

1 Reasons or Rationalisations: The Role of Principles in the Moral Dumbfounding Paradigm

2 Cillian McHugh¹, Marek McGann², Eric R. Igou¹, & Elaine L. Kinsella¹

3 ¹ University of Limerick

4 ² Mary Immaculate College ~ University of Limerick

5 Author Note

6 All procedures performed in studies involving human participants were approved by
7 institutional research ethics committee and conducted in accordance with the Code of
8 Professional Ethics of the Psychological Society of Ireland, and with the 1964 Helsinki
9 declaration and its later amendments or comparable ethical standards. Informed consent was
10 obtained from all individual participants included in the study. The authors declare that
11 there are no potential conflicts of interest with respect to the research, authorship, and/or
12 publication of this article. All authors consented to the submission of this manuscript.

13 Correspondence concerning this article should be addressed to Cillian McHugh,
14 Department of Psychology, University of Limerick, Limerick, Ireland, V94 T9PX. E-mail:
15 cillian.mchugh@ul.ie

Abstract

16

17 Moral dumbfounding occurs when people maintain a moral judgment even though they
18 cannot provide reasons for it. Recently, questions have been raised about whether
19 dumbfounding is a real phenomenon. Two reasons have been proposed as guiding the
20 judgments of dumbfounded participants: harm-based reasons (believing an action may cause
21 harm) or norm-based reasons (breaking a moral norm is inherently wrong). Participants who
22 endorsed either reason were excluded from analysis, and instances of moral dumbfounding
23 seemingly reduced to non-significance. We argue that endorsing a reason is not sufficient
24 evidence that a judgment is grounded in that reason. Stronger evidence should additionally
25 account for (a) articulating a given reason, and (b) consistently applying the reason in
26 different situations. Building on this, we develop revised exclusion criteria across 2 studies.
27 Study 1 included an open-ended response option immediately after the presentation of a
28 moral scenario. Responses were coded for mention of harm-based or norm-based reasons.
29 Participants were excluded from analysis if they both articulated and endorsed a given
30 reason. Using these revised criteria for exclusion, we found evidence for dumbfounding, as
31 measured by the selecting of an admission of not having reasons. Study 2 included a further
32 three questions relating to harm-based reasons specifically, assessing the consistency with
33 which people apply harm-based reasons across differing contexts. As predicted, few
34 participants consistently applied, articulated, and endorsed harm-based reasons, and
35 evidence for dumbfounding was found.

36

Keywords: morality, moral judgment, dumbfounding, intuition, rationalism, reasons

37

Word count: 11,298

38 Reasons or Rationalisations: The Role of Principles in the Moral Dumbfounding Paradigm

39

List of Tables

40 *Table 1.* Multinomial logistic regression predicting responses to the critical slide
41 where providing reasons is the referent in each case.

42 *Table 2.* Multinomial logistic regression predicting responses to the critical slide
43 where providing reasons is the referent in each case.

44

List of Figures

45 *Figure 1.* Responses to critical slide for the entire sample, and for each measure
46 of convergence: (i) endorsing only, and (ii), endorsing and articulating;
47 percentages of full sample displayed within plot, percentages of relevant
48 sample displayed in parenthesis below the count.

49 *Figure 2.* Responses to critical slide for the entire sample, and for each measure
50 of convergence: (i) endorsing only, (ii) endorsing and articulating, and
51 (iii), endorsing, articulating, and applying; percentages of full sample
52 displayed within plot, percentages of relevant sample displayed in paren-
53 thesis below the count.

54 *Figure 3.* Probability of selecting each response to the critical slide depending on
55 Religiosity

56 *Figure 4.* Probability of selecting each response to the critical slide depending on
57 MLQ: Presence

58 *Figure 5.* Probability of selecting each response to the critical slide depending on
59 MLQ: Search

60 *Figure 6.* Probability of selecting each response to the critical slide depending on
61 Initial Judgement.

1 | Introduction

62

63 Moral dumbfounding occurs when people maintain a moral judgment even though they
64 cannot provide a reason in support of this judgment (Haidt, 2001; Haidt, Björklund, &
65 Murphy, 2000). It is typically evoked when people encounter taboo behaviors that do not
66 result in any harm (Haidt, 2001; Haidt et al., 2000; see also McHugh et al., 2017). One
67 example of such a behavior can be found in the widely discussed *Incest* scenario, which reads
68 as follows:

69

Julie and Mark, who are brother and sister are traveling together in France.

70

They are both on summer vacation from college. One night they are staying

71

alone in a cabin near the beach. They decide that it would be interesting and fun

72

if they tried making love. At very least it would be a new experience for each of

73

them. Julie was already taking birth control pills, but Mark uses a condom too,

74

just to be safe. They both enjoy it, but they decide not to do it again. They

75

keep that night as a special secret between them, which makes them feel even

76

closer to each other. (Haidt et al., 2000, p. 22)

77

Incest is considered taboo in most cultures, and in violating this taboo, Julie and

78

Mark's actions are typically judged as wrong. However, the consensual and harmless nature

79

of their actions means that the reasons people generally provide do not apply in this case.

80

People who maintain their judgment in the absence of reasons are identified as morally

81

dumbfounded. McHugh et al. (2017), building on the original work by Haidt et al. (2000),

82

identified two measurable responses that may be taken as indicators of moral dumbfounding.

83

Firstly, people may explicitly admit to not having reasons for their judgment. Secondly,

84

people may use unsupported declarations ("it's just wrong") or tautological reasons

85

("because it's incest") as justifications for a judgment.

86 1.1 | The Influence of Moral Dumbfounding

87 The discovery of moral dumbfounding (Haidt et al., 2000; see also Haidt, Koller, &
88 Dias, 1993) coincided with, and arguably contributed to, some of the key developments in
89 moral psychology over the past two decades. It had a clear influence on the development of
90 Haidt's social intuitionist model of moral judgment (SIM, Haidt, 2001), and by extension
91 may be seen as contributing to the growth of intuitionist theories of moral judgment that
92 followed (e.g., Cushman, Young, & Greene, 2010; Haidt, 2001; Prinz, 2005).

93 Haidt proposed the SIM in opposition to the perceived dominance of rationalist
94 approaches (Kohlberg, 1969, 1971; Narvaez, 2005; Topolski, Weaver, Martin, & McCoy,
95 2013). According to rationalist approaches our moral judgments are grounded in reason,
96 informed by discernible moral principles (Fine, 2006; Kennett & Fine, 2009; Kohlberg, 1969,
97 1971; Royzman, Kim, & Leeman, 2015); Haidt (2001). Moral dumbfounding is presented by
98 Haidt (2001) and by Prinz (2005) as evidence against this rationalist perspective, in that, if
99 moral judgments were grounded in reason, people would be able to provide reasons for their
100 judgments (and moral dumbfounding would not occur). Intuitionist theorists propose that
101 moral judgments are grounded in an emotional or intuitive automatic response rather than
102 slow deliberate reasoning (Cameron, Payne, & Doris, 2013; Haidt, 2001; Prinz, 2005). In
103 recent years the joint role of reason/deliberation and intuition in the making of moral
104 judgments has been emphasised in dual-process theories (Brand, 2016; Crockett, 2013;
105 Cushman, 2013a; Cushman et al., 2010; Greene, 2008). The dumbfounding paradigm may be
106 useful in developing and extending these theories; developing an understanding of moral
107 dumbfounding and the processes that lead to it, may inform the further development of
108 theories of moral judgment, leading to a greater understanding of the processes that underlie
109 moral judgment more generally.

110 The influence of dumbfounding may be observed in everyday discourse, particularly in
111 relation to highly sensitive and divisive social issues. Real-world interactions differ from a

112 laboratory study designed to elicit a dumbfounded response, and as such, in the absence of
113 explicit and consistent refuting of arguments, it is unlikely that people in everyday life would
114 admit to not having reasons for their moral judgments. Despite this, it is not uncommon to
115 hear unsupported declarations/tautological statements as arguments in support of a position
116 with no further justification (e.g., Mustonen, Paakkonen, Ryökäs, & Nieminen, 2017;
117 Stepniak, 1995). Similarly, moral positions are often justified by appealing to emotions (e.g.,
118 Mustonen et al., 2017; Stepniak, 1995; see also Rozin, Haidt, MacCauley, McKay, & Olatunji,
119 2008; Rozin, Lowery, Imada, & Haidt, 1999). This type of appeal to emotion has previously
120 been discussed as similar/equivalent to dumbfounding (see Prinz, 2005, p. 101; see also
121 Haidt & Hersh, 2001). These responses may not clearly demonstrate dumbfounding, however
122 they illustrate the way in which discussions of reasons for moral positions are occasionally
123 absent from the public debate.

124 That people may defend a judgment in the absence of articulated reasons, and
125 maintain it even in the knowledge of their own inconsistencies poses a challenge for the type
126 of rational debate that is supposed to form the basis of public discourse and inform the
127 development of public policy. The study of moral dumbfounding, as an extreme case, may
128 lead to a better understanding of the underlying cognitive processes that lead to these types
129 of problematic practices that have no place in public debate. Identifying these processes and
130 explaining moral dumbfounding is beyond the scope of the current research. Here, in light of
131 recent critiques, here we test whether or not dumbfounding is a real phenomenon, worthy of
132 further study.

133 **1.2 | Challenging the Dumbfounding Paradigm**

134 A key concern regarding the dumbfounding paradigm is that the eliciting scenarios
135 have been artificially construed to remove potentially harmful consequences to the point that
136 they become unrealistic or otherwise not credible (e.g., Jacobson, 2012). It could be argued
137 that studying such idiosyncratic scenarios does little to inform our understanding of

138 everyday moral decision making; similar criticisms have been made regarding the widely
139 used trolley-type sacrificial dilemmas (e.g., Bauman, McGraw, Bartels, & Warren, 2014;
140 Bostyn, Sevenhant, & Roets, 2018). However, responses to hypothetical trolley dilemmas
141 have been found to predict behaviour in a money burning game with real pay-off
142 consequences (Dickinson & Masclot, 2018), and the study of trolley-type dilemmas arguably
143 contributed to key theoretical advancements of the past two decades (e.g., Plunkett &
144 Greene, 2019; see also Greene, 2008; Christensen, Flexas, Calabrese, Gut, & Gomila, 2014;
145 Christensen & Gomila, 2012; Greene, Sommerville, Nystrom, Darley, & Cohen, 2001). If
146 moral dumbfounding is a real phenomenon it may prove a useful paradigm to further
147 advance theories of moral judgment, and examine the mechanisms and cognitive processes
148 that underlie the making of moral judgments (e.g., the relative roles of emotion versus
149 deliberation). It may be possible to identify specific contextual features that may lead people
150 to change their mind rather than provide a dumbfounded response (or vice-versa).
151 Experimental manipulations that may increase dumbfounded responding (e.g., cognitive
152 load) or reduce dumbfounded responding (e.g., distancing) could be investigated. There may
153 also be individual difference variables that predict susceptibility to dumbfounding.

