.org/10.1037/xan0000213
485	Learning and Cognition, 43(3), 330-333. https://doi.org/10.1037/xdii/0000213
486	Aparecida Martins, R., Ribeiro Caldara, F., Crone, C., Markiy Odakura, A., Bevilacqua,
487	A., Oliveira dos Santos Nieto, V. M., Aparecida Felix, G., Pereira dos Santos, A., Sousa dos
488	Santos, L., Garófallo Garcia, R., & de Castro Lippi, I. C. (2021). Strategic use of straw as
489	environmental enrichment for prepartum sows in farrowing crates. <i>Applied Animal Behaviour</i>
490	Science, 234, 105194. https://doi.org/10.1016/j.applanim.2020.105194
491	
492	Beran, M. J., French, K., Smith, T. R., & Parrish, A. E. (2019). Limited evidence of
493	number-space mapping in rhesus monkeys (Macaca mulatta) and capuchin monkeys (Sapajus
494	apella). Journal of Comparative Psychology, 133(3), 281–293.
495	https://doi.org/10.1037/com0000177
496	
497	Brecht, K. F., Müller, J., & Nieder, A. (2020). Carrion crows (Corvus corone corone) fail
498	the mirror mark test yet again. Journal of Comparative Psychology, 134(4), 372–378.
499	https://doi.org/10.1037/com0000231
500	
501	Cimarelli, G., Schoesswender, J., Vitiello, R., Huber, L., & Virányi, Z. (2021). Partial
502	rewarding during clicker training does not improve naïve dogs' learning speed and induces a
503	pessimistic-like affective state. Animal Cognition, 24(1), 107–119.
504 505	https://doi.org/10.1007/s10071-020-01425-9
1111	



506	Cohen, J. (1994). The Earth is round (p $< .05$ ). American Psychologist, 49(12), 997–
507	1003. https://doi.org/10.1037/0003-066X.49.12.997
508	
509	Cunningham, P. J., & Shahan, T. A. (2020). Delays to food-predictive stimuli do not
510	affect suboptimal choice in rats. Journal of Experimental Psychology: Animal Learning and
511	Cognition, 46(4), 385–397. https://doi.org/10.1037/xan0000245
512	
513	DeVries, M. S., Winters, C. P., & Jawor, J. M. (2020). Similarities in expression of
514	territorial aggression in breeding pairs of northern cardinals, Cardinalis cardinalis. Journal of
515	Ethology, 38(3), 377–382. https://doi.org/10.1007/s10164-020-00659-x
516	
517	Farrar, B. G., Boeckle, M., & Clayton, N., S. (2020). Replications in comparative
518	cognition: What should we expect and how can we improve? Animal Behavior and Cognition,
519	7(1), 1–22. https://doi.org/10.26451/abc.07.01.02.2020
520	
521	Fidler, F., Burgman, M. A., Cumming, G., Buttrose, R., & Thomason, N. (2006). Impact
522	of Criticism of Null-Hypothesis Significance Testing on Statistical Reporting Practices in
523	Conservation Biology. Conservation Biology, 20(5), 1539–1544.
524	
525	Fiedler, K., Kutzner, F., & Krueger, J. I. (2012). The long way from α-error control to
526	validity proper: Problems with a short-sighted false-positive debate. Perspectives on
527	Psychological Science, 7(6), 661–669. https://doi.org/10.1177/1745691612462587
528	
529	Fritz, A., Scherndl, T., & Kühberger, A. (2013). A comprehensive review of reporting
530	practices in psychological journals: Are effect sizes really enough? <i>Theory &amp; Psychology</i> , 23(1),
531	98–122. https://doi.org/10.1177/0959354312436870
532	
533	Gelman, A., & Carlin, J. (2014). Beyond power calculations. <i>Perspectives on</i>
534	Psychological Science, 9(6), 641–651. https://doi.org/10.1177/1745691614551642
535	
536	Gigerenzer, G., Krauss, S., & Vitouch, O. (2004). The null ritual: What you always
537	wanted to know about null hypothesis testing but were afraid to ask. Handbook on Quantitative
538	Methods in the Social Sciences. Sage, Thousand Oaks, CA, 389–406.
539	
540	Goodman, S. (2008). A Dirty Dozen: Twelve p-value misconceptions. Seminars in
541	Hematology, 45(3), 135–140. https://doi.org/10.1053/j.seminhematol.2008.04.003
542	
543	Guadarrama, S. S., Domínguez-Vega, H., Díaz-Albiter, H. M., Quijano, A., Bastiaans, E.,
544	Carrillo-Castilla, P., Manjarrez, J., Gómez-Ortíz, Y., & Fajardo, V. (2020). Hypoxia by altitude
545	and welfare of captive beaded lizards (Heloderma Horridum) in Mexico: Hematological



