

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11 **Psychology is a Feature of Persons, Not Averages or Distributions:**  
12  
13 **The Group-to-Person Generalizability Problem in Social Cognition Research**  
14  
15  
16  
17

18 Ryan M. McManus<sup>1\*</sup>

19 Liane Young<sup>1</sup>

20 Joseph Sweetman<sup>2</sup>

21  
22  
23  
24  
25  
26  
27  
28 <sup>1</sup>Department of Psychology and Neuroscience, Boston College, Boston, MA, USA

29  
30  
31 <sup>2</sup>Department of Psychology, University of Exeter, Exeter, Devon, UK

32  
33  
34 \* Corresponding author email:

35  
36 [mcmnurrd@bc.edu](mailto:mcmnurrd@bc.edu)

37  
38  
39  
40 *Acknowledgements:* We would like to thank Stefano Anzellotti, Hiram Brownell, Richard  
41 Morey, Ehri Ryu, and Jordan Theriault for helpful feedback at the beginning of this project. We  
42 would also like to thank Adam Bear, Tony Chen, Minjae Kim, Aditi Kodipady, Gordon Kraft-  
43 Todd, Matthew Leitao, Shangzan (Sunny) Liu, Michael (Mookie) Manalili, Julia Marshall, Isaac  
44 Handley-Miner, Joshua Rottman, and Abraham Rutchick for helpful conversations at various  
45 stages of this project, as well as providing feedback on an early draft of the manuscript.  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## Abstract

When experimental psychologists make a claim (e.g., “Participants judged X as morally worse than Y”), how many participants are represented? Such claims are often based exclusively on group-level analyses; here, psychologists often fail to report, or perhaps even investigate, how many participants judged X as morally worse than Y. More troubling, group-level analyses do not necessarily generalize to the person-level: “the group-to-person generalizability problem.” We first argue for the necessity of designing experiments that allow investigation of whether claims represent most participants. Second, building on prior approaches, we document claims in the social cognition literature, derived from sets of typical group-level analyses, that describe only a (sometimes small) minority of participants. Third, we reason through an example to illustrate this group-to-person generalizability problem. Additionally, we show how severe the problem can be, demonstrating how claims from sets of simulated group-level effects can describe zero participants. Fourth, we conduct four experiments that rule out several methodology-based noise explanations of the problem. Fifth, we survey psychology researchers (and laypeople), finding that most interpret claims based on group-level effects as being intended to represent most participants in a study. Importantly, most believe this ought to be the case if a claim is used to support a general psychological theory. Finally, we propose that, if experimental psychologists are indeed interested in person-level psychology, then they should deploy different analytic strategies from those typically used. Overall, our approach offers a simple and flexible method to help researchers begin to engage with person-level analysis.

## Introduction

Francis Galton attended the 1906 “West of England Fat Stock and Poultry Exhibition” where attendees, hoping to win a prize, estimated an ox’s weight. Galton calculated that the crowd’s average estimate was 1,197 pounds, a perfect match to the ox’s true weight (Galton, 1907; Wallis, 2014). In this case, we might reasonably say that “people judged the ox’s weight perfectly.” Though this impressive example suggests the “wisdom of crowds” (Surowiecki, 2005), it is worth noting the considerable variability in person-to-person estimates, ranging below 1,000 pounds to above 1,400 pounds. In fact, the person-level data reveals that only one person guessed the correct weight of 1,197 pounds (Wallis, 2014). Consequently, we might question whether “people judged the ox’s weight perfectly” in truth describes what happened, as the group-level average represented only one person. Due to the ubiquity of aggregation approaches, this problem of group-to-person generalizability plagues modern-day experimental psychology. Psychologists average sets of person-level responses—largely ignoring person-to-person variability—and then use these averages to make claims about the mind. However, if psychology aims to understand and describe *persons*—to uncover the uniqueness or universality of certain cognitive processes—person-level responses ought to be the explananda.

In this paper, we argue that although experimental psychologists often strive to describe person-level phenomena, they sometimes fail to do so. First, we argue for closely matching experimental designs and analytic methods to precise research questions. Second, we document instances in published literature in which a person-level analytic approach yields different conclusions than typical group-level approaches. Third, in a tutorial, we show readers how this can occur, and how to conduct person-level analyses on their own data. Fourth, we conduct four pre-registered experiments to rule out several methodology-based explanations of group-to-

1  
2  
3 person generalizability failures. Fifth, we survey laypeople and psychology researchers to  
4 understand what is inferred about person-level phenomena from group-level analyses. Finally,  
5 we argue that experimental psychologists, if interested in person-level cognition, ought to deploy  
6 different design and analytic strategies than those typically used. We note here that, although we  
7 believe that our arguments apply to all areas of psychology, we focus on moral and social  
8 cognition in this paper.  
9  
10  
11  
12  
13  
14  
15

### 16 17 **Psychology as the Study of Person-Level (Not Group-Level) Phenomena**

18  
19 Psychology is often defined as “the study of the mind and behavior.” Therefore, its  
20 essential goals are describing cognitive functions and uncovering their antecedents and  
21 consequences. We contend that researchers intend to apply these goals to the study of individual  
22 persons, as it is minds that possess psychological processes, and minds reside in individuals. To  
23 strengthen this argument, we ask readers to engage in a thought exercise. Recall your most recent  
24 meeting with collaborators in which you discussed hypotheses and experimental designs to test  
25 them. At any point in that meeting, did you reason about possible patterns in a way that reflected  
26 how *individuals* may respond to different stimuli, or did you exclusively reason in a way that  
27 reflected how different stimuli would affect *averages or locations of distributions*? Furthermore,  
28 given the seeming frequency with which studied phenomena are described as applying to people  
29 generally, we also contend that most experimental psychologists intend to uncover phenomena  
30 that describe a *majority* of individuals (i.e., “general psychological laws”; Hamaker, 2012).  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45

46 Therefore, what follows are the most important takeaways from this paper:  
47  
48  
49

- 50 1. Psychologists sometimes fail to design experiments that permit investigation of  
51 person-level hypotheses.  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 2. Even when appropriate experimental designs are used, psychologists often report  
4 *only* their group-level analyses and interpret them *as if* they support or falsify person-  
5  
6 level hypotheses.  
7  
8  
9

10  
11  
12 Because it is possible for the above statements to be misinterpreted or overgeneralized,  
13  
14 we next communicate what we mean and do not mean.  
15  
16

17 Between-subjects experiments do not permit tests of person-level hypotheses (Speelman  
18 & McGann, 2020; Whitsett & Shoda, 2014). These common designs make it impossible to ask  
19 the simple question, “How many people respond this way?” (see Grice et al., 2020; Speelman &  
20 McGann, 2020), and they prohibit examination of unfolding person-level processes (e.g., Brandt  
21 & Morgan, 2022; Fisher, Medaglia, & Jeronimus, 2018; Moeller, 2022). For example, in our  
22 own recent moral cognition research, we assessed moral character judgments to test their  
23 sensitivity to social relationship information in the context of helping behavior (McManus,  
24 Mason, & Young, 2021). Among other variations, participants in our experiments were given  
25 two scenarios: one in which someone helps a total stranger, and another in which someone helps  
26 a distant family member. Group-level analyses suggested that participants—*on average*—judged  
27 agents who helped strangers as more morally good than agents who helped family members,  
28 presumably because people believe that there is no obligation to help strangers. Importantly, this  
29 was tested using a within-subjects design. Therefore, although it was not reported, our design  
30 permitted investigation of the question, “How many people respond this way?” A between-  
31 subjects design would have disallowed such investigation.  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50

51 However, having within-subjects designs does not automatically prevent inference errors  
52 from occurring. Researchers can still commit ecological or ergodic fallacies (Kuppens & Pollet,  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 2014; Speelman & McGann, 2020), due to special instances of Simpson’s paradox—when  
4 group-level patterns poorly represent lower-level units constituting the group (Simpson, 1951;  
5 Kievit, Frankenhuys, Waldorp, & Borsboom, 2013). A popular example of this is the relation  
6 between typing speed and mistake frequency (Hamaker, 2012), where a group-level correlation  
7 suggests a negative relation (i.e., faster typers make fewer mistakes), but person-level  
8 correlations suggest a positive relation (i.e., within individuals, typing faster results in *more*  
9 mistakes). These are not novel insights (see Hamaker, 2013; Speelman & McGann, 2020), but  
10 we communicate them here because we believe they have either not reached most research-  
11 practicing psychologists, or if they have, they have seemingly been ignored.  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23

24 To reiterate, even when psychologists deploy appropriate experimental designs, they  
25 often, if not always, only report their group-level analyses. Continuing to use our moral  
26 cognition research as an example, imagine that we had instead reported a null effect, generating a  
27 reasonable explanation for it (e.g., people perceive distant family members as otherwise stranger-  
28 like, so our manipulation was not strong enough to induce a change in people’s moral character  
29 judgments). While this sounds plausible, it is based on an average null effect. Perhaps, though,  
30 one third of participants’ responses supported the social relationship hypothesis, whereas another  
31 third’s were in the opposite direction, and yet another third showed no difference. Without  
32 investigating person-level responses, this distribution could not be uncovered.  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43

