



Assessing theoretical conclusions with blinded inference to investigate a potential inference crisis

Journal:	<i>Advances in Methods and Practices in Psychological Science</i>
Manuscript ID	AMPPS-18-0147.R2
Manuscript Type:	Empirical Article
Date Submitted by the Author:	n/a
Complete List of Authors:	Starns, Jeffrey; University of Massachusetts Amherst, Psychological and Brain Sciences Cataldo, Andrea; University of Massachusetts, Amherst, Psychological & Brain Sciences Rotello, Caren; University of Massachusetts, Psychology
Substance Keywords:	memory
Method and Stats :	modeling < Bayesian, bootstrap
Additional Keywords:	metascience, blinded inference

Assessing theoretical conclusions with blinded inference to investigate a potential inference crisis

Running Head: Blinded inference

Jeffrey J. Starns^{*1}, Andrea M. Cataldo¹, Caren M. Rotello¹, Jeffrey Annis², Andrew Aschenbrenner², Arndt Bröder², Gregory Cox², Amy Criss², Ryan A. Curl², Ian G. Dobbins², John Dunn², Tasnuva Enam², Nathan J. Evans², Simon Farrell², Scott H. Fraundorf², Scott D. Gronlund², Andrew Heathcote², Daniel W. Heck², Jason L. Hicks², Mark J. Huff², David Kellen², Kylie N. Key², Asli Kilic², Karl Christoph Klauer², Kyle R. Kraemer², Fábio P. Leite², Marianne E. Lloyd², Simone Malejka², Alice Mason², Ryan M. McAdoo², Ian M. McDonough², Robert B. Michael², Laura Mickes², Eda Mizrak², David P. Morgan², Shane T. Mueller², Adam Osth², Angus Reynolds², Travis M. Seale-Carlisle², Henrik Singmann², Jennifer F. Sloane², Andrew M. Smith², Gabriel Tillman², Don van Ravenzwaaij², Christoph T. Weidemann², Gary L. Wells², Corey N. White², Jack Wilson²

¹Organizing Authors (University of Massachusetts, Amherst)

²Contributing Authors (Multiple Institutions)

*Correspondence to: jstarns@psych.umass.edu

Abstract: Scientific advances across a range of disciplines hinge on our ability to make inferences about unobservable theoretical entities based on empirical data patterns. Accurate inferences rely on both a) discovering valid, replicable data patterns, and b) accurately interpreting those patterns in terms of their implications for theoretical constructs. The replication crisis in science has led to widespread efforts to improve the reliability of research findings, but comparatively little attention has been devoted to the validity of inferences based on those findings. Using an example from cognitive psychology, we demonstrate a blinded inference paradigm for assessing the quality of theoretical inferences from data. Our results reveal substantial variability in expert judgements on the very same data, hinting at a possible *inference crisis*.

Data and materials availability: Data and analyses are available at https://osf.io/92ahy/?view_only=2f6d9b285c2d4e279f144b6fed363142.

1
2
3 Assessing theoretical conclusions with blinded inference to investigate a potential inference
4
5
6 crisis
7

8
9 At the most fundamental level, science is the process of creating, testing, and refining
10 ideas that explain and predict natural phenomena. Two core components are necessary for this
11 process to be effective: First, researchers must be able to produce reliable data patterns. Second,
12
13 5 researchers must be able to reach sound theoretical conclusions based on those patterns.
14
15 Scientists in a variety of fields have developed techniques to minimize failure in the first
16
17 component, that is, to correct the surprisingly high rate of unreliable data patterns reported in the
18
19 scientific literature, often referred to as the *replication crisis* (Open Science Collaboration,
20
21 2015). These techniques, including pre-registration (Miguel et al., 2014), an increased emphasis
22
23 on direct replication (Open Science Collaboration, 2015), and blinded analysis (MacCoun &
24
25 10 Perlmutter, 2015), are crucial for promoting reliable scientific findings. However, we suggest
26
27 that researchers looking to reform the scientific process should broaden the scope of their
28
29 investigation to assess whether researchers can make valid theoretical conclusions by analyzing
30
31 empirical outcomes. This broader perspective could reveal whether some fields suffer from an
32
33 *inference crisis*; that is, a situation in which researchers have a surprisingly high likelihood of
34
35 making incorrect theoretical conclusions even if they are working with reliable, replicable data
36
37 15 patterns (Rotello, Heit, & Dubé, 2015).
38
39
40
41
42
43
44

45
46 The most direct way to assess inference quality is to create data sets for which the correct
47
48 20 inferences are known and to determine whether researchers can discover these correct inferences
49
50 through blinded data analysis. This *blinded inference* procedure represents an extension of
51
52 blinding techniques already in common practice. As outlined in 1, blinding techniques applied
53
54 during data collection and analysis are used routinely to reduce the tendency of researchers
55
56
57
58
59
60

1
2
3 and/or participants to promote desired outcomes. Specifically, “blinded data collection” refers to
4
5 experimental designs that blind the experimental participant, the researcher, or both to the
6
7 assigned condition (e.g., placebo v. drug), minimizing the ability of these agents to change their
8
9 behavior according to their beliefs about the assigned condition. “Blinded analysis” techniques,
10
11 increasingly common in physics (MacCoun & Perlmutter, 2015), hide from the data analyst
12 5
13 either the true experimental condition from which each observation is drawn (e.g., scrambled
14
15 conditions) or the true value of the observation itself (e.g., addition of removable random noise),
16
17 thereby limiting the ability of analysts to promote desired outcomes with their analysis choices,
18
19 such as in the well-documented practice of *p*-hacking (Simmons, Nelson, & Simonsohn, 2011).
20
21
22
23
24 10 These blinding procedures are valuable tools to limit the malign effects of “researcher degrees of
25
26 freedom (*df*),” a term that describes the wide range of design and analysis choices researchers
27
28 can use to address the same research question (Simmons et al., 2011). A recent study (Silberzahn
29
30 et al., 2018) highlighted the influence of researcher degrees of freedom by sending the same data
31
32 set to 29 teams of researchers and asking each team to determine whether soccer referees
33
34 disproportionately “red-card” darker-skinners players. The results showed substantial variability
35 15
36
37 in analysis techniques and conclusions across the research teams.
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

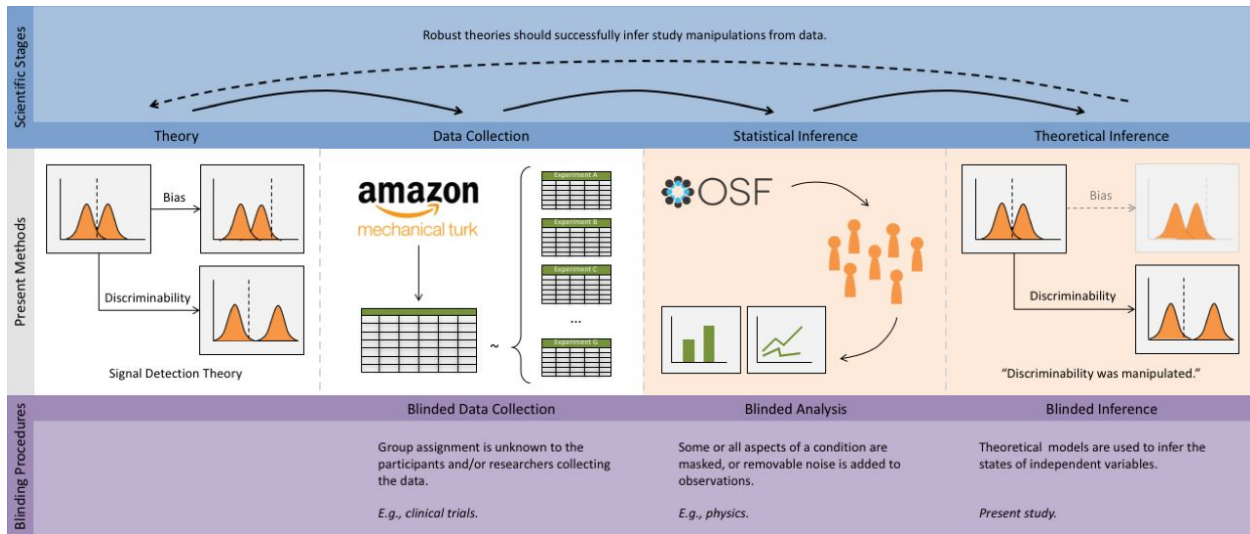


Fig. 1. Diagram of the scientific process. The top panel denotes the main stages. The middle panel outlines the methods used in the present paper. The bottom panel denotes common blinding techniques applied in each of the scientific stages, including the blinded inference paradigm advocated for in the present paper.