546	approaches. Journal of Applied Animal Welfare Science, 23(1), 74–82.
547	https://doi.org/10.1080/10888705.2018.1562350
548	
549	Harris, J. A., & Bouton, M. E. (2020). Pavlovian conditioning under partial
550	reinforcement: The effects of nonreinforced trials versus cumulative conditioned stimulus
551	duration. Journal of Experimental Psychology: Animal Learning and Cognition, 46(3), 256-272.
552	https://doi.org/10.1037/xan0000242
553	
554	Hashmi, A., & Sullivan, M. (2020). The visitor effect in zoo-housed apes: The variable
555	effect on behaviour of visitor number and noise. Journal of Zoo and Aquarium Research, 8(4),
556	268–282. https://doi.org/10.19227/jzar.v8i4.523
557	
558	Hoekstra, R., Finch, S., Kiers, H. A. L., & Johnson, A. (2006). Probability as certainty:
559	Dichotomous thinking and the misuse of p values. Psychonomic Bulletin & Review, 13(6), 1033-
560	1037. https://doi.org/10.3758/BF03213921
561	
562	Kawaguchi, M., & Kuriwada, T. (2020). Effect of predator cue on escape and oviposition
563	behaviour of freshwater snail. <i>Behaviour</i> , 157(7), 683–697. https://doi.org/10.1163/1568539X-
564	bja10018
565	
566	Kawai, N., Nakagami, A., Yasue, M., Koda, H., & Ichinohe, N. (2019). Common
567	marmosets (Callithrix jacchus) evaluate third-party social interactions of human actors but
568	Japanese monkeys (Macaca fuscata) do not. Journal of Comparative Psychology, 133(4), 488-
569	495. https://doi.org/10.1037/com0000182
570	
571	Koczura, M., Martin, B., Musci, M., Massimo, M. D., Bouchon, M., Turille, G., Kreuzer,
572	M., Berard, J., & Coppa, M. (2021). Little difference in milk fatty acid and terpene composition
573	among three contrasting dairy breeds when grazing a biodiverse mountain pasture. Frontiers in
574	Veterinary Science, 7, 612504. https://doi.org/10.3389/fvets.2020.612504
575	
576	Kvarnemo, C., Andersson, S. E., Elisson, J., Moore, G. I., & Jones, A. G. (2021). Home
577	range use in the West Australian seahorse Hippocampus subelongatus is influenced by sex and
578	partner's home range but not by body size or paired status. <i>Journal of Ethology</i> , 39(2), 235–248.
579	https://doi.org/10.1007/s10164-021-00698-y
580	
581	Lakens, D. (2017). Equivalence tests: A practical primer for t tests, correlations, and
582	meta-analyses. Social Psychological and Personality Science, 8(4), 355–362.
583	https://doi.org/10.1177/1948550617697177
584	