44 We are not suggesting, though, that there are only certain kinds of experimental designs  
45 that can yield knowledge. For example, it would be foolish to discard insights from Milgram’s  
46 obedience experiments simply because no participant engaged in shocking behavior across  
47 multiple conditions. Relatedly, all intervention-like insights (e.g., assessments of reading  
48 comprehension programs) would need to be discarded if we claimed that other designs did not  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 yield useful information. We are not arguing for this extreme position. However, even in this  
4  
5 class of research, a person-level analytic approach may provide valuable insights. For example,  
6  
7 in Milgram's obedience research, some participants declined to "continue the experiment."  
8  
9 Assuming these participants believed the ruse, they must have had reasons for declining. Probing  
10  
11 these reasons, as Milgram did, could help to explain the person-to-person variability. Similarly,  
12  
13 an assessment of a reading comprehension intervention might yield an average positive effect,  
14  
15 but whether this is due to a majority of the experimental students doing better than the control  
16  
17 students, or a subset of the experimental students doing better than all the rest, is not knowable  
18  
19 on the basis of *only* a group-level analysis. Therefore, even in cases where a person-level  
20  
21 approach might at first seem misguided, it may provide important insights not gained from  
22  
23 typical analysis strategies.  
24  
25  
26  
27

28 Overall, we are suggesting that, if a theory or research question is a person-level one, and  
29  
30 the goal of a study is to make a general claim (Hamaker, 2012), then researchers ought to choose  
31  
32 appropriate designs and analytic procedures. However, such careful matching does not always  
33  
34 occur in practice. The rest of this paper focuses on instances in which within-subjects group-  
35  
36 level effects fail to describe the majority of sampled individuals. From here on, we refer to this  
37  
38 as the "group-to-person generalizability problem."  
39  
40  
41

### 42 **Group-to-Person Generalizability Problems in the Wild**

43

44 We examined open data from the past five years in social cognition research (2016–  
45  
46 2021), looking for the group-to-person generalizability problem. Due to the larger reform  
47  
48 movements in psychology, publications from this era should be relatively more rigorous than  
49  
50 prior eras (e.g., larger samples, better statistical inferences). Our investigation was not systematic  
51  
52 in the sense that we can say, "X% of publications contain the problem." Rather, using a person-  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 level analytic approach, we re-analyzed open data with the goal of finding five instances of the  
4 problem from moral cognition—as we ourselves are moral psychologists—and five instances  
5 from social cognition generally (see Table 1)<sup>1</sup>.  
6  
7

8  
9  
10 To accomplish person-level analysis, we adopted “persons-as-effect-sizes” or  
11 “pervasiveness” approaches (see Grice et al., 2020; Speelman & McGann, 2020). Put simply, we  
12 created variables in each dataset that distinguished participants based on whether their response  
13 patterns supported the reported group-level patterns. If a participant’s responses had at least *some*  
14 distance between experimental conditions (e.g., 1-point on a Likert/sliding scale) and were  
15 directionally consistent with a group-level pattern, then that participant was categorized as  
16 supporting generalizability. Therefore, we used an extremely liberal threshold. An important  
17 nuance is that all of the investigated claims are based on *sets* of group-level tests (e.g., multiple  
18 paired t-tests). We therefore extended the person-level approaches to accommodate such claims.  
19 Specifically, we categorized participants as supporting generalizability if their full set of  
20 responses matched the full set of group-level patterns. For example, if a 2x2 interaction pattern  
21 underlied the claim, we counted person-level responses as supporting generalizability if a  
22 participant’s simple effects’ directions and differential magnitudes reflected the group-level  
23 pattern. But not all four person-level judgments had to be ordered identically to the group-level  
24 pattern. Readers can imagine (and if they wish, investigate) what these analyses look like under  
25 stricter constraints (see our OSF page:

26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47 [https://osf.io/xyse4/?view\\_only=dc61885b85b74681b10c116361c0cdad](https://osf.io/xyse4/?view_only=dc61885b85b74681b10c116361c0cdad)).  
48

49 We statistically defined claims as being unsupported at the person-level if we could not  
50 rule out the possibility that fewer than a simple majority of people would show the group-level  
51 pattern, conducting a binomial test against 0.50. If the 95% CIs contained a proportion of 0.50 or  
52  
53  
54  
55  
56  
57  
58  
59  
60



1  
2  
3 lower, then the claim was unsupported at the person-level. This, too, is a liberal threshold (i.e., a  
4 simple majority) for categorizing claims as favoring generalizability. As Table 1 shows,  
5 proportions of participants favoring generalizability varied across these publications but was low  
6 overall (3%-50%, with most proportions ranging between 20%-40%). Importantly, this occurred  
7 across a variety of dependent variables (e.g., sliding scales, Likert scales, reaction times, error  
8 rates) and pattern types (crossover interactions, attenuation interactions, ordinal patterns,  
9 conjunctive differences). We note that authors of these studies may not believe that their claims  
10 are describing most of their participants (although see *An Important Objection*). However,  
11 descriptions of results are, at the least, ambiguous enough to warrant uncertainty. We next break  
12 down a moral cognition example showing how the group-to-person generalizability problem can  
13 occur.  
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

**Table 1.** Quotes, relevant tests, and person-level statistics for instances of the group-to-person generalizability problem

Publication	Exact Quote(s)	Group-Level Test(s)	Person-Level Proportions
McManus, Mason, & Young (2021)	“On the one hand, people judged agents who helped a stranger as more morally good than agents who helped a family member. On the other hand, people judged agents who helped a stranger instead of a family member as less morally good than agents who helped a family member instead of a stranger.”	<u>Experiments 1a-b</u> -2 x 2 interactions -Set of paired t-tests -See Figure 2	E1a: <b>31% [34%-37%]</b> (62 / 203)  E1b: <b>29% [23%-36%]</b> (59 / 203)
Law, Campbell, & Gaesser (2021)	“People consistently view socially distant altruism as less morally acceptable as the person not receiving help becomes closer to the agent helping.”	<u>Experiments 1 &amp; 4</u> -Set of paired t-tests -See Figures 1 & 7b (Country vs Town vs Friend vs Family)	E1: <b>3% [1%-9%]</b> (3 / 97)  E4: <b>8% [5%-11%]</b> (30 / 397)
Fowler, Law, & Gaesser (2021)	“The results showed that moral judgments of empathy are biased toward preferring more empathy for a socially close over a socially distant individual. Despite this bias in moral judgments, however, people consistently judged feeling equal empathy as the most morally right perspective.”	<u>Experiment 2</u> -Set of paired t-tests -See Figure 3 (More For Distant vs More For Close vs Equal)	<b>32% [27%-37%]</b> (97 / 304)
Soter, Berg, Gelman, & Kross (2021)	“Participants said they should protect close others more than distant others. However, the effect of relationship was consistently weaker for “should” judgments than “would” judgments, revealing that people show <i>relatively less</i> partiality in their judgments of what is morally right, compared to judgments of how they would act.”	<u>Experiment 2</u> -2 x 2 interaction -Simple comparisons -See Figure 2	<b>29% [25%-34%]</b> (104 / 356)
Rottman & Young (2019)	“In three studies, adult participants judged the moral wrongness of harm and purity transgressions that varied in frequency (e.g., occasionally vs. regularly) or magnitude (e.g., small vs large) with the same sets of modifiers or the same quantities (e.g., a single drop vs. a teaspoon) repeated across content domains. All studies found that evaluations of purity violations were considerably less sensitive to variations in scope than evaluations of harms, yielding robust statistical interactions between domain and dosage.”	<u>Experiments 1-3</u> -2x2 interactions -Simple comparisons -See Figures 1-3	E1: <b>29% [22%-36%]</b> (51 / 177)  E2: <b>46% [35%-57%]</b> (37 / 81)  E3: <b>22% [16%-29%]</b> (37 / 168)

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

Deska et al. (2020)	“We also observed an interaction between target race and target gender for life hardship. As with social pain, it was clear that participants generally agreed that Black targets experience greater life hardship than White targets; however, this seemed to be especially true for male targets.”	<u>Experiment 4</u> -2x2 interaction -Simple comparisons	<b>50% [42%-59%]</b> (66 / 131)
Stroessner et al. (2020)	“An association between a gender category and a shape would be revealed by faster categorization speeds following compatible (masculine-square and feminine-circle) compared with incompatible (masculine-circle and feminine-square) prime-target pairings.”  “Along with the results of Studies 3a–3c, these data demonstrate that gender categorization of basic squares and circles occurs without intention.”	<u>Experiments 2 &amp; 4</u> -2x2 interaction -Sets of paired t-tests -See Figure 3	E2: <b>38% [26%-50%]</b> (26 / 69)  E4: <b>41% [33%-49%]</b> (61 / 150)
Craig, Nelson, & Dixon (2019)	“We found that the presence of a beard increased the speed and accuracy with which participants recognized displays of anger but not happiness.”  “In Experiment 1, facial hair facilitated recognition of anger, and the advantage in response times cannot be attributed to a shift toward responding “angry.” Recognition of facial expressions of happiness, which are positive and nonthreatening, was slowed by the presence of a beard in this task.”	<u>Experiment 1</u> -2x2 interactions -Sets of paired t-tests -See Figure 2	Speed: <b>45% [38%-52%]</b> (99 / 219)  Accuracy: <b>25% [20%-31%]</b> (55 / 219)  Both: <b>13% [9%-18%]</b> (29 / 219)
Decelles, Adams, Lowe, & John (2021)	“Using a sample of working professionals, including fraud investigators and auditors, we found in Study 4 that an angry response to an accusation was interpreted as a sign of guilt, relative to remaining calm. Moreover, compared with remaining calm and with angrily denying an accusation, remaining silent was also perceived as a cue of guilt and therefore does not appear to be a viable solution for the accused to avoid the negative effects of anger.”	<u>Experiment 4</u> -Set of paired t-tests (Anger vs Calm & Silent vs Calm)	<b>38% [30%-47%]</b> (52 / 136)
Thai, Borgella, & Sanchez (2019)	“Study 3 demonstrated that it was deemed most acceptable for a person to make jokes about a particular social group if they themselves were a part of that social group. This remained true for both minority-directed and majority-directed humor. This pattern emerged consistently for all three categories of humor studied, including race-based, sexual orientation-based, or gender-based humor.”	<u>Experiment 3</u> -2x2 interaction -Simple comparisons -See Figure 4 (Gender-based Jokes)	<b>45% [33%-57%]</b> (31 / 70)

