



Lagisz, M., Zidar, J., Nakagawa, S., Neville, V. M., Sorato, E., Paul, E. S., Bateson, M., Mendl, M. T., & Løvlie, H. (2020). Optimism, pessimism and judgement bias in animals: a systematic review and meta-analysis. *Neuroscience and Biobehavioral Reviews*, 118, 3-17. <https://doi.org/10.1016/j.neubiorev.2020.07.012>

Peer reviewed version

License (if available):
CC BY-NC-ND

Link to published version (if available):
[10.1016/j.neubiorev.2020.07.012](https://doi.org/10.1016/j.neubiorev.2020.07.012)

[Link to publication record in Explore Bristol Research](#)
PDF-document

This is the author accepted manuscript (AAM). The final published version (version of record) is available online via Elsevier at <https://doi.org/10.1016/j.neubiorev.2020.07.012>. Please refer to any applicable terms of use of the publisher.

University of Bristol - Explore Bristol Research

General rights

This document is made available in accordance with publisher policies. Please cite only the published version using the reference above. Full terms of use are available: <http://www.bristol.ac.uk/red/research-policy/pure/user-guides/ebr-terms/>

1 **Optimism, pessimism and judgement bias in animals: a systematic**
2 **review and meta-analysis**

3

4 Lagisz, Malgorzata ^{1†}, Zidar, Josefina^{2†}, Nakagawa, Shinichi^{1,†*}, Neville, Vikki³, Sorato, Enrico²,
5 Paul, Elizabeth S.³, Bateson, Melissa⁴, Mendl, Michael^{3#}, Løvlie, Hanne^{2#}

6

7 **Addresses:**

8 ¹ Evolution and Ecology Research Centre, School of Biological, Earth and Environmental
9 Sciences, University of New South Wales, Sydney, New South Wales, Sydney, NSW 2052,
10 Australia

11 ² The Department of Physics, Chemistry and Biology, IFM Biology, Linköping University, SE-581
12 83 Linköping, Sweden

13 ³ Centre for Behavioural Biology, Bristol Veterinary School, University of Bristol, Langford, BS40
14 5DU, United Kingdom

15 ⁴ Centre for Behaviour and Evolution, Biosciences Institute, Newcastle University, Newcastle
16 upon Tyne, NE2 4HH, United Kingdom

17

18 ***Correspondence:** Shinichi Nakagawa s.nakagawa@unsw.edu.au and Michael Mendl

19 Mike.Mendl@bristol.ac.uk

20

21 † These authors contributed equally to this work

22 # These authors supervised this work equally and are joint senior authors

23

24 **Author contributions:** two groups of the authors, SN & HL and MM, EP & MB conceived the idea
25 independently, SN developed study design and methods from the inputs from others. JZ, ES, VN

1 and EP collected the data with inputs from SN, ML, MB, MM and HL. ML, JZ and SN conducted the
2 analysis with inputs from MM. JZ, ML, SN, VN and MM co-wrote the first draft and all contributed
3 to revisions of the manuscript.

4
5 **Declarations of interest:** none

6
7 **Abstract**

8 Just as happy people see the proverbial glass as half-full, 'optimistic' or 'pessimistic' responses to
9 ambiguity might also reflect affective states in animals. Judgement bias tests, designed to
10 measure these responses, are an increasingly popular way of assessing animal affect and there is
11 now a substantial, but heterogeneous, literature on their use across different species, affect
12 manipulations, and study designs. By conducting a systematic review and meta-analysis of 459
13 effect sizes from 71 studies of non-pharmacological affect manipulations on 22 non-human
14 species, we show that animals in relatively better conditions, assumed to generate more positive
15 affect, show more 'optimistic' judgements of ambiguity than those in relatively worse conditions.
16 Overall effects are small when considering responses to all cues, but become more pronounced
17 when non-ambiguous training cues are excluded from analyses or when focusing only on the
18 most divergent responses between treatment groups. Task type (go/no-go; go/go active choice),
19 training cue reinforcement (reward-punishment; reward-null; reward-reward) and sex of
20 animals emerge as potential moderators of effect sizes in judgement bias tests.

21 **Keywords:** research synthesis, affective state, cognitive bias, animal welfare

1 **Introduction**

2 Accurate assessment of affect (emotion) in non-human animals is an important goal in
3 disciplines including animal welfare science, neuroscience, psychopharmacology and drug
4 development. A prevailing view in the study of human emotion is that affective states comprise
5 subjective, behavioural, neural and physiological components (Paul et al., 2020; Scherer, 2005).
6 Whilst the subjective component of animal affective states (feelings) is not currently accessible
7 to direct measurement and we cannot be certain which species consciously experience such
8 states (see Paul et al., 2020), we can objectively assess the other components. In his book *The*
9 *Expression of Emotions in Man and Animals*, Darwin (1872) focused on behavioural
10 manifestations of animal emotion, namely “expressive movements of the face and body”, and
11 such measures continue to be used as indicators of animal affect today (e.g. Girard & Bellone,
12 2020). But other measures focus more directly on the role of affect in behavioural control and
13 decision-making. A relatively new and promising approach is to measure biases in decision-
14 making under ambiguity as indicators of animal affect (Harding et al., 2004; Mendl et al., 2009).
15 This is because there are empirical and theoretical reasons to expect that responses to such
16 ambiguity reflect affective valence (positivity or negativity of an affective state). For example,
17 people in negative states are more likely to make negative (‘pessimistic’) judgements about
18 ambiguous events or stimuli than people in more positive states (Blanchette and Richards, 2010;
19 Paul et al., 2005). Such assessments could reflect an adaptive use of background affect (or mood)
20 as a Bayesian prior over the likelihood of future positive or negative outcomes (Mendl et al.,
21 2010; see Mendl & Paul, 2020 for a fuller discussion).

22 In line with these findings and ideas, a generic assay for measuring these so-called ‘*judgement*
23 *biases*’ has been developed for animals and has now been used in a large number of studies
24 across a range of species. The original assay (Harding et al., 2004) involves training subjects to
25 make one response (positive response) to a ‘positive’ cue (a single frequency tone) in order to

1 achieve a positive outcome (e.g. food) and a different response (negative response) to a
2 'negative' cue (a tone of a different frequency) in order to avoid a negative outcome (e.g. white
3 noise) (Figure 1a). Once subjects have learnt this conditional discrimination task, training
4 continues but includes occasional ambiguous cues (tones of intermediate frequency) designed to
5 assess whether subjects would make the positive response indicating anticipation of a positive
6 outcome, or the negative response indicating anticipation of a negative outcome. This allows one
7 to test whether, for example, animals in a putative negative affective state (e.g. as a result of
8 some sort of experimental treatment, Figure 1b) are more likely to make the negative response,
9 as predicted (Figure 1c,d). Making the positive or negative response under ambiguity can be
10 operationally defined as 'optimistic' or 'pessimistic' (Bateson, 2016) without implying that
11 animals experience optimism or pessimism as humans do.

12 Published studies using this judgement bias task (also referred to as an 'ambiguous cue
13 interpretation' task (Rygula et al., 2013)), and variants of it, have supported the general
14 prediction, but also generated null and opposite results. These findings have been summarised
15 narratively in a number of review papers that have also identified various methodological and
16 theoretical questions regarding the task and approach (Baciadonna and McElligott, 2015;
17 Bethell, 2015; Gyax, 2014; Hales et al., 2014; Mendl et al., 2009; Mendl & Paul 2020; Roelofs et
18 al., 2016). What has been lacking, and much needed, is a systematic review and meta-analysis of
19 the findings to date to evaluate whether the general predictions behind the approach are
20 supported, and how results may be influenced by a variety of moderators, including aspects of
21 task design, methods used to manipulate affect, species studied, and age and sex of subjects.
22 Recently, we published the first such meta-analysis focusing on the effects of pharmacological
23 manipulations of affective state on judgement biases (Neville et al., 2020). Here we
24 systematically review and meta-analyse the much larger number of studies that have used non-
25 pharmacological affect manipulations.

1 We focus on judgement bias tasks based on the Harding et al. (2004) method, since these have
2 been more widely studied in animals than other cognitive biases such as attention (Bethell et al.,
3 2016; Crump et al., 2018) and memory biases (Burman and Mendl, 2018). Although details of
4 the procedures and criteria used to select suitable studies and extract appropriate data for the
5 meta-analysis are explained in the Methods section, three points should be noted here.

6 First, a major challenge in any study of animal affect is to establish a 'ground truth' for the
7 affective state that the animal is in when under study. This is necessary, for example, if an aim of
8 a study is to determine what behavioural, physiological or neural changes occur in animals in
9 particular affective states, and hence to develop reliable indicators of such states. Therefore, in
10 studies which seek to evaluate whether judgement bias is a valid indicator of affective valence,
11 we need to know whether the animal is in a relatively positive or negative state, so that we can
12 test whether animals that are in a more positive state do indeed show more optimistic decisions
13 under ambiguity, than those in a more negative state. In most judgement bias studies,
14 researchers attempt to use an experimental treatment to induce a relatively positive or negative
15 affective state compared to a control or 'benign' treatment group, or they impose both a positive
16 and a negative treatment, and compare these. Because we cannot know for certain where the
17 intermediate 'neutral' state lies, we use terminology that emphasises the relative nature of these
18 manipulations. Thus, we refer to 'better' (more positive), 'benign/control', and 'worse' (more
19 negative) treatments, and assign them to either 'relatively better' or 'relatively worse' groups for
20 pair-wise comparison in the meta-analysis.

21 Second, there are two main types of task used in judgement bias trials: active choice (go/go) and
22 go/no-go (Figure 1a). In go/go active choice tasks, the animal has to choose between two
23 alternative responses (e.g. press the left or right lever), while in go/no-go tasks the animal's
24 options are to perform a response (e.g. approach a location or press a lever) or suppress it. The
25 response of animals can be reported as a proportion (e.g. proportion of trials in which the

1 subject pressed the left lever or approached the location), or a latency (e.g. time taken to press
2 the lever or approach the location). Latency and proportion data have different statistical
3 distributions; they require different transformations and use of different formulae to calculate
4 effect sizes. There are also biological reasons for separating latency and proportion data.
5 Different measures may represent different aspects of cognitive processes and their utility
6 depends on the type of cognitive bias task used. In go/no-go tasks, latency to perform the
7 response under ambiguity is a direct measure of judgement bias. For example, if the positive
8 response is to approach the cue, then quick approach to an ambiguous cue indicates an
9 'optimistic' response. In contrast, go/go active choice tasks require responses to both cues (e.g.
10 press left or right lever), meaning that the latency to perform whichever response the animal
11 selects is more difficult to interpret in terms of 'optimism' or 'pessimism'. Rather, the proportion
12 of positive or negative responses provides more definitive information about 'optimistic' or
13 'pessimistic' decisions. This measure is also of use in go/no-go tasks. Therefore, proportion of
14 positive vs. negative responses is preferable to latency as a measure of judgment bias for go/go
15 active choice tasks, whereas for the go/no-go tasks both measures are, in principle, suitable.