585	Lakens, D. (2021). Sample Size Justification [Preprint]. PsyArXiv.
586	https://doi.org/10.31234/osf.io/9d3yf
587	
588	Lakens, D., Adolfi, F. G., Albers, C. J., Anvari, F., Apps, M. A. J., Argamon, S. E.,
589	Baguley, T., Becker, R. B., Benning, S. D., Bradford, D. E., Buchanan, E. M., Caldwell, A. R.,
590	Van Calster, B., Carlsson, R., Chen, SC., Chung, B., Colling, L. J., Collins, G. S., Crook, Z.,
591	Zwaan, R. A. (2018). Justify your alpha. Nature Human Behaviour, 2(3), 168–171.
592	https://doi.org/10.1038/s41562-018-0311-x
593	
594	Lambdin, C. (2012). Significance tests as sorcery: Science is empirical—significance
595	tests are not. Theory & Psychology, 22(1), 67–90. https://doi.org/10.1177/0959354311429854
596	
597	Lazarowski, L., Thompkins, A., Krichbaum, S., Waggoner, L. P., Deshpande, G., &
598	Katz, J. S. (2020). Comparing pet and detection dogs (Canis familiaris) on two aspects of social
599	cognition. Learning & Behavior, 48(4), 432–443. https://doi.org/10.3758/s13420-020-00431-8
600	
601	Lilley, M. K., de Vere, A. J., & Yeater, D. B. (2020). Laterality of eye use by bottlenose
602	(Tursiops truncatus) and rough-toothed (Steno bredanensis) dolphins while viewing predictable
603	and unpredictable stimuli. International Journal of Comparative Psychology, 33.
604	https://doi.org/10.46867/ijcp.2020.33.03.01
605	
606	Mayo, D. G. (2018). Statistical inference as severe testing: How to get beyond the
607	statistics wars. Cambridge University Press.
608	
609	Meza, P., Elias, D. O., & Rosenthal, M. F. (2021). The effect of substrate on prey capture
610	does not match natural substrate use in a wolf spider. <i>Animal Behaviour</i> , 176, 17–21.
611	https://doi.org/10.1016/j.anbehav.2021.03.014
612	
613	Neyman, J. (1976). Tests of statistical hypotheses and their use in studies of natural
614	phenomena. Communications in Statistics - Theory and Methods, 5(8), 737–751.
615	https://doi.org/10.1080/03610927608827392
616	
617	Neyman, J., & Pearson, E. S. (1933). On the problem of the most efficient tests of
618	statistical hypotheses. <i>Philosophical Transactions of the Royal Society A: Mathematical</i> ,
619	Physical and Engineering Sciences, 231(694–706), 289–337.
620	https://doi.org/10.1098/rsta.1933.0009
621	N 1 C I W I W W 1 I O D ' I D (2017) D A I ' C
622	Nord, C. L., Valton, V., Wood, J., & Roiser, J. P. (2017). Power-up: A reanalysis of
623	"Power Failure" in neuroscience using mixture modeling. <i>Journal of Neuroscience</i> , 37(34),
624	8051–8061. https://doi.org/10.1523/JNEUROSCI.3592-16.2017



625	
626	O'Donoghue, E. M., Broschard, M. B., & Wasserman, E. A. (2020). Pigeons exhibit
627	flexibility but not rule formation in dimensional learning, stimulus generalization, and task
628	switching. Journal of Experimental Psychology: Animal Learning and Cognition, 46(2), 107-
629	123. https://doi.org/10.1037/xan0000234
630	
631	Paijmans, K. C., Booth, D. J., & Wong, M. Y. L. (2021). Odd one in: Oddity within
632	mixed-species shoals does not affect shoal preference by vagrant tropical damselfish in the
633	presence or absence of a predator. Ethology, 127(2), 125-134. https://doi.org/10.1111/eth.13110
634	
635	Pereira, F. C., Teixeira, D. L., Boyle, L. A., Pinheiro Machado Filho, L. C., Williams, S.
636	R. O., & Enriquez-Hidalgo, D. (2021). The equipment used in the SF6 technique to estimate
637	methane emissions has no major effect on dairy cow behavior. Frontiers in Veterinary Science,
638	7, 620810. https://doi.org/10.3389/fvets.2020.620810
639	
640	Piefke, T. J., Bonnell, T. R., DeOliveira, G. M., Border, S. E., & Dijkstra, P. D. (2021).
641	Social network stability is impacted by removing a dominant male in replicate dominance
642	hierarchies of a cichlid fish. Animal Behaviour, 175, 7–20.
643	https://doi.org/10.1016/j.anbehav.2021.02.012
644	
645	Pinto, P., & Hirata, S. (2020). Does size matter? Examining the possible mechanisms of
646	multi-stallion groups in horse societies. Behavioural Processes, 181, 104277.
647	https://doi.org/10.1016/j.beproc.2020.104277
648	
649	Ribes-Iñesta, E., Hernández, V., & Serrano, M. (2020). Temporal contingencies are
650	dependent on space location: Distal and proximal concurrent water schedules. Behavioural
651	Processes, 181, 104256. https://doi.org/10.1016/j.beproc.2020.104256
652	
653	Rose, E. M., Mathew, T., Coss, D. A., Lohr, B., & Omland, K. E. (2018). A new
654	statistical method to test equivalence: An application in male and female eastern bluebird song.
655	Animal Behaviour, 145, 77-85. https://doi.org/10.1016/j.anbehav.2018.09.004
656	
657	Rouder, J. N., Speckman, P. L., Sun, D., Morey, R. D., & Iverson, G. (2009). Bayesian t
658	tests for accepting and rejecting the null hypothesis. Psychonomic Bulletin & Review, 16(2),
659	225–237. https://doi.org/10.3758/PBR.16.2.225
660	
661	Schino, G., Boggiani, L., Mortelliti, A., Pinzaglia, M., & Addessi, E. (2021). Testing the
662	two sides of indirect reciprocity in tufted capuchin monkeys. Behavioural Processes, 182,
663	104290. https://doi.org/10.1016/j.beproc.2020.104290
664	