*Note:* Across publications, it was sometimes difficult to find specific claims which could be connected back to specific hypothesis tests. For some publications, there was not a specific, insulated claim which clearly referenced a specific hypothesis test (e.g., Stroessner et al., 2020), which is why some quoted sections are taken from multiple places of the publication. In Law, Campbell, & Gaesser (2021), the verbal claim was not an accurate representation of the set of group-level patterns (some necessary group-level patterns did not emerge). However, re-analysis of their data was based on the claim rather than the group-level patterns.

## Tutorial for the Group-to-Person Generalizability Problem (McManus et al., 2021)

Here we demonstrate how researchers can understand and perform analyses to make person-level inferences. For relevant background, consider the two earlier moral cognition scenarios: someone helps an unrelated stranger, and someone helps their cousin. We predicted that agents who helped strangers should be judged as more morally good than agents who helped their cousin, due to stranger-helping agents lacking an obligation to help but doing so anyway. Now consider these two scenarios in a slightly different context: someone helps an unrelated stranger *instead of* their cousin, and someone helps their cousin instead of an unrelated stranger. We predicted the opposite pattern here, as stranger-helping agents would be violating their family obligation. These two contexts were described as “No Choice” and “Choice” contexts, respectively. Indeed, this interaction and context-based reversal of simple effects emerged at the group-level.

In the general discussion, we communicated this effect as follows: “On the one hand, people judged agents who helped a stranger as more morally good than agents who helped a family member. On the other hand, people judged agents who helped a stranger instead of a family member as less morally good than agents who helped a family member instead of a stranger.” As two of the three authors of the current paper were authors, we can say, honestly, that we intended to communicate this effect as applying to a majority of participants. Therefore, our claim is interesting, and arguably, accurate, if *and only if* the interaction describes most participants’ psychology.

To investigate this at the person-level, each simple effect and the interaction can be described by a set of directional patterns. The No Choice simple effect can be computed by

1  
2  
3 subtracting the “helped a family member” ratings from the “helped a stranger” ratings, whereas  
4  
5 the Choice simple effect can be computed by subtracting the “helped a family member instead of  
6  
7 a stranger” ratings from the “helped a stranger instead of a family member ratings.” An  
8  
9 interaction effect can then be computed by subtracting the Choice effect from the No Choice  
10  
11 effect. See Table 2 for an example of 13 hypothetical participants who reflect all possible  
12  
13 qualitative patterns, and Table 3 for example R code to create generalizable 2x2 person-level  
14  
15 patterns, investigate their frequencies, and conduct a binomial test. The person-level combination  
16  
17 which matches the published claim is the “Positive, Negative, Positive” pattern (No Choice  
18  
19 simple effect, Choice simple effect, Interaction effect). As shown in Figure 1, less than 30% of  
20  
21 our participants show the group-level effect.  
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

**Table 2.** Example hypothetical participants, showing all possible patterns in McManus et al. (2021)

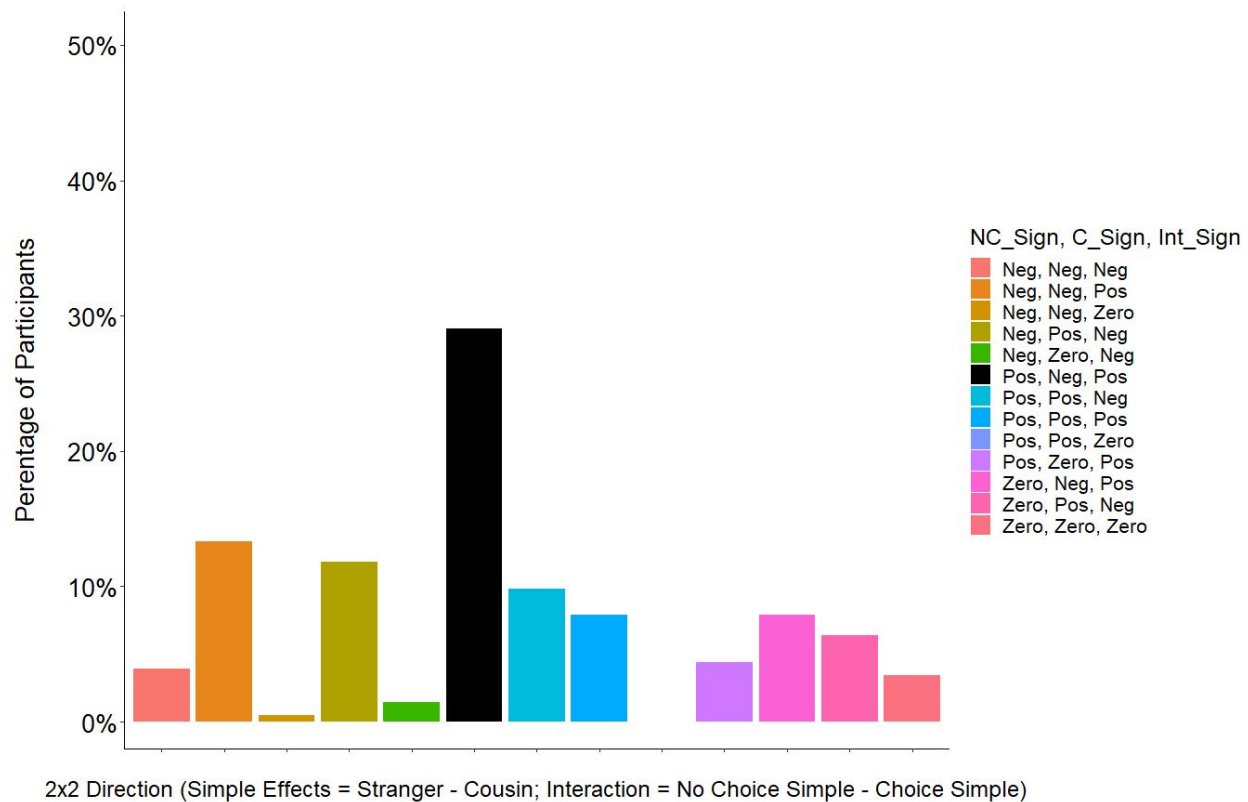
Subj	NC_Stranger	NC_Cousin	C_Stranger	C_Cousin	NC_Diff	C_Diff	Interaction	NC_Direction	C_Direction	Int_Direction
1	1	3	2	3	-2	-1	-1	Negative	Negative	Negative
2	2	3	1	3	-1	-2	1	Negative	Negative	Positive
3	2	3	2	3	-1	-1	0	Negative	Negative	Zero
4	2	3	2	1	-1	1	-2	Negative	Positive	Negative
5	2	3	2	2	-1	0	-1	Negative	Zero	Negative
<b>6</b>	<b>3</b>	<b>2</b>	<b>1</b>	<b>2</b>	<b>1</b>	<b>-1</b>	<b>2</b>	<b>Positive</b>	<b>Negative</b>	<b>Positive</b>
7	3	2	3	1	1	2	-1	Positive	Positive	Negative
8	3	1	3	2	2	1	1	Positive	Positive	Positive
9	3	2	3	2	1	1	0	Positive	Positive	Zero
<b>10</b>	3	2	2	2	1	0	1	Positive	Zero	Positive
11	3	3	1	2	0	-1	1	Zero	Negative	Positive
12	3	3	2	1	0	1	-1	Zero	Positive	Negative
13	3	3	2	2	0	0	0	Zero	Zero	Zero

*Note:* Each of these hypothetical person-level patterns constitute all possible combinations of two simple effects directions, leading to 13 possible interaction patterns. “NC” and “C,” denote No Choice and Choice, respectively, as communicated in McManus et al., (2021). Subject row 6 is bolded to highlight the pattern that matches the claimed effect. The first four non-subject columns are hypothetical raw scores in each within-subjects condition. The next two columns are hypothetical difference scores which constitute the simple effects of interest. Simple effects (NC\_Diff and C\_Diff) are calculated by subtracting “Cousin” scores from “Stranger” scores. The “Int” column contains the interaction values which are computed by subtracting the second simple effect from the first simple effect. The last three columns are directional labels to communicate the full person-level pattern for each subject. For ease of calculation and communication, this table assumes that hypothetical participants used a simple three-point scale. In principle, the number of scale points are irrelevant so long as the scale has more than two points (otherwise, there could not be differential magnitudes of simple effects). Importantly, these patterns do not consider other features of interaction patterns, such as the rank-ordering of all four conditions on the numerical response scale.