16 Third, many judgement bias studies use more than one ambiguous cue during test trials (Figure
17 1c). Often three such cues are used; one (MID) which is assumed to be perceived by the animal
18 as being at the mid-point of the sensory scale (e.g. sound frequency) between the positive (P)
19 and negative (N) training cues, one (near positive: NP) which is half way between MID and the
20 positive (P) cue, and one (near negative: NN) which is halfway between MID and the negative
21 (N) cue. There are theoretical and methodological reasons for why an affect manipulation
22 treatment might have an influence at one ambiguous cue but not at others in the same study. For
23 example, non-midpoint ambiguous cues (NP, NN) may be perceptually too similar to the P and N
24 training cues for animals to moderate their responses to them, whilst the midpoint (MID) cue is
25 usually ambiguous enough for background affect to influence responses to it. In some studies

1 MID could be perceived as closer to P or N and the most ambiguous cue becomes either NN or
2 NP, respectively. Moreover, the perceived payoff of the positive and negative response
3 outcomes, and hence associated decisions, may be asymmetrical. For example, if the perceived
4 negative value of a foot-shock outcome is much stronger than the perceived positive value of a
5 food pellet, animals may be strongly motivated to avoid shock risk and thus respond negatively
6 to both MID and NN ambiguous cues, with variation in response limited to the 'safest' NP
7 ambiguous cue (Mendl et al., 2009). Conversely, in a test variant where negative cues are simply
8 lacking a reward instead of bearing a punishment, animals may respond positively even to
9 negative cues, because the cost of doing so is negligible. Because it is likely that biased responses
10 are unevenly spread across ambiguous cues – in fact some studies report effects which are
11 strongest or only statistically significant at one ambiguous cue location (e.g. Bethell and Koyama,
12 2015; Zidar et al., 2018) – we investigate the effect of relative cue position and also conduct
13 sensitivity analyses. These additional analyses use data subsets with different decision rules for
14 selecting the most representative data points from response curves (e.g. using only the
15 ambiguous cue with the largest absolute between-treatment effect size; more details in the
16 Methods section).

17 This systematic review and meta-analysis aims to: (i) quantify the overall effect size that affect
18 manipulations have on measures of judgement bias in animals; (ii) estimate heterogeneity of the
19 results among different studies; (iii) explore the influences of different biological and
20 methodological moderators (explanatory variables for variation in effect sizes).

21 **Methods**

22 *Literature search*

23 We conducted a systematic literature search and recorded relevant information required in the

1 Preferred Reporting Items for Systematic reviews and Meta-Analyses statement (PRISMA;
2 Moher et al., 2009; see Supplementary Materials for additional search details). We ran the first
3 online database search on 29 October 2015, a second search in December 2017, and a final
4 database search to update the dataset again on 27 March 2019. For these searches, we used the
5 broad-coverage interdisciplinary databases *Scopus* and *Web of Science*, covering the titles,
6 abstracts and keywords of academic publications.

7 The initial search string used in *Scopus* was: TITLE-ABS-KEY (("cognitive bias*" OR "judgment
8 bias*" OR "judgement bias*" OR "cognitive affective bias*") AND (pessimis* OR optimis* OR
9 valence OR mood* OR emotion* OR "affective state*" OR "emotional state*" ambig* OR animal*
10 OR "animal welfare")) AND PUBYEAR > 2003 and in *Web of Science*: TS=(("cognitive bias*" OR
11 "judgment bias*" OR "judgement bias*" OR "cognitive affective bias*") AND (pessimis* OR
12 optimis* OR valence OR mood* OR emotion* OR "affective state*" OR "emotional state*" ambig*
13 OR animal* OR "animal welfare")) AND LANGUAGE: (English) AND DOCUMENT TYPES:
14 (Article), Indexes=SCI-EXPANDED, SSCI Timespan=2004-2015. We restricted the publication
15 years to those following the seminal paper on animal judgement bias (Harding et al., 2004). We
16 restricted the subsequent updates of the literature search to the years since the previous search
17 update (i.e. 2015-2017 and 2017-2019, respectively) and otherwise used the same search
18 strings. We collected additional relevant studies from the authors whom we contacted to
19 request data or other additional information that was missing from their publications. We also
20 performed searches of reference lists of relevant review articles and research articles citing the
21 seminal study by Harding et al. (2004).

22 The searches of the online databases generated over 900 potential article references and
23 searches of other sources generated almost 500 additional references for screening (Figure 2).
24 We removed duplicated results from these separate search paths. Two authors (J.Z. and E.S.)
25 independently screened 482 abstracts from the articles identified in the 2015 search using the

1 software AbstrackR (Wallace et al., 2012). M.L. performed two updates of literature searches in
2 2017 and 2019, following the same methodology as in the first search. Overall, we identified 74
3 published studies as potentially suitable for inclusion in our meta-analysis after screening of full
4 texts and removal of duplicated studies. We excluded three studies during the data extraction
5 stage (due to missing data), resulting in data from 71 studies being included in the meta-
6 analysis.

7 *Inclusion and exclusion criteria*

8 We screened titles and abstracts from bibliometric records to identify empirical studies on
9 judgement bias in animals in which subjects were exposed to an affect manipulation aimed at
10 inducing either a relatively positive or negative state. We then screened full text versions of the
11 articles that passed this initial screening stage. At the full-text screening stage, the following six
12 criteria had to be met for the study to be included in the meta-analysis: i) study had to be
13 experimental and designed to investigate variation in judgement bias (i.e. 'optimistic' or
14 'pessimistic' interpretation of stimuli) in non-human animals; ii) experiments had to include at
15 least two treatment groups (or control/'benign' and treatment groups); iii) experimental
16 treatments had to be designed to induce 'relatively better' or 'relatively worse' affective states
17 (see decision-tree in Supplementary Materials Figure S1); iv) for go/no-go tasks studies had to
18 report either latency to make a response to ambiguous cues, or proportion of go or no-go
19 responses towards ambiguous cues; for active choice tasks, studies had to report proportion of
20 positive or negative responses; if the data available could be translated into such latencies or
21 proportions, they were included; v) studies had to present data usable for effect size calculation;
22 if suitable data could not be retrieved by contacting the authors, the study was excluded from
23 the meta-analysis; vi) studies had to be published in peer-reviewed journals, but student reports
24 and data from unpublished work, as well as articles that were written in languages other than
25 English, could have been included if they met the above criteria.

1 We also excluded studies for the following additional reasons. We only considered data from
2 studies investigating judgement bias, i.e. we excluded studies investigating other cognitive
3 biases, such as attention bias and memory bias. We also excluded studies only describing
4 judgement bias theory or methods or reviewing previous findings and studies that used the
5 generic judgement bias task for humans, because our focus was on non-human animals. As
6 studies investigating effects of drugs on judgement bias often include several doses that cannot
7 easily be assigned into relatively better and relatively worse treatment groups, we also excluded
8 all drug studies from this meta-analysis. As mentioned earlier, the drug studies were recently
9 subjected to a separate meta-analysis by our group (Neville et al., 2020).

10 *Data extraction*

11 After compiling a final list of included studies, we extracted measurements representing
12 behavioural responses to cues in the judgement bias tests. Each pairwise comparison consisted
13 of a pair of outcome measures comparing behaviour of animals from 'relatively better' to
14 'relatively worse' affect manipulation groups. Our classification of treatments as inducing
15 'relatively better' or 'relatively worse' affective states was based on a decision tree involving
16 screening articles and assessing treatments based on the following three criteria. First, if stated,
17 we used the *a priori* hypothesis and reasoning outlined in the research article. Second, where
18 possible, we employed Rolls' (2005, p.11) operational definition of emotion as "states elicited by
19 rewards and punishers", where "a reward is anything for which an animal will work" and "a
20 punisher is anything that an animal will work to escape or avoid". Thus, if a treatment involved
21 stimuli that the subject animal is known to actively avoid, we deemed it to induce a relatively
22 worse affective state than one which involved neutral or preferred stimuli. Third, we considered
23 evidence from previous studies on the effects of the treatments in question on affective state
24 (e.g. their effects on other putative indicators of affective state, such as abnormal repetitive
25 behaviour or physiological stress indicators).

1 The decision tree for assigning affect treatments to relative affect manipulation categories is
2 presented in Supplementary Materials Figure S1. If the first criterion in the decision tree was
3 fulfilled (i.e. the authors of the original paper explicitly stated whether the treatment is expected
4 to have positive/negative effect on animals' affective state), we ignored the subsequent decision
5 criteria. If not, we evaluated the subsequent decision criteria. We classified all extracted
6 treatment groups within a study relative to each other. For example, in a study with a control
7 (benign/unmanipulated) and enriched housing group, the enriched group would be considered
8 'relatively better' and the control/benign group 'relatively worse'. Conversely, in a study with
9 control/benign and stress-induction groups, the stress group would be considered 'relatively
10 worse', and the control/benign group would be considered 'relatively better'.

11 We tackled variation in study design and outcome measurement as follows. First, for the go/no-
12 go judgement bias tasks we extracted either or both (depending on which was reported) latency
13 and proportion outcome measures (the signs of the effect sizes calculated from latency
14 measurements were later inverted, so that interpretation of the effect direction was consistent
15 with that for the proportion data). For active choice go/go judgement bias tasks, we extracted
16 only proportion outcome measures (as explained earlier, latency measures in active choice tasks
17 cannot be clearly linked to more 'optimistic' or 'pessimistic' responding). The extracted mean
18 and standard error (or standard deviation) of responses to ambiguous and non-ambiguous cues
19 during the tests were used to calculate values of effect sizes (and their variances) for each
20 pairwise comparison of the relatively better and relatively worse treatment groups at the same
21 cue. Relevant sample sizes were also recorded representing the number of animals from each
22 group participating in the judgement bias test.