154 In defending the claim that moral judgments are not caused by reasoning, Haidt (2001)
155 presents moral dumbfounding as a demonstration of inconsistency between judgment and
156 reasons available. The implicit alternative to this argument is that the absence of reasons
157 would lead a moral judgment to change or to be revised; i.e., the presence or absence of
158 reasons can cause a judgment to change. Haidt does not clearly distinguish between
159 *reasoning* as a cause versus *reasons* as a cause of judgments (e.g., 2001, p. 822). Despite
160 being inconsistent with approaches beyond the moral domain (e.g., Mercier, 2016;
161 Johnson-Laird, 2006; Mercier & Sperber, 2011, 2017; Todd & Gigerenzer, 2012), this
162 ambiguity can still be seen in discussions of moral judgment (and moral dumbfounding),
163 such that, for the rationalist perspective (see Haidt, 2001), reasons appear to play a causal
164 role, (e.g., Jacobson, 2012, p. 17; Flanagan, Sarkissian, & Wong, 2008, p. 7; Triskiel, 2016, p.

165 93). Furthermore, this assumption is implicit in challenges to the dumbfounding narrative,
166 whereby these challenges attempt to demonstrate that people do have “warrantable reasons”
167 for their judgments (Royzman et al., 2015, p. 309). Here we identify and address
168 methodological limitations of one example of this type of challenge to the dumbfounding
169 paradigm (Royzman et al., 2015).

170 Gray, Schein, and Ward (2014) argue that people’s moral judgments are grounded in
171 harm-based reasons, suggesting that when judging moral scenarios, people implicitly perceive
172 harm even in scenarios that are construed as objectively harmless. If people perceive harm in
173 the scenarios, then, even when the experimenter claims that they are harm free, this
174 perception of harm still serves as a reason to condemn the behavior. They conducted a series
175 of experiments demonstrating that people do implicitly perceive harm in supposedly
176 victim-less scenarios; e.g., “masturbating to a picture of one’s dead sister, watching animals
177 have sex to become sexually aroused, having sex with a corpse, covering a Bible with feces”
178 (Gray et al., 2014, p. 1063). This suggests that in studies of moral dumbfounding people
179 may also be making judgments based on an implicit perception of harm.

180 Jacobson (2012) makes specific reference to the scenarios used in the study of moral
181 dumbfounding, and presents a number of plausible reasons why a person may condemn the
182 actions of the characters in these scenarios. In the case of the *Incest* scenario, he suggests
183 that the behavior of Julie and Mark was risky, “reckless and licentious” (Jacobson, 2012, p.
184 25). Jacobson also discusses another scenario, *Cannibal*, that has been used in studies of
185 moral dumbfounding. This scenario describes an act of cannibalism by a researcher in a
186 pathology lab (Jennifer) on a cadaver from the lab. Jacobson argues that if Jennifer’s
187 behavior became known, people would be less willing to donate their bodies to the lab. In
188 addition to providing reasons that may explain the judgments of participants, Jacobson
189 suggests that when participants appear to be dumbfounded they have simply given up on
190 the argument and conceded to the experimenter who is in a position of authority. While this

191 claim is not directly tested empirically by Jacobson, it has been studied by Royzman et al.
192 (2015), as discussed in the following section.

193 **1.3 | Evidence for Judgments Based on Reasons or Principles**

194 A recent series of studies by Royzman et al. (2015), investigating the *Incest* scenario
195 specifically, aimed to identify if participants presenting as dumbfounded genuinely had no
196 reasons to support their judgments. In line with Jacobson (2012), they claim that
197 dumbfounding occurs as a result of social pressure to adhere to conversational norms,
198 arguing that dumbfounded participants do have reasons for their judgments and that these
199 reasons are incorrectly dismissed as invalid by the experimenter. They argue that
200 dumbfounded responding occurs as a result of social pressure to avoid appearing
201 “uncooperative” (Royzman et al., 2015, p. 299), “inattentive” or “stubborn” (2015, p. 300).
202 In addition to this claim, Royzman et al. (2015) identify two justifying principles that may
203 be guiding participants’ judgments: the harm principle and the norm principle. They argue
204 that when excluding from analysis participants who endorse either of these principles,
205 incidences of dumbfounding are negligible.

206 In identifying the *harm principle*, Royzman et al. (2015) draw on the work of Gray et
207 al. (2014). They hypothesised that participants may not believe the scenario to be harm free
208 even in the face of repeated assurances from the experimenter that it is harm free. If a
209 participant does not believe that an act is truly harm free then this provides them with a
210 perfectly valid reason to judge it as morally wrong (Gray et al., 2014; Royzman et al., 2015).
211 They devised two questions which served as a “credulity check” (Royzman et al., 2015, p.
212 309), to assess whether or not participants believed that the *Incest* scenario was harm-free.
213 The questions read as follows: (i) “Having read the story and considering the arguments
214 presented, are you able to believe that Julie and Mark’s having sex with each other will not
215 negatively affect the quality of their relationship or how they feel about each other later
216 on?”; (ii) “Having read the story and considering the arguments presented, are you able to

217 believe that Julie and Mark’s having sex with each other will have no bad consequences for
218 them personally and/or for those close to them?” (Royzman et al., 2015, p. 302–303). If
219 participants responded “No” to either of these questions, their judgments were attributed to
220 harm-based reasons, and therefore they could not be identified as dumbfounded.

221 The second principle identified by Royzman et al. (2015) is the *norm principle*. They
222 argue that if people believe that committing a particular act is wrong, regardless of the
223 circumstances, then, for these people, this belief may be sufficient to serve as a reason to
224 condemn the behavior of the characters in the scenario. Royzman et al. (2015) presented
225 participants with two statements: (a) “violating an established moral norm just for fun or
226 personal enjoyment is wrong only in situations where someone is harmed as a result, but is
227 acceptable otherwise”; (b) “violating an established moral norm just for fun or personal
228 enjoyment is inherently wrong even in situations where no one is harmed as a result”
229 (Royzman et al., 2015, p. 305). If participants endorsed (b) over (a) they reasoned that a
230 judgment could be legitimately defended using a normative statement. They suggest that
231 the “unsupported declarations” (Haidt et al., 2000, p. 12) identified by Haidt et al. (2000)
232 are statements of a normative position, and that, rather than being viewed as a
233 dumbfounded response, they may be viewed as reasons for judgments.

234 Royzman et al. (2015) used the credulity check to assess if participants’ judgments
235 could be attributed to the harm principle, while attributing judgments to the norm principle
236 was based on the norm statements. Royzman et al. (2015) use the phrase “fully convergent”
237 to describe participants who, in their view, are eligible for analysis (Royzman et al., 2015, p.
238 306). According to Royzman et al. (2015), a participant is fully convergent if their judgment
239 cannot be attributed to either the harm principle or the norm principle. Using these stricter
240 criteria for dumbfounding, Royzman et al. (2015) initially identified 4 participants, from a
241 sample of 53, who presented as dumbfounded. Each of these participants was then
242 interviewed and the inconsistencies in their responses pointed out to them. During these

243 interviews 2 participants changed their judgment of the behavior and 1 participant changed
244 her position on the normative statements. This left just 1 fully convergent, dumbfounded
245 participant. This participant did not resolve the inconsistency in his responses to the
246 questions, and, following post-experiment interviews, Royzman and colleagues found
247 dumbfounding to occur once in a sample of 53. This was found to be not significantly greater
248 than 0 (Royzman et al., 2015, p. 309), supporting the claim that moral dumbfounding is
249 “highly irregular” or even “non-existent” (Royzman et al., 2015, p. 300; see also Guglielmo,
250 2018).

251 **1.4 | Reasons or Rationalisations**

252 The studies conducted by Royzman et al. (2015) introduce an additional level of
253 methodological rigor to the study of moral dumbfounding. They clearly demonstrate that
254 people will endorse a reason for a judgment if it is available to them. This undermines the
255 dumbfounding narrative, that people defend a judgment in the absence of reasons, and poses
256 a strong challenge to the existence of moral dumbfounding.

257 We (McHugh et al., 2017) have previously outlined some limitations with the
258 conclusions presented by Royzman et al. (2015). Firstly, Royzman et al. (2015) suggest that
259 people who present as morally dumbfounded do so in an attempt to avoid appearing
260 “stubborn” or “inattentive” (2015, p. 310). However, Royzman et al. (2015) also employ the
261 original Haidt et al. (2000) definition of moral dumbfounding, which defines moral
262 dumbfounding as “the stubborn and puzzled maintenance of a judgment without supporting
263 reasons” (Haidt et al., 2000, p. 2; see also Haidt & Björklund, 2008, p. 197; Haidt & Hersh,
264 2001, p. 194). This means that according to Royzman et al. (2015), people who present as
265 dumbfounded, paradoxically present as stubborn in an attempt to avoid appearing stubborn.

266 Secondly, the means by which Royzman et al. (2015) arrive at their estimate of 1
267 instance of moral dumbfounding out of a sample of 53 is problematic for the claim that

268 moral dumbfounding occurs as a result of social pressure. They present their estimate of
269 1/53 as not significantly greater than 0/53 ($z = 1, p = .315$).¹ However their original
270 estimate of instances of moral dumbfounding was 4/53, which is significantly greater than
271 0/53 ($z = 2.04, p = .041$). These participants were invited back into the lab and the
272 “inconsistencies” in their “responses were pointed out directly” to them (Royzman et al.,
273 2015, p. 308). Furthermore they were then “advised to carefully review and, if appropriate,
274 revise” their responses (Royzman et al., 2015, p. 308). This procedure subjected participants
275 to social pressure to appear consistent in their responding. This illustrates that
276 dumbfounded responding can be influenced by social pressure, however it does not support
277 the stronger claim (by Royzman et al., 2015) that dumbfounded responding can be
278 attributed to social pressure (McHugh et al., 2017). The role of social pressure in eliminating
279 instances of dumbfounded responding is not acknowledged by Royzman et al. (2015).

280 Finally, demonstrating that people endorse principles that are consistent with their
281 judgments does not provide evidence that these principles are guiding their judgments. In
282 relying on participants’ endorsing of a given principle to attribute their judgment to that
283 principle, Royzman et al. (2015) may have falsely excluded some participants from analysis.
284 Consider the following scenario to illustrate this point:

285 Two friends (John and Pat) are bored one afternoon and trying to think of
286 something to do. John suggests they go for a swim. Pat declines stating that it’s
287 too much effort – to get changed, and then to get dried and then washed and
288 dried again after; he says he’d rather do something that requires less effort. John
289 agrees and adds “Oh yeah, and there’s that surfing competition on today so the
290 place will be mobbed”. To which Pat replies “Yeah exactly!” (McHugh et al.,
291 2017, p. 20)

¹ No explanation for the responding of this participant is offered. Neither can this participant’s response be explained by the theoretical position adopted by Royzman et al. (2015).