**Table 3.** Instructions and Example R Code to Investigate Person-Level Patterns in a 2x2 Design

<b>Step 1</b>	Use wide-formatted data (i.e. 1 row per participant) to create simple effects of interest.	<pre>data_wide &lt;- data_wide %&gt;%   mutate(SimpleEff1 = A1 - A2) %&gt;%   mutate(SimpleEff2 = B1 - B2)</pre>
<b>Step 2</b>	Create variables which constitute person-level pattern possibilities.	<pre>data_wide &lt;- data_wide %&gt;%   mutate(`2x2_Pattern` = case_when(     (SimpleEff1 == 0 &amp; SimpleEff2 == 0) ~ "Zero, Zero, Zero",     (SimpleEff1 == 0 &amp; SimpleEff2 &lt; 0) ~ "Zero, Neg, Pos",     (SimpleEff1 == 0 &amp; SimpleEff2 &gt; 0) ~ "Zero, Pos, Neg",     (SimpleEff1 &lt; 0 &amp; SimpleEff2 == 0) ~ "Neg, Zero, Neg",     (SimpleEff1 &lt; 0 &amp; SimpleEff2 &lt; 0 &amp; SimpleEff1 == SimpleEff2) ~ "Neg, Neg, Zero",     (SimpleEff1 &lt; 0 &amp; SimpleEff2 &gt; 0) ~ "Neg, Pos, Neg",     (SimpleEff1 &lt; 0 &amp; SimpleEff2 &lt; 0 &amp; SimpleEff1 &gt; SimpleEff2) ~ "Neg, Neg, Pos",     (SimpleEff1 &lt; 0 &amp; SimpleEff2 &lt; 0 &amp; SimpleEff1 &lt; SimpleEff2) ~ "Neg, Neg, Neg",     (SimpleEff1 &gt; 0 &amp; SimpleEff2 == 0) ~ "Pos, Zero, Pos",     (SimpleEff1 &gt; 0 &amp; SimpleEff2 &lt; 0) ~ "Pos, Neg, Pos", # predicted effect     (SimpleEff1 &gt; 0 &amp; SimpleEff2 &gt; 0 &amp; SimpleEff1 == SimpleEff2) ~ "Pos, Pos, Zero",     (SimpleEff1 &gt; 0 &amp; SimpleEff2 &gt; 0 &amp; SimpleEff1 &lt; SimpleEff2) ~ "Pos, Pos, Neg",     (SimpleEff1 &gt; 0 &amp; SimpleEff2 &gt; 0 &amp; SimpleEff1 &gt; SimpleEff2) ~ "Pos, Pos, Pos"))</pre>
<b>Step 3</b>	Investigate frequencies of all person-level patterns.	<pre>data_wide %&gt;%   group_by(`2x2_Pattern`) %&gt;%   summarize(freq = n())</pre>
<b>Step 4</b>	Test the predicted effect's frequency against 0.50, using a binomial test.	<pre>binom.test(x = Predicted Effect Freq, n = Total N, p = 0.50, alternative = "two.sided")</pre>

Note: The above R code was created using functions from the “tidyverse” package. In Step 2, all text-based patterns reflect the direction of the first simple effect, the second simple effect, and the interaction (e.g., “Zero, Zero, Zero”), in that order.



33  
34  
35  
36

**Figure 1.** Person-Level Patterns from McManus, Mason, & Young (2021). The black bar represents the group-level pattern.

37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

How can this happen? Consider first the crossover interaction. This interaction is typically tested for using a 2x2 repeated-measures ANOVA, as we did. Importantly, the interaction can be assessed using t-tests, which can help to explain the discrepancy. To use the t-test methods, the analyst first creates difference score variables by subtracting the second response from the first response within each simple effect of interest. The paired-samples t-test method is completed by conducting a t-test on the two difference scores. The one-sample t-test method involves an extra step, creating a third difference score variable—the interaction score—by subtracting the second simple effect's difference score from the first simple effect's difference score. The one-sample t-test method is completed by conducting a t-test (against zero)



1  
2  
3 on the interaction scores. If either t-test returns a below-alpha p-value, then an interaction effect  
4 exists. Importantly, in this context, the p-value from both t-test methods would be identical to  
5 one another and to the p-value of the ANOVA's interaction F-test, as all methods are testing for  
6 a difference in differences (see SOM for a demonstration).  
7  
8  
9  
10

11  
12 Why does this matter? As shown in Table 2, there are five patterns which yield a positive  
13 interaction value, only one of which is the claimed pattern. This is problematic considering that  
14 the interaction test is simply assessing whether the interaction scores' average differs from zero,  
15 nothing more. Therefore, it is possible that more participants had a positive interaction value  
16 constituted by the "incorrect" set of simple effects than had a positive interaction value  
17 constituted by the "correct" set of simple effects. Indeed, more than 60% of our sample had a  
18 positive interaction value that contributed to the group-level interaction test.  
19  
20  
21  
22  
23  
24  
25  
26  
27

28 Now consider the opposite-signed simple effects. It is an obvious but crucial point that a  
29 person-level claim about the full interaction pattern requires that participants show *both* simple  
30 effects. However, what seems non-obvious is that *sets* of typical inferential tests cannot provide  
31 this evidence. Because the units of analysis for a single paired-samples t-test are the person-level  
32 difference scores, two separate paired-samples t-tests cannot connect units across analyses. The  
33 only way to ensure that a particular proportion of participants show both group-level patterns is  
34 to first count how many show each individual pattern. Tabulations of within-person differences  
35 showed that the first simple effect described 51% of participants, whereas the second simple  
36 effect described 55% of participants. Consequently, the *maximum* proportion of participants who  
37 could have shown both patterns was 51%. As has been established, however, fewer than 30% of  
38 participants showed both patterns.  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52

53  
54 **The Problem Worsens (and is Difficult to Fix)!**  
55  
56  
57  
58  
59  
60

1  
2  
3 We believe that we have provided compelling reasoning that person-level claims need to  
4 be tested using persons-as-effect-sizes or pervasiveness approaches— tabulating the proportion  
5 of participants whose responses match predictions (Grice et al., 2020; Speelman & McGann,  
6 2020). To provide further supporting evidence, we generated hypothetical datasets in which sets  
7 of group-level analyses are extremely poor representations of person-level cognition. In these  
8 datasets, we created 2x2 crossover interactions, 2x2 attenuation interactions, and three-level  
9 ordinal effects, all of which yield group-level effects (and survive non-parametric tests) but  
10 describe *zero* participants (see SOM). Although we are unaware of real-world instances, the  
11 theoretical possibility of group-level patterns being perfectly unrepresentative of persons should  
12 warrant caution.  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25

26 Despite these existence proofs, it could be argued that most discrepancies between group-  
27 level and person-level analyses are due to methodological features of experiments which can be  
28 remedied. That is, most experiments may not be designed to minimize noise and therefore  
29 maximize the probability of individuals exhibiting the group-level pattern. If such barriers could  
30 be addressed, then group-level patterns may better represent person-level patterns. To address  
31 this, using our moral cognition paradigm described earlier, we conducted four pre-registered  
32 experiments which systematically varied methodological features hypothesized as partial, noise-  
33 inducing causes of the group-to-person generalizability problem. Within each experiment, we  
34 replicated our original group-level effects, as well as the low proportions of participants  
35 represented by them (17%-27%; see also Devezer et al. [2021] for a discussion on how  
36 replicability need not imply “true”). However, none of our experiments were successful in  
37 explaining the problem and therefore shifting person-level patterns to be better aligned with the  
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  
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

group-level pattern (see Table 4 for a summary of the experiments' logic and results, and SOM for full details).