23 Second, included studies used varying numbers of ambiguous cues (range 1-13, mean 2.99,
24 mode 3). We only extracted data for a maximum of three ambiguous cues per measurement
25 (response curve). We always extracted data for the middle cue (midpoint between the positive

1 and negative cues, MID) and, if available, two intermediate cues between the middle cue and
2 positive and negative cues (near-positive NP and near-negative NN, respectively). If response
3 data to positive (P) and negative (N) cues were reported for judgement bias tests, these were
4 also extracted.

5 Third, when judgement bias was measured on several consecutive days following a treatment,
6 we extracted the first measure only as it was usually closest in time to the acute affect
7 manipulation treatment (Destrez et al., 2013; Doyle et al., 2011). In a few studies, animals were
8 exposed to several judgement bias tests during a long-term treatment (Douglas et al., 2012;
9 Hales et al., 2016; Rygula et al., 2013). In these cases, we extracted the last test occurring during
10 each treatment, thus maximising the time available for it to exert its effects. We assumed that
11 the cumulative impact of chronic exposure to the treatment likely out-weighed any potential
12 effect of learning about repeated tests that were spaced out across time.

13 Fourth, some studies with a within-subject design measured judgement bias before, during and
14 after an affect manipulation (e.g. pre-stress, stress, post-stress), or repeated the 'baseline'
15 treatment (e.g. enriched, barren, enriched) (Brilot et al., 2010; da Cunha Nogueira et al., 2015;
16 Hales et al., 2016; Murphy et al., 2013). In these studies, we compared measures taken before
17 treatment ('baseline') to those taken during it and did not include measures taken after it.

18 Fifth, studies using a between-subject design sometimes tested both control/benign and
19 treatment groups before, during and after a manipulation. In these cases, we compared the
20 control/benign group to the treatment group during treatment and ignored the pre- and post-
21 treatment measurements (Hales et al., 2016; Oliveira et al., 2016; Rygula et al., 2013).

22 Finally, if several treatments were applied where one or more treatments were hypothesized to
23 be intermediate in effect to the two most extreme treatments, only the two extreme treatments

1 were included (Ash and Buchanan-Smith, 2016; Burman et al., 2009; Keen et al., 2014; Wheeler
2 et al., 2015).

3 For each experiment, we gathered information on the potential moderator variables to
4 characterise our dataset and explain potential heterogeneity in the data. Detailed descriptions of
5 all the originally extracted moderators are included in Table S1. In brief, the three key groups of
6 extracted moderators considered information about the article, biological variables, and test
7 design. Paper-specific information included authors, title, journal, and publication year. For each
8 data point (i.e. comparison between two groups of animals), we extracted the following
9 biological variables: taxa studied (mammals, birds, insects), sex (female, male, mixed-sex), age
10 class (juvenile, adult) and source of animals (captive, wild-caught). Test-specific information
11 included affect manipulation category (enrichment, stress, other), affect manipulation timing
12 (before/during test, long-term), comparison category (Better-Worse, Benign-Worse, Better-
13 Benign), type of cue used in judgement bias test (spatial, visual, auditory, tactile, olfactory),
14 whether animals were food deprived prior to behavioural trials (yes, no/no information),
15 automation of response measurement (yes, no/no information), blinding of personnel
16 performing trials (yes, no/no information), combination of reinforcement used during training
17 (Reward Vs. Null, Reward Vs. Punishment; Reward Vs. Smaller Reward), task type (active choice
18 go/go, go/no-go), whether ambiguous cues were reinforced (yes, no/no information),
19 measurement type (latency, proportion), location of ambiguous cues relative to positive and
20 negative cues (P – positive, NP – near-positive, MID – midpoint, NN – near-negative, N –
21 negative). We also noted any pertinent additional details about study designs (between-subjects,
22 within-subjects), affect manipulations, source of the data in the original studies, and any
23 associated comments. When data were provided in a graph instead of a table or text, we
24 extracted the values using GraphClick 3.0.3 (<http://www.arizona-software.ch/graphclick/>).
25 Data extraction was performed by J.Z., M.L. and V.N. and was checked by M.L., V.N. and E.S.

1 *Effect-size calculation*

2 We used Hedges' unbiased standardized mean difference (Hedges' g) as the measure of effect
3 size. Because latency and proportion data are bounded (i.e. latencies start at 0 and are often
4 censored, and percentages are bounded between 0 and 100, proportions between 0 and 1), we
5 used natural log (for latencies) or logit-transformed data (for proportions, and percentages
6 expressed as proportions) to calculate Hedges' g (details provided in Supplementary Materials
7 Methods and Figure S2). In brief, to calculate Hedges' g , we focused on positive responses (i.e.
8 those which indicated that the subject was anticipating a more rewarding outcome) and
9 subtracted the mean value of the relatively worse treatment from the mean of the relatively
10 better treatment, and divided the difference by the pooled standard deviation (SD) with
11 correction for small sample sizes (Hedges and Olkin, 1985). Thus, if animals from the relatively
12 better treatment group were making a higher proportion of positive responses than animals
13 from the relatively worse treatment, the difference between the means would be positive and
14 the effect size too. However, the expected pattern would be reversed when latencies to make the
15 positive response were measured in go/no-go tasks: if animals from the relatively better
16 treatment group were quicker to make the positive response (i.e. had lower latencies), than
17 animals from the relatively worse treatment, the difference between the means (and the
18 resulting effect size) would be negative. To allow for easier comparison and interpretation of the
19 effect sizes from latency and proportion measures, we reversed the sign of the effect sizes based
20 on latency measures. Thus, after the sign adjustment, across all data positive values of Hedges' g
21 can be interpreted as optimistic responses of animals exposed to relatively better treatments
22 compared to those exposed to relatively worse treatments. For the go/no-go tests that reported
23 the outcomes as both latency and proportion, we calculated Pearson's correlation between these
24 two measures.

25 *Meta-analysis and meta-regression models*

1 We ran all statistical analyses in *R* version 3.6.0 (R Development Core Team, 2019); we created
2 main forest-like (orchard) plots of effects using *orchaRd* package (Nakagawa et al., 2020). For
3 multilevel meta-analysis and meta-regression we used the *rma.mv* function from the package
4 *metafor* (Viechtbauer, 2010).

5 To estimate the overall mean of the effect sizes we constructed intercept-only models (i.e. meta-
6 analysis) with study ID, experiment ID, cue ID, and effect size ID as random effects. To explore
7 effect of species identity and phylogenetic relatedness, we also evaluated meta-analytic models
8 with phylogeny and species ID added to the random effects list. We calculated I^2 values for each
9 random factor and the overall heterogeneity, I^2_{Total} , in the meta-analytic models (Nakagawa and
10 Santos, 2012).

11 To evaluate the effects of moderators of interest (e.g. subject sex or age class, test task type, test
12 cue type and level of cue ambiguity), we ran univariate multilevel phylogenetic meta-regression
13 models with moderators as fixed effects, and the same random effects as in the meta-analytic
14 models (except species ID). In the multivariate meta-regression models (i.e. models with
15 multiple moderators), we included only moderators that were significant in the univariate meta-
16 regression models. We then performed AICc-based model selection using MuMIn package
17 (Barton, 2009) to infer relative contributions of included moderators. To assess the fit of meta-
18 regression models, we calculated marginal R^2 values (*sensu* Nakagawa and Schielzeth, 2013;
19 Nakagawa et al., 2017).

20 *Publication bias*

21 Statistically significant results are more likely to be published, resulting in a non-random sample
22 of data available for meta-analysis (Rosenthal, 1979). To examine publication bias in our data
23 set, we visually inspected a funnel plot for asymmetry in the distribution of the residuals of
24 effect sizes (which are the sum of effect size level effects and sampling variance effects; i.e. meta-

1 analytic residuals: *sensu* Nakagawa and Santos, 2012). We also performed Egger's regression on
2 the residuals and measurement errors from the full meta-regression model (multilevel version
3 of the publication bias test; Nakagawa and Santos, 2012). Egger's regression indicates
4 publication bias if the regression intercept is significantly different from 0 (Egger et al., 1997).
5 Finally, we tested for a special type of publication bias, a time-lag bias, i.e. a tendency for studies
6 with larger effects to be published earlier (Jennions and Møller, 2002).

7 *Sensitivity analyses (robustness of results)*

8 To test robustness of our results to the estimation method, we ran a meta-regression model and
9 a multilevel mixed-effect full meta-regression model (with subject sex, task type, cue type, and
10 reinforcement type as moderators), using a Bayesian approach, as implemented in the
11 *MCMCglmm* package (Hadfield, 2010). These models were run with 110,000 iterations, 10,000
12 burn-in periods, and thinning by every 100 resulting in an effective sample size of 1000. We
13 used a parameter-expanded prior ($V = 1$, $nu = 1$, $alpha.mu = 0$, $alpha.V = 1000$), with EffectID
14 (units) fixed at one.

15 We also ran the meta-analytic models using four additional data configurations representing
16 different ways of interpreting results from pairs of response curves with multiple cues tested.
17 First, we used a dataset with positive and negative test cues excluded, so that only responses to
18 ambiguous cues were used (maximum of 3 effect sizes per comparison of pair of response
19 curves: for near-positive, midpoint, near-negative cues). In the remaining data subsets, we
20 selected only one cue per response curve comparison. Thus, to create the second data subset, we
21 only included data from the mid-point ambiguous cue location (MID data points and effect
22 sizes). In the third data subset, we selected the effect sizes data from the cue location with the
23 largest absolute value within each response curve comparison; notably, in 71.3% of the
24 comparisons, the largest absolute effect size was not located at the mid-point ambiguous cue. In

- 1 the fourth data subset, we used effect sizes with the biggest absolute value in the direction of the
- 2 mean value, within each response curve comparison, as in Neville et al. (2020).
- 3

1 **Results**

2 *Description of data set*

3 The workflow and outcomes of our systematic literature searches are presented in a PRISMA
4 diagram (Figure 2). The list of included studies is provided in Supplementary Table S2. Excluded
5 studies, with reasons for exclusion, are listed in Supplementary Table S3. To retrieve missing
6 data, or additional information, we contacted 39 authors about 35 studies. We attained raw data
7 for 18 studies and additional information for 10 studies. Ultimately, we extracted 459 effect
8 sizes, representing 91 experiments published in 71 articles. These studies were performed on 22
9 species, ranging from bees to monkeys. The main characteristics of the included studies are
10 summarised in Figure 3, showcasing significant variation in study subjects and methodologies.
11 Individual studies contributed between 1 and 30 effect sizes to our final data set.