292 It is clear from reading this scenario that even though he endorsed it to support or to
293 rationalise his decision, the surfing competition was not the reason for John's decision not to
294 go to the beach. It would be incorrect to attribute his decision to this reason. The studies
295 conducted by Royzman et al. (2015) do not guard against the possibility of this type of false
296 attribution, and it is likely that some participants were incorrectly excluded from analysis on
297 this basis. This possibility of false exclusion presents a key limitation Royzman et al. (2015)
298 that casts doubt on their findings.

299 We suggest that attributing people's judgments to principles requires stronger evidence
300 than endorsing alone. We propose two measures that may be useful in establishing whether
301 or not a given principle may truly be identified as a reason for the judgments made by
302 participants. Firstly, participants should be given the opportunity to provide the reason(s)
303 that they based their judgment on, and the reasons provided should inform decisions of
304 inclusion or exclusion.² Attributing participants' judgments to particular reasons/principles
305 should account for both the endorsing and the articulating of the reason/principle. Secondly,
306 if a principle is guiding the judgments of participants, this principle should be applied
307 consistently across different contexts. We predict that when these two measures are applied
308 evidence for dumbfounding will be found.

309 1.5 | The Current Studies

310 The aim of the current studies was to investigate whether or not people's moral
311 judgments can be attributed to moral principles based on their endorsing of these principles.
312 Specifically, aim to address the concerns raised by McHugh et al. (2017) and test the claim
313 by Royzman et al. (2015) that participants' judgments in the *Incest* scenario can be
314 attributed to the harm principle or the norm principle. Firstly, the degree to which
315 participants articulate either the harm principle or the norm principle as informing their

² Participants in Royzman et al. (2015) provided reasons however these reasons did not inform their exclusion criteria.

316 judgment is examined (Study 1). Secondly, the consistency with which participants apply
317 the harm principle across differing contexts is additionally assessed (Studies 2 and 3). We
318 hypothesise that by developing more rigorous exclusion criteria the rates of false exclusion of
319 participants would be reduced and that evidence for moral dumbfounding would be found,
320 posing a challenge to the type of rationalist perspective described by Haidt (2001). The
321 failure to identify dumbfounded responding would serve as support for these alternative
322 perspectives (e.g., Gray et al., 2014; Guglielmo, 2018; Jacobson, 2012; Royzman et al., 2015;
323 Sneddon, 2007; Wielenberg, 2014) and pose a challenge to SIM as described by Haidt (2001).
324 Given that the exclusion criteria used by Royzman et al. (2015) were developed for the
325 *Incest* dilemma, the studies reported here similarly focus on the *Incest* dilemma specifically.

326 **2 | Study 1: Articulating and Endorsing**

327 In Study 1 we use an existing method for the evoking of dumbfounded responding
328 (McHugh et al., 2017), however, we incorporate to additional materials taken from Royzman
329 et al. (2015) as a more stringent set of criteria for inclusion in analysis. This serves two
330 purposes. If effective, it reduces the likelihood of false inclusions for analysis to identify rates
331 of dumbfounded responding, and also allows us to assess rates at which participants will
332 explicitly articulate or endorse the principles when given the opportunity to do so. In
333 addition to the stricter measure of inclusion proposed by Royzman et al. (2015), we
334 introduce an additional change designed to reduce the possibility of false exclusions. Study 1
335 was an extension the work of Royzman et al. (2015), using largely the same materials. One
336 moral judgment vignette (*Incest*) was taken from Haidt et al. (2000, Appendix A). Targeted
337 questions, designed to assess participants endorsements of the harm principle or the norm
338 principle, were taken directly from Royzman et al. (2015).

339 As noted above, if a participant endorses a principle this does not necessarily provide
340 evidence that this principle was guiding their judgment. Relying on the endorsing of
341 principles to determine participants' eligibility for analysis may result in some participants

342 being falsely excluded from analysis, and any resulting estimate of the prevalence of
343 dumbfounded responding would be inaccurate. In an attempt to control for the possibility of
344 falsely attributing participants' judgments to principles based on endorsing alone, we
345 included an open-ended response option to assess whether or not participants could also
346 articulate these principles. This was presented to participants immediately after the
347 presenting of the vignette. The inclusion or exclusion of participants from analysis, depended
348 on both endorsing and articulating either principle. Participants' judgments were only
349 attributed to a given principle if they both articulated and endorsed that principle. It was
350 hypothesised that participants' endorsing of a principle would not be predictive of their
351 ability to articulate this principle, and that by accounting for this, rates of false attribution
352 and false exclusion would be reduced. We hypothesised that in reducing rates of false
353 exclusion, dumbfounded responding would be observed.

354 **2.1 | Method**

355 **2.1.1 | Participants and design.** Study 1 was a frequency based extension of
356 Royzman et al. (2015). A combined sample of 110 (60 female, 49 male, 1 other; $M_{\text{age}} =$
357 32.44, min = 18, max = 69, $SD = 11.28$) took part. Fifty-eight (25 female, 32 male, 1 other;
358 $M_{\text{age}} = 38.47$, min = 19, max = 69, $SD = 12.34$) were recruited through MTurk.³
359 Participation was voluntary and participants were paid 0.50 US dollars for their
360 participation. Participants were recruited from English speaking countries or from countries
361 where residents generally have a high level of English (e.g., The Netherlands, Denmark,
362 Sweden). Fifty-two (35 female, 17 male; $M_{\text{age}} = 25.71$, min = 18, max = 38, $SD = 3.80$)
363 were recruited through direct electronic correspondence. Participants in this sample were

³ In order to prevent repeat participation from MTurk workers, this study and all remaining studies conducted on MTurk, were included as part of the same MTurk project as Study 3b from McHugh et al. (2017). In addition, a probe question was included to check if participants had encountered the scenario before. This probe included a follow-up question to determine the nature of participants' previous experience with the scenario.

364 undergraduate students, postgraduate students, and alumni from Mary Immaculate College
365 (MIC), and University of Limerick (UL). Participation was voluntary and participants did
366 not receive a reward for their participation. Previous research on moral dumbfounding found
367 responses from an MTurk sample and a College sample are largely comparable (see McHugh
368 et al., 2017 Study 3a and 3b).

369 **2.1.2 | Procedure and materials.** Data were collected using an online
370 questionnaire generated using Questback (Unipark, 2013). The questionnaire opened with
371 the information sheet and consent form. The main questionnaire was only accessible once
372 consent had been provided. Following the consent form, participants were presented with
373 questions relating to basic demographics. Participants were then presented with two
374 statements to assess if participants' judgments may be grounded in the norm principle.
375 These were taken directly from Royzman et al. (2015): (i) "violating an established moral
376 norm just for fun or personal enjoyment is wrong only in situations where someone is
377 harmed as a result, but is acceptable otherwise."; (ii) "violating an established moral norm
378 just for fun or personal enjoyment is inherently wrong even in situations where no one is
379 harmed as a result.". Participants read both statements and were asked to select the
380 statement they "identify with the most". The order of these statements was randomised.
381 Participants who selected (ii) were then asked to elaborate on their position through an
382 open-ended response question. The purpose of these statements was to assess participants'
383 own *prior beliefs* regarding moral judgment and justifications (see Royzman et al., 2015, p.
384 331). In order to prevent the potentially confounding influence of a salient example moral
385 scenario, these statements were presented before the moral judgment task.

386 Participants were then presented with the *Incest* vignette (Appendix A) from the
387 original moral dumbfounding study (Haidt et al., 2000). They were asked to rate on a
388 seven-point Likert scale how right or wrong they would rate the behavior of Julie and Mark
389 (where, 1 = *Morally wrong*; 4 = *Neutral*; 7 = *Morally right*). They were asked to provide a
390 reason for their judgment through open-ended response, and, rated their confidence in their

391 judgment. Participants were then presented with a series of prepared counter-arguments
392 designed to refute commonly used justifications for rating the behavior as “wrong”
393 (Appendix B).

394 Dumbfounding was measured using a “critical slide” (developed by McHugh et al.,
395 2017). The critical slide is a page in an online or computer based questionnaire specifically
396 designed to measure dumbfounded responding. It contains a statement defending the
397 behavior and a question as to how the behavior could be wrong (“Julie and Mark’s behavior
398 did not harm anyone, how can there be anything wrong with what they did?”). There are
399 three possible answer options: (a) “There is nothing wrong”; (b) an admission of not having
400 reasons (“It’s wrong but I can’t think of a reason”); and finally a judgment with
401 accompanying justification (c) “It’s wrong and I can provide a valid reason”. The order of
402 these response options is randomised. Participants who select (c) are prompted on a
403 following slide to type a reason. In line with McHugh et al. (2017), the selecting of option
404 (b), the admission of not having reasons, was taken to be a dumbfounded response.

405 Following the critical slide, participants rated the behavior, and rated their confidence
406 in their judgment again. They also indicated, on a 7-point Likert scale, how much they
407 changed their mind. A post-discussion questionnaire containing self-report reaction to the
408 scenario across various dimensions (confidence, confusion, irritation, etc.) taken from Haidt
409 et al. (2000) was administered after these revised judgments had been made (Appendix C).

410 Two targeted questions were taken directly from Royzman et al. (2015) to assess
411 whether or not participants’ judgments may be grounded in the harm principle: (i) “Having
412 read the story and considering the arguments presented, are you able to believe that Julie
413 and Mark’s having sex with each other will not negatively affect the quality of their
414 relationship or how they feel about each other later on?”; (ii) “Having read the story and
415 considering the arguments presented, are you able to believe that Julie and Mark’s having
416 sex with each other will have no bad consequences for them personally and/or for those close

417 to them?”. Participants responded “Yes” or “No” to each of these statements. The order of
418 these questions was randomised.

419 Two other measures were also taken for exploratory purposes: Meaning in Life
420 questionnaire (MLQ; Steger, Kashdan, Sullivan, & Lorentz, 2008). This ten item scale is
421 made up of two five item sub scales: presence (e.g., “I understand my life’s meaning.”) and
422 search (e.g., “I am looking for something that makes my life feel meaningful.”). Responses
423 were recorded using a 7-point Likert scale ranging from 1 (*strongly disagree*) to 7 (*strongly*
424 *agree*); and CRSi7 a seven item scale taken from The Centrality of Religiosity Scale (Huber
425 & Huber, 2012). Participants responded to questions relating to the frequency with which
426 they engage in religious or spiritual activity (e.g., “How often do you think about religious
427 issues?”). Responses were recorded using a 5-point Likert scale ranging from 1 (*never*) to 5
428 (*very often*). The seven item inter-religious version of the scale was selected because some
429 non-religious activities (such as meditation) may also have a bearing on a person’s ability to
430 reason about moral issues.