For Review Only

**Table 4. Underlying Logic and Results for Methodology-Based Experiments (see SOM for full details)**

Manipulation	Underlying Logic	Results
Absence/Presence of Calibration Trials	<p><b>Problem 1:</b> If participants do not engage in calibration trials or get feedback about their scale use, then different participants may have different interpretations of identical points along the scale.</p> <p><b>Problem 2:</b> If participants do not engage in calibration trials which are designed to elicit responses along the entire range of the scale, then, when the main task starts, some participants may use extreme ends of the scale for the first stimulus they see, disallowing them from distinguishing between the first stimulus and a later stimulus which they truly wish to judge as more extreme.</p> <p><b>Solution:</b> Before the main experimental task, give participants calibration trials and normative feedback about how most other people use the scale.</p> <p><b>Hypothesis:</b> If the group- versus person-level discrepancy is due to noise of this kind, then participants in an experimental condition (i.e., those who engage in pre-task calibration trials) should be more likely to show the person-level response pattern that matches the group-level pattern, compared to participants in a control condition (i.e., those who do not engage in pre-task calibration trials).</p>	<p><b>N per Condition</b>  <i>N</i>Control: 658  <i>N</i>Experimental: 589</p> <p><b>Predicted Interaction</b>  Control: 24%  Experimental: 27%</p> <p><b>Eq of Proportions Test</b>  <math>\chi^2 = 1.17, p = .280</math></p> <p><b>Hypothesis Decision</b>  Unsupported</p>
Inability/Ability to Respond to Stimuli Simultaneously	<p><b>Problem 1:</b> If participants cannot consider all stimuli simultaneously, then some participants may fail to distinguish between stimuli that they truly wish to distinguish between.</p> <p><b>Problem 2:</b> If participants cannot consider all stimuli simultaneously (and they instead encounter stimuli sequentially), then some participants may use the extreme end of a scale for an early stimulus and be unable to distinguish between it and a later stimulus which they believe is more extreme.</p> <p><b>Solution:</b> Give participants the opportunity to see all stimuli before making any judgments. Then, re-present the important details of all stimuli simultaneously, requesting that participants make any single judgment while considering how they would make their other judgments.</p> <p><b>Hypothesis:</b> If the group- versus person-level discrepancy is due to noise of this kind, then participants in an experimental condition (i.e., those who can see all stimuli and make judgments simultaneously) should</p>	<p><b>N per Condition</b>  <i>N</i>Control: 628  <i>N</i>Experimental: 609</p> <p><b>Predicted Interaction</b>  Control: 24%  Experimental: 19%</p> <p><b>Eq of Proportions Test</b>  <math>\chi^2 = 4.65, p = .031</math></p>

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

	<p>be more likely to show the person-level response pattern that matches the group-level pattern, compared to participants in a control condition (i.e., those who see stimuli and make judgments sequentially).</p>	<p><b>Hypothesis Decision</b> Unsupported</p>
<p>Absence/Presence of Matched Stimuli</p>	<p><b>Problem:</b> If participants respond to stimuli which differ in content across experimental conditions (even if all stimuli variants appear in each condition across the entire sample), then some participants may attend to non-experimental features of stimuli when responding.</p> <p><b>Solution:</b> Give participants matched-in-content stimuli across experimental conditions, varying only the experimental features of interest.</p> <p><b>Hypothesis:</b> If the group- versus person-level discrepancy is due to noise of this kind, then participants in an experimental condition (i.e., those who see perfectly matched stimuli) should be more likely to show the person-level response pattern that matches the group-level pattern, compared to participants in a control condition (i.e., those who see different-in-content stimuli).</p>	<p><b>N per Condition</b> NControl: 638 NExperimental: 641</p> <p><b>Predicted Interaction</b> Control: 24% Experimental: 17%</p> <p><b>Eq of Proportions Test</b> <math>\chi^2 = 10.94, p &lt; .001</math></p> <p><b>Hypothesis Decision</b> Unsupported</p>
<p>Inability/Ability to “Opt Out” of using Measures/Scales</p>	<p><b>Problem:</b> If participants do not have the opportunity to “opt out” of using a measurement scale, then some participants’ responses may not reflect the construct of interest in exactly the way that researchers intend. For example, participants may not believe a measurement scale captures how they think; therefore, they may actively transform the scale or respond completely randomly.</p> <p><b>Solution:</b> Give participants the ability to opt out of using a measurement scale.</p> <p><b>Hypothesis:</b> If the group- versus person-level discrepancy is due to noise of this kind, then participants in an experimental condition (i.e., of those who have an opportunity to opt out, those who do not) should be more likely to show the person-level response pattern that matches the group-level pattern, compared to participants in a control condition (i.e., those who cannot opt out).</p>	<p><b>N per Condition</b> NControl: 746 NExperimental: 691</p> <p><b>Predicted Interaction</b> Control: 22% Experimental: 23%</p> <p><b>Eq of Proportions Test</b> <math>\chi^2 = 0.09, p = .779</math></p> <p><b>Hypothesis Decision</b> Unsupported</p>

### (Empirically Addressing) An Important Objection

We have argued that there is a group-to-person generalizability problem in social cognition, documenting published instances of it, showing how it can occur, demonstrating its potential severity, and its resistance to obvious method-based remedies. However, there is obvious subjectivity involved when deciding what should count as person-level evidence for a claim. For example, many claims which we viewed as instances of the group-to-person generalizability problem (see Table 1) may seem unproblematic to other researchers. It could be argued that percentages of participants in the 20-40% range, who show the group-level patterns, are quite high. Moreover, perhaps readers of psychology research (laypeople and psychology researchers themselves) do not interpret authors as intending to make claims that represent at least a majority of participants. We therefore set out to answer two questions empirically:

1. Do a majority of people who read psychology research believe that authors intend to communicate claims as representing most participants in their data?
2. Do a majority of people who read psychology research believe that claims ought to represent most participants if the authors use their data to claim support for a general theory of individual psychology?

To answer these questions, we surveyed laypeople and social psychology researchers by presenting modified excerpts of “results” and “general discussion” sections from publications which contain unrepresentative group-level patterns. We report how we determined our sample sizes, all data exclusions, all manipulations, and all measures.

#### **Method**

**Participants.** All laypeople were U.S. residents recruited and compensated via CloudResearch’s “approved participants” list. Participants from McManus et al. (2021) were unable to access the

1  
2  
3 current study. Additionally, participants from our methods experiments could not participate.  
4  
5 Researchers were affiliated with the Society for Personality and Social Psychology (SPSP),  
6  
7 recruited via SPSP's Open Forum listserv and compensated with Amazon gift cards. As pre-  
8  
9 registered ([https://osf.io/6qay8/?view\\_only=9aca048fe58042d1a6835c5cc84cc293](https://osf.io/6qay8/?view_only=9aca048fe58042d1a6835c5cc84cc293) and  
10  
11 [https://osf.io/nucbf/?view\\_only=220229726d1249978df2e342c2af2098](https://osf.io/nucbf/?view_only=220229726d1249978df2e342c2af2098)), we aimed to collect at  
12  
13 least 642 analyzable laypeople and 280 analyzable researchers. In total, we were able to collect  
14  
15 705 and 256 unique responses, respectively. After applying the pre-registered exclusion criterion  
16  
17 (failing a comprehension check), this resulted in  $N_{Laypeople}=588$  and  $N_{Researchers}=244$ . We did not  
18  
19 resample due to still having high statistical power for our focal hypothesis tests (see *Statistical*  
20  
21 *Power & Hypotheses*).  
22  
23  
24  
25

26 **Design.** Participants were randomly assigned to one of two conditions. Half of participants  
27  
28 learned about a simple effect comparison, whereas the other half of participants learned about a  
29  
30 more complex, two-way interaction effect. We used both simple and complex effect examples to  
31  
32 test the generality of our hypotheses.  
33  
34

35 **Materials and Procedure.** At the beginning of the study, all participants were informed that they  
36  
37 would be answering questions about a moral cognition experiment. For the simple effect  
38  
39 condition, participants learned about an effect from the supplemental materials of Law,  
40  
41 Campbell, & Gaesser (2021). For the complex effect condition, participants learned about the  
42  
43 interaction effect from McManus et al. (2021).  
44  
45  
46

47 Participants first read text communicating results in typical journal article format (with  
48  
49 means, SDs, t-values, p-values, within-subject standardized effect sizes for comparisons of  
50  
51 interest [ $d_z$ ], and a barplot; see OSF for full materials). After learning the results, they then read  
52  
53 text that simulated how data-based claims are made in a general discussion section (e.g., "People  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 judged fictional agents who helped a stranger as more morally good than fictional agents who  
4 helped a cousin, but they judged fictional agents who helped a stranger instead of a cousin as less  
5 morally good than fictional agents who helped a cousin instead of a stranger”).  
6  
7  
8  
9

10 After learning about the claim, participants were then asked to respond to a series of true-  
11 false questions about what the reported results suggested. However, these questions were not of  
12 primary interest (see OSF for Rmarkdown results). Participants were then again shown the claim  
13 in general discussion format, and asked “By *people*, approximately what percentage of the  
14 study’s participants do you think the researchers mean?” We call this measure the “empirical  
15 proportion estimate.” Responses ranged from 0-100% on a sliding scale, with the starting  
16 position (0, 50, 100) counterbalanced across participants. This measure allows categorization of  
17 responses into two categories: less than a simple majority (50% or less), and equal to or greater  
18 than a simple majority (51% or more).  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29

30  
31 Next, participants learned about a (fictional) general, individual-level theory that the  
32 authors had developed pre-study. Participants were then asked to respond to a series of true-false  
33 questions about how the reported results informed the theory (see OSF). Participants were again  
34 shown the claim in general discussion format and told that, later in the paper, the authors used  
35 their study’s results to claim support for their theory. Participants were then asked, “In order for  
36 the study’s results to support the researchers’ theory/model, approximately what percentage of  
37 the study’s participants do you think need to respond in the way described by [the general  
38 discussion’s language]?” We call this measure the “theoretical proportion estimate.” Responses  
39 were measured identically to the empirical estimate. Finally, participants could write an open-  
40 ended response to communicate anything that they were unable to communicate thus far. After  
41 the main task, participants answered several demographic questions.  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