12 Mammals were the best-represented taxonomic group (56 out of 71 studies; 330 out of 459
13 effect sizes), and almost all studies were performed on captive animals (65 studies; 414 effect
14 sizes). Females were more frequently used in experiments than males or mixed-sex groups (225,
15 118, 116 effect sizes, respectively; for the numbers of studies see Figure 3), and adults were
16 more commonly used than juveniles (333 and 126 effect sizes, respectively). Most often, affect
17 manipulation was a form of stress induction compared to standard/benign conditions (benign-
18 worse comparison: 230 effect sizes). Enrichment compared to control/benign conditions was
19 the next most common manipulation (better-benign comparison: 135 effect sizes), and a few
20 studies compared positive treatments (e.g. enrichment) to negative treatments (e.g. handling)
21 (better-worse comparison: 94 effect sizes). Manipulations were usually long-term (292 effect
22 sizes), lasting for days or weeks before affect was measured.

23 Between-subject designs (independent groups of animals exposed to manipulation or
24 control/benign treatment) accounted for 302 effect sizes and within-subject designs accounted

1 for 157 effect sizes. Go/no-go tasks dominated over active choice go/go tasks (389 and 70 effect
 2 sizes, respectively). Spatial and visual cues were most commonly used in judgement bias tests
 3 (177 and 167 effect sizes respectively), and reward-punishment training schemes were more
 4 common than reward-null (283 and 132 effect sizes, respectively), with the remaining studies
 5 using different reward strengths (44 effect sizes). Most studies did not report whether the
 6 personnel performing measurements of animal behaviour were blinded to treatments (only 113
 7 effect sizes came from blinded trials), or whether the measurements were automated (only 71
 8 effect sizes came from automated trials). Finally, latency and proportion outcome measures
 9 were reported at similar levels (258 and 201 effect sizes, respectively). Only 5 studies using
 10 go/no-go tasks reported outcome measures as both latency and proportion, and these were
 11 moderately correlated ($r = 0.578$, $t = 3.085$, $df = 19$, p -value = 0.006), although not for the data
 12 subset using only the largest effect sizes from each experiment to remove non-independence
 13 ($r = 0.443$, $t = 0.857$, $df = 3$, p -value = 0.455).

14 *An overall effect and heterogeneity among effect sizes*

15 Overall, we found a statistically significant effect of experimental treatments on judgement bias
 16 in animals (phylogenetic multilevel meta-analysis: Hedges' g (H_g)_[overall mean] = 0.201, 95%
 17 Confidence Interval (CI) = 0.028 to 0.374; Figure 4, Table S4). A similar model, but without
 18 controlling for phylogeny, also showed a statistically significant overall effect (multilevel meta-
 19 analysis: H_g _[overall mean] = 0.204, 95% CI = 0.087 to 0.320, Table S5). Therefore, animals in a
 20 relatively better treatment usually behaved in a more 'optimistic' way than animals in a
 21 relatively worse treatment, whereas animals in a relatively worse treatment were more
 22 'pessimistic'. Notably, this overall effect is comparable to a small effect, as suggested by the
 23 benchmark values (0.2, 0.5 and 0.8 as small, medium and large effects; Cohen, 1969). The total
 24 heterogeneity in the whole data set was high ($I^2_{total} = 76.4\%$; according to Higgins' benchmark
 25 25, 50 and 75% can be interpreted as low, moderate and high heterogeneity, respectively;

1 Higgins and Thomson, 2002). About 68.1% of the variability across studies was due to sampling
2 error, while phylogeny contributed little to account for this heterogeneity (2.0%), suggesting a
3 weak phylogenetic signal (see Nakagawa & Santos, 2012).

4 High observed total heterogeneity in the data set warrants investigation of potential moderators
5 of heterogeneity. We, thus, present findings of the univariate multilevel phylogenetic meta-
6 regression models examining the effects of different moderators (see Figures 5, Figure 6, Figure
7 S3).

8 *Species-specific effects*

9 A meta-regression model estimating mean effect for each included species did not show a clear
10 pattern of differences among species (Figure 5; $R^2 = 0.070$, Table S6). Some of the species-
11 specific point estimates were medium or large, but they were accompanied by wide confidence
12 intervals crossing zero (no-effect) line. We note that the distribution of studies among species
13 was not balanced, with the data set being dominated by studies on rats, cattle, and pig (15, 11
14 and 8 studies, respectively), while most of remaining species are each represented by a single
15 study (Figure 5).

16 *Sex-effects*

17 Effects of judgement bias manipulations on males were small-to-medium and statistically
18 different from zero ($H_{g[\text{males}]} = 0.365$, 95% CI = 0.155 to 0.575), while effects on females were, on
19 average, close to zero ($H_{g[\text{females}]} = 0.104$, 95% CI = -0.063 to 0.271; Figure 6a). The difference
20 between mean effects in males and females was small ($H_{g[\text{male vs. female difference}]} = 0.261$, 95% CI = -
21 0.001 to 0.522; $R^2 = 0.024$, Table S7), indicating that affect manipulations on judgement bias
22 measurements tend to be more pronounced in studies on males than females.

23 *Tasks type effects*

1 Effects of judgement bias manipulations tended to be larger in studies using active choice tasks
 2 in comparison to studies using go/no-go tasks ($H_{g[\text{go/no-go vs. active choice difference}]} = -0.277$, 95% CI = -
 3 0.567 to 0.012; Figure 6b). On average, tasks with active choice had a medium effect size and
 4 were statistically different from zero ($H_{g[\text{active choice}]} = 0.432$, 95% CI = 0.151 to 0.712), while the
 5 average effect size in in go/no-go tasks was small, but still statistically different from zero
 6 ($H_{g[\text{go/no-go}]} = 0.154$, 95% CI = 0.005 to 0.304; $R^2 = 0.021$, Table S8).

7 *Cue types used during judgement bias tests*

8 Across the five categories of cues used during judgement bias tests, only tests using auditory and
 9 tactile cues consistently revealed differences between control and affect-manipulated groups of
 10 animals ($H_{g[\text{auditory cues}]} = 0.393$, 95% CI = 0.136 to 0.651; $H_{g[\text{tactile cues}]} = 0.658$, 95% CI = 0.136 to
 11 1.118; Figure 6c). These two categories of cues were only significantly different from the results
 12 from studies using visual cues, which on averaged had the weakest effect ($H_{g[\text{visual cues}]} = 0.067$,
 13 95% CI = -0.133 to 0.268; $R^2 = 0.044$, Table S9).

14 *Reinforcement scheme during judgement bias tests*

15 Studies using Reward-Punishment and Reward-Reward training cue reinforcement schemes
 16 usually generated small-medium statistically significant effect sizes in the predicted direction
 17 ($H_{g[\text{Reward-Punishment}]} = 0.216$, 95% CI = 0.036 to 0.396; $H_{g[\text{Reward-Reward}]} = 0.488$, 95% CI = 0.137 to
 18 0.839), but not the Reward-Null reinforcement scheme (Figure 6d). Reward-Reward studies
 19 generally showed significantly larger judgement bias than those that used a Reward-Null
 20 reinforcement scheme ($H_{g[\text{Reward-Reward vs. Reward-Null}]} = 0.436$, 95% CI = 0.045 to 0.827; $R^2 = 0.030$,
 21 Table S10). Studies using non-reinforced ambiguous cues (which was the vast majority of
 22 included studies) generated effect sizes in the predicted direction ($H_{g[\text{ambig. cue not reinforced}]} = 0.204$,
 23 95% CI = 0.026 to 0.382), although not statistically different from studies in which ambiguous

1 cues were reinforced ($H_g[\text{ambig. cue not reinforced vs. reinforced}] = 0.080$, 95% CI = -0.527 to 0.686; $R^2 =$
 2 0.001), whose effect sizes were close to zero (Table S11).

3 *Cue ambiguity level*

4 Ambiguous cues that were halfway between the positive and negative cues, as well as cues that
 5 were closer to the negative cues, were most likely to reveal judgement bias in tested animals
 6 ($H_g[\text{mid-point cue}] = 0.250$, 95% CI = 0.042 to 0.458; $H_g[\text{near-negative cue}] = 0.303$, 95% CI = 0.075 to 0.530;
 7 $R^2 = 0.014$, Figure 6e). Ambiguous near-negative cues were also significantly different from the
 8 effects of positive training cues, with the latter on average being least likely to show judgement
 9 bias effect ($H_g[\text{positive cues}] = 0.063$, 95% CI = -0.153 to 0.278, Table S12).

10 *Other moderators in univariate models*

11 Variation in the other considered moderators did not appear to significantly influence the
 12 magnitude of judgement bias effects. These moderators were: source of animals (captive
 13 vs. wild-caught), animal age, type of affect manipulation (stress vs. enrichment), timing of affect
 14 manipulation (short vs. long-term), whether manipulation was compared to benign or worse
 15 reference condition, type of study design (within-individual vs. between-individuals), food
 16 deprivation during judgement bias tests, measurement type of behavioural response (latency vs.
 17 proportion), automation and blinding of measurements of animal responses (Figure S3; Tables
 18 S13 – S22; $R^2 = 0$ to 0.010).

19 *Multivariate (full) meta-regression models and model selection*

20 The full meta-regression model included four moderators that were significant or close to
 21 statistical significance in univariate models (after confirming they were not co-linear with each
 22 other): sex of test animals, task type (go/no-go vs. active choice go/go), type of cue used in the
 23 test, and type of reinforcement for positive and negative training cues. In the multivariate meta-
 24 regression, none of the considered moderators was significant (Table S23). These moderators

1 can jointly explain only about 7% of variation in the data ($R^2 = 0.072$). Model selection analysis
2 indicated that type of the task and type of reinforcement used could be the most influential
3 moderators, followed by the sex of animals (Table S24).