431 2.2 | Results and Discussion

432 Eighty-seven of the total sample ($N = 110$; 79.09%) initially rated the behavior of Julie
433 and Mark as wrong; no difference in initial rating between the MTurk sample ($M = 1.98$, SD
434 $= 1.52$), and the MIC sample, ($M = 2.10$, $SD = 1.39$), $t(107.94) = -0.41$, $p = .683$, $d = 0.08$.
435 Eighty-six of the total sample, ($N = 110$; 78.18%) rated the behavior as wrong after viewing
436 the counter-arguments and the critical slide; no difference in revised rating between the
437 MTurk sample, ($M = 2$, $SD = 1.53$), and the MIC sample, ($M = 2.33$, $SD = 1.54$), $t(106.55)$
438 $= -1.11$, $p = .268$, $d = 0.21$. A paired samples t-test revealed a significant difference in
439 rating of behavior from time one, initial rating, ($M = 2.04$, $SD = 1.45$), to time two, revised
440 rating, ($M = 2.15$, $SD = 1.54$), $t(109) = -2.38$, $p = .019$, $d = 0.08$. This result may be due
441 to changes in the severity of the judgments as opposed to changing the judgment. Further
442 analysis revealed that only eight (7.27%) participants changed their judgment: two

443 participants changed their judgment from “wrong” to “neutral”; one participant changed
444 their judgment from “right” to “neutral”; four changed their judgment from “neutral” to
445 “right”; and one participant changed their judgment from “neutral” to “wrong”. A chi-square
446 test for independence revealed no significant association between time of judgment and
447 valence of judgment made, $\chi^2(2, N = 220) = 0.73, p = .694, V = 0.06$. This rate of
448 changing judgments is lower than the 12% reported in Haidt et al. (2000), however, as noted
449 above, social pressure appears to influence responses in the dumbfounding paradigm. It is
450 likely that the lower rates of changing judgments can be attributed to the reduced social
451 pressure in a computerized task.

452 Ten participants (9%) indicated that they had encountered the scenario before. When
453 asked to elaborate, participants provided anecdotes, or referred to previous readings (either
454 fiction or philosophy). Two participants (2%) indicated that they had encountered it in a
455 previous survey. The low numbers mean that any potential influence of previous experience
456 on the results is negligible and these participants were not excluded from the analyses.

457 **2.2.1 | Measuring dumbfounding.** Participants who selected the admission of not
458 having reasons on the critical slide were identified as dumbfounded. Rates of of each
459 response to the critical slide are for the entire sample ($N = 110$) are displayed in Figure 1.
460 Twenty participants (18.18%) were initially identified as dumbfounded.⁴ The exclusion
461 criteria developed by Royzman et al. (2015) were applied, all participants who endorsed
462 either the harm principle or the norm principle were excluded from analysis. This left a
463 sample of 14 participants who were eligible for analysis. None of these 14 selected the
464 dumbfounded response.

⁴ Unsupported declarations and tautological responses provided in the open-ended responses resulted in an additional six participants presenting as potentially dumbfounded; given that Royzman et al. (2015) argue that these responses are an articulation of a norm/principle, these participants are not identified as dumbfounded here.

465 The purpose of the Study 1 was to assess if participants could articulate the principles
466 identified by Royzman et al. (2015), independently of the targeted statements/questions, as
467 these may serve as a prompt. A revised measure of convergence is developed here. A
468 participant's endorsement of either principle should lead to their exclusion from analysis,
469 only if the participant also articulated this principle when given the opportunity. The
470 open-ended responses were analysed and coded for any mention of either the harm principle
471 or the norm principle. Participants were only excluded from analysis if they both endorsed
472 and articulated either principle. For the purposes of consistency with Royzman et al. (2015),
473 unsupported declarations and tautological responses (identified as dumbfounded responses by
474 McHugh et al., 2017) were coded as an articulation of the norm principle here.⁵ As
475 predicted, the number of participants who both articulated and endorsed either principle was
476 much lower than the number of participants who only endorsed either principle. Fifty two
477 participants were eligible for analysis according to the revised exclusion criteria. Eight of
478 these participants (15.38%) selected the dumbfounded response, providing some evidence for
479 moral dumbfounding. Figure 1 shows the responses to the critical slide for the entire sample
480 and for participants eligible for analysis according to each measure of convergence.

481 **2.2.2 | Consistency between endorsed principles and expressed judgments.**

482 The exclusion criteria developed by Royzman et al. (2015) (endorsing only), led to a large
483 proportion of participants who selected "There is nothing wrong" to be excluded from
484 analysis (12 participants; 54.55% of the 22 participants who selected this option). Both the
485 harm principle and the norm principle provide legitimate reasons for participants to judge
486 the behavior as wrong (Royzman et al., 2015). It follows that if a participant endorsed either
487 principle, they would also judge the behavior as wrong. It is surprising then that, 12 of the

⁵ By only identifying participants who explicitly admitted to not having a reason as dumbfounded we also reduced the potential risk of "false inclusions", where people provide a dumbfounded response through laziness or inattentiveness. While the motivations for selecting various responses cannot be known, previous research has identified the selecting of an admission of not having reasons as a conservative indicator of moral dumbfounding (McHugh et al., 2017, p. 16).

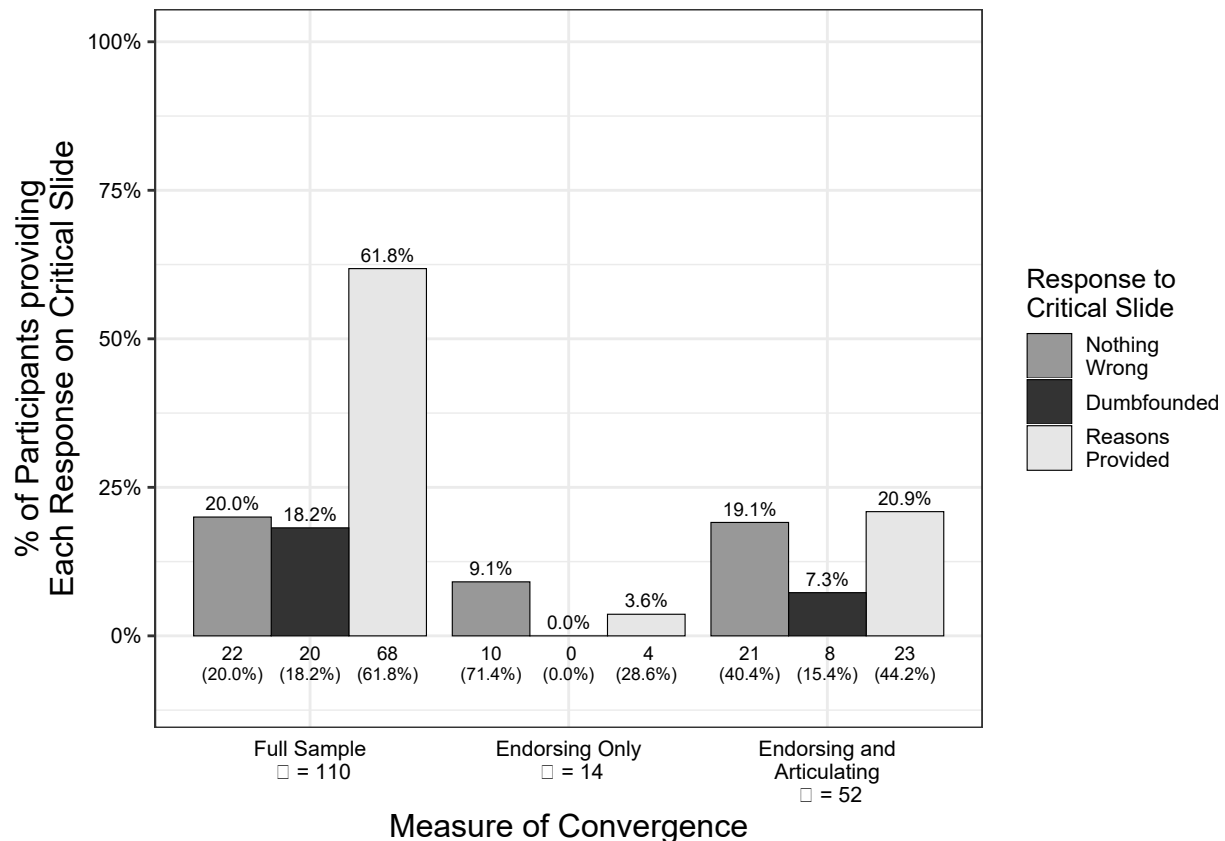


Figure 1. Responses to critical slide for the entire sample, and for each measure of convergence: (i) endorsing only, and (ii), endorsing and articulating; percentages of full sample displayed within plot, percentages of relevant sample displayed in parenthesis below the count.

488 22 participants who selected “There is nothing wrong” on the critical slide, also endorsed
 489 either the harm principle or the norm principle. The endorsing of these principles meant that
 490 these participants were excluded from analysis on the grounds they had a legitimate reason
 491 to rate the behavior as wrong. However, these participants did not rate the behavior as
 492 wrong. This demonstrates an inconsistency between the endorsing of the principles through
 493 targeted questions and statements and the apparent use of these principles as reasons
 494 guiding the participants’ judgments. The endorsing only measure of convergence, using the
 495 targeted questions and statements developed by Royzman et al. (2015) led to participants
 496 being falsely excluded from analysis.

497 According to the revised criteria for exclusion, in which participants are only excluded
498 from analysis if they were also able to articulate the principle that they endorsed, only one of
499 the 22 participants (4.55%) who selected “There is nothing wrong” was excluded from
500 analysis. The revised measure of convergence developed in Study 1 shows a reduced
501 incidence of false exclusion of participants who selected “There is nothing wrong”. This
502 suggests that accounting for both the articulating and the endorsing of principles provides
503 more accurate (though still not quite perfect) exclusion criteria.

504 The aim of Study 1 was to extend previous research by Royzman et al. (2015). They
505 excluded participants from analysis based on their endorsing of either the harm principle or
506 the norm principle through targeted questions/statements. Using these criteria for exclusion,
507 they found minimal dumbfounded responding (1 participant from a sample of 53 (Royzman
508 et al., 2015, p. 309)). It was hypothesised that their exclusion criteria were too broad, and
509 that participants’ endorsing of either principle does imply that participants can articulate
510 the given principle. Revised criteria for exclusion were developed which accounted for both
511 the endorsing and the articulation of either the harm principle or the norm principle. Our
512 initial analysis replicated the findings of Royzman et al. (2015).

513 Further analysis, using the revised measure of convergence demonstrated considerably
514 more consistency in the exclusion/inclusion of participants who selected “There is nothing
515 wrong”. These revised criteria identified eight (7.27% of the total sample of $N = 110$)
516 participants as dumbfounded. Study 1 demonstrated inconsistency in the endorsing and
517 articulation of the harm principle and the norm principle, and provided evidence for moral
518 dumbfounding, however rates of dumbfounded responding were low, with the majority of
519 participants (68; 61.82%) providing reasons for their judgments. A second study was devised
520 to assess the consistency in the application of the harm principle across differing contexts,
521 along with the endorsing, and articulation of the each principle.