1  
2  
3 **Statistical Power.** As pre-registered, we aimed for at least 321 participants per condition for the  
4 laypeople sample, and 140 participants per condition for the researcher sample. The pre-  
5 registered laypeople sample size yielded 95% power to detect a 10-point proportion difference  
6 from 50% (e.g., 60%) using a two-tailed binomial test and assuming an alpha level = 0.05, the  
7 focal test to examine whether a majority of empirical/theoretical proportion estimates reflect  
8 inferences being made about a majority of a study's participants. As explained in our pre-  
9 registrations, we planned the researcher sample based on the results of the laypeople sample. For  
10 the researcher sample, the pre-registered sample size yielded 95% power to detect a 15-point  
11 proportion difference from 50% using identical test specifications as the laypeople sample.  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23

24 In the laypeople sample, applying the pre-registered exclusion criterion (i.e., missing a  
25 comprehension check question) led to  $N_{Simple}=303$  and  $N_{Complex}=285$ . In the researcher sample, we  
26 were unable to successfully recruit our entire desired sample size. After one attempt to get more  
27 responses (via reposting to SPSP's Open Forum listserv), we decided to close the survey once  
28 incoming responses completely stalled, which occurred after two weeks. Applying the same  
29 exclusion criterion led to  $N_{Simple}=123$  and  $N_{Complex}=121$ . We did not resample for either  
30 population because sensitivity analyses revealed that we still had more than 90% power to detect  
31 our pre-registered minimal effect sizes.  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41

## 42 **Hypotheses**

- 43  
44  
45 1) Empirical Proportion: The majority of people (i.e., 51% or more) within each sample  
46 (laypeople and researchers) will believe authors' claims are intended to describe at least a  
47 simple majority (i.e., 51% or more) of their study's participants.  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 2) Theoretical Proportion: The majority of people within each sample will believe at least a  
4 simple majority of a study's participants ought to be described by the authors' claims in order  
5  
6 for the results to support a general theory of individual psychology.  
7  
8  
9

## 10 **Results**

11  
12 ***Empirical Proportion Estimate.*** The majority of laypeople believed authors intended to describe  
13 at least a simple majority of their study's participants, for both simple (81%) and complex (88%)  
14 effects. The majority of researchers agreed for both simple (73%) and complex (80%) effects<sup>2</sup>  
15 (see Table 5 for additional statistics). Strikingly, as shown in Figure 2, there is no discernible  
16 pattern as a function of being relatively inexperienced (e.g., layperson or undergraduate) and  
17 relatively experienced with academic research (e.g., professor).  
18  
19  
20  
21  
22  
23  
24  
25

26 ***Theoretical Proportion Estimate.*** The majority of laypeople believed that at least a simple  
27 majority of a study's participants ought to be described by authors' claims for the results to  
28 support an individual-level psychological theory, for both simple (93%) and complex (92%)  
29 effects. The majority of researchers agreed for both simple and (80%) and complex (90%) effects  
30 (see Table 6 for additional statistics). As shown in Figure 2, there again is no discernible pattern  
31 as a function of research experience.  
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

**Table 5. Empirical Estimate Tests within Each Effect Type (split by Population)**

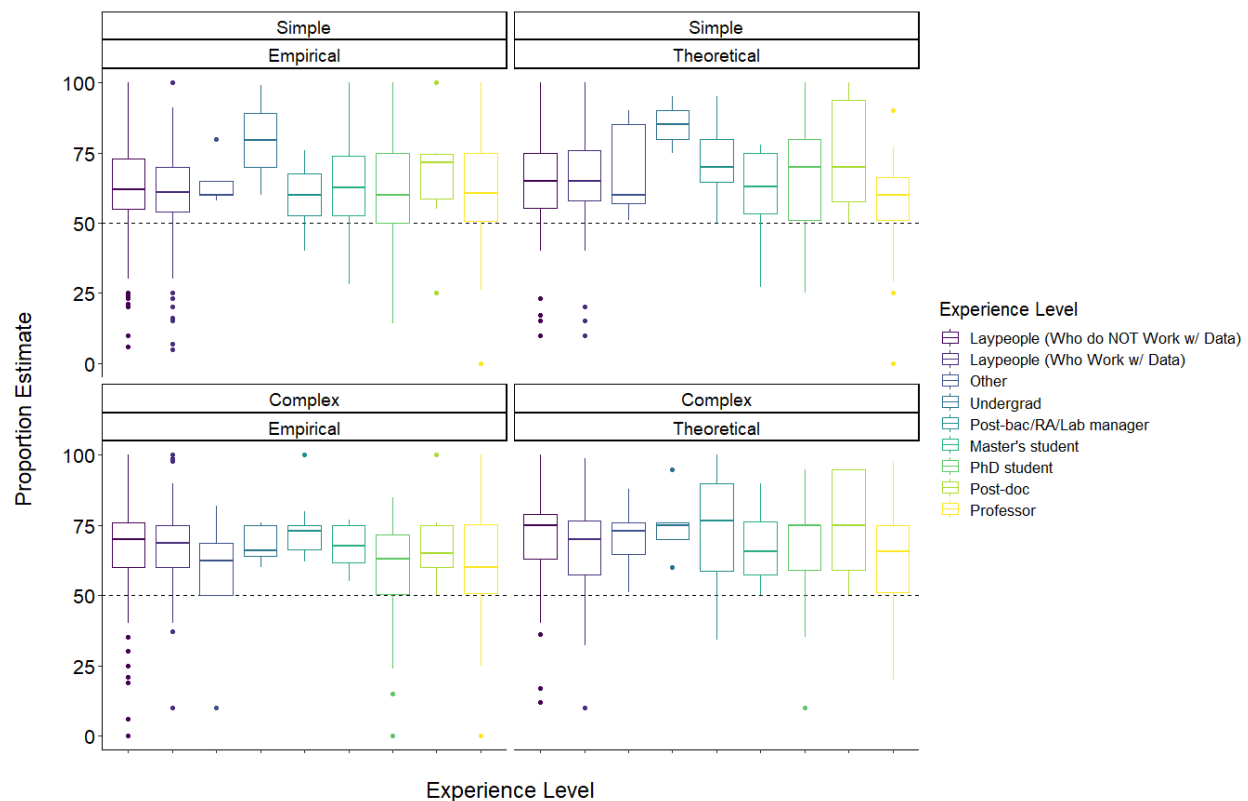
Effect Type	Population	Proportion	<i>p</i> -value
<b>Simple</b>	<i>Laypeople</i>	81% [76% - 85%]	< .001
	<i>Researchers</i>	73% [64% - 81%]	< .001
<b>Complex</b>	<i>Laypeople</i>	88% [84% - 92%]	< .001
	<i>Researchers</i>	80% [72% - 87%]	< .001

*Note:* Proportions of laypeople/researchers who indicated that the empirical proportion of the study's participants who matched the claim was at least a simple majority. Brackets underneath proportions indicate 95% CIs for the proportion estimate. P-values were computed via binomial tests against 0.50.

**Table 6. Theoretical Estimate Tests within Each Effect Type (split by Population)**

Effect Type	Population	Proportion	<i>p</i> -value
<b>Simple</b>	<i>Laypeople</i>	93% [90% - 96%]	< .001
	<i>Researchers</i>	80% [72% - 87%]	< .001
<b>Complex</b>	<i>Laypeople</i>	92% [89% - 95%]	< .001
	<i>Researchers</i>	90% [83% - 95%]	< .001

*Note:* Proportions of laypeople/researchers who indicated that the proportion of the study's participants who needed to match the claim was at least a simple majority if the results were to be used to support an individual-level psychological theory. Brackets underneath proportions indicate 95% CIs for the proportion estimate. P-values were computed via binomial tests against 0.50.



**Figure 2.** Boxplots of empirical/theoretical proportion estimates by effect type (simple versus complex), and by participants' level of experience. Note that "Other" ( $n = 13$ ) refers to people involved in academic research in some way (via SPSP) but who indicated that they have never held an academic position.