These blinding methods are excellent strategies to limit the influence of researcher degrees of freedom and/or to assess the consistency of inferences across researchers, but they do not address the *validity* of those inferences. This extra step is crucial because researchers might make inference errors even if they are not promoting a desired outcome with their analysis choices, and these errors could be consistent across researchers who make similar choices (for examples, see Rotello et al., 2015). To assess the validity of theoretical inference, we advocate widespread use of a blinded inference design to supplement traditional approaches. In such a design, researchers who are blinded to condition assignment make inferences about the state of independent variables that are linked to theoretical constructs. Our characterization of the blinded inference technique is heavily influenced by a recent study by Dutilh et al. (2018) in which condition-blinded data sets were sent to response-time modelers who were asked to infer whether the conditions differed in terms of psychological constructs such as response caution and evidence strength. Our general characterization of the blinded inference approach relies on

1
2
3 Dutilh et al.'s innovative design with two modifications: (1) analysts should be asked to make
4 inferences about empirically manipulated factors rather than latent constructs so that the correct
5 inferences can be unambiguously defined, and (2) analysts should be required to communicate
6 the level of uncertainty associated with their inferences in terms of a probability distribution.
7
8
9

10
11
12 5 As characterized here, blinded inference can be used in any scenario in which researchers
13 claim that they can (a) measure a theoretical construct based on data patterns and (b) manipulate
14 that theoretical construct with independent variables. If both of these claims are true, then
15 researchers should be able to make accurate inferences about the state of independent variables
16 specifically linked to the theoretical construct by analyzing data. If researchers fail in this task,
17 then it suggests that at least one of the claims is false, i.e., researchers either lack valid
18 techniques for measuring the theoretical construct, lack valid ways to manipulate it, or both. In
19 turn, failures to validly measure theoretical constructs could arise from a variety of problems.
20 One class of problems applies to the process of selecting a measurement model to map patterns
21 of data to underlying processes. Different models might suggest different inferences even if they
22 have a similar ability to match observed data patterns. Another class of problems applies to the
23 process of applying the model, and includes malign factors like parameter estimation biases and
24 mishandling of data.
25
26
27
28
29
30
31
32
33
34
35 15
36
37
38
39
40
41

42 Concretely, consider a famous example: Mendel and his peas. Mendel recorded
43 systematic patterns of variables, i.e., the relationship between the traits of parents and offspring,
44 and linked them to unobservable theoretical constructs, i.e., hereditary "factors" that obeyed
45 certain laws. His data have been described as being too clean, with too few extreme observations,
46 which may be a result of "unconscious bias in classifying ambiguous phenotypes, stopping the
47 counts when satisfied with the results, recounting when results seem suspicious, and repeating
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 experiments whose outcome is mistrusted” (Hartl & Fairbanks, 2007). Thus, Mendel’s
4
5 conclusions might represent the first documented case of *p(ea)*-hacking. Clearly, Mendel would
6
7 have benefitted from using blinded analyses to eliminate researcher biases, but we wish to
8
9 demonstrate how he could have gone further.
10
11

12 5 By applying his theory of genetics, Mendel claimed to be able to (a) measure underlying
13
14 heritable factors by evaluating the phenotype of a plant and (b) manipulate heritable factors
15
16 in offspring by selecting parents with certain phenotypes. These are precisely the claims related
17
18 to the validity of theoretical inference that can be tested in a blinded inference paradigm. For
19
20 example, someone could have given Mendel a number of plants produced by mating parents with
21
22 certain traits (unknown to Mendel) and asked him to use his laws of heritability to predict the
23
24 10 likely traits of the *parent* plants by interpreting the traits of the offspring. Mendel would not have
25
26 been able to make perfect inferences, of course, given that some phenotypes can be produced by
27
28 multiple genotypes, but he should have been able to make substantially more accurate inferences
29
30 than someone without a valid theory linking the phenotypes of parents and offspring. We claim
31
32 that a procedure like this one would have provided a more compelling demonstration of the
33
34 15 predictive value of Mendel's laws than unblinded data that could be "massaged." Moreover, by
35
36 revealing specific offspring phenotypes for which the parents’ phenotypes were particularly
37
38 difficult to predict accurately, it might have allowed the limitations in Mendel’s basic theory to
39
40 be identified more quickly.
41
42
43
44
45
46

47 20 Many modern scientists share with Mendel the challenge of making inferences about
48
49 theoretical constructs on the basis of indirect evidence. For example, modern geoscientists infer
50
51 the composition and dynamics of Earth’s interior from a variety of indirect methods, including
52
53 radar and magnetic fields. Likewise, cosmologists have inferred that dark matter exists in the
54
55
56
57
58
59
60

1
2
3 absence of direct observation. In the authors' discipline, cognitive processes are inferred from
4
5 observable behaviors such as decision accuracy or response times. Thus, a critical step in
6
7 establishing the validity of many scientific claims is to test the inferential power of the data, and
8
9 this is precisely what the blinded inference procedure achieves: If the researcher is blind to the
10
11 nature of the manipulation(s), conclusions about what experimental factor was manipulated
12 5
13
14 depend entirely on the data and not on the expectations or unconscious biases of the researcher.
15
16

17 In what follows, we demonstrate the blinded inference paradigm using an example study
18
19 from recognition memory research. The scheme in the middle of 1 summarizes the design. We
20
21 conducted a study in which we sent recognition memory researchers ("contributors") seven data
22
23 sets generated with common experimental manipulations and asked them to make inferences
24 10
25
26 about memory performance. In a recognition memory task, participants are asked to indicate
27
28 whether they previously encountered a stimulus (often a word) in a certain context (typically a
29
30 study list). A common question is whether, and to what extent, an independent variable produces
31
32 changes in discriminability (the ability to distinguish stimuli that were and were not seen in the
33
34 target context), and in many cases this determination is obscured by differences in response bias
35 15
36
37 (the overall predilection for saying "studied"). Signal detection theory (SDT; Macmillan &
38
39 Creelman, 2005) was developed in the 1950s with the goal of separating discriminability and
40
41 bias, and SDT-based measures have been in common use throughout psychology and other
42
43 disciplines ever since. Several other models or measurement techniques have been developed as
44
45 alternatives to SDT (Ratcliff, 1978; Riefer & Batchelder, 1988), and some of these also achieved
46
47 20
48
49 wide popularity throughout psychology (e.g., Erdfelder et al., 2009). Thus, researchers have had
50
51 nearly seven decades to hone their ability to distinguish discriminability and bias as theoretical
52
53 constructs, and thousands of papers have been published using models and measures that claim
54
55
56
57
58
59
60

1
2
3 to be able to do so. We tested published memory researchers on their ability to detect whether
4
5 memory discriminability varied between experimental conditions that might have also varied in
6
7 terms of response biases.
8
9