665	Seitz, B. M., Flaim, M. E., & Blaisdell, A. P. (2020). Evidence that novel flavors
666	unconditionally suppress weight gain in the absence of flavor-calorie associations. Learning &
667	Behavior, 48(3), 351-363. https://doi.org/10.3758/s13420-020-00419-4
668	
669	Stevens, A., Doneley, R., Cogny, A., & Phillips, C. J. C. (2021). The effects of
670	environmental enrichment on the behaviour of cockatiels (Nymphicus hollandicus) in aviaries.
671	Applied Animal Behaviour Science, 235, 105154.
672	https://doi.org/10.1016/j.applanim.2020.105154
673	
674	Vadillo, M. A., Konstantinidis, E., & Shanks, D. R. (2016). Underpowered samples, false
675	negatives, and unconscious learning. Psychonomic Bulletin & Review, 23(1), 87-102.
676	https://doi.org/10.3758/s13423-015-0892-6
677	
678	Vernouillet, A., Clary, D., & Kelly, D. M. (2021). Highly social pinyon jays, but not less
679	social Clark's nutcrackers, modify their food-storing behaviour when observed by a
680	heterospecific. BioRxiv, 2021.02.28.433225. https://doi.org/10.1101/2021.02.28.433225
681	
682	Wu, Y., Petrosky, A. L., Hazzi, N. A., Woodward, R. L., & Sandoval, L. (2021). The role
683	of learning, acoustic similarity and phylogenetic relatedness in the recognition of distress calls in
684	birds. Animal Behaviour, 175, 111–121. https://doi.org/10.1016/j.anbehav.2021.02.015
685	
686	Yang, C., Tsedan, G., Fan, Q., Wang, S., Wang, Z., Chang, S., & Hou, F. (2021).
687	Behavioral patterns of yaks (Bos grunniens) grazing on alpine shrub meadows of the Qinghai-
688	Tibetan Plateau. Applied Animal Behaviour Science, 234, 105182.
689	https://doi.org/10.1016/j.applanim.2020.105182
690	



Table 1(on next page)

Sources of articles containing negative results in their abstracts



Source	N articles
Animal Behaviour	13
Animal Behavior and Cognition	14
Animal Cognition	17
Animals	15
Applied Animal Behaviour Science	15
Behaviour	14
Behavioural Processes	15
Ethology	16
Frontiers in Psychology: Comparative Psychology	14
Frontiers in Veterinary Science: Animal Behaviour and Welfare	15
International Journal of Comparative Psychology	13
Journal of Applied Animal Welfare Science	15
Journal of Comparative Psychology	15
Journal of Ethology	15
Journal of Experimental Psychology: Animal Learning and Cognition	16
Journal of Zoo and Aquarium Research	15
Learning and Behavior	15
PeerJ: Animal Behaviour	15
bioRxiv: Animal Behaviour and Cognition	14
PCI: Animal Science	2



Table 2(on next page)

Example categorization of sample-level statements

Category	Non-Significant	No Effect	Ambiguous, Similar, or Small Effect Size
Description	Reports that there was no <i>significant</i> difference between two conditions, or words to that effect.	A statement that there was not a difference within the sample, when in fact there was – it was just not significant in their analysis.	A statement about the results that neither suggests they were the same, nor that there was no significant difference.
	There was no significant/detectable difference between X and Y.	There was no difference between X and Y.  There was no effect.	X and Y were similar.
	We did not detect a difference between X and Y (or any other statement implying	There was no evidence of an effect.  There was no relationship between X	There was no large/clear difference between X and Y.
Examples	failing to find a signal within noise).  We did not find a significant effect.  X was not significantly	and Y.  We did not find/observe/see a difference between X and Y.	There was no large effect of X on Y.
	related to Y.  X did not perform significantly above chance.	We did not find an effect.  We found no evidence of an effect.	
	X performed significantly above chance, but Y did not.	X performed at chance levels.	
	There were no significant differences between X and Y's performance.		