3 | Study 2: Applying Moral Principles Across Contexts

In Study 1, we tested if participants could articulate the harm principle and the norm principle as identified by Royzman et al. (2015). In Study 2, we investigated the role of the harm principle in the making of judgments. Specifically, we examined if the harm principle can legitimately be said to be guiding the judgments of participants. This was done by assessing whether or not the harm principle is applied consistently across different contexts

Drawing on the research by Royzman et al. (2015), the harm principle may summarised as follows “it is wrong for two people to engage in an activity whereby harm may occur”. Royzman et al. (2015) do not offer clarification on specific types of harm that may fall under this principle, it is therefore assumed that this is a generalised principle concerning any form of harm. According to the argument proposed by Royzman et al. (2015), participants’ moral judgments are grounded in this principle, such that applying this principle to the *Incest* dilemma gives people a good reason to judge the behavior of Julie and Mark as wrong. If this general harm principle is to be considered as guiding participants’ judgments, it should be consistently applied across differing contexts.

Study 2 tested if this was the case by including a set of targeted questions relating to the generalisation and application of the harm principle across different contexts (the rest of the materials were largely the same as those used in Study 1). We hypothesised that participants’ responses to these targeted questions would reveal inconsistency in the application of the harm principle across differing contexts. Any exclusion criteria based on the harm principle should account for the endorsing of the principle (Royzman et al., 2015), articulating the principle (Study 1), and the application of the principle (Study 2).

3.1 | Method

3.1.1 | Participants and design. Study 2 was a frequency-based extension of Study 1. The aim was to investigate the prevalence of moral dumbfounding when controlling

547 for (a) the consistency with which people articulate and endorse the norm principle and the
548 harm principle, and (b) the consistency with which people apply the norm principle principle.
549 A combined sample of 111 (67 female, 44 male; $M_{\text{age}} = 34.23$, min = 19, max = 74, $SD =$
550 11.42) took part.

551 Sixty-one (36 female, 25 male; $M_{\text{age}} = 39.08$, min = 20, max = 74, $SD = 12.25$) were
552 recruited through MTurk. Participation was voluntary and participants were paid 0.50 US
553 dollars for their participation. Participants were recruited from English speaking countries or
554 from countries where residents generally have a high level of English (e.g., The Netherlands,
555 Denmark, Sweden). Fifty (31 female, 19 male; $M_{\text{age}} = 28.32$, min = 19, max = 48, $SD =$
556 6.65) were recruited through direct electronic correspondence. Participants in this sample
557 were undergraduate students, postgraduate students, and alumni from Mary Immaculate
558 College (MIC), and University of Limerick (UL). Participation was voluntary and
559 participants were not reimbursed for their participation.

560 **3.1.2 | Procedure and materials.** Data were collected using an online
561 questionnaire generated using Questback (Unipark, 2013). The questionnaire in Study 2 was
562 the same as that presented in Study 1, with the inclusion of three additional targeted
563 questions which aimed to assess the consistency with which participants generalise and apply
564 the harm principle. The questions were: (a) “How would you rate the behavior of two people
565 who engage in an activity that could potentially result in harmful consequences for either of
566 them?”; (b) “Do you think boxing is wrong?”; (c) “Do you think playing contact team sports
567 (e.g. rugby; ice-hockey; American football) is wrong?”. Responses to (a) were recorded on a
568 7-point Likert scale (where, 1 = *Morally wrong*; 4 = *Neutral*; 7 = *Morally right*). Responses
569 to (b) and (c) were recorded using a binary “Yes/No” option. These questions were
570 presented sequentially, in randomised order. The randomised sequence was grouped as Block
571 A. Similarly all slides and questions directly relating the moral scenario were grouped as
572 Block B. Block B also included the targeted questions relating to the endorsing of the harm
573 principle. The order of presentation of these blocks was randomised.

574 As with Study 1, the questionnaire opened with the information sheet, and the main
575 body of the questionnaire could not be accessed until participants consented to continue.
576 Once consent was given participants were asked a number of questions relating to basic
577 demographics. They were then presented with the two targeted statements relating to the
578 norm principle (in randomised order) and asked to select the statement they “identify with
579 the most”. Participants were then presented with either Block A (containing the targeted
580 questions relating to the application of the harm principle) or Block B (containing the moral
581 scenario, related questions, and targeted questions relating to the endorsing of the harm
582 principle). Following this participants were presented with the second block. As in Study 1,
583 the questionnaire ended with the MLQ (Steger et al., 2008); and CRSi7 (Huber & Huber,
584 2012).

585 3.2 | Results and Discussion

586 Seventy-nine of the total sample ($N = 111$; 71.17%) initially rated the behavior of Julie
587 and Mark as wrong. An independent samples t-test revealed no difference in initial rating
588 between the MTurk sample ($M = 2.08$, $SD = 1.48$), and the MIC sample, ($M = 2.68$, $SD =$
589 1.83), $t(93.31) = 1.86$, $p = .066$, $d = 0.36$. Sixty seven of the total sample, ($N = 111$;
590 60.36%) rated the behavior as wrong after viewing the counter-arguments and the critical
591 slide. An independent samples t-test revealed a significant difference in revised rating
592 between the MTurk sample, ($M = 2.31$, $SD = 1.53$), and the MIC sample, ($M = 3$, $SD =$
593 1.84), $t(95.40) = 2.11$, $p = .037$, $d = 0.41$. A paired samples t-test revealed a significant
594 difference in rating of behavior from time one, initial rating, ($M = 2.35$, $SD = 1.67$), to time
595 two, revised rating, ($M = 2.62$, $SD = 1.54$), $t(110) = -3.47$, $p < .001$, $d = 0.16$. Further
596 analysis revealed that although 15 participants changed their judgment, only two
597 participants changed fully the valence of their judgment, changing their judgment from
598 “wrong” to “right”. Of the other changes in judgment, ten participants changed their
599 judgment from “wrong” to “neutral”; two participants changed their judgment from “right”

600 to “neutral”; and one changed their judgment from “neutral” to “right”. A chi-square test for
601 independence revealed no significant association between time of judgment and valence of
602 judgment made, $\chi^2(2, N = 222) = 3.40, p = .183, V = 0.12$.

603 Eighteen participants (16%) indicated that they had encountered the scenario before.
604 As in Study 1, when asked to elaborate, participants provided anecdotes, or referred to
605 previous readings/TV (either fiction or philosophy), 8 participants (7%) indicated that they
606 had encountered it in a previous survey. The number of participants indicating previous
607 experience with the scenario was higher than in Study 1 and as such the possibility that it
608 may have confounded the results was investigated. An independent samples t-test revealed
609 no difference in judgment between participants who had previously seen the scenario, ($M =$
610 $2.83, SD = 1.86$), and participants who had not previously seen the scenario, ($M = 2.26, SD$
611 $= 1.62$), $t(22.31) = 1.23, p = .232, d = 0.35$. Furthermore, a chi-squared test for
612 independence revealed no significant association between previous experience with the
613 scenario and response to the critical slide, $\chi^2(2, N = 111) = 3.16, p = .206, V = 0.17$.
614 These participants were not excluded from the analyses.

615 **3.2.1 | Testing for order effects.** The order of the blocks had no influence on the
616 any of the responses of interest (see supplementary materials for details of analysis). Of the
617 questions relating to the application of the harm principle, there were differences in
618 responding to general question only (“How would you rate the behavior of two people who
619 engage in an activity that could potentially result in harmful consequences for either of
620 them?”). This question was more abstract than the two questions it appeared with, in which
621 participants were asked to judge a named behavior (boxing or contact team sports). The
622 description in the general question could apply to either of the named behaviors.
623 Participants who responded to this question first rated the behavior as more wrong than
624 participants who responded to it after reading one or both of the named behaviors. It seems
625 likely that the named behaviors provided an example of a situation in which the behavior
626 described in the general question may be acceptable, leading participants to respond more

627 favorably to the general question.

628 **3.2.2 | Measuring dumbfounding.** As in Study 1, participants who selected the
629 admission of not having reasons on the critical slide were identified as dumbfounded. Rates
630 of each response to the critical slide are for the entire sample ($N = 111$) are displayed in
631 Figure 1. Twenty one participants (18.92%) were initially identified as dumbfounded.⁶ The
632 exclusion criteria developed by Royzman et al. (2015; the endorsing of either principle) were
633 applied, and this left a sample of 20 who were eligible for analysis. Two of these fully
634 convergent participants selected the dumbfounded response. We then applied the revised
635 criteria for exclusion (both articulating and endorsing either principle) developed in Study 1,
636 and the number of participants eligible for analysis increased to 61. Of these, nine (14.75%)
637 selected the dumbfounded response. Again this also led to a reduction in false exclusions,
638 three 3 of the 36 (8.33) participants who selected “There is nothing wrong” were excluded by
639 this measure.

640 The responses to the three targeted questions relating the application of the harm
641 principle were analysed together. Only one participant was consistent in their application of
642 the harm principle across all three targeted questions and this meant that only one
643 participant was consistent in the application, articulation, and, endorsing of the harm
644 principle (as measured by the open-ended responses and the targeted questions taken from
645 Royzman et al. (2015)). This was combined with the exclusion criteria developed in Study 1
646 leaving a sample of 73 participants who were eligible for analysis. Ten (9.01% of the total
647 sample) of these participants selected the dumbfounded response. The responses to the
648 critical slide across all measures of convergence used are displayed in Figure 2.

649 **3.2.3 | Consistency between endorsed principles and expressed judgments.**
650 As in Study 1, the initial criteria for exclusion (endorsing only) excluded a large proportion

⁶ Unsupported declarations and tautological responses provided in the open-ended responses resulted in an additional six participants presenting as potentially dumbfounded; again, these participants are not identified as dumbfounded here.

651 of the participants who selected “There is nothing wrong”; 20 of the 36 participants (55.56%)
652 who selected “There is nothing wrong” were excluded. When articulation of the principles
653 was accounted for, only three (8.33%) of these 36 participants were excluded. This is higher
654 than in Study 1 (one participant, 4.55% of those who selected “There is nothing wrong”),
655 however in reducing the obvious false exclusion of participants who selected “There is nothing
656 wrong” it remains an improvement on the original criteria. This suggests that accounting for
657 participants’ ability to articulate the principles endorsed provides a more accurate criteria
658 for exclusion than accounting only for the endorsing of a given principle. Furthermore, when
659 the applying of the harm principle was also accounted for, only one of the 36 participants
660 who selected “There is nothing wrong” was excluded. The criteria for convergence developed
661 here lead to greater consistency between a participant’s eligibility for analysis and their
662 judgment made than the original criteria described by Royzman et al. (2015).