10
11 We have two primary research questions: First, how variable are inferences across
12 5 researchers? Finding high variability across researchers would be unsettling, given that they all
13
14 analyzed the same data. Second, and more importantly, how accurate are researcher inferences?
15
16 If recognition memory researchers have effective methods for manipulating and measuring
17
18 discriminability and bias based on seven decades of investigating these constructs, then they
19
20 should be able to make accurate inferences about whether conditions come from the same level
21
22 or from different levels of a discriminability manipulation.
23
24 10

25
26 To preview, we found surprisingly high variability in the inferences of memory
27
28 researchers asked to interpret the same data, and we also found that many researchers made more
29
30 inferential errors than would be expected from sampling variability in the data. Given that our
31
32 task required a relatively simple inference, we suspect that this pattern of surprisingly low
33
34 inferential accuracy is likely to be found in other research areas. Broadly, however, we
35 15
36
37 emphasize key positive outcomes of this study. Our study exemplifies scientists' commitment to
38
39 improving the research process, in that many respected memory researchers had the courage to
40
41 put their conclusions to a public test. Moreover, despite the troubling error rate of the group, our
42
43 framework identified multiple researchers as having made highly accurate inferences. We
44
45 therefore believe that our study demonstrates a promising methodology for the future goal of
46
47 20
48
49 improving inference quality by identifying best practices.
50
51
52
53
54
55
56
57
58
59
60

Methods

Experimental Design

There were two main phases of data collection. In Phase 1, we collected experimental data in a large-scale recognition memory experiment that used standard study materials and included orthogonally-varied factors known to influence memory discriminability and response bias.¹ The between-subjects design of Phase 1 is analogous to any comparison of memory performance between a special population (e.g., Alzheimer's patients) and a control group, except that our participants were randomly assigned to conditions. In Phase 2, subsets of the full data set were selected to generate seven two-condition experiments in which only the factor affecting discriminability varied (2 experiments), only the factor affecting response bias varied (2 experiments), both factors varied (2 experiments), or neither varied (1 experiment). The conditions in these seven experiments were masked and the data were shared with researchers who had published papers investigating recognition memory, and these experts (or "contributors") were asked to rate the probability that each experiment had only a memory discriminability manipulation, only a response bias manipulation, both, or neither. Contributors were not told how many experiments of each type were included in the data sets, and they were free to select their preferred strategy for distinguishing memory discriminability and response bias.

Phase 1

¹ All study procedures were approved by the Institutional Review Board at the University of Massachusetts Amherst.

1
2
3 **Participants.** A total of 459 participants were recruited through Amazon’s Mechanical
4 Turk (Buhrmester, Kwang, & Gosling, 2011) using psiTurk (Gureckis et al., 2016). Participants
5 earned \$1.00 for completing the experiment.
6
7

8
9
10 **Materials.** The experiment utilized 104 high-frequency (at least 100 occurrences/million
11 in Kučera & Francis, 1967) English nouns that were 3-7 letters long. Four words were used in
12 5 the practice block, and the remaining 100 were equally divided into two study lists, A and B.
13
14 the practice block, and the remaining 100 were equally divided into two study lists, A and B.
15
16 Participants were randomly assigned to study either list A or list B. All participants were tested
17 on the combined list of all 100 words, resulting in complete counterbalancing of stimulus status
18 (studied or unstudied) across participants.
19
20
21
22

23
24 10 **Procedure.** The experiment was coded in javascript using the jsPsych library (de Leeuw,
25 2015). Participants were given detailed instructions that included comprehension checks for key
26 components, and they completed a brief practice block before beginning the main task. Word
27 order in the study and test phases was independently randomized for each participant. On each
28 trial of the study phase, participants were asked to report whether the presented word represented
29 an animate object. All of the stimulus words represented clearly animate or inanimate objects, as
30 judged by four independent raters. Each word remained on the screen until the participant
31 entered a response for the animacy question. On each trial of the test phase, participants were
32 first asked to report whether or not they had seen the presented word in the study phase.
33
34
35 15 Participants were then asked to report how confident they were in their response on a 1-3 scale,
36 in which a “1” meant “Not Sure” and a “3” meant “Very Sure”. All responses were made via key
37 press, and participants were asked to balance speed and accuracy throughout the experiment.
38
39
40
41
42
43
44
45
46
47 20
48
49
50

51 Memory discriminability and bias were manipulated between participants.
52
53
54 Discriminability was manipulated by varying the number of times each word was presented in
55
56
57

1
2
3 the study phase (1, 2, or 3). Bias was manipulated by instructing participants to avoid making
4 particular kinds of errors in the test phase. Specifically, conservative participants were told to
5 particularly avoid false alarms (“old” responses to unstudied items), liberal participants were told
6 to particularly avoid misses (“new” responses to studied items), and neutral participants were
7 told to avoid both errors equally. This manipulation was reinforced by varying the quality of the
8 error feedback in the test phase, such that conservative participants saw a “BAD ERROR!”
9 message after false alarms and standard “ERROR” message after misses, liberal participants saw
10 a standard “ERROR” message after false alarms and a “BAD ERROR!” message after misses,
11 and neutral participants saw a standard “ERROR” message in both cases. The “BAD ERROR!”
12 message was accompanied by a reminder of the type of error to particularly avoid and was
13 presented longer than the standard message (2500ms vs. 500ms).
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28

29 **Phase 1 results.** Complete data are available at the OSF site. A summary of the data
30 analyzed in each of the seven experiments appears in Table 1. We offer no statistical
31 interpretation of these data, given our goal of crowd-sourcing that interpretation in Phase 2
32 (described next). However, we note that the outcome of this experiment is very consistent with
33 decades of recognition memory literature. For example, hit rates increased and false alarm rates
34 decreased with repeated learning opportunities (as in, e.g., Lachman & Field, 1965; Ratcliff,
35 Clark, & Shiffrin, 1990; Stretch & Wixted, 1998; Verde & Rotello, 2007). We also observed
36 typical effects of response bias manipulations: both hit and false alarm rates tended to increase as
37 increasingly liberal responding was encouraged (e.g., Dube, Starns, Rotello, & Ratcliff, 2012;
38 Han & Dobbins, 2009; Starns, Hicks, Brown, & Martin, 2008; Swets, Tanner, & Birdsall, 1961)
39 and the effects of bias appeared weaker when encoding strength was greater (e.g., Ratcliff, Sheu,
40 & Gronlund, 1992).
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Experiment	Condition	Discriminability	Bias	N	Hit Rate	False Alarm Rate
A	1	3	Liberal	24	.873	.239
	2	3	Conservative	24	.875	.126
B	1	1	Liberal	27	.865	.266
	2	2	Conservative	25	.840	.191
C	1	2	Neutral	27	.861	.205
	2	3	Neutral	24	.911	.174
D	1	1	Neutral	27	.781	.256
	2	1	Conservative	26	.739	.195
E	1	1	Conservative	26	.742	.192
	2	3	Neutral	24	.815	.190
F	1	1	Liberal	26	.812	.287
	2	3	Liberal	26	.935	.164
G	1	2	Liberal	26	.847	.208
	2	2	Liberal	26	.913	.208

Notes: Discriminability represents the number of times each target word was presented in the study phase (1, 2, or 3). Liberal and conservative biases refer to instructions to particularly avoid missing studied items and false alarms to unstudied memory probes, respectively, in the test phase; neutral bias emphasized both errors equally. N indicates sample size, and hit and false alarm rates indicate the proportion of correct and erroneous “old” judgments.

Table 1. Definition and summary statistics of the seven experiments sent to contestants.