4 *Publication bias*

5 We conducted 3 kinds of publication bias analyses: 1) contour-enhanced funnel plots of
6 residuals, 2) a variant of Egger's regression, and 3) a regression-based time-lag bias test. Visual
7 inspection of enhanced-contour funnel plots of residuals did not reveal skewness indicative of
8 publication bias (Figure S4). Further, the intercept of Egger's multivariate regression, controlling
9 for potentially important moderators from univariate models, was not significantly different
10 from zero ($t = 0.017$, $df = 457$, $p = 0.986$), confirming lack of publication bias in the full data set.
11 Finally, we found no evidence for time-lag bias, as the slope of linear regression between
12 publication year and effect size was not significantly different from zero (Slope_[Year] = -0.002,
13 95% CI = -0.121 to 0.118, $p = 0.980$, Table S25).

14 *Sensitivity analyses (robustness of results)*

15 The estimates from Bayesian models run on full data set gave qualitatively identical results to
16 the REML models used in the main data analyses. Namely, the overall effect was small and
17 statistically significant ($H_{g[\text{overall mean}]} = 0.206$, 95% CI = 0.041 to 0.383; $I^2_{\text{total}} = 76.8\%$; Table S26).

18 In the Bayesian multivariate meta-regression, none of the moderators significantly influenced
19 judgement bias test outcomes, as in the equivalent log-likelihood model.

20 Finally, we ran meta-analytic models on four data subsets, representing different ways of
21 looking at the results from response curves with multiple cues: i) including only data from
22 ambiguous cues (81 NP, 108 MID, and 80 NN effect sizes for cue locations included in this data
23 subset), ii) including only data from mid-point ambiguous cues (108 MID effect sizes included),
24 iii) including only data for maximum response, in absolute terms (26 P, 13 NP, 31 MID, 22 NN,

1 and 16 N effect sizes included), iv) including only data for maximum response in the overall
2 direction of response (19 P, 12 NP, 34 MID, 28 NN, and 15 N effect sizes included). All these data
3 subsets tended to have larger overall effect size estimates, than in the full data set meta-analyses
4 (Figure 4, Tables S4 and S5). Univariate and multivariate meta-regression models usually
5 showed similar patterns to these observed in the analyses on the full dataset (Tables S6-S23).

6

1 **Discussion**

2 Our meta-analysis revealed that non-pharmacological affect manipulations generally influenced
3 judgement bias in the predicted direction (i.e. manipulations assumed to generate a relatively
4 positive state were likely to generate an 'optimistic' response to cues). However, effects were
5 usually small to large (average Hedges' g of 0.2 – 0.6), and they were highly variable, with total
6 observed heterogeneity (I^2) over 75%. The moderators that potentially influenced magnitude of
7 effects included cue type, type of task used in judgement bias trials, reinforcement combination
8 used for training positive and negative cues, cue ambiguity level, and sex of tested animals.
9 However, small R^2 values (1.4 to 4.4%) indicated that these moderators explained only a small
10 proportion of variance. We discuss these findings in detail below.

11 *Validity and efficacy of judgement bias tests*

12 Our main finding generally supports judgement bias tests as a valid approach to measure affect
13 in non-human animals. This is in line with conclusions of a narrative cross-species review
14 (Bethell, 2015) and a recent systematic review of 20 rodent studies on judgement bias (Nguyen,
15 et al. 2020). However, the latter considered both pharmacological and non-pharmacological
16 manipulations and only conducted a qualitative synthesis of their rodent data set. Effects of
17 pharmacological manipulations across species were recently quantitatively synthesised by our
18 team (Neville, et al. 2020) and our current work provides the first quantification of non-
19 pharmacological manipulations across different taxa.

20 Our quantitative results show that the observed behavioural effect of the affect manipulations
21 investigated is, on average, small (Hedges' g of 0.2) and highly heterogeneous. However, we base
22 this conclusion on the analyses of the full dataset, which included mean latency and/or
23 proportion data from all cues used in the judgement bias tests. Thus, we likely underestimated
24 the overall effect size, due to the inclusion of positive and negative training (unambiguous) cues,

1 in the analysis. Weak discrimination performance for the training cues could decrease the
2 likelihood of detecting a judgement bias (Roelofs et al. 2016), warranting further investigations
3 of how selected learning criteria, or actual discrimination levels, affect the sensitivity of
4 judgement bias tests. It is also possible that the affect manipulations used in many of the
5 included studies were rather "mild" – not many authors used severe stressors or pain stimuli,
6 either for welfare reasons and/or because pain stimuli might exert a general suppressive effect
7 on responses to cues. It is thus possible that some manipulations failed to influence animal
8 affect.

9 As noted earlier, there are theoretical and empirical reasons for why judgement biases may not
10 occur at training cues, and also for why they may not occur at all ambiguous cues. When we
11 restricted analysis to the cue with the largest absolute effect size in the direction of the overall
12 mean effect size from each response curve – the estimated overall effect sizes were between
13 moderate to large (Hedges' g of 0.6). The overall effect sizes were moderate when we used the
14 other three data subsets: (i) ambiguous cues only; (ii) middle cue only; (iii) cue with the largest
15 absolute effect size. Yet, analyses of the full dataset are most powerful, given that they include
16 data points representing the whole response curve (Gygax, 2014).

17 The high data heterogeneity is congruent with the levels observed in most ecological and
18 evolutionary meta-analyses (70 – 95%; Senior et al., 2016). High heterogeneity (> 75%) of the
19 effect sizes in our data set indicates variability in the influences of non-pharmacological
20 manipulations of affective state on judgement bias in animals, but is perhaps not surprising
21 given how diverse the studies were in terms of, for example, species used (22 diverse species;
22 Figure 5), task variants, affect manipulations, and other methodological specifics. Accordingly,
23 the lack of phylogenetic effects in our data set is consistent with the observation that meta-
24 analyses on phylogenetically diverse sets of species are unlikely to show a strong phylogenetic
25 signal (Chamberlain et al., 2012).

1 *Key moderators of judgement bias tests*

2 We also revealed five important moderators of responses in the judgement bias task. Four are
3 related to methodology and one is a biological factor. First, active choice go/go tasks tended to
4 yield larger effects than go/no-go tasks. It is possible that the former are more cognitively
5 challenging given that the response needs to be deployed to different stimuli. Such a potential
6 cognitive load might render go/go tasks less susceptible to habitual responding and thereby
7 more sensitive to affect manipulations. Furthermore, go/no-go tasks are likely to be vulnerable
8 to the influence of Pavlovian action predispositions (e.g. go-for-reward; no-go to avoid
9 punishment; Guitart-Masip et al., 2014; Jones et al., 2017), that could inadvertently bias
10 responding (Mendl & Paul, 2020) and obscure affect manipulation effects. Additionally, subjects
11 may sometimes perform no-go responses for reasons unrelated to affect manipulations (e.g.
12 failing to detect or attend to a cue; Bethell, 2015; Jones et al., 2018), making these tests less
13 dependable. Still, we observed that go/no-go tasks are more commonly used in judgement bias
14 studies (in 57 vs. 14 studies; Figure 3), probably because they are easier and quicker to train.

15 Second, Reward-Reward tasks usually generated larger effect sizes than Reward-Null tasks. Part
16 of the reason for this may be that Reward-Reward tasks usually involve a go/go active choice
17 response and this itself predisposes stronger effects, as just discussed. The most frequently used
18 Reward-Punishment tasks had the largest observed average effect size. It is possible that the
19 Reward-Punishment design, providing a more affectively-laden task (i.e. decision outcomes can
20 range from a desired reward to an aversive punisher), is more sensitive to manipulations of
21 affective state (see Mendl et al. 2009).

22 Third, the use of auditory and tactile cues tended to reveal the largest effects compared to when
23 spatial, visual and olfactory cues were employed. There may be a number of reasons for this,
24 some of which may be linked to differences in species biology (Bethell, 2015). For example,
25 whilst people are strongly visually focused when information gathering, many other animal

1 species are not, and may not readily exhibit human-like processing of visual cues. Conversely,
2 olfactory sensitivity in humans is poor, relative to many other species, and this may impair the
3 ability of researchers to design or use meaningful cues in this sensory dimension. It is also
4 possible that cue modality and presentation method can influence the uncertainty of
5 information provided by 'ambiguous' cues. For example, there may be greater uncertainty about
6 the information provided by a single tone intermediate between two training tones, than by a
7 spatial location situated between two training locations. Such differences in uncertainty may
8 have knock-on effects on animal's decisions.

9 Fourth, cue ambiguity level (P, NP, MID, NN, N) was important. We found predicted judgement
10 bias only at ambiguous cues in the full dataset analysis, and not at positive or negative training
11 cues, on average. Still, some individual studies in the dataset yielded large effects at positive or
12 negative training cues (e.g. Deakin, 2018; Horváth et al., 2016; Zidar et al., 2018). In line with
13 this, Neville et al. (2020) noted that pharmacological manipulations of affect altered judgement
14 bias principally at ambiguous cues, but also at the negative training cue. Large effects at non-
15 ambiguous cues could occur in at least two ways. First, if affect manipulations altered valuation
16 of decision outcomes (e.g. by decreasing food valuation and hence generating a weaker response
17 to the positive cue), the manipulations could change propensities to perform specific responses
18 (e.g. go vs. no-go) and interfere with memory of training cue-outcome associations. Second, large
19 effects at non-ambiguous cues might occur if training was brief or ineffective such that there was
20 considerable ambiguity about the training cue-outcome association during testing (see Mendl et
21 al., 2009; Bateson et al. 2011; Bethell, 2015; Mendl and Paul, 2020). As mentioned earlier,
22 further research of the effects of variation in discrimination training criteria on test performance
23 would shed light on this issue.

24 Finally, in all analyses, larger predicted effect sizes tended to be reported for male subjects than
25 for females or mixed sex groups. This pattern could be due to existence of sex differences in

1 neurobiology of learning and memory (Jonasson et al., 2005) or sex differences in stress effects
2 on memory, with different patterns for acute and prolonged stress (Andreano and Cahill, 2009).
3 Effects of enrichment may also be sex-specific (Lin et al., 2011; ter Horst et al., 2012).