663 Study 2 investigated the consistency with which people apply, articulate, and endorse
664 the harm principle. Only one participant consistently applied, articulated, and endorsed the
665 harm principle. As such, the harm principle as a basis for exclusion from analysis becomes
666 practically redundant, and it seems unlikely that there is a generalised harm principle that
667 underlies moral judgments (though does not rule out the possibility of more focused, content
668 specific harm principles). The endorsing and articulation of the norm principle resulted in
669 the exclusion of 37 participants. The degree to which the articulation or the endorsing of the
670 norm principle may render participants ineligible for consideration as dumbfounded is
671 unclear, this is discussed in more detail below. However, even if participants are excluded
672 from analysis based on the norm principle, dumbfounded responding is still observed, with
673 ten participants (13.70% of sample eligible for analysis; 9.01% of the total sample) selecting
674 the admission of having no reason on the critical slide. As in Study 1, rates of observed
675 dumbfounding are low, and providing reasons appears to be the preferred response, with
676 more participants (54; 48.65%) providing reasons than selecting either of the other responses
677 to the critical slide.

4 | Study 3: Replication and Extension

678

679 Studies 1 and 2 demonstrated that people do not consistently articulate and endorse
680 the norm principle, or consistently articulate, endorse and apply the harm principle. Both
681 studies found evidence of dumbfounding, however the exclusion of participants resulted in
682 relatively small numbers of participants being eligible for analysis. As such we conducted a
683 third study, an attempt to replicate Study 2, with a larger sample.

684 4.1 | Method

685 **4.1.1 | Participants and design.** Study 3 was a frequency-based replication of
686 Study 2. The aim was to investigate the prevalence of moral dumbfounding when controlling
687 for (a) the consistency with which people articulate and endorse the norm principle and the
688 harm principle, and (b) the consistency with which people apply the norm principle principle.
689 A total sample of 502 (287 female, 212 male; $M_{\text{age}} = 39.05$, min = 18, max = 81, $SD =$
690 12.46) took part. All participants were recruited through MTurk. Participation was
691 voluntary and participants were paid 0.50 US dollars for their participation. Participants
692 were recruited from English speaking countries or from countries where residents generally
693 have a high level of English (e.g., The Netherlands, Denmark, Sweden).

694 **4.1.2 | Procedure and materials.** The materials and procedure were identical to
695 Study 2.

696 4.2 | Results and Discussion

697 Three-hundred-and-seventy-nine of the total sample ($N = 502$; 75.50%) rated the
698 behavior of Julie and Mark as wrong initially; and 357 participants, ($N = 502$; 71.12%) rated
699 the behavior as wrong after viewing the counter-arguments and the critical slide. A paired
700 samples t-test revealed a significant difference in rating of behavior from time one, initial
701 rating, ($M = 2.21$, $SD = 1.72$), to time two, revised rating, ($M = 2.38$, $SD = 1.79$), $t(501) =$
702 -4.74 , $p < .001$, $d = 0.10$. However a chi-square test for independence revealed no significant

703 association between time of judgment and valence of judgment made, $\chi^2(2, N = 1004) =$
704 $3.59, p = .166, V = 0.08$.⁷

705 **4.2.1 | Testing for order effects.** As in Study 2, the order of the blocks did
706 influence on the any of the responses of interest, and the general harm question was the only
707 question relating to the application of the harm principle that varied significantly with order
708 (see supplementary materials for details of analysis). Again, it is likely that encountering a
709 behaviour where harm may be acceptable (through the content of the other two questions),
710 led participants to respond to the general question more favourably.

711 **4.2.2 | Measuring dumbfounding.** Participants who selected the admission of not
712 having reasons on the critical slide were identified as dumbfounded. This option was selected
713 by 88 participants (17.53% of the entire sample $N = 502$).⁸

714 The exclusion criteria developed by Royzman et al. (2015; the endorsing of either
715 principle) were applied, and this left a sample of 84 who were eligible for analysis. Of these,
716 9 participants selected the dumbfounded response.

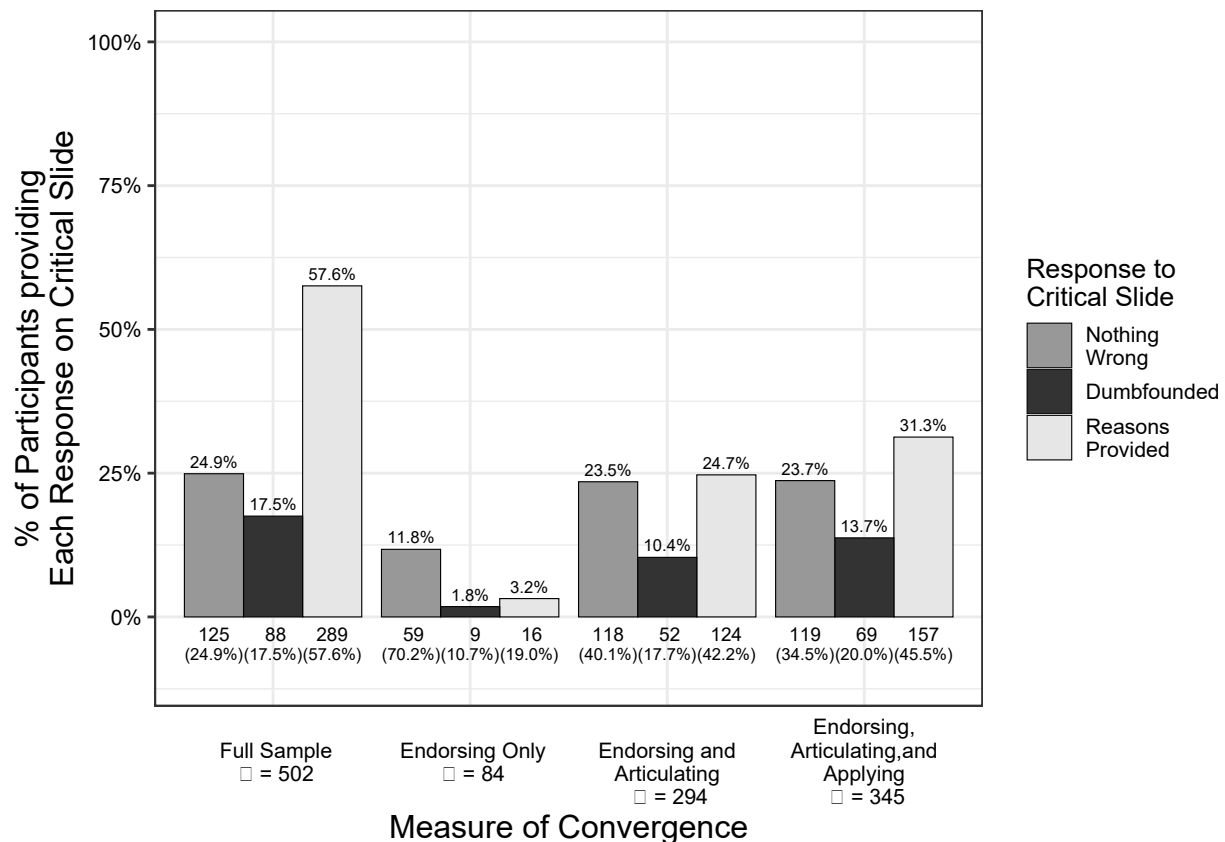
717 We then applied the exclusion criteria developed in Study 1 (both articulating and
718 endorsing either principle), and the number of participants eligible for analysis increased to
719 294. Of these, 52 (17.69%) selected the dumbfounded response.

720 Finally, the exclusion criteria developed in Study 2 were applied, leaving a sample of
721 345 participants who were eligible for analysis; Sixty nine of whom (13.75% of the total

⁷ Further analysis revealed that 42 participants changed their judgment, only seven participants changed fully the valence of their judgment, with five changing their judgment from “wrong” to “right”, and two changing their judgement from “right” to “wrong”. Of the other changes in judgment, twenty two participants changed their judgment from “wrong” to “neutral”; six participants changed their judgment from “right” to “neutral”; and four changed their judgment from “neutral” to “right”.

⁸ Unsupported declarations and tautological responses provided in the open-ended responses resulted in an additional 50 participants presenting as potentially dumbfounded; again, these participants are not identified as dumbfounded here.

722 sample) selected the dumbfounded response. The responses to the critical slide for the entire
 723 sample, and for each measure of convergence used are displayed in Figure 2.



724 *Figure 2.* Responses to critical slide for the entire sample, and for each measure of convergence:
 725 (i) endorsing only, (ii) endorsing and articulating, and (iii), endorsing, articulating, and
 726 applying; percentages of full sample displayed within plot, percentages of relevant sample
 727 displayed in parenthesis below the count.

724 4.2.3 | Consistency between endorsed principles and expressed judgments.

725 As in Studies 1 and 2, the exclusion criteria developed here resulted in fewer false exclusions.
 726 In the current study, the exclusion criteria developed by Royzman et al. (2015, endorsing
 727 only), led to 66 of the 125 participants who selected “There is nothing wrong” being
 728 excluded from analysis (52.80%). Conversely, applying the exclusion criteria developed in
 729 Study 1 resulted in seven of these 125 participants being excluded (5.60%); and the exclusion
 730 criteria from Study 2 resulted in six of these 125 participants being excluded (4.80%).

731 Further analysis, using the revised measure of convergence demonstrated considerably
732 more consistency in the exclusion/inclusion of participants who selected “There is nothing
733 wrong”. These revised criteria identified sixty-nine (20% of the total eligible sample of $N =$
734 345) participants as dumbfounded. Study 1 provided evidence for moral dumbfounding and
735 demonstrated inconsistency in the endorsing and articulation of the harm principle and the
736 norm principle, a second study was devised to assess the consistency in the application of the
737 harm principle across differing contexts, along with the endorsing, and articulation of the
738 each principle. Study 3 replicated the findings of both Studies 1 and 2 with a larger sample.
739 By applying our revised exclusion criteria, we found clear evidence for the existence of moral
740 dumbfounding, though observed rates of dumbfounding were low, with the majority of
741 participants (157; 45.51%) providing reasons.

742 The analyses of the individual difference variables are reported in the Supplementary
743 Materials (Appendix D).

744 5 | General Discussion

745 The overarching goal of Studies 1, 2, and 3 was to re-assess the occurrence of moral
746 dumbfounding. That is, we examined whether the judgments of dumbfounded participants
747 can be attributed to moral principles based on their endorsing of these principles. This was
748 done by assessing the consistency with which participants articulate and apply these moral
749 principles. Royzman et al. (2015) argue that, if participants endorse a principle, their
750 judgment can be attributed to that principle. They claimed that by attributing participants’
751 judgments to particular principles in this way, moral dumbfounding can be eliminated.
752 However, attributing judgments to reasons based on the endorsing of a related principle is
753 problematic. Stronger evidence that a participant’s judgment may be attributed to a given
754 principle should account for (a) the participant’s ability to articulate this principle,
755 independent of a prompt; or (b) the consistency with with the participant applies the
756 principle across differing contexts. Three studies were conducted to address these issues.