Phase 2

Participants. Contributors were recruited through targeted e-mails to researchers with a background in recognition memory and/or models of memory and decision making. These individuals were encouraged to forward our invitation to other experts. Out of the 121 researchers who were initially contacted, a total of 46 contributors (comprising 27 PIs and 19 members of their labs) submitted analyses. The data were available in two phases, one for which the confidence-rating data were withheld and another that included the confidence ratings. The purpose of the phases was to investigate whether or not confidence ratings improved inference

1
2
3 quality. Of the 27 groups of contributors, 14 also submitted new analyses when the confidence
4 rating data were released. Two contributors declined authorship, and their inferences are de-
5 identified. Of the 44 contributors who accepted authorship, 33 (representing 19 labs) opted to
6 have their inferences associated with their identities; the others chose to remain anonymous. The
7
8
9
10
11
12 5 27 PIs had an average of 14.7 years of post-Ph.D. experience.

13
14
15 **Materials.** Subsets of data collected in Phase 1 were sampled to form seven
16
17 “experiments” for the contributors to analyze, summarized in Table 1. Each experiment was
18
19 designed to have two between-participant conditions that differed in terms of either a memory
20
21 discriminability manipulation, a response bias manipulation, both, or neither. The data for each
22
23
24 10 condition were created by taking separate random samples of participants who studied list A and
25
26 participants who studied list B and combining them. Each condition had either an equal number
27
28 of participants from the two lists or very close to equal (off by one). The data sets that
29
30 contributors received for the binary analyses included data from the test phase with variables for
31
32 participant ID, condition (1 or 2), study list (A or B), trial (1-100), test word, whether or not the
33
34
35 15 tested word had been studied (target or lure), the participant’s binary response (“old” or “new”),
36
37 and response time for the binary response. The data sets that contributors received for the
38
39 confidence rating analyses additionally included the participant’s confidence rating, both on the
40
41 original 1-3 scale and on a recoded 1-6 scale that ranged from “Very Sure New” to “Very Sure
42
43
44
45 Old”, and response time for the confidence rating response.

46
47 20 Each contributor completed a submission template summarizing their analyses (see OSF
48
49 site for an example). The template asked contributors to report the authors collaborating on the
50
51 submission, accept or decline authorship, and indicate whether they would prefer their
52
53 conclusions be de-identified. Contributors were then asked to provide a description of their
54
55
56
57
58
59
60

1
2
3 process for analyzing the data in sufficient detail for external replication, a description of any
4 exclusion criteria that were applied, and any code that they were comfortable sharing. All shared
5 code is available at the OSF site. Contributors were lastly asked to report four probabilities for
6 the four possible types of experiment; namely, experiments for which the two conditions were
7 from (1) different levels of a memory strength (discriminability) manipulation but not different
8 levels of a bias manipulation, (2) different levels of a bias manipulation but not different levels
9 of a memory strength manipulation, (3) different levels of both a memory strength and a bias
10 manipulation, or (4) the same levels of memory strength and bias (i.e., null data sets).

11
12 **Procedure.** Materials for the binary and confidence rating data analyses were posted to
13 separate private OSF pages. The materials for the binary data analyses were made accessible to
14 contributors on July 7, 2017 and analyses were due August 31, 2017. The materials for the
15 confidence rating data analyses were made accessible on September 9, 2017 and analyses were
16 due on November 1, 2017. No changes to the binary data contributions were allowed after the
17 confidence rating data were released. To support the independence of contributors' inferences,
18 all communication of the coordinating team with contributors was conducted via individually-
19 generated emails, contributors' identities were not shared until mid-November of 2017, and
20 contributors were strongly discouraged from discussing their interpretations of the data with one
21 another in case they accidentally discovered their common participation.

22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 **Results**

49
50 Our response format was designed to highlight the fact that contributors needed to
51 distinguish the effects of discriminability and bias, but we are primarily interested in conclusions
52 about whether there was a discriminability manipulation. A wide range of research questions in
53
54
55
56
57
58
59
60

1
2
3 the recognition memory literature require conclusions about discriminability, whereas bias is
4 more often considered a “nuisance” process. Moreover, focusing on discriminability gives our
5 contributors the best chance to succeed because discriminability is better understood and less
6 theoretically contentious than bias (Macmillan & Creelman, 2005). To isolate discriminability
7 inferences, we collapsed the “memory alone” and “both” categories to represent the reported
8 probability of a discriminability manipulation and the “bias alone” and “neither” categories to
9 represent the reported probability of no discriminability manipulation (see OSF for bias results,
10 which unsurprisingly showed poorer inference performance than the discriminability results).

11
12 Fig. 2A shows histograms of the reported probability of a discriminability manipulation
13 across contributors for each of the seven experiments, with regions reflecting correct and
14 incorrect inferences marked in green and red, respectively. The most striking finding shown in
15 Fig. 2A is the extremely high variability across contributors, with responses spanning a wide
16 range of probabilities for all experiments. For example, some contributors reported a 0% chance
17 that the conditions in Experiment A came from different levels of a memory discriminability
18 manipulation, some reported a 100% chance, and the rest follow an essentially uniform
19 distribution of probability estimates between these two extremes. Responses were concentrated
20 on the correct side for some experiments (e.g., D, F), but not for others (A, B). The high level of
21 variability is surprising given that all researchers received the same data sets. Note that Fig. 2B,
22 addressed in greater detail below, shows the data that informed the researchers’ inferences,
23 namely the proportion of studied and non-studied items called “studied” (or the “hit rate” and
24 “false alarm rate” in signal detection terms). The dark symbols show results with no participants
25 or trials excluded and grey symbols show results of applying the exclusion criteria used by each
26 contributor. A priori, some experiments seemed likely to be easier to interpret, for example,

1
2
3 when both the hit and false alarm rate effects were large and consistent with the same theoretical
4 inference (e.g., in Exp. F, the higher hit rate and lower false alarm rate for Cond. 2 both indicate
5 higher memory discriminability in this condition).
6
7
8
9

10 The variability in inferences was matched by high variability in the analysis methods
11 selected by our contributors. These methods, identified on the y-axis of Fig. 2D and described in
12 5 the Supplemental Materials, are purportedly capable of distinguishing memory discriminability
13 and response bias. Within most of these techniques, some contributors used traditional
14 frequentist statistical methods (e.g., maximum likelihood estimation, significance tests) and
15 others used Bayesian methods (e.g., posterior distributions of parameters or model selection via
16 Bayes Factors). When all analysis choices were considered, no two contributors used exactly the
17 same analysis approach (e.g., same exclusion criteria, measurement technique, and statistical
18 approach).
19
20
21
22
23
24 10
25
26
27
28
29

30 To summarize inferential accuracy, we counted the number of times across experiments
31 that each contributor reported the true discriminability effect status as the most likely outcome,
32 that is, reported a greater than 50% chance of a discriminability manipulation when
33 15 that is, reported a greater than 50% chance of a discriminability manipulation when
34 discriminability was in fact manipulated or reported a less than 50% chance of a discriminability
35 manipulation when it was not. A histogram of these results appears in Fig. 2C. Slightly over half
36 of the contributors performed well by this measure, correctly describing five or six of the seven
37 data sets, but the other contributors performed more poorly. We note that the contributor with
38 zero correct inferences estimated a 50% chance of a discriminability manipulation for every
39 experiment, so in fairness, this contributor did not make any *incorrect* inferences either.
40
41
42
43
44
45
46
47 20
48
49
50
51
52
53
54
55
56
57
58
59
60