## General Discussion

Drawing on recent "persons-as-effect-sizes" and "pervasiveness" approaches (Grice et al., 2020; Speelman & McGann, 2020), the current work documents instances of social psychological claims, derived from typical sets of group-level statistical tests, that upon re-analysis are quite poor representations of person-level psychology. Our work extends these approaches, showing that group-level effects in multi-factor experiments (and single-factor experiments with more than two levels) cannot provide the person-level evidence that psychologists likely desire (i.e., "Do people respond this way, and if so, how many?"). Additionally, the current research experimentally tested multiple method-based noise explanations for this group-to-person generalizability problem in a moral judgment paradigm,

1  
2  
3 with obvious remedies failing to address the problem. Finally, our research shows that a majority  
4 of laypeople and social psychology researchers interpret authors of psychology articles as  
5 intending to make claims that represent a majority of their study's participants. Moreover, a  
6 majority of laypeople and researchers believe that this ought to be the case if authors are using a  
7 study's results to claim support for a general, individual-level psychological theory.  
8  
9  
10  
11  
12  
13

14 Our research is consistent with recent critiques put forth, in which some researchers (e.g.,  
15 Richters, 2021; Speelman & McGann, 2020) have argued that there is a pervasive mismatch  
16 between psychological theorizing and the analytic procedures used for testing—typical  
17 theorizing occurs at the person-level but analytic procedures operate at the group-level. Over the  
18 past decade, much effort has gone toward correcting, and promoting better, statistical inferences  
19 (e.g., Lakens, 2021), but relatively fewer reform efforts have been aimed at appropriate  
20 psychological inference (e.g., Moeller et al., *preprint*; Navarro, 2019; Liew, Howe, & Little,  
21 2016) and proper theory development (e.g., van Rooij & Baggio, 2021). The current research  
22 suggests that even if theorizing indeed improves, inference can still go wrong if familiar  
23 statistical methods are privileged over ones that address specific psychological questions. Put  
24 simply, psychologists seem to have put the statistical cart ahead of the psychological horse. This  
25 problem, however, should not be judged as just another instance of “psychology in crisis.”  
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

### *Responses to Potential Objections*

One objection is that, if the modal person-level pattern matches the group-level pattern,  
then there is no group-to-person generalizability problem. We suggest, however, that if the  
modal person-level pattern matches the group-level pattern but is not *also* the pattern describing

1  
2  
3 most participants, then the majority of responses would be unexplained. Moreover, most polled  
4  
5 laypeople and social psychology researchers agreed that for a claim to support a theory, it ought  
6  
7 to represent a majority of participants. Extrapolating from this, it stands to reason that *any*  
8  
9 reported effect that does not describe a majority of participants (regardless of its use for theory)  
10  
11 ought to be tagged as unrepresentative.  
12  
13

14  
15 A related objection is that the proportion of participants showing the group-level pattern  
16  
17 may be above chance given how many possible qualitative patterns exist, meaning that there is  
18  
19 no group-to-person generalizability problem. For example, in a simple two-cell design, there are  
20  
21 only three possible person-level effects (i.e., positive, negative, or zero). Based on this objection,  
22  
23 if any of those effects described more than 33% of the sample and was consistent with the group-  
24  
25 level pattern, then the study would not contain the group-to-person generalizability problem. Our  
26  
27 reasoning against the first objection also applies here: if the modal pattern does not also describe  
28  
29 most participants, then the group-level pattern represents only a minority of responses. In  
30  
31 addition, there is unfortunately no ground truth to a “better than chance” criterion. That is,  
32  
33 designs could be modified to include additional conditions which are likely to engender uniform  
34  
35 responses across participants, creating more possible qualitative patterns while artificially  
36  
37 reducing the chance of any single pattern.  
38  
39  
40  
41

42  
43 A final objection is that there are other sources of noise accounting for the group-to-  
44  
45 person generalizability problem, beyond those tested here (see SOM). For example, some  
46  
47 participants are distracted, leading to frequencies of person-level patterns which do not represent  
48  
49 the “true” frequencies. First, consistent with our experimental results, there is no reason to  
50  
51 believe, if such noise was reduced, that most person-level patterns would conveniently shift to  
52  
53 the group-level pattern. Second, as our tutorial and hypothetical datasets show, there are simple  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 non-method explanations for how group-level patterns can be (even perfectly) unrepresentative  
4 of persons. Therefore, rather than assuming that there are solvable methodological issues  
5 underlying the problem, it should be accepted that person-level patterns cannot be inferred from  
6 group-level analyses.  
7  
8  
9

### 10 11 12 *Recommendations* 13

14  
15 Given the group-to-person generalizability problem, what should experimental  
16 psychologists do? Our recommendations are consistent with those in a recent critique (Yarkoni,  
17 2020). We note that our generalizability critique refers to generalizing across levels of analysis,  
18 whereas the recent critique refers to generalizing across stimuli, tasks, etc. Specifically, we  
19 propose predicting specific orderings of observations based on theory (e.g., A1 higher than A2 in  
20 B1, but A2 higher than A1 in B2), while specifying a proportion of participants whose responses  
21 should match predictions for the theory-derived hypothesis to survive. To conduct a “severe test”  
22 (Mayo, 2018) or corroborate a “risky prediction,” (Meehl, 1990a, 1990b), the empirical  
23 proportion should be close to the theory-predicted proportion, and other theories should not  
24 predict this proportion.  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36

37  
38 Unfortunately, there may not be many psychological theories (especially in social  
39 cognition which typically lacks formal models) that can make such predictions. However,  
40 theoretical progress can still be made. To test face validity, researchers can make minimum  
41 proportion predictions. To do this, researchers can tabulate the proportion of participants whose  
42 responses match a predicted pattern and test this proportion against 50%, using a binomial test.  
43  
44 To be clear, we do not advocate 50% as the benchmark against which psychologists should test  
45 theory; we are simply suggesting a method to discard those which fail to support general  
46 psychological regularities (Hamaker, 2012). Others have suggested even higher proportions as  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 convincing evidence (e.g., 80%; Speelman & McGann, 2020), though responses from the current  
4  
5 work suggest that psychologists disagree about the appropriate cutoff. We also do not  
6  
7 recommend ignoring theory-inconsistent patterns, or patterns represented by a small minority of  
8  
9 participants. Understanding if and why other patterns exist allows refinement of theory by  
10  
11 postulating and testing whether there are substantive moderating variables (e.g., individual  
12  
13 differences), or simple violations of auxiliary assumptions (e.g., divergent interpretations of  
14  
15 measures; see Quintana, 2021, for a discussion).  
16  
17

### 18 19 *Limitations and Future Directions*

20  
21 One constraint of this person-level approach is that it ignores magnitude information  
22  
23 (e.g., participants who use two extreme ends of a measure are treated identically to participants  
24  
25 who use two close points of a measure). However, magnitude information can be incorporated  
26  
27 into this approach. Researchers can choose an “imprecision value” (Grice et al., 2020), allowing  
28  
29 only certain magnitudes to support a qualitative pattern. Additionally, researchers can plot  
30  
31 frequencies of qualitative patterns by different imprecision values, allowing discernment  
32  
33 between participants who show small versus large effects (see Speelman & McGann, 2020,  
34  
35 Figure 4).  
36  
37

38  
39 Relatedly, there are other (potentially better) methods for evaluating person-level effects  
40  
41 in high-repetition studies which yield magnitude information, such as person-level standardized  
42  
43 effect sizes and confidence intervals (e.g., Kurz, Johnson, Kellum, & Willer, 2019). However,  
44  
45 the intensive sampling methods these require are often infeasible in social cognition research.  
46  
47 The current method can be used for both high-repetition and single-response-per-condition  
48  
49 studies. In high-repetition designs, an analyst can simply average over an individuals’ multiple  
50  
51 responses within each condition and then apply the person-level analysis to those averages. This  
52  
53  
54  
55  
56  
57  
58  
59  
60



1  
2  
3 suggestion is not inconsistent with our argument to avoid making inferences from averages, as  
4 the suggested averaging would occur within rather than across persons. Other strengths of the  
5 current approach are clear: it requires no advanced statistical knowledge, is easy to implement  
6 and interpret, and therefore, is easy to communicate.  
7  
8  
9  
10

11  
12 Another limitation is that we used only one moral judgment paradigm to test method-  
13 based noise explanations for the group-to-person generalizability problem. Additionally, much  
14 research in moral cognition—including our current experiments (see SOM)—utilizes on-the-fly  
15 measurement practices (see Flake & Fried, 2020). Future research is needed to determine  
16 whether method manipulations fail to remedy the problem in other paradigms and areas of  
17 psychology with better measurement practices. However, as shown earlier, there are obvious  
18 non-method (and non-measurement) explanations for the problem. Therefore, a person-level  
19 approach should still be used in more rigorous disciplines to ensure generalizability.  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30

31 Finally, we did not assess the ubiquity of the group-to-person generalizability problem.  
32 We simply documented (and replicated) existence proofs. Documenting its ubiquity is a  
33 necessary next step to examine its generalizability (see Simons, Shoda, & Lindsay, 2017). We  
34 hypothesize that the problem will be least frequent in low-level research, such as visual and  
35 auditory cognition, due to assumptions that most people share basic physiological similarities,  
36 and some lack of conscious control, which underlie these disciplines' studied effects. As the  
37 content of study becomes higher-level, the problem should be more prevalent, as responses will  
38 rely more on individual differences (e.g., values and knowledge). It may even be the rule rather  
39 than the exception in some areas, such as social cognition.  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50

## 51 **Conclusion**

52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 Psychological scientists often make claims about, and interpret others' claims as being  
4 about, person-level cognition. Sometimes, however, these claims are made from experiments  
5 which disallow investigation of person-level phenomena. Even when such investigation is  
6 possible, these claims are typically derived from group-level patterns, interpreted *as if* they  
7 reveal the pervasiveness of some phenomenon. The current work confirms and builds upon  
8 previous warnings that this practice can lead to serious errors in inference, as (sets of) group-  
9 level patterns need not reflect even a simple majority of sampled individuals. Put simply,  
10 psychology is a feature of persons, not averages or distributions. Therefore, person-level design  
11 and analytic approaches should be customary in psychological science.  
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

## Footnotes

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

1. Anecdotally, three points should be noted: (1) the first author spent less than a week finding and reanalyzing these instances, (2) instances reported here were the *only* publications tested, (3) many seemingly person-level claims (not reported here) were made from between-subjects designs that fundamentally disallowed reanalysis.