2 .



Table 3(on next page)

Example categorization of population-level or title claims

Category	Justified	Caveated, Ambiguous or Similar	No Effect
Description	Comments on statistical power, uses equivalence tests or otherwise justifies why a non-significant result suggests that there is no theoretically important difference in the population, or that the study provides no strong evidence of a difference.	Interprets the non- significant results as suggesting/indicating etc. that X and Y do not differ in the population, or are similar.	Interprets the non- significant result as showing that X and Y do not differ in the population.
Examples	Because the test was high-powered to detect a meaningful difference, this nonsignificant result suggests that A is not related to Y in a theoretically important way.  In addition to being not statistically different to each other, X and Y were also statistically equivalent (if a frequentist equivalence or non-inferiority test was performed), suggesting that X is not meaningfully related to Y.	suggesting that X is not related to Yindicating that X is not related to Ysuggesting/indicating that there is no difference between X and Ysuggesting that X has not changed Y. Our results provide no strong evidence that X and Y are differentsuggesting that X and Y are similar.	meaning that X is not related to Y. showing that X is not related to Y.  There is no difference between X and Y.  X and Y do not differ.  X and Y are similar.  X and Y are the same (show the same effect, etc).  X does not change Y.  Our results provide no evidence that X and Y are different.







### Table 4(on next page)

Examples of claims in the titles of articles following non-significant NHST classified as "No Effect" and "Caveated or Similar"

#### No Effect

N = 37 (84%)

"Home range use in the West Australian seahorse Hippocampus subelongatus is influenced by sex and partner's home range but not by body size or paired status"

Kvarnemo et al., 2021

"Delays to food-predictive stimuli do not affect suboptimal choice in rats."

Cunningham & Shahan, 2020

"Common Marmosets (*Callithrix jacchus*) Evaluate Third-Party Social Interactions of Human Actors But Japanese Monkeys (*Macaca fuscata*) Do Not"

Kawai et al., 2019

#### Caveated, Ambiguous, or Similar

N = 7 (16%)

"Limited Evidence of Number-Space Mapping in Rhesus Monkeys (Macaca mulatta) and Capuchin Monkeys (Sapajus apella)"

Beran et al., 2019

"Little Difference in Milk Fatty Acid and Terpene Composition Among Three Contrasting Dairy Breeds When Grazing a Biodiverse Mountain Pasture"

Koczura et al., 2021

"The Equipment Used in the SF6 Technique to Estimate Methane Emissions Has No Major Effect on Dairy Cow Behavior"

Pereira et al., 2021







### Table 5(on next page)

Examples of claims about the sample in the abstracts of papers following non-significant NHST classified as "No Effect", "Similar or Small Effect Size", "Non-Significant" or "Ambiguous"

#### No Effect

N = 174, 63%

"Levels of individuals sitting with their back to the window was unaffected by visitor number or noise."

Hashmi & Sullivan, 2020

"The groups did not differ in their ability to follow human signals"

Lazarowski et al., 2020

#### **Similar or Small Effect Size**

N = 17, 6%

"Pair members demonstrated comparable responses towards a male 'intruder', as latency to respond and proximity scores were very similar between pair members in the majority of pairs examined"

DeVries et al., 2020

"We found that individuals called back to sympatric and allopatric calls within similar amounts of time,"

Wu et al., 2021

#### Non-Significant

N = 71, 26%

"Nutcrackers... did not significantly change their caching behaviour when observed by a pinyon jay."