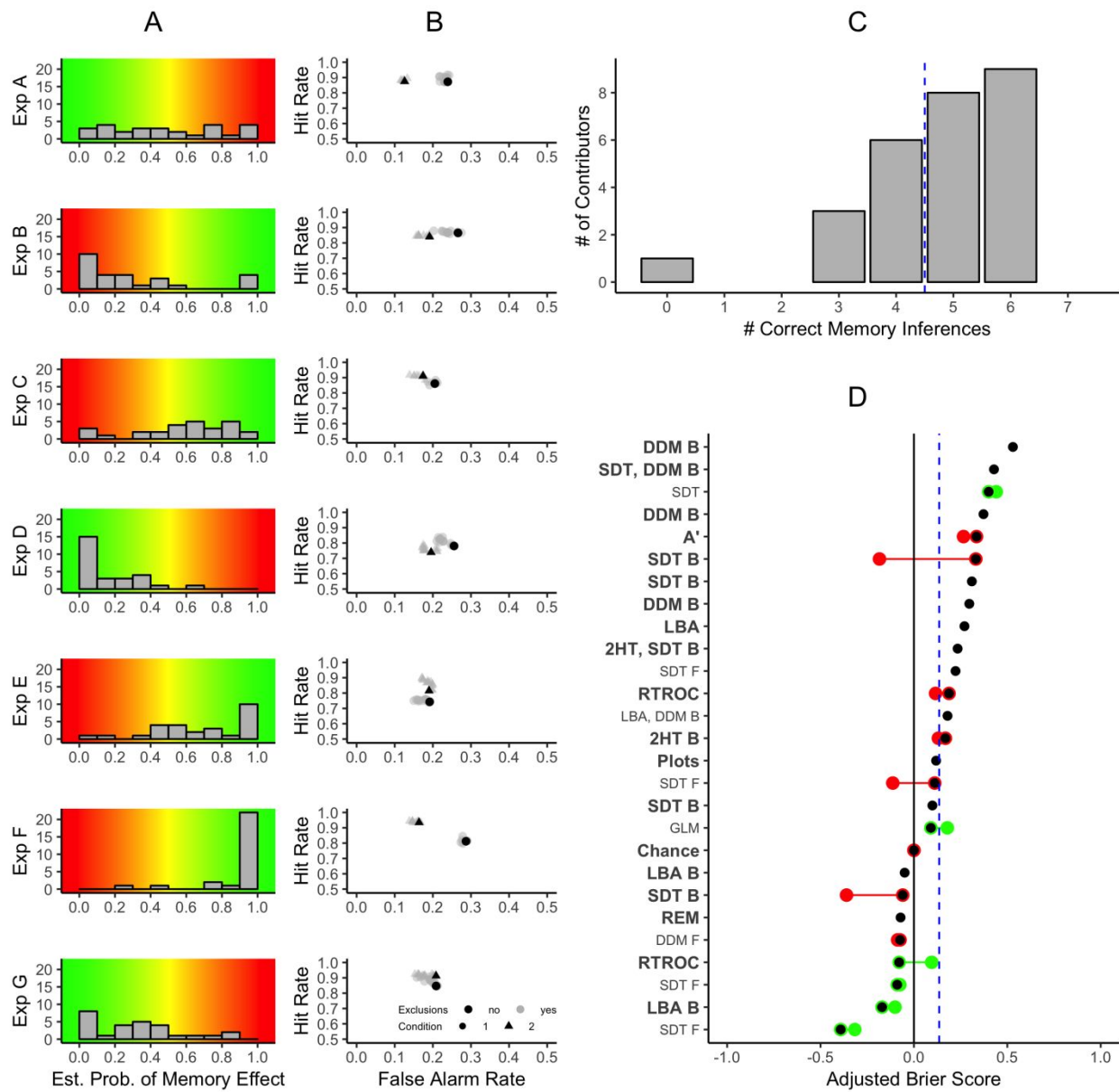


Fig. 2. Discriminability inference performance. Panel A: Histograms of contributors' estimated probabilities of an effect in each of the seven experiments. Red denotes incorrect estimates whereas green denotes correct estimates. Panel B: Hit and false alarm rates for each of the seven experiments. Black points represent original values. Grey points represent values after applying each contributor's specified exclusion criteria. Panel C: Histogram of the number of correct inferences out of the seven experiments analyzed for each contributor. The blue dashed line denotes simulation-based benchmark for reasonable performance. Panel D: The adjusted Brier score for each contributor, labelled by their chosen method of analysis. (Note that 19 contributor groups, highlighted in bold, were willing to have their names associated with their responses. The OSF page includes a figure that identifies these contributors.) Black points represent scores for the binary data analysis. Red and green points represent scores for the data analysis with confidence ratings where performance decreased or increased, respectively. The black vertical line denotes chance performance; The blue dashed line denotes simulation-based benchmark for

1
2
3 reasonable performance. Labels on the y-axis denote analysis strategies (defined in the
4 Supplemental Materials) and statistical choices (B = Bayesian; F = frequentist).
5
6
7

8 Even a valid inference procedure will sometimes reach inaccurate conclusions due to
9
10 5 sampling variability, so we needed to identify a benchmark accuracy level below which it would
11
12 be reasonable to conclude that an invalid inference technique had been applied. We performed
13
14 model simulations to identify this benchmark. In the simulations, we generated data sets by
15
16 randomly sampling data from a signal detection model and analyzing those data sets with
17
18 measures derived from the same model (see the Supplementary Materials for details). Each
19
20 simulated data set contained the same type of information as the data sets sent to contributors
21 10
22 with no labeling to identify the experimental manipulation. Thus, the simulation code performed
23
24 blinded inference just like our contributors. The key difference between the simulation code and
25
26 the contributors' analyses is that the former uses an inference procedure that is known to be valid
27
28 (i.e., consistent with the process that generated the data), so the results represent expected
29
30 performance levels when sampling variability is the only source of inaccuracy. We set
31 15
32 performance benchmarks such that only 10% of the simulated studies fell below the value,
33
34 meaning that performance is rarely that bad when a valid inference method is applied.
35
36
37
38
39

40 The benchmark for number correct is indicated with a dashed line in Fig. 2C. Nearly half
41
42 of the contributors fell below this benchmark, suggesting that some aspect of their inference
43
44 method was ineffective. To assess whether our empirical data sets were a particularly misleading
45 20
46 sample (like the 10% of simulated data sets that produced accuracy below our benchmark even
47
48 when a valid inference technique was applied), we used the analysis technique from the
49
50 simulation on the actual data sets sent to contributors and obtained correct inferences for 6 of the
51
52 7 data sets. Thus, the empirical data sets do not seem to be a "bad" or misleading sample.
53
54
55
56
57
58
59
60

1
2
3 Scientists should be able to express appropriate degrees of certainty in their conclusions,
4
5 so we also assessed accuracy with a measure that is sensitive to the contributors' probability
6
7 estimates: the Brier score (Brier, 1950). Brier scores compute the variance between the predicted
8
9 probability that an outcome will occur and the actual outcome (coded as a 0 or 1). In our case,
10
11 the outcome is whether or not the two conditions in an experiment come from different levels of
12 5
13 a discriminability manipulation. Therefore, the best possible performance is produced by
14
15 reporting a 0% predicted chance of a discriminability manipulation for all data sets without a
16
17 discriminability manipulation and a 100% predicted chance for all data sets with a
18
19 discriminability manipulation, the worst possible performance is the converse, and "chance"
20
21 performance means reporting a 50% chance for all data sets (meaning that estimates provide no
22
23 information about which data sets have discriminability manipulations). We adjusted our Brier
24 10
25 scores such that 0 represents chance performance, 1 represents the best possible performance,
26
27 and -1 represents the worst possible performance (see the Supplementary Materials for details).
28
29 In our simulations to explore performance levels for a valid inference technique, the median
30
31 adjusted Brier score was .44 and 10% of scores fell below .13, which will thus serve as our
32
33 benchmark for problematic inferences. Applying the analysis technique from the simulations to
34
35 15
36 the empirical data sets sent to contributors produced a Brier score of .38, which is well above our
37
38 benchmark.
39
40
41
42
43