4 *Potential limitations and recommendations*

5 The results of our meta-analysis come with six caveats, which we list here alongside
6 recommendations for future studies of judgement bias. First, captive and domesticated
7 mammals dominate the dataset making our conclusions particularly relevant to research on
8 welfare of such animals. Conversely, the analyses are less informative for wild animals,
9 vertebrates other than mammals, and invertebrates. Indeed, Bethell's narrative review (2015)
10 highlighted biased taxonomic representation in empirical evidence. Thus, future work in this
11 area could aim to increase the representation of non-domesticated species, such as those kept in
12 zoos and for research (where animal welfare is of concern; Baumans, 2005; Bethell, 2015;
13 Wolfensohn et al., 2018) and invertebrates (where welfare is an emerging issue; Drinkwater et
14 al., 2019).

15 Second, we had limited statistical power to detect clear differences between the levels of a
16 number of the tested moderators. Also, the small sample sizes at some levels of the considered
17 moderators might have introduced some spurious findings. For example, relatively few studies
18 used tactile or olfactory cues (e.g. Barker et al., 2017; Novak et al., 2016), and very few used
19 reinforced ambiguous cues during tests (e.g. Bailoo et al., 2018; Keen et al., 2014). To address
20 this limitation, future studies of commonly used laboratory and domesticated species should
21 systematically investigate the role of different cue types. Researchers should also attempt to
22 make cue types relevant for a given species, and vary the perceptual closeness of training cues
23 and hence the difficulty of the task and uncertainty of ambiguous cues.

1 Third, for some moderators, especially those related to study quality, poor reporting might have
2 obscured statistical relationships. Very few of the included studies explicitly stated that they
3 used automation or blinding, and we had to assume that the remaining studies did not use these.
4 Thus, automation of measurements could be used more often, and/or their use should be clearly
5 reported. Notably, Nguyen et al. (2020) in their systematic review of 20 rodent studies
6 highlighted limited information on the details of experimental procedures and analyses in 65%
7 of assessed studies, undermining confidence in the findings. Nevertheless, we found no
8 statistical evidence for publication bias in our meta-analytic data set. The lack of publication bias
9 is potentially due to our full data set containing data points across the whole response curve,
10 which are usually a mixture of small and large positive effects (in the expected direction) and
11 even some negative ones (not in the expected direction). Also related to reporting, mixed-sex
12 groups of animals comprised almost one-third of the data in our meta-analysis, potentially
13 obscuring sex-specific effects. Providing sex-disaggregated data in research is absolutely
14 essential for improving our understanding of animal behaviour and cognition (Shansky and
15 Woolley, 2016; Palanza and Parmigiani, 2017).

16 Fourth, we were also not able to include strength of manipulation in our analyses (there is no
17 common scale for the diverse types of manipulations included in our data set). To overcome this
18 problem, in future studies it would be valuable to test and synthesize relationships between
19 measures of cognitive bias and different biomarkers of stress, such as cortisol, adrenaline, alpha-
20 amylase, testosterone, leucocyte profiles (Keay et al., 2006; Davis et al., 2008).

21 Fifth, we also noted some outliers in the dataset, which usually came from studies with severe
22 manipulations and/or small sample sizes. We, however, conducted extensive sensitivity analyses
23 to test robustness of our conclusions, with the results generally conforming to our predictions
24 and being robust across different statistical approaches. Further, in individual empirical studies
25 comparing two means, to achieve power of 0.8 at alpha of 0.05, it is necessary to have sample

1 sizes of at least 50 animals per group for detecting moderate effect sizes (Hedges' $g = 0.4$). Best-
2 case scenario, when effect is large (Hedges' $g = 0.8$), would require only 13 animals per group to
3 achieve the same power. Conducting power analyses to determine suitable sample sizes for
4 planned experiments can help reducing animal use and also prevent wasting animals on
5 underpowered studies.

6 Finally, the largest responses often do not appear at the most intermediate/ambiguous cue.
7 Because of this, we suggest that multiple ambiguous (probe) cues (at least 3) are needed for
8 robust and comprehensive judgement bias tests although 25% of response curves in our data set
9 included only one ambiguous cue.

10 Taken together, while it is unlikely that a single "perfect" version of a judgement bias test exists,
11 our analyses suggest that the most sensitive setup would comprise a go/go active choice task
12 employing a reward-punishment or reward-reward reinforcement contingency, and using at
13 least three ambiguous cues of a sensory modality appropriate to the species of interest. Such
14 tasks may require more lengthy training than conventional go/no-go tasks, but promising new
15 go/no-go variants that achieve many of the benefits of go/go tasks by incorporating active trial
16 initiation are being developed and can be trained relatively rapidly (Hintze et al. 2018; Jones et
17 al. 2018). For more detailed guidance on how tests should be designed and conducted, and what
18 types of adjustments may be needed for different organisms, we refer readers to the works of
19 Bethell (2015) and Roelofs et al. (2016).

20 **Conclusions**

21 In summary, judgement bias tests are a valid method of measuring animal affective state.
22 However, high heterogeneity among studies, which can be only partially explained by simple
23 influences of considered moderators, warrants care in designing and interpreting judgement

1 bias manipulations and tests. We call for better reporting of experimental designs, especially
2 blinding and automation, disaggregation of data by sex of subjects, and other experimental
3 details that might influence study results. Also, there is a need for more empirical studies that
4 compare different experimental designs and setups, including using different types of tasks,
5 cues, and cue ambiguity levels.

6 **Acknowledgments**

7 We are thankful to all authors that have provided us with additional information or data, as
8 indicated in Table S2. M.L. and S.N. were supported by the Australian Research Council
9 Discovery Project (DP200100367). E.S. and J.Z. were funded by Carl Trygger's Foundation and
10 the Swedish research council Formas, respectively, awarded to H.L., M.M., E.S.P. and V.N. thank
11 the UK Biotechnology and Biological Sciences Research Council (BBSRC grants BB/P019218/1,
12 BB/T002654/1 and BBSRC SWBio DTP grant BB/M009122/1), and the UK National Centre for
13 the Replacement, Refinement and Reduction of Animals in Research (NC3Rs grant
14 NC/K00008X/1) for supporting their work in this area. The authors declare no conflicts of
15 interests.

16 **Data and Code Availability**

17 All data and code are available from the following online OSF repository: <https://osf.io/anfhm/>

18 **References**

- 19 Andreano, J.M., Cahill, L., 2009. Sex influences on the neurobiology of learning and memory.
20 Learn. Mem. 16, 248–266. <https://doi.org/10.1101/lm.918309>
- 21 Ash, H., Buchanan-Smith, H.M., 2016. The long-term impact of infant rearing background on the
22 affective state of adult common marmosets (*Callithrix jacchus*). Appl. Anim. Behav. Sci. 174,
23 128–136. <https://doi.org/10.1016/j.applanim.2015.10.009>

- 1 Baciadonna, L., McElligott, A.G., 2015. The use of judgement bias to assess welfare in farm
2 livestock. *Anim. Welf.* 24, 81–91. <https://doi.org/10.7120/09627286.24.1.081>
- 3 Bailoo, J.D., Murphy, E., Boada-Saña, M., Varholick, J.A., Hintze, S., Baussière, C., Hahn, K.C.,
4 Göpfert, C., Palme, R., Voelkl, B., Würbel, H., 2018. Effects of cage enrichment on behavior,
5 welfare and outcome variability in female mice. *Front. Behav. Neurosci.* 12, 232.
6 <https://doi.org/10.3389/fnbeh.2018.00232>.
- 7 Barker, T.H., Bobrovskaya, L., Howarth, G.S., Whittaker, A.L., 2017. Female rats display fewer
8 optimistic responses in a judgment bias test in the absence of a physiological stress
9 response. *Physiol. Behav.* 173, 124–131. <https://doi.org/10.1016/j.physbeh.2017.02.006>
- 10 Barton, K., 2009. MuMIn: Multi-model inference. R Package Version 0.12.2/r18. [http://R-](http://R-Forge.R-project.org/projects/mumin/)
11 [Forge.R-project.org/projects/mumin/](http://R-Forge.R-project.org/projects/mumin/)
- 12 Bateson, M., 2016. Optimistic and pessimistic biases: a primer for behavioural ecologists. *Curr.*
13 *Opin. Behav. Sci.* 12, 115–121. <https://doi.org/10.1016/j.cobeha.2016.09.013>
- 14 Bateson, M., Desire, S., Gartside, S. E., Wright, G. A., 2011. Agitated honeybees exhibit pessimistic
15 cognitive biases. *Curr. Biol.* 21, 1070–1073. doi: 10.1016/j.cub.2011.05.017
- 16 Baumans, V., 2005. Science-based assessment of animal welfare: laboratory animals. *Rev. Sci.*
17 *Tech.* 24, 503-513. <http://dx.doi.org/10.20506/rst.24.2.1585>
- 18 Bethell, E.J., 2015. A “how-to” guide for designing judgment bias studies to assess captive animal
19 welfare. *J. Appl. Anim. Welf. Sci.* 18, S18–S42. [https://doi.org/](https://doi.org/10.1080/10888705.2015.1075833)
20 [10.1080/10888705.2015.1075833](https://doi.org/10.1080/10888705.2015.1075833).
- 21 Bethell, E.J., Holmes, A., MacLarnon, A., Semple, S., 2016. Emotion evaluation and response
22 slowing in a non-human primate: new directions for cognitive bias measures of animal
23 emotion? *Behav. Sci. (Basel)*. 6, 2. <https://doi.org/10.3390/bs6010002>
- 24 Bethell, E.J., Koyama, N.F., 2015. Happy hamsters? Enrichment induces positive judgement bias
25 for mildly (but not truly) ambiguous cues to reward and punishment in *Mesocricetus*
26 *auratus*. *R. Soc. Open Sci.* 2, 140399. <https://doi.org/10.1098/rsos.140399>
- 27 Blanchette, I., Richards, A., 2010. The influence of affect on higher level cognition: a review of
28 research on interpretation, judgement, decision making and reasoning. *Cogn. Emot.*
29 <https://doi.org/10.1080/02699930903132496>
- 30 Brilot, B.O., Asher, L., Bateson, M., 2010. Stereotyping starlings are more “pessimistic.” *Anim.*
31 *Cogn.* 13, 721–731. <https://doi.org/10.1007/s10071-010-0323-z>