2. In the researcher sample, a small minority used the open-ended question to *correctly* communicate that inferences about percentages cannot be derived from average differences ( $n = 17$ ). Therefore, some of the empirical estimates were not true beliefs, as the researchers simply had no other option but to respond. (Despite our methods experiments, we ironically did not think to allow participants to opt out). To conduct the most stringent test of our hypothesis, we recoded all of these hypothesis-consistent slider responses ( $n = 6$ ) as being hypothesis-inconsistent. We did not remove any of the 17 responses to ensure that, even accounting for some researchers understanding the problem, a majority still responded in a hypothesis-consistent way. This resulted in similar proportions for both simple (70%) and complex (79%) effects.

## References

- Brandt, M.J., & Morgan, G.S. (2022). Between-person methods provide limited insight about within-person belief systems. *Journal of Personality and Social Psychology*.
- Craig, B.M., Nelson, N.L., & Dixson, B.J.W. (2019). Sexual selection, agnostic signaling, and the effect of beards on recognition of men's anger displays. *Psychological Science, 30*(5), 728-738.
- Decelles, K.A., Adamas, G.S., Howe, H.S., & John, L.K. (2021). Anger damns the innocent. *Psychological Science, 32*(8), 1214-1226.
- Deska, J.C., Kuntsman, J., Lloyd, P.E., Almaraz, S.M., Bernstein, M.J., Gonzales, J.P., & Hugenberg, K. (2020). Race-based biases in judgments of social pain. *Journal of Experimental Social Psychology, 88*, 103964.
- Devezer, B., Navarro, D.J., Vandekerckhove, J., & Buzbas, E.O. (2021). The case for formal methodology in scientific reform. *Royal Society Open Science, 8*, 200805.
- Fisher, A.J., Medaglia, J.D., & Jeronimus, B.F. (2018). Lack of group-to-individual generalizability is a threat to human subjects research. *Proceedings of the National Academy of Sciences, 115*(27), E6106-E6115.
- Flake, J.K., & Fried, E.I. (2020). Measurement schmeasurement: Questionable measurement practices and how to avoid them. *Advances in Methods and Practices in Psychological Science, 3*(4), 456-465.
- Fowler, Z., Law, K.F., & Gaesser, B. (2021). Against empathy bias: The moral value of equitable empathy. *Psychological Science, 32*(5), 766-779.
- Galton, F. (1907). Vox populi. *Nature, 75*, 450-451.

- 1  
2  
3 Grice, J.W., Medellin, E., Jones, I., Horvath, S., McDaniel, H., O'lansen, C., & Baker, M.  
4  
5 (2020). Persons as effect sizes. *Advances in Methods and Practices in Psychological*  
6  
7 *Science*, 3(4), 443-455.  
8  
9
- 10 Hamaker, E. (2012). Why researchers should think “within-person”: A paradigmatic rationale. In  
11  
12 M.R. Mehl & T.S. Conner (Eds.). *Handbook of Research Methods for Studying Daily*  
13  
14 *Life*, 43-61, NY, NY: Guilford.  
15  
16
- 17 Lakens, D. (2021). The practical alternative to the p-value is the correctly used p-value.  
18  
19 *Perspectives on Psychological Science*, 16(3), 639-648.  
20  
21
- 22 Law, K.F., Campbell, D., & Gaesser, B. (2021). Biased benevolence: The perceived morality of  
23  
24 effective altruism across social distance. *Personality and Social Psychological Bulletin*,  
25  
26 48(3), 426-444.  
27  
28
- 29 Liew, S.H., Howe, P.D.L., & Little, D.R. (2016). The appropriacy of averaging in the study of  
30  
31 context effects. *Psychonomic Bulletin and Review*, 23(5), 1639-1646.  
32  
33
- 34 Kievit, R.A., Frankenhuys, W.E., Waldorp, L.J., & Borsboom, D. (2013). Simpson’s paradox in  
35  
36 psychological science: A practical guide. *Frontiers in Psychology*, 4, 513.  
37  
38
- 39 Kuppens, T. Pollet, T.V. (2014). Mind the level: Problems with two recent national-level  
40  
41 analyses in psychology. *Frontiers in Psychology*, 5, 1110.  
42  
43
- 44 Kurz, A.S., Johnson, Y.L., Kellum, K.K., & Wilson, K.G. (2019). How can process-based  
45  
46 researchers bridge the gap between individuals and groups? Discover the dynamic p-  
47  
48 technique. *Journal of Contextual Behavioral Science*, 13, 60-65.  
49  
50
- 51 Mayo, D.G. (2018). *Statistical inference as severe testing*.  
52  
53  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3 McManus, R.M., Mason, J.E., Young, L. (2021). Re-examining the role of family relationships  
4 in structuring perceived helping obligations, and their impact on moral evaluation.  
5  
6 *Journal of Experimental Social Psychology*, 96, 104182.  
7  
8  
9
- 10 Meehl, P.E. (1990a). Appraising and amending theories: The strategy of Lakatosian defense and  
11 two principles that warrant it. *Psychological Inquiry*, 1(2), 108-141.  
12  
13  
14
- 15 Meehl, P.E. (1990b). Why summaries of research on psychological theories are often  
16 uninterpretable. *Psychological Reports*, 66(1), 195-244.  
17  
18
- 19 Moeller, J. (2022). Averting the next credibility crisis in psychological science. Within-person  
20 methods for personalized diagnostic and intervention. *Journal for Person-Oriented*  
21 *Research*, 7(2), 53-77.  
22  
23  
24
- 25  
26 Moeller, J. et al. (*preprint*). Generalizability crisis meets heterogeneity revolution: Determining  
27 under which boundary conditions findings replicate and generalize.  
28  
29
- 30  
31 Navarro, D.J. (2019). Between the Devil and the Deep Blue Sea: Tensions between scientific  
32 judgment and statistical model selection. *Computational Brain and Behavior*, 2(1), 28-34.  
33  
34
- 35 Quintana, D.S. (2021). Towards better hypothesis tests in oxytocin research: Evaluating the  
36 validity of auxiliary assumptions. *Psychoneuroendocrinology*, 105642.  
37  
38  
39
- 40 Richters, J.E. (2021). Incredible utility: The lost causes and causal debris of psychological  
41 science. *Basic and Applied Social Psychology*, 43(6), 366-405.  
42  
43  
44
- 45 Rottman, J., & Young, L. (2019). Specks of dirt and tons of pain: Dosage distinguishes impurity  
46 from harm. *Psychological Science*, 30(8), 1151-1160.  
47  
48
- 49 Simons, D.J., Shoda, D.Y., & Lindsay, D.S. (2017). Constraints on generality (COG): A  
50 proposed addition to all empirical papers. *Perspectives on Psychological Science*, 12(6),  
51 1123-1128.  
52  
53  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3 Simpson, E.H. (1951). The interpretation of interaction in contingency tables. *Journal of the*  
4  
5 *Royal Statistical Society. Series B (Methodological)*, 13(2), 238-241.  
6  
7  
8 Soter, L.K., Berg, M.K., Gelman, S.A., & Kross, E. (2021). What we would (but shouldn't) do  
9  
10 for those we love: Universalism versus partiality in responding to others' moral  
11  
12 transgressions. *Cognition*, 217, 104886.  
13  
14  
15 Speelman, C.P., & McGann, M. (2020). Statements about the pervasiveness of behavior require  
16  
17 data about the pervasiveness of behavior. *Frontiers in Psychology*, 11, 1-16.  
18  
19  
20 Stroessner, S.J., Benitez, J., Perez, M.A., Wyman, A.B., Carpinella, C., Johnson, K.L. (2020).  
21  
22 What's in a shape? Evidence of gender category associations with basic forms. *Journal of*  
23  
24 *Experimental Social Psychology*, 87, 103915.  
25  
26  
27 Surowiecki, J. (2005). *The wisdom of crowds*.  
28  
29 Thai, M., Borgella, A.M., & Sanchez, M.S. (2019). It's only funny if we say it: Disparagement  
30  
31 humor is better if it originates from a member of the group being disparaged. *Journal of*  
32  
33 *Experimental Social Psychology*, 85, 103838.  
34  
35  
36 Van Rooij, I., & Baggio, G. (2021). Theory before the test. How to build high-verisimilitude  
37  
38 explanatory theories in psychological science. *Perspectives on Psychological Science*,  
39  
40 16(4), 682-697.  
41  
42  
43 Wallis, K.F. (2014). Revisiting Francis Galton's forecasting competition. *Statistical Science*,  
44  
45 29(3), 420-424.  
46  
47  
48 Whitsett, D.D., & Shoda, Y. (2014). An approach to test for individual differences in the effects  
49  
50 of situations without using moderator variables. *Journal of Experimental Social*  
51  
52 *Psychology*, 50(1), 94-104.  
53  
54  
55 Yarkoni, T. (2020). The generalizability crisis. *Behavioral and Brain Sciences*, 45, E1.  
56  
57  
58  
59  
60