757 All three studies showed that participants do not consistently articulate principles that
758 they may endorse. This inconsistency between the endorsing and articulation of principles
759 that are purported to be governing moral judgments suggests that endorsing alone provides a
760 poor measure of whether these principles directly underpin a given judgment. In these cases
761 participants' judgments were not attributed to these principles, and evidence for
762 dumbfounding was found, though rates of dumbfounding were quite low. Studies 2 and 3
763 demonstrated that people do not consistently apply the harm principle across different
764 contexts. This poses a challenge to the argument that the judgments of dumbfounded
765 participants can be attributed to the harm principle (e.g., Royzman et al., 2015; see also
766 Gray et al., 2014; Jacobson, 2012). Our studies showed evidence for dumbfounding. Despite
767 the low rates of dumbfounding observed, the consistency across all three studies provides
768 some evidence that dumbfounded responding may indeed be indicative of a state of
769 dumbfoundedness, rather than being entirely attributed to features of the experimental
770 design.

771 **5.1 | The Norm Principle and Unsupported Declarations**

772 In all three studies, unsupported declarations were coded as an articulation of the
773 norm principle, and therefore not taken as dumbfounded responses. However, in previous
774 work, we identified parallels between the providing of unsupported declarations and the
775 providing of admissions of not having reasons (similar proportion of time spent (a)
776 smiling/laughing, (b) in silence; see McHugh et al., 2017). There is also a strong theoretical
777 case for the inclusion of unsupported declarations as dumbfounded responses. Propositional
778 beliefs/deontological judgments may be viewed as habitual/model-free intuitions (e.g.,
779 Crockett, 2013; Cushman, 2013b). The reasons for these judgments are independent of the
780 intuition. Stating the content of the intuition, is not the same as providing a reason for the
781 intuition. Royzman et al. (2015) argue that endorsing the propositional belief is sufficient
782 evidence of that belief playing an influential role in relevant judgments, however, this is

783 holding participants to a different standard. There is a difference between having a reason
784 for an intuition/propositional belief and claiming the direct basis for a judgment is an
785 associated propositional belief. In view of this, it is possible that by not including
786 unsupported declarations or tautological reasons as dumbfounded responses, the rates of
787 dumbfounding reported here are not representative of the phenomenon, providing instead an
788 overly conservative estimate. However, even according to this stricter measure adopted here,
789 evidence for dumbfounding was found.

790 **5.2 | Consistency Between Endorsed Principles and Expressed Judgments**

791 The most convincing evidence that the exclusion criteria developed in these studies are
792 more accurate than the criteria proposed by Royzman et al. (2015) is the greater consistency
793 between valence of judgment and eligibility for analysis. Participants' eligibility for analysis
794 is determined by whether or not their judgment can be attributed to either the harm
795 principle or the norm principle. If a participant's judgment can be attributed to a given
796 principle, this participant is deemed to have a reason for their judgment and they cannot be
797 identified as dumbfounded (rendering them ineligible for analysis). In order for a judgment
798 to legitimately be attributed to a particular principle, it is necessary that the valence of the
799 judgment is consistent with what is predicted by the application of that principle. In the
800 case of both principles, applying either the harm principle or the norm principle (as
801 described by Royzman et al., 2015) results in the behavior being judged as wrong. This
802 means that the judgments of participants who selected "There is nothing wrong" cannot be
803 attributed to either principle. Any participants who are excluded from analysis but selected
804 "There is nothing wrong", are clearly identifiable as being falsely excluded from analysis such
805 that this may be used as a measure of the relative accuracy of the different exclusion criteria
806 employed.

807 According to Royzman et al. (2015), a participant's judgment can be attributed to a
808 given principle if they endorse this principle. However, in each of the studies reported here,

809 excluding participants based on the endorsing of a principle resulted in over half of the
810 participants who selected “There is nothing wrong” to be falsely excluded from analysis;
811 participants’ judgments were incorrectly attributed to either the harm principle or the norm
812 principle (12 of the 22 participants who selected “There is nothing wrong” in Study 1 were
813 falsely excluded 54.55%; 20 of the 36 participants who selected “There is nothing wrong” in
814 Study 2 were falsely excluded 55.56%; and 66 of the 125 participants who selected “There is
815 nothing wrong” in Study 3 were falsely excluded 52.80%). This suggests that the endorsing
816 of a principle is a flawed indicator of the degree to which the principle is guiding participants’
817 judgments.

818 We made two changes to the exclusion criteria that aimed to reduce the numbers of
819 participants being falsely excluded from analysis. We hypothesised that providing
820 participants with an opportunity to articulate the reasons for their judgment would more
821 accurately identify the principles that guided participants’ judgments than their endorsing of
822 particular principles. This was found to be the case; in Study 1, only one of the 22
823 participants who selected “There is nothing wrong” was falsely excluded from analysis; in
824 Study 2 only three of the 36 participants who selected “There is nothing wrong” were falsely
825 excluded from analysis; and in Study 3 seven of the 125 participants who selected “There is
826 nothing wrong” were falsely excluded from analysis. Taking participants’ articulating of the
827 reasons for their judgments into account reduced measurable rate of false exclusion from
828 54.55% to 4.55% in Study 1; 55.56% to 8.33% in Study 2; and 52.80% to 5.60% in Study 3.
829 Furthermore, in Studies 2 and 3, with specific reference to the harm principle, we
830 hypothesised that assessing the degree to which people’s judgments could be attributed to
831 the harm principle would be related to whether or not they apply the harm principle across
832 different contexts. Again this was found to be the case, as evidenced by a further reduction
833 in the measurable rate of false exclusion from 8.33% (3/36) to 2.78% (1/36) in Study 2, and
834 from 5.60% (7/125) to 4.80% (6/125) in Study 3.

835 **5.3 | Implications**

836 The existence of moral dumbfounding and the associated support for intuitionist
837 theories of moral judgment (e.g. Cushman et al., 2010; Haidt, 2001; Hauser, Young, &
838 Cushman, 2008; Prinz, 2005; see also Crockett, 2013; Cushman, 2013b; Greene, 2008, 2013)
839 has been questioned in recent years. The majority of these challenges are theoretical (e.g.,
840 Jacobson, 2012; Sneddon, 2007; Wielenberg, 2014). The work of Gray et al. (2014),
841 appeared to give some empirical weight to these challenges, while Royzman et al. (2015)
842 extended these challenges to the dumbfounding paradigm specifically. We conducted three
843 studies addressing specific methodological limitations associated with the work by Royzman
844 et al. (2015). Their criteria for exclusion were found to be overly liberal, as evidenced by the
845 high rates of false exclusion of participants who selected “There is nothing wrong”. and
846 evidence for dumbfounding was found. Adopting the more rigorous exclusion criteria
847 developed here led to a reduction in the false exclusion of participants. In using these
848 criteria, evidence for dumbfounding was found, and the explanation of dumbfounded
849 responding proposed by Royzman et al. (2015) was not supported.

850 Our findings provide further evidence that the distinction between implicit and explicit
851 cognition (e.g., Bonner & Newell, 2010; Evans, 2003, 2006, 2008; Evans & Over, 2013; Reber,
852 1989) extends to the moral domain. It has long been known that people have poor
853 introspective awareness of how judgments are made (e.g., Nisbett & Wilson, 1977) and it
854 appears that in some cases this may also be true for moral judgments.

855 **5.4 | Limitations and Future Directions**

856 The research we present here consists of three studies with a combined sample of $N =$
857 723, from MTurk ($N = 621$) and third level institutions ($N = 102$). Follow-up studies should
858 investigate the phenomenon with larger and more diverse samples. Such follow-up work may
859 inform investigations into the influence of cultural and societal norms on the prevalence of
860 moral dumbfounding. Previous work by Haidt and Hersh (2001) provides suggestive evidence

861 that political orientation may influence a person's susceptibility to moral dumbfounding;
862 furthermore, there is some evidence to indicate that cultural and socio-economic factors may
863 also play a role (Haidt et al., 1993). Future research should draw on the methods developed
864 here and by both McHugh et al. (2017) and Royzman et al. (2015) to investigate these
865 influences further.

866 The procedures we used were very similar across both studies. They were also very
867 similar to those used by McHugh et al. (2017) and by Royzman et al. (2015). A more
868 rigorous test of moral dumbfounding should employ a variety of methods. We recommend
869 that future research develops a broader selection of "dumbfounding scenarios", and
870 investigate the feasibility of alternative procedures that may elicit dumbfounding.

871 The role of social pressure and conversational norms in the emergence of moral
872 dumbfounding is not well understood. The studies described here were conducted using
873 online surveys and therefore there was no immediate social pressure on participants to either
874 appear consistent or to conform to conversational norms. Furthermore, the argument
875 proposed by Royzman et al. (2015), that participants' judgment are grounded in reasons
876 (harm-based/norm-based) and that they drop these reasons in response to social pressure is
877 not supported by the evidence presented here; harm-based/norm based reasons were not
878 consistently articulated or applied by participants in these studies. It is apparent then that
879 dumbfounded responding cannot be attributed to social pressure alone. The processes by
880 which we make moral judgments also give rise to moral dumbfounding. This means that
881 isolating the underlying mechanisms that give rise to moral dumbfounding may contribute to
882 our overall understanding of the making of moral judgments.

883

6 | Conclusion

884 Based on three studies we conclude: moral dumbfounding seems to be real, if not as
885 widespread as initial reports might suggest (Haidt, 2001; Haidt et al., 2000; Haidt & Hersh,

2001). By reconsidering approaches of earlier research, our procedures found clear evidence for this phenomenon. People are not always able to justify their moral judgments. Indeed, in our studies, between 13% and 18% of people showed dumbfounding. Gaining insights into the occurrence and underlying processes equips society with the tools to confront and reduce dumbfounding. Further research in the area may inform improvements in the conduct of public debate, particularly in relation to polarizing issues. Perhaps in the future, the influence dumbfounding in public discourse and public policy (e.g., MacNab, 2016; Sim, 2016) will be reduced or even eliminated.

7 | Data Accessibility Statement

All participant data, and analysis scripts can be found on this paper's project page on the Open Science Framework at <https://osf.io/m4ce7/>.