Vernouillet et al., 2021

"No significant correlations between degree of laterality and behavioral interest in the stimuli were found"



Lilley et al., 2020
Ambiguous
N = 16 (6%)
"We also found no conclusive evidence that either the visual or the vibratory sensory modalities are critical for prey capture."
Meza et al., 2021
"No systematic variations on space allocation were observed in neither experiment"
Ribes-Iñesta et al., 2020



### Table 6(on next page)

Examples of claims about populations in the abstracts of papers following non-significant NHST classified as "No Effect" and "Caveated, Ambiguous or Similar"

#### No Effect

N = 18 (29%)

"Partial rewarding does not improve training efficacy"

Cimarelli et al., 2021

"Our findings show that H. horridum does not respond to hypoxic environments"

Guadarrama et al., 2020

"Oviposition site choice is not by-product of escape response"

Kawaguchi & Kuriwada, 2020

#### Caveated, Ambiguous, or Similar

N = 45 (71%)

"These results suggest capuchin monkeys do not engage in indirect reciprocity"

Schino et al., 2021

"These results suggest that shoal composition may not be an important driver of shoal choice in this system"

Paijmans et al., 2021

"...suggesting that size is not a determinant factor for feral horse society."

Pinto & Hirata, 2020



### Table 7(on next page)

Examples of statement reporting the results in the results sections of papers using non-significant NHST classified as "No Effect", "Similar or Small Effect Size", "Non-Significant" or "Ambiguous"

#### No Effect

N = 113 (41%)

During farrowing, No Effect of the treatments was seen on the percentage of time spent (3.22 % vs. 1.90 %, P = 0.372) on the nest-building behaviour"

Aparecida Martins et al., 2021

"There were no differences between treatments in the frequency or duration of birds flying between walls"

Stevens et al., 2021

#### **Similar or Small Effect Size**

 $\overline{N} = \overline{12} (4\%)$ 

"The average time yaks spent grazing was similar among shrub coverage groups (P = 0.663)"

Yang et al., 2021

"The number of sessions required to reach criterion didn't reliably differ between groups"

O'Donoghue et al., 2020

#### Non-Significant

N = 140 (52%)

"Comparing the pooled data of all crows, no significant increase in the number of mark-directed behaviors during the mirror mark condition was found compared with the no-mirror sham condition."

Brecht et al., 2020

"There was no significant effect of removal type on changes in display strength in either dominant males or subordinate males."



Piefke et al., 2021	
Ambiguous	
N = 10 (4%)	

"As can be seen in Figure 1D, there was no difference in response rates after R and NR trials across days for rats under reward uncertainty." [where in Figure 1D the bars on the graph look almost identical)

Anselme & Robinson, 2019

"It showed that there was a significant main effect of session, but no main effect of CS"

Harris & Bouton, 2020

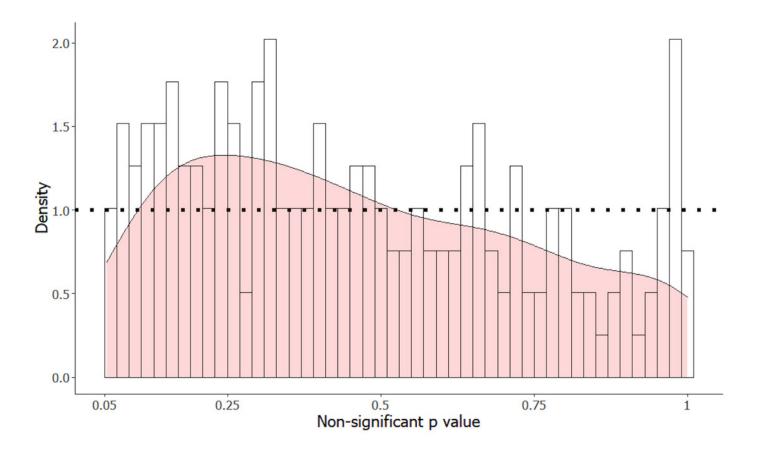
2



# Figure 1

Distribution of non-significant p-values from result sections of 198 articles in animal cognition and related fields, with a density distribution overlaid in pink.

The dotted line shows the average density.





## Figure 2

The observed p-value distribution of 198 p-values > .05, compared to 3 simulated distributions where 80% of alternative hypotheses were correct.

The observed p-value distribution was manually extracted from results corresponding to negative claims present in the abstracts of animal cognition articles. The observed p-value distribution was compared to 3 simulated distributions where 80% of alternative hypotheses were correct, with studies performed at either 10%, 33%, 50% or 80% statistical power.