- 1 Burman, O.H.P., Mendl, M.T., 2018. A novel task to assess mood congruent memory bias in non-
2 human animals. *J. Neurosci. Methods* 308, 269–275.
3 <https://doi.org/10.1016/j.jneumeth.2018.07.003>
- 4 Burman, O.H.P., Parker, R.M.A., Paul, E.S., Mendl, M.T., 2009. Anxiety-induced cognitive bias in
5 non-human animals. *Physiol. Behav.* 98, 345–350.
6 <https://doi.org/10.1016/j.physbeh.2009.06.012>
- 7 Chamberlain, S.A., Hovick, S.M., Dibble, C.J., Rasmussen, N.L., Van Allen, B.G., Maitner, B.S., Ahern,
8 J.R., Bell - Dereske, L.P., Roy, C.L., Meza - Lopez, M., Carrillo, J., Siemann, E., Lajeunesse, M.J.
9 and Whitney, K.D., 2012. Does phylogeny matter? Assessing the impact of phylogenetic
10 information in ecological meta - analysis. *Ecol. Lett.* 15, 627-636. doi:10.1111/j.1461-
11 0248.2012.01776.x
- 12 Crump, A., Arnott, G., Bethell, E.J., 2018. Affect-Driven Attention Biases as Animal Welfare
13 Indicators: Review and Methods. *Animals* 8. <https://doi.org/10.3390/ani8080136>
- 14 da Cunha Nogueira, S.S., Fernandes, I.K., Oliveira Costa, T.S., Gama Nogueira-Filho, S.L., Mendl, M.,
15 2015. Does trapping influence decision-making under ambiguity in white-lipped peccary
16 (*Tayassu pecari*)? *PLOS One* 10. <https://doi.org/10.1371/journal.pone.0127868>
- 17 Cohen, J. (1969). *Statistical power analysis for the behavioral sciences* (1st ed.). New York, NY:
18 Academic Press.
- 19 Darwin, C., 1872. *The expression of the emotions in man and animals*, 3rd ed., *The expression of*
20 *the emotions in man and animals*, 3rd ed. Oxford University Press, New York, NY, US.
- 21 Davis, A.K., Maney, D.L., Maerz, J.C., 2008. The use of leukocyte profiles to measure stress in
22 vertebrates: a review for ecologists. *Funct. Ecol.* 22: 760-772.
23 <https://doi.org/10.1111/j.1365-2435.2008.01467.x>
- 24 Destrez, A., Deiss, V., Lévy, F., Calandreau, L., Lee, C., Chaillou-Sagon, E., Boissy, A., 2013. Chronic
25 stress induces pessimistic-like judgment and learning deficits in sheep. *Appl. Anim. Behav.*
26 *Sci.* 148, 28–36. <https://doi.org/10.1016/j.applanim.2013.07.016>
- 27 Douglas, C., Bateson, M., Walsh, C., Bédoué, A., Edwards, S.A., 2012. Environmental enrichment
28 induces optimistic cognitive biases in pigs. *Appl. Anim. Behav. Sci.* 139, 65–73.
29 <https://doi.org/10.1016/j.applanim.2012.02.018>
- 30 Doyle, R.E., Lee, C., Deiss, V., Fisher, A.D., Hinch, G.N., Boissy, A., 2011. Measuring judgement bias
31 and emotional reactivity in sheep following long-term exposure to unpredictable and

- 1 aversive events. *Physiol. Behav.* 102, 503–510.
2 <https://doi.org/10.1016/j.physbeh.2011.01.001>
- 3 Drinkwater, E., Robinson, E.J.H., Hart, A.G., 2019. Keeping invertebrate research ethical in a
4 landscape of shifting public opinion. *Methods Ecol Evol.* 10, 1265–1273.
5 <https://doi.org/10.1111/2041-210X.13208>
- 6 Egger, M., Smith, G.D., Schneider, M., Minder, C., 1997. Bias in meta-analysis detected by a simple
7 , graphical test measures of funnel plot asymmetry. *BMJ* 315, 629–634.
8 <https://doi.org/10.1136/bmj.315.7109.629>
- 9 Girard, B., Bellone, C., 2020. Revealing animal emotions. *Science* 368, 33–34. doi:
10 10.1126/science.abb2796
- 11 Guitart-Masip M., Duzel, E., Dolan, R., Dayan, P., 2014. Action versus valence in decision making.
12 *Trends Cogn. Sci.* 18,194-202. <https://doi.org/10.1016/j.tics.2014.01.003>
- 13 Gygax, L., 2014. The A to Z of statistics for testing cognitive judgement bias. *Anim. Behav.* 95, 59–
14 69. <https://doi.org/10.1016/j.anbehav.2014.06.013>
- 15 Hadfield, J.D., 2010. MCMC methods for multi-response generalized linear mixed models: The
16 MCMCglmm R package. *J. Stat. Softw.* 33, 1–22. <https://doi.org/10.18637/jss.v033.i02>
- 17 Hales, C.A., Robinson, E.J., Houghton, C.J., 2016. Diffusion modelling reveals the decision making
18 processes underlying negative judgement bias in rats. *PLOS One* 11.
19 <https://doi.org/10.1371/journal.pone.0152592>
- 20 Hales, C.A., Stuart, S.A., Anderson, M.H., Robinson, E.S., 2014. Modelling cognitive affective biases
21 in major depressive disorder using rodents. *Br. J. Pharmacol.* 171, 4524–4538.
22 <https://doi.org/10.1111/bph.12603>
- 23 Harding, E.J., Paul, E.S., Mendl, M., 2004. Cognitive bias and affective state. *Nature* 427, 312.
- 24 Hedges, L. V., Olkin, I., 1985. *Statistical Methods for Meta-Analysis*, *Statistical Methods for Meta-*
25 *Analysis*. New York: Academic Press. <https://doi.org/10.1016/c2009-0-03396-0>
- 26 Hintze, S., Melotti, L., Colosio, S., Bailoo, J.D., Boada-Sana, M., Wurbel, H., Murphy, E., 2018. A
27 cross-species judgement bias task: integrating active trial initiation into a spatial Go/No-go
28 task. *Sci. Rep.* 8, 5104. doi:10.1038/s41598-018-23459-3.
- 29 Jennions, M.D., Møller, A.P., 2002. Relationships fade with time: a meta-analysis of temporal
30 trends in publication in ecology and evolution. *Proceedings. Biol. Sci.* 269, 43–8.
31 <https://doi.org/10.1098/rspb.2001.1832>

- 1 Jonasson, Z., 2005. Meta-analysis of sex differences in rodent models of learning and memory: a
2 review of behavioral and biological data. *Neurosci. Biobehav. Rev.*
3 <https://doi.org/10.1016/j.neubiorev.2004.10.006>
- 4 Jones, S., Neville, V., Higgs, L., Paul, E.S., Dayan, P., Robinson, E.S.J., Mendl, M., 2018. Assessing
5 animal affect: an automated and self-initiated judgement bias task based on natural
6 investigative behaviour. *Scientific Reports* 8, 12400. doi: 10.1038/s41598-018-30571-x
- 7 Jones, S., Paul, E.S., Dayan, P., Robinson, E.S.J., Mendl, M., 2017. Pavlovian influences on learning
8 differ between rats and mice in a counter-balanced Go/NoGo judgement bias task.
9 *Behavioural Brain Research* 331:214-224. doi: 10.1016/j.bbr.2017.05.044
- 10 Keay, J.M., Singh, J., Gaunt, M.C., Kaur, T., 2006. Fecal glucocorticoids and their metabolites as
11 indicators of stress in various mammalian species: a literature review. *J. Zoo Wildlife Med.*,
12 37, 234-244. <https://doi.org/10.1638/05-050.1>
- 13 Keen, H.A., Nelson, O.L., Robbins, C.T., Evans, M., Shepherdson, D.J., Newberry, R.C., 2014.
14 Validation of a novel cognitive bias task based on difference in quantity of reinforcement
15 for assessing environmental enrichment. *Anim. Cogn.* 17, 529–541.
16 <https://doi.org/10.1007/s10071-013-0684-1>
- 17 Lin, E.J., Choi, E., Liu, X., Martin, A., & Dusing, M.J. (2011). Environmental enrichment exerts sex-
18 specific effects on emotionality in C57BL/6J mice. *Behav. Brain Res.* 216, 349–357.
19 <https://doi.org/10.1016/j.bbr.2010.08.019>
- 20 Mendl, M., Burman, O.H.P., Parker, R.M.A., Paul, E.S., 2009. Cognitive bias as an indicator of
21 animal emotion and welfare: emerging evidence and underlying mechanisms. *Appl. Anim.*
22 *Behav. Sci.* 118, 161–181. <https://doi.org/10.1016/j.applanim.2009.02.023>
- 23 Mendl, M., Burman, O.H.P., Paul, E.S., 2010. An integrative and functional framework for the
24 study of animal emotion and mood. *Proc. R. Soc. B Biol. Sci.* 277, 2895–2904.
25 <https://doi.org/10.1098/rspb.2010.0303>
- 26 Mendl, M., Paul, E.S. 2020. Animal affect and decision-making. *Neurosci. Biobehav. Rev.* 112, 144-
27 163. <https://doi.org/10.1016/j.neubiorev.2020.01.025>
- 28 Moher, D., Liberati, A., Tetzlaff, J., Altman, D.G., Altman, D., Antes, G., Atkins, D., Barbour, V.,
29 Barrowman, N., Berlin, J.A., Clark, J., Clarke, M., Cook, D., D'Amico, R., Deeks, J.J., Devereaux,
30 P.J., Dickersin, K., Egger, M., Ernst, E., Gøtzsche, P.C., Grimshaw, J., Guyatt, G., Higgins, J.,
31 Ioannidis, J.P.A., Kleijnen, J., Lang, T., Magrini, N., McNamee, D., Moja, L., Mulrow, C., Napoli,
32 M., Oxman, A., Pham, B., Rennie, D., Sampson, M., Schulz, K.F., Shekelle, P.G., Tovey, D.,