44 Fig. 2D shows ranked Brier scores for our contributors (contributions are labeled by their
45
46 inference technique). The contributor who reported 50% for every data set is on the chance line.
47 20
48 Although this contributor returned no correct inferences, their probability estimates
49
50 outperformed about one-third of contributors in terms of Brier scores. The contributors who are
51
52 below chance made multiple incorrect inferences with high confidence levels; in other words,
53
54
55
56
57
58
59
60

1
2
3 their reported probabilities provided *misinformation* as to which data sets were likely to have a
4 discriminability manipulation. Roughly half of contributors were below the benchmark for
5 problematic inferences, shown by the dashed vertical line, demonstrating that researchers fairly
6 commonly made the mistake of being inappropriately confident in their incorrect inferences.
7
8 Reassuringly, some contributors achieved Brier scores that are basically as high as can be
9 expected given sampling variability in the data, suggesting that they applied appropriate
10 inference methods. Given the poor overall performance, one might wonder whether these high-
11 performing contributors were simply lucky, indicating that none of our contributors truly
12 succeeded in the inference task. The Supplementary material includes analyses that strongly
13 support the conclusion that at least some of our contributors applied valid inference procedures.
14
15

16
17 Inference errors were not associated with the choice of any particular analysis technique.
18 The *y*-axis of Fig. 2D reveals no clear pattern. Methods used by multiple contributors tend to be
19 distributed among the top, middle, and bottom rankings, as are techniques relying on frequentist
20 and Bayesian approaches. Our simulation results also showed that inferences about
21 discriminability are generally robust to different measurement methods, at least for data patterns
22 similar to those in our experiments. Specifically, we reanalyzed all of the simulated data sets
23 using a different measure of discriminability ($P_r = \text{hit rate} - \text{false alarm rate}$) that is consistent
24 with a different class of models (Pazzaglia, Dube, & Rotello, 2013; Snodgrass & Corwin, 1988)
25 than the data-generating signal detection model. The P_r analyses achieved accuracy levels that
26 were well above our benchmarks for problematic inferences in terms of number correct and Brier
27 scores (see the Supplementary materials for details). P_r depends on different processing
28 assumptions than the signal-detection model used to sample the simulated data sets, but the two
29 models often make similar discriminability inferences for data set like the ones we sent to
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 contributors (inferences start to diverge for data sets that have large bias effects, but our bias
4 effects were moderate). Thus, it is possible to make appropriate conclusions about
5
6 discriminability when using a measurement model that does not exactly match the processes
7
8 generating the data, and selecting an incorrect measurement model cannot entirely explain the
9
10
11
12 5 poor inference performance revealed in Figure 2.

13
14
15 Variability in inferences was not predictable from contributors' rules for censoring data.
16
17 Recall that the grey symbols in Fig. 2B show the mean hit and false alarm rates for each
18
19 condition with the exclusion criteria used by each contributor. Although these censoring rules
20
21 clearly resulted in different hit and false alarm rates, we were unable to identify any systematic
22
23 relationship between these rules and inference accuracy. Moreover, seven contributors did not
24 10 exclude any data, yet they used different analytic tools and reached different conclusions about
25
26 the probability of a discriminability effect.

27
28
29
30
31 Theoretically, discriminability and bias effects are more easily distinguished with
32
33 receiver operating characteristics (ROCs) formed from confidence-rating data than with binary
34
35 15 old/new response data (Rotello et al., 2015). In a second round of blinded inference, we re-sent
36
37 the data sets with an addition column for the reported confidence level on each trial, and 14
38
39 contributors offered new probability ratings based on the ROCs in each experiment. The
40
41 resulting Brier scores appear in Fig. 2D with lines to mark the difference from the corresponding
42
43 Brier scores based on the binary-response data. The largest changes were actually negative,
44
45 reflecting reduced inferential accuracy with ROC data.
46
47 20
48
49
50

51 Discussion

52
53
54
55
56
57
58
59
60

1
2
3 Distinguishing memory discriminability effects from bias effects is a common empirical
4 issue for recognition memory researchers that has important theoretical and practical
5 implications; for example, understanding memory processes in a special population (e.g., older
6 adults) hinges on the ability to determine if differences from a control group reflect a memory
7 discriminability effect. The available tools to interpret discriminability are well-established, and
8 some have been in use for nearly 70 years (Macmillan & Creelman, 2005). Despite these truths,
9 our expert contributors had mixed success when faced with the task of inferring whether
10 discriminability had been manipulated across conditions that might have also had different levels
11 of response bias. Strikingly, the reported probability of a discriminability effect was highly
12 variable across contributors even though they all received the same data sets. One natural
13 interpretation of these results is that the data themselves were too noisy to allow clear inference.
14 Our simulations are inconsistent with that conclusion: 90% of simulated sets of experiments
15 yielded five or more (of seven possible) correct inferences about discriminability. Thus, we view
16 the outcome of this blinded inference study as a challenge to recognition memory researchers;
17 one which should result in a re-evaluation of our methods, and in humbler presentation of future
18 conclusions that rely on the ability to distinguish discriminability and bias effects. The fact that
19 we found generally low inference quality when researchers used decades-old analysis tools
20 shows that the normal practice of science is not sufficient to ensure effective analysis techniques.
21 Indeed, some examples of systematically problematic inferences have survived decades of
22 scientific review, to the detriment of theoretical progress in those domains (see, e.g., Dube,
23 Rotello, & Heit, 2010, for a specific example and Rotello et al., 2015, for a more general
24 treatment). Widespread use of the blinded inference procedure will help to quickly identify these
25 inference problems and refine analysis methods to optimize inference quality.
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 **Constraints on generality.** Our study only provides information about a single research
4 scenario – assessing discriminability changes based on recognition memory data – but the fact
5 that we found surprisingly low accuracy for this relatively simple inferential task suggests that
6 problematic inference procedures may plague a broad range of research domains. However, these
7 different domains must be assessed individually in future work, and our results should not be
8 used to make general conclusions about general validity of scientific research. Even within the
9 field of recognition memory, our results are only directly troubling for studies that attempt to
10 make conclusions about discriminability and bias when both processes can potentially vary.
11 Although this is an unavoidable situation for some research questions (e.g., comparing memory
12 across different populations), for other questions memory researchers can substantially simplify
13 the inferential process by experimentally controlling bias when evaluating discriminability, or
14 vice versa. Moreover, memory researchers use a wide range of different types of paradigms and
15 data beyond the recognition tasks that we investigated.

16
17 The blinded inference paradigm demonstrated here is also not a substitute for good theory
18 testing and development. A theory that makes correct assumptions could perform poorly in
19 blinded inference based on limitations in the analysis tools available to implement the
20 measurement properties of the theory, and a theory that makes incorrect assumptions might
21 nevertheless serve as a useful tool in some situations (e.g., Newton's Laws are sufficient for
22 many applications despite being incomplete). Our results show that inference problems are not
23 limited to particular theoretical approaches in recognition memory: even researchers who relied
24 on the same measurement model were highly variable in their inferences. Good theory
25 development should run on several parallel tracks simultaneously – empirical assessment,

1
2
3 quantitative modeling or analysis, and, we argue, blinded inference studies – to establish that
4 applications of the theory can truly measure what they are intended to measure.
5
6