- 1 Tugwell, P., 2009. Preferred reporting items for systematic reviews and meta-analyses: the
2 PRISMA statement. *PLoS Med.* <https://doi.org/10.1371/journal.pmed.1000097>
- 3 Murphy, E., Nordquist, R.E., van der Staay, F.J., 2013. Responses of conventional pigs and
4 Gottingen miniature pigs in an active choice judgement bias task. *Appl. Anim. Behav. Sci.*
5 148, 64–76. <https://doi.org/10.1016/j.applanim.2013.07.011>
- 6 Nakagawa, S., Johnson, P.C.D., Schielzeth, H., 2017. The coefficient of determination R^2 and intra-
7 class correlation coefficient from generalized linear mixed-effects models revisited and
8 expanded. *J. R. Soc. Interface* 14. <https://doi.org/10.1098/rsif.2017.0213>
- 9 Nakagawa, S., Lagisz, M., O’Dea, R.E., Rutkowska, J., Yang, Y., Noble, D., Senior, A.M., 2020. The
10 orchard plot: cultivating forest plots for use in ecology, evolution and beyond. *Res. Syn.*
11 *Meth.* 1– 9. <https://doi.org/10.1002/jrsm.1424>
- 12 Nakagawa, S., Santos, E.S.A., 2012. Methodological issues and advances in biological meta-
13 analysis. *Evol. Ecol.* <https://doi.org/10.1007/s10682-012-9555-5>
- 14 Nakagawa, S., Schielzeth, H., 2013. A general and simple method for obtaining R^2 from
15 generalized linear mixed-effects models. *Meth. Ecolol. Evol.* 4, 133–142
16 <https://doi.org/10.1111/j.2041-210x.2012.00261.x>
- 17 Neville, V., Nakagawa, S., Zidar, J., Paul, E.S., Lagisz, M., Bateson, M., Løvlie, H., Mendl, M., 2020.
18 Pharmacological manipulations of judgement bias: a systematic review and meta-analysis.
19 *Neurosci. Biobehav. Rev.* 108, 269–286. <https://doi.org/10.1016/j.neubiorev.2019.11.008>
- 20 Nguyen, H.A.T., Guo, C., Homberg, J.R., 2020. Cognitive bias under adverse and rewarding
21 conditions: a systematic review of rodent studies. *Front. Behav. Neurosci.* 14:14. doi:
22 10.3389/fnbeh.2020.00014
- 23 Novak, J., Stojanovski, K., Melotti, L., Reichlin, T.S., Palme, R., Würbel, H., 2016. Effects of
24 stereotypic behaviour and chronic mild stress on judgement bias in laboratory mice. *Appl.*
25 *Anim. Behav. Sci.* 174, 162–172. <https://doi.org/10.1016/j.applanim.2015.10.004>
- 26 Oliveira, F.R.M., Nogueira, S.L.G., Sousa, M.B.C., Dias, C.T.S., Mendl, M., Nogueira, S.S.C., 2016.
27 Measurement of cognitive bias and cortisol levels to evaluate the effects of space restriction
28 on captive collared peccary (Mammalia, Tayassuidae). *Appl. Anim. Behav. Sci.* 181, 76–82.
29 <https://doi.org/10.1016/j.applanim.2016.05.021>
- 30 Palanza, P., Parmigiani, S., 2017. How does sex matter? Behavior, stress and animal models of
31 neurobehavioral disorders. 76A: 134-143. *Neurosci. Biobehav. Rev.*

- 1 <https://doi.org/10.1016/j.neubiorev.2017.01.037>
- 2 Paul, E.S., Harding, E.J., Mendl, M., 2005. Measuring emotional processes in animals: the utility of
3 a cognitive approach. *Neurosci. Biobehav. Rev.* 29, 469–491.
4 <https://doi.org/10.1016/j.neubiorev.2005.01.002>
- 5 Paul, E.S., Sher, S., Tamietto, M., Winkielman, P., Mendl, M.T., 2020. Towards a comparative
6 science of emotion: affect and consciousness in humans and animals. *Neurosci. Biobehav.*
7 *Rev.* <https://doi.org/10.1016/j.neubiorev.2019.11.014>
- 8 R Development Core Team (2018) R: a language and environment for statistical computing. R
9 Foundation for Statistical Computing, Vienna.
- 10 Roelofs, S., Boleij, H., Nordquist, R.E., van der Staay, F.J., 2016. Making decisions under ambiguity:
11 Judgment bias tasks for assessing emotional state in animals. *Front. Behav. Neurosci.* 10.
- 12 Rolls, E.T., 2009. *Emotion Explained*. Oxford University Press.
13 <https://doi.org/10.1093/acprof:oso/9780198570035.001.0001>
- 14 Rosenthal, R., 1979. The file drawer problem and tolerance for null results. *Psychol. Bull.* 86,
15 638–641. <https://doi.org/10.1037/0033-2909.86.3.638>
- 16 Rygula, R., Papciak, J., Popik, P., 2013. Trait pessimism predicts vulnerability to stress-induced
17 anhedonia in rats. *Neuropsychopharmacology* 38, 2188–2196.
18 <https://doi.org/10.1038/npp.2013.116>
- 19 Scherer, K.R., 2005. What are emotions? And how can they be measured? *Soc. Sci. Inf.* 44, 695–
20 729. <https://doi.org/10.1177/0539018405058216>
- 21 Senior, A.M., Grueber, C.E., Kamiya, T., Lagisz, M., O'Dwyer, K., Santos, E.S.A., Nakagawa, S., 2016.
22 Heterogeneity in ecological and evolutionary meta-analyses: its magnitude and
23 implications. *Ecology* 97. <https://doi.org/10.1002/ecy.1591>
- 24 Shansky, R.M., Woolley, C.S., 2016. Considering sex as a biological variable will be valuable for
25 neuroscience research. *J. Neurosci.* 36, 11817–11822.
26 <https://doi.org/10.1523/JNEUROSCI.1390-16.2016>
- 27 ter Horst, J.P., de Kloet, E.R., Schächinger, H., & Oitzl, M S. (2012). Relevance of stress and female
28 sex hormones for emotion and cognition. *Cell. Mol. Neurobiol.* 32, 725–735.
29 <https://doi.org/10.1007/s10571-011-9774-2>
- 30 Viechtbauer, W., 2010. Conducting meta-analyses in R with the metafor package. *J. Stat. Softw.*
31 <https://doi.org/10.1103/PhysRevB.91.121108>

- 1 Wallace, B.C., Small, K., Brodley, C.E., Lau, J., Trikalinos, T.A., 2012. Deploying an interactive
2 machine learning system in an Evidence-based Practice Center: Abstrackr, in: IHI'12 -
3 Proceedings of the 2nd ACM SIGHIT International Health Informatics Symposium. pp. 819-
4 823. <https://doi.org/10.1145/2110363.2110464>
- 5 Wheeler, R.R., Swan, M.P., Hickman, D.L., 2015. Effect of multilevel laboratory rat caging system
6 on the well-being of the singly-housed Sprague Dawley rat. *Lab. Anim.* 49, 10-19.
7 <https://doi.org/10.1177/0023677214547404>
- 8 Wolfensohn, S., Shotton, J., Bowley, H., Davies, S., Thompson, S., Justice, W.S.M., 2018.
9 Assessment of welfare in Zoo animals: towards optimum quality of life. *Animals* 8, 110.
10 <https://doi.org/10.3390/ani8070110>
- 11 Zidar, J., Campderrich, I., Jansson, E., Wichman, A., Winberg, S., Keeling, L., Løvlie, H., 2018.
12 Environmental complexity buffers against stress-induced negative judgement bias in
13 female chickens. *Sci. Rep.* 8, 5404. <https://doi.org/10.1038/s41598-018-23545-6>
- 14

1 **Figure captions**

2 **Figure 1**

3 Conceptual diagram presenting the main elements of a typical judgement bias study. a) The
4 basic task is trained using either a go/no-go, or active choice (go/go) design. b) Manipulations of
5 affective state usually, but not always, occur after training of the task and may be acute or
6 longer-term. c) Tests involve the standard training protocol plus the addition of occasionally
7 presented ambiguous cues whose properties are usually intermediate between the trained
8 positive and negative cues (NP = near positive cue, MID = intermediate between positive and
9 negative cue, NN = near negative cue). d) 'Optimistic' and 'pessimistic' responding to the cues is
10 inferred from the proportion of positive responses and/or the latency to make positive
11 responses, which are inversely related.

12 **Figure 2**

13 PRISMA flow diagram. Articles identified and number of articles included and excluded during
14 each screening stage.

15 **Figure 3**

16 Main characteristics of the included studies. Blue bars represent numbers of studies represented
17 in each level of categorical variables. Between one and 30 effect sizes were extracted per study
18 and the distribution of effect sizes generally follows the pattern of the presented data
19 aggregated to the study level (e.g. some studies reported data for only one sex, others reported
20 data for both sexes together, and 3 studies that included both sexes reported data for females
21 and males separately, shown here as 'female and male sep.'). Numbers do not add up to 71 for
22 some of the variables due to multiple experiments being present within some studies, or
23 complex experimental designs being used.

24 **Figure 4**

25 Forest-like (orchard) plots showing effect size (Hedges' g) estimates from meta-analyses on: a)
26 whole data set (all cues reported for judgement bias tests), and b-e) four subsets of this data set,
27 representing different ways of interpreting the judgement bias test results. Positive effect sizes

1 indicate a positive effect of affect manipulation treatments on judgement bias in a relatively
2 better condition compared to a relatively worse condition, i.e. affect manipulations working in
3 the expected direction. The effects are statistically significant when the thick horizontal error
4 bars (95% confidence intervals) do not cross zero. Thin horizontal whiskers indicate prediction
5 intervals. k is number of effect sizes. Dots represent individual effect sizes scaled proportionally
6 to their precision.

7 **Figure 5**

8 Forest plot showing mean effect size (Hedges' g) estimates from meta-regression analysis using
9 species identity as a moderator. Positive effect sizes indicate a positive effect of affect
10 manipulation treatments on judgment bias in a relatively better condition compared to a
11 relatively worse condition, i.e. affect manipulations working in the expected direction. The
12 effects are statistically significant when the horizontal error bars (95% confidence intervals) do
13 not cross zero. k is number of effect sizes, K is number of studies.

14 **Figure 6**

15 Forest plots showing effect size (Hedges' g) estimates from the univariate meta-regression
16 analyses (one moderator at a time) with potentially influential moderators. Effect sizes with
17 positive values indicate a positive effect of affect manipulations on judgement bias in a relatively
18 better condition compared to a relatively worse condition, i.e. affect manipulation treatment
19 working in the expected direction. The mean effects (black unfilled circles) for each group of
20 individual effect sizes (grey filled circles) are statistically different from zero when their
21 horizontal error bars (95% confidence intervals) do not cross zero. Thin horizontal whiskers
22 indicate prediction intervals. k is number of effect sizes. Dots represent individual effect sizes
23 scaled proportionally to their precision.