7
8 Another potential limitation of our results is that contributors might have applied
9
10 different analysis standards for our project than they would in a “real” study conducted in their
11
12 5 labs. We cannot rule out the possibility that our contributors might have made better inferences if
13
14 they were analyzing their own data for their own purposes, but there are many good reasons to
15
16 consider this unlikely. The vast majority of our contributors elected to be co-authors on this
17
18 manuscript, and a majority (19/27) agreed to have their name directly linked to their performance
19
20 level in presentations and publications (note that while inference methods were used as labels in
21
22 10 Figure 2, results identified by contributor are available on OSF). Thus, one could argue that our
23
24 contributors had a stronger incentive for rigor compared to typical studies in which no one is
25
26 likely to re-run the analyses and conclusions are never compared to an “answer key.” Indeed, our
27
28 contributors generally displayed a remarkable level of motivation and dedication to the project,
29
30 with some applying state-of-the-art techniques like hierarchical Bayesian modeling and/or
31
32 15 analyzing the data with multiple measurement models to inform their conclusions. Moreover, the
33
34 majority of contributors (14/27) agreed to make their analysis code publicly available (see OSF).
35
36 Thus, we are confident that the inference problems that we observed are not based on a simple
37
38 lack of effort, and although we cannot rule out the possibility that some contributors made
39
40 careless, easily correctable mistakes, we seriously doubt that these mistakes can fully explain the
41
42 inference problems that we observed.
43
44
45
46
47 20

48
49 **Comparison to similar studies.** Our results are similar to those of Silberzahn et al.
50
51 (2018) in that both reveal high variability in inferences across contributors who all received the
52
53 same data. In many ways, though, the high variability in our contributors’ inferences is even
54
55
56
57
58
59
60

1
2
3 more surprising – and troubling – given that our inference task represented a fairly common
4
5 research scenario. Whereas Silberzahn et al. (2018) asked contributors to address the novel
6
7 research question of whether referees are biased against darker-skinned players by analyzing
8
9 real-world data that lacked an experimental control, we asked our contributors to address a
10
11 research question that has been a focus of recognition memory research for decades and to do so
12 5
13
14 with data from controlled experiments.
15
16

17 Our results are also similar in some respects to the previous blinded inference study
18
19 reported by Dutilh et al. (2018), but direct comparisons are difficult based on procedural
20
21 differences between the two studies. In that study, response-time (RT) modelers analyzed
22
23 unlabeled data sets with the goal of inferring whether the conditions differed in psychological
24 10
25 constructs represented in RT models. Unfortunately, contributors disagreed about which
26
27 cognitive processes should theoretically vary as a function of certain experimental
28
29 manipulations; in other words, they had different views about what the “answer key” should be.
30
31 Different scoring rules were developed in light of this disagreement, making it difficult to
32
33 characterize overall performance. Using the originally planned scoring, at least, the proportion of
34
35 15
36 correct inferences (71%) was similar to our overall accuracy rate (68%). We recommend that
37
38 future blinded inference studies adopt our strategy of asking contributors to make inferences
39
40 about experimental manipulations as opposed to underlying theoretical processes to avoid
41
42 scoring ambiguities. A second difference between our study and Dutilh et al. (2018) also limits
43
44 our ability to compare the results: Their contributors were not required to express their
45
46 20
47 uncertainty with probability distributions. As a result, we do not know if their contributors’
48
49 inferences varied as dramatically as ours, with contributors reporting effect probabilities ranging
50
51
52
53
54
55
56
57
58
59
60

1
2
3 from 0% to 100% for some data sets, and we cannot compare Brier score results between the two
4
5 studies.
6

7
8 **Refining analysis quality.** Blinded inference can be a method to not only assess
9
10 inference quality, but also to improve it. Many of our contributors expressed surprise when they
11
12 learned of their performance level and conveyed that they would carefully re-evaluate their
13 5
14 chosen analysis techniques. Our results show that inference problems in recognition memory are
15
16 not a simple matter of choosing poor measurement techniques, as there are many instances of the
17
18 same technique being used by both high- and low-performing contributors. Defining the
19
20 characteristics of effective inference will require additional research, but for now we recommend
21
22 that analysts try a variety of analysis techniques and, ideally, have multiple researchers
23
24 10
25 independently analyze the data, reserving high confidence for consistent inferences.
26
27

28
29 **Conclusion.** We will end by again emphasizing that all of our contributors drew
30
31 inferences about the same data. Thus, the disparate conclusions reached by our contributors are
32
33 not another example of the replication crisis. Contributors were allowed to use any analysis and
34
35 15
36 any data censoring criteria they preferred, but those researcher degrees of freedom could not
37
38 systematically influence their conclusions because contributors were blind to the nature of the
39
40 experimental manipulation. Thus, our findings suggest that current efforts to improve research
41
42 quality are incomplete, in that they largely focus on limiting researchers' ability to bias results by
43
44 promoting desired outcomes (whether implicitly or explicitly). Even unbiased analysis
45
46 techniques can be *ineffective*, so it is critical for scientists to put their skills as analysts to direct
47 20
48 (and public) tests. The blinded inference paradigm is a promising method of assessing inference
49
50 quality and improving analysis procedures, so any field that uses analysis techniques to link data
51
52 patterns to unobserved theoretical constructs will benefit from applying this method. Our results
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

suggest that even well-established areas of research may be facing an inference crisis that warrants equal consideration with the replication crisis.

5

For Review Only

References

- 1
2
3
4
5 Aarts, A., & Open Science Collaboration. (2015). Estimating the reproducibility of psychological
6 science. *Science*, 349(6251), aac4716-aac4716. <https://doi.org/10.1126/science.aac4716>
7
8 Brier, G. W. (1950). Verification of forecasts expressed in terms of probability. *Monthly Weather*
9 5 *Review*, 78(1), 1–3. [https://doi.org/10.1175/1520-0493\(1950\)078<0001:VOFEIT>2.0.CO;2](https://doi.org/10.1175/1520-0493(1950)078<0001:VOFEIT>2.0.CO;2)
10
11 Buhrmester, M., Kwang, T., & Gosling, S. D. (2011). Amazon’s Mechanical Turk: A new source
12 of inexpensive, yet high-quality, data? *Perspectives on Psychological Science*, 6(1), 3–5.
13 <https://doi.org/10.1177/1745691610393980>
14
15 de Leeuw, J. R. (2015). jsPsych: A JavaScript library for creating behavioral experiments in a
16 10 Web browser. *Behavior Research Methods*, 47(1), 1–12. <https://doi.org/10.3758/s13428-014-0458-y>
17
18 Dube, C., Rotello, C. M., & Heit, E. (2010). Assessing the belief bias effect with ROCs: It’s a
19 response bias effect. *Psychological Review*, 117(3), 831–863.
20 <https://doi.org/10.1037/a0019634>
21
22 15 Dube, C., Starns, J. J., Rotello, C. M., & Ratcliff, R. (2012). Beyond ROC curvature: Strength
23 effects and response time data support continuous-evidence models of recognition memory.
24 *Journal of Memory and Language*, 67(3), 389–406.
25 <https://doi.org/10.1016/J.JML.2012.06.002>
26
27 Dutilh, G., Annis, J., Brown, S. D., Cassey, P., Evans, N. J., Grasman, R. P. P. P., ... Donkin, C.
28 20 (2018). The Quality of Response Time Data Inference: A Blinded, Collaborative
29 Assessment of the Validity of Cognitive Models. *Psychonomic Bulletin & Review*, 1–19.
30 <https://doi.org/10.3758/s13423-017-1417-2>
31
32 Erdfelder, E., Auer, T.-S., Hilbig, B. E., Abfal, A., Moshagen, M., & Nadarevic, L. (2009).
33 Multinomial Processing Tree Models. *Zeitschrift Für Psychologie / Journal of Psychology*,
34 217(3), 108–124. <https://doi.org/10.1027/0044-3409.217.3.108>
35 25
36 Gureckis, T. M., Martin, J., McDonnell, J., Rich, A. S., Markant, D., Coenen, A., ... Chan, P.
37 (2016). psiTurk: An open-source framework for conducting replicable behavioral
38 experiments online. *Behavior Research Methods*, 48(3), 829–842.
39 <https://doi.org/10.3758/s13428-015-0642-8>
40
41 30 Han, S., & Dobbins, I. G. (2009). Regulating recognition decisions through incremental
42 reinforcement learning. *Psychonomic Bulletin & Review*, 16(3), 469–474.
43 <https://doi.org/10.3758/PBR.16.3.469>
44
45 Hartl, D. L., & Fairbanks, D. J. (2007). Mud sticks: on the alleged falsification of Mendel’s data.
46 *Genetics*, 175(3), 975–979. Retrieved from <http://www.ncbi.nlm.nih.gov/pubmed/17384156>
47
48 35 Kučera, H., & Francis, W. N. (1967). *Computational analysis of present-day American English*.
49 Dartmouth Publishing Group.
50
51 Lachman, R., & Field, W. H. (1965). Recognition and recall of verbal material as a function of
52 degree of training. *Psychonomic Science*, 2(1–12), 225–226.
53 <https://doi.org/10.3758/BF03343418>
54
55 40 MacCoun, R., & Perlmutter, S. (2015). Blind analysis: Hide results to seek the truth. *Nature*,

1
2
3 526(7572), 187–189. <https://doi.org/10.1038/526187a>

4
5 Macmillan, N. A., & Creelman, C. D. (2005). *Detection Theory: A User's Guide*. Lawrence
6 Erlbaum Associates.

7
8 Miguel, E., Camerer, C., Casey, K., Cohen, J., Esterling, K. M., Gerber, A., ... Van der Laan, M.
9 5 (2014). Promoting Transparency in Social Science Research. *Science*, 343(6166), 30–31.
10 <https://doi.org/10.1126/science.1245317>

11
12 Pazzaglia, A. M., Dube, C., & Rotello, C. M. (2013). A critical comparison of discrete-state and
13 continuous models of recognition memory: Implications for recognition and beyond.
14 *Psychological Bulletin*, 139(6), 1173–1203. <https://doi.org/10.1037/a0033044>

15
16 10 Ratcliff, R. (1978). Theory of Memory Retrieval. *Psychological Review*, 85(2), 59–108.
17 <https://doi.org/10.1037//0033-295X.85.2.59>

18
19 Ratcliff, R., Clark, S. E., & Shiffrin, R. M. (1990). List-strength effect: I. Data and discussion.
20 *Journal of Experimental Psychology. Learning, Memory, and Cognition*, 16(2), 163–178.
21 <https://doi.org/10.1037//0278-7393.16.2.179>

22
23 15 Ratcliff, R., Sheu, C. F., & Gronlund, S. D. (1992). Testing global memory models using ROC
24 curves. *Psychological Review*, 99(3), 518–535. <https://doi.org/10.1037/0033-295X.99.3.518>

25
26 Riefer, D. M., & Batchelder, W. H. (1988). Multinomial modeling and the measurement of
27 cognitive processes. *Psychological Review*, 95(3), 318–339. <https://doi.org/10.1037//0033-295X.95.3.318>

28
29 20 Rotello, C. M., Heit, E., & Dubé, C. (2015). When more data steer us wrong: replications with
30 the wrong dependent measure perpetuate erroneous conclusions. *Psychonomic Bulletin &*
31 *Review*, 22(4), 944–954. <https://doi.org/10.3758/s13423-014-0759-2>

32
33 Silberzahn, R., Uhlmann, E. L., Martin, D., Anselmi, P., Aust, F., Awtrey, E. C., ... Nosek, B. A.
34 (2018). Many analysts, one dataset: Making transparent how variations in analytical choices
35 25 affect results. *Advances in Methods and Practices in Psychological Science*, 1(3), 337–356.
36 <https://doi.org/10.17605/OSF.IO/QKWST>

37
38 Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2011). False-Positive Psychology.
39 *Psychological Science*, 22(11), 1359–1366. <https://doi.org/10.1177/0956797611417632>

40
41 Snodgrass, J. G., & Corwin, J. (1988). Pragmatics of measuring recognition memory:
42 30 applications to dementia and amnesia. *Journal of Experimental Psychology. General*,
43 117(1), 34–50. Retrieved from <http://www.ncbi.nlm.nih.gov/pubmed/2966230>

44
45 Starns, J. J., Hicks, J. L., Brown, N. L., & Martin, B. A. (2008). Source memory for
46 unrecognized items: Predictions from multivariate signal detection theory. *Memory &*
47 *Cognition*, 36(1), 1–8. Retrieved from
48 35 [http://silk.library.umass.edu/login?url=http://search.ebscohost.com/login.aspx?direct=true&](http://silk.library.umass.edu/login?url=http://search.ebscohost.com/login.aspx?direct=true&db=psyh&AN=2008-03226-001&site=ehost-live&scope=site)
49 [db=psyh&AN=2008-03226-001&site=ehost-live&scope=site](http://silk.library.umass.edu/login?url=http://search.ebscohost.com/login.aspx?direct=true&db=psyh&AN=2008-03226-001&site=ehost-live&scope=site)

50
51 Stretch, V., & Wixted, J. T. (1998). On the difference between strength-based and frequency-
52 based mirror effects in recognition memory. *Journal of Experimental Psychology:*
53 *Learning, Memory, and Cognition*, 24(6), 1379–1396. [https://doi.org/10.1037/0278-](https://doi.org/10.1037/0278-7393.24.6.1379)
54 40 [7393.24.6.1379](https://doi.org/10.1037/0278-7393.24.6.1379)

1
2
3 Swets, J. A., Tanner, W. P., & Birdsall, T. G. (1961). Decision processes in perception.
4 *Psychological Review*, 68(5), 301–340. <https://doi.org/10.1037/h0040547>
5

6 Verde, M. F., & Rotello, C. M. (2007). Memory strength and the decision process in recognition
7 memory. *Memory & Cognition*, 35(2), 254–262. <https://doi.org/10.3758/BF03193446>
8

9 5

10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Review Only

Author contributions:

Conceptualization: Starns, Rotello

Data Curation: Cataldo

Formal Analysis: Starns, Cataldo, Measuring Memory Project contributors

Investigation: Cataldo

Methodology: Starns, Rotello, Cataldo

Project Administration: Starns, Cataldo, Rotello

Software: All authors

Supervision: Starns, Rotello

Visualization: Cataldo

Writing: Starns, Cataldo, Rotello

Inference Contributions: Jeffrey Annis, Andrew Aschenbrenner, Arndt Bröder, Gregory Cox, Amy Criss, Ryan A. Curl, Ian G. Dobbins, John Dunn, Tasnuva Enam, Nathan J. Evans, Simon Farrell, Scott H. Fraundorf, Scott D. Gronlund, Andrew Heathcote, Daniel W. Heck, Jason L. Hicks, Mark J. Huff, David Kellen, Kylie N. Key, Asli Kilic, Karl Christoph Klauer, Kyle R. Kraemer, Fábio P. Leite, Marianne E. Lloyd, Simone Malejka, Alice Mason, Ryan M. McAdoo, Ian M. McDonough, Robert B. Michael, Laura Mickes, Eda Mizrak, David P. Morgan, Shane T. Mueller, Adam Osth, Angus Reynolds, Travis M. Seale-Carlisle, Henrik Singmann, Jennifer F. Sloane, Andrew M. Smith, Gabriel Tillman, Don van Ravenzwaaij, Christoph T. Weidemann, Gary L. Wells, Corey N. White, Jack Wilson

Competing interests: Authors declare no competing interests.

Data and materials availability: Data and analyses are available at https://osf.io/92ahy/?view_only=2f6d9b285c2d4e279f144b6fed363142.

List of Supplementary Materials:

Supplementary Text

Figures S1-S5

Table S1