All statistical analysis was conducted using R (Version 3.6.0; R Core Team, 2017) and the R-packages *afex* (Version 0.23.0; Singmann, Bolker, & Westfall, 2015), *boot* (Version 1.3.23; Davison & Hinkley, 1997), *Cairo* (Version 1.5.10; Urbanek & Horner, 2019), *car* (Version 3.0.3; Fox & Weisberg, 2011; Fox, Weisberg, & Price, 2018), *carData* (Version 3.0.2; Fox et al., 2018), *citr* (Version 0.3.2; Aust, 2016), *DescTools* (Version 0.99.28; et mult. al., 2019), *desnum* (Version 0.1.1; McHugh, 2017), *devtools* (Version 2.1.0; Wickham & Chang, 2017), *emmeans* (Version 1.4; Lenth, 2019), *extrafont* (Version 0.17; Chang, 2014), *foreign* (Version 0.8.72; R Core Team, 2018), *Formula* (Version 1.2.3; Zeileis & Croissant, 2010), *ggplot2* (Version 3.2.1; Wickham, 2009), *koRpus* (Version 0.11.5; Michalke, 2018a, 2019), *koRpus.lang.en* (Version 0.1.3; Michalke, 2019), *lme4* (Version 1.1.21; Bates, Mächler, Bolker, & Walker, 2015), *lmtest* (Version 0.9.37; Zeileis & Hothorn, 2002), *lsmeans* (Lenth, 2016), *lsr* (Version 0.5; Navarro, 2015), *MASS* (Version 7.3.51.4; Venables & Ripley, 2002a), *Matrix* (Version 1.2.17; Bates & Maechler, 2017), *metap* (Version 1.1; Dewey, 2017), *mlogit* (Version 1.0.1; Croissant, 2013), *nnet* (Version 7.3.12; Venables & Ripley, 2002b), *papaja* (Version 0.1.0.9842; Aust & Barth, 2018), *plyr* (Version 1.8.4; Wickham, 2011), *powerMediation*

912 (Version 0.2.9; Qiu, 2018), *pwr* (Version 1.2.2; Champely, 2018), *QuantPsyc* (Version 1.5;
913 Fletcher, 2012), *reshape2* (Version 1.4.3; Wickham, 2007), *scales* (Version 1.0.0; Wickham,
914 2016), *sjstats* (Version 0.17.5; Lüdecke, 2018), *syllly* (Version 0.1.5; Michalke, 2018b), *tibble*
915 (Version 2.1.3; Müller & Wickham, 2017), *usethis* (Version 1.5.1; Wickham & Bryan, 2019),
916 *VGAM* (Version 1.1.1; Yee & Wild, 1996; Yee, 2010, 2013; Yee & Hadi, 2014; Yee, Stoklosa,
917 & Huggins, 2015), *wordcountaddin* (Version 0.3.0.9000; Marwick, 2019), and *zoo* (Version
918 1.8.6; Zeileis & Grothendieck, 2005).

References

- 919
- 920 Aust, F. (2016). *Citr: 'RStudio' Add-in to Insert Markdown Citations*. Retrieved from
921 <https://CRAN.R-project.org/package=citr>
- 922 Aust, F., & Barth, M. (2018). *Papaja: Create APA manuscripts with R Markdown*.
923 Retrieved from <https://github.com/crsh/papaja>
- 924 Bates, D., & Maechler, M. (2017). *Matrix: Sparse and Dense Matrix Classes and Methods*.
925 Retrieved from <https://CRAN.R-project.org/package=Matrix>
- 926 Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting Linear Mixed-Effects Models
927 Using lme4. *Journal of Statistical Software*, *67*(1), 1–48. doi:10.18637/jss.v067.i01
- 928 Bauman, C. W., McGraw, A. P., Bartels, D. M., & Warren, C. (2014). Revisiting External
929 Validity: Concerns about Trolley Problems and Other Sacrificial Dilemmas in Moral
930 Psychology. *Social and Personality Psychology Compass*, *8*(9), 536–554.
931 doi:10.1111/spc3.12131
- 932 Bonner, C., & Newell, B. R. (2010). In conflict with ourselves? An investigation of heuristic
933 and analytic processes in decision making. *Memory & Cognition*, *38*(2), 186–196.
934 doi:10.3758/MC.38.2.186
- 935 Bostyn, D. H., Sevenhant, S., & Roets, A. (2018). Of Mice, Men, and Trolleys: Hypothetical
936 Judgment Versus Real-Life Behavior in Trolley-Style Moral Dilemmas. *Psychological
937 Science*, *29*(7), 1084–1093. doi:10.1177/0956797617752640
- 938 Brand, C. (2016). *Dual-Process Theories in Moral Psychology: Interdisciplinary Approaches
939 to Theoretical, Empirical and Practical Considerations*. Springer.
- 940 Cameron, C. D., Payne, B. K., & Doris, J. M. (2013). Morality in high definition: Emotion
941 differentiation calibrates the influence of incidental disgust on moral judgments.

- 942 *Journal of Experimental Social Psychology*, 49(4), 719–725.
943 doi:10.1016/j.jesp.2013.02.014
- 944 Champely, S. (2018). *Pwr: Basic Functions for Power Analysis*. Retrieved from
945 <https://CRAN.R-project.org/package=pwr>
- 946 Chang, W. (2014). *Extrafont: Tools for using fonts*. Retrieved from
947 <https://CRAN.R-project.org/package=extrafont>
- 948 Christensen, J. F., Flexas, A., Calabrese, M., Gut, N. K., & Gomila, A. (2014). Moral
949 judgment reloaded: A moral dilemma validation study. *Emotion Science*, 5, 607.
950 doi:10.3389/fpsyg.2014.00607
- 951 Christensen, J. F., & Gomila, A. (2012). Moral dilemmas in cognitive neuroscience of moral
952 decision-making: A principled review. *Neuroscience & Biobehavioral Reviews*, 36(4),
953 1249–1264. doi:10.1016/j.neubiorev.2012.02.008
- 954 Crockett, M. J. (2013). Models of morality. *Trends in Cognitive Sciences*, 17(8), 363–366.
955 doi:10.1016/j.tics.2013.06.005
- 956 Croissant, Y. (2013). *Mlogit: Multinomial logit model*. Retrieved from
957 <https://CRAN.R-project.org/package=mlogit>
- 958 Cushman, F. A. (2013a). Action, Outcome, and Value A Dual-System Framework for
959 Morality. *Personality and Social Psychology Review*, 17(3), 273–292.
960 doi:10.1177/1088868313495594
- 961 Cushman, F. A. (2013b). The role of learning in punishment, prosociality, and human
962 uniqueness. In K. Sterelny, B. Calcott, & B. Fraser (Eds.), *Signaling, Commitment*
963 *and Emotion, Vol. 2: Psychological and Environmental Foundations of Cooperation*.
964 MIT Press.

- 965 Cushman, F. A., Young, L., & Greene, J. D. (2010). Multi-system Moral Psychology. In J.
966 M. Doris (Ed.), *The Moral Psychology Handbook* (pp. 47–71). Oxford; New York:
967 Oxford University Press.
- 968 Davison, A. C., & Hinkley, D. V. (1997). *Bootstrap Methods and Their Applications*.
969 Cambridge: Cambridge University Press. Retrieved from
970 <http://statwww.epfl.ch/davison/BMA/>
- 971 Dewey, M. (2017). *Metap: Meta-analysis of significance values*.
- 972 Dickinson, D. L., & Maslet, D. (2018). *Using Ethical Dilemmas to Predict Antisocial*
973 *Choices with Real Payoff Consequences: An Experimental Study* (SSRN Scholarly
974 Paper No. ID 3205879). Rochester, NY: Social Science Research Network. Retrieved
975 from <https://papers.ssrn.com/abstract=3205879>
- 976 et mult. al., A. S. (2019). *DescTools: Tools for Descriptive Statistics*. Retrieved from
977 <https://cran.r-project.org/package=DescTools>
- 978 Evans, J. S. B. T. (2003). In two minds: Dual-process accounts of reasoning. *Trends in*
979 *Cognitive Sciences*, 7(10), 454–459. doi:10.1016/j.tics.2003.08.012
- 980 Evans, J. S. B. T. (2006). The heuristic-analytic theory of reasoning: Extension and
981 evaluation. *Psychonomic Bulletin & Review*, 13(3), 378–395.
982 doi:10.3758/BF03193858
- 983 Evans, J. S. B. T. (2008). Dual-Processing Accounts of Reasoning, Judgment, and Social
984 Cognition. *Annual Review of Psychology*, 59(1), 255–278.
985 doi:10.1146/annurev.psych.59.103006.093629
- 986 Evans, J. S. B. T., & Over, D. E. (2013). *Rationality and Reasoning*. Psychology Press.
- 987 Fine, C. (2006). Is the emotional dog wagging its rational tail, or chasing it? *Philosophical*

- 988 *Explorations*, 9(1), 83–98. doi:10.1080/13869790500492680
- 989 Flanagan, O., Sarkissian, H., & Wong, D. (2008). Naturalizing Ethics. In W.
990 Sinnott-Armstrong (Ed.), *Moral Psychology Volume 1: The evolution of morality*
991 *adaptations and innateness* (pp. 1–26). Cambridge, Mass.; London, England: The
992 MIT press.
- 993 Fletcher, T. D. (2012). *QuantPsyc: Quantitative Psychology Tools*. Retrieved from
994 <https://CRAN.R-project.org/package=QuantPsyc>
- 995 Fox, J., & Weisberg, S. (2011). *An R Companion to Applied Regression* (Second.). Thousand
996 Oaks CA: Sage. Retrieved from
997 <http://socserv.socsci.mcmaster.ca/jfox/Books/Companion>
- 998 Fox, J., Weisberg, S., & Price, B. (2018). *carData: Companion to Applied Regression Data*
999 *Sets*. Retrieved from <https://CRAN.R-project.org/package=carData>
- 1000 Gray, K. J., Schein, C., & Ward, A. F. (2014). The myth of harmless wrongs in moral
1001 cognition: Automatic dyadic completion from sin to suffering. *Journal of*
1002 *Experimental Psychology: General*, 143(4), 1600–1615. doi:10.1037/a0036149
- 1003 Greene, J. D. (2008). The Secret Joke of Kant’s Soul. In W. Sinnott-Armstrong, *Moral*
1004 *Psychology Volume 3: The neurosciences of morality: Emotion, brain disorders, and*
1005 *development* (pp. 35–79). Cambridge (Mass.): the MIT press.
- 1006 Greene, J. D. (2013). *Moral tribes: Emotion, reason, and the gap between us and them*.
- 1007 Greene, J. D., Sommerville, R. B., Nystrom, L. E., Darley, J. M., & Cohen, J. D. (2001). An
1008 fMRI investigation of emotional engagement in moral judgment. *Science (New York,*
1009 *N.Y.)*, 293(5537), 2105–2108. doi:10.1126/science.1062872
- 1010 Guglielmo, S. (2018). Unfounded dumbfounding: How harm and purity undermine evidence

- 1011 for moral dumbfounding. *Cognition*, 170, 334–337.
1012 doi:10.1016/j.cognition.2017.08.002
- 1013 Haidt, J. (2001). The emotional dog and its rational tail: A social intuitionist approach to
1014 moral judgment. *Psychological Review*, 108(4), 814–834.
1015 doi:10.1037/0033-295X.108.4.814
- 1016 Haidt, J., & Björklund, F. (2008). Social Intuitionists Answer Six Questions about Moral
1017 Psychology. In W. Sinnott-Armstrong (Ed.), *Moral psychology Volume 2, The*
1018 *cognitive science of morality: Intuition and diversity* (pp. 181–217). London: MIT.
- 1019 Haidt, J., Björklund, F., & Murphy, S. (2000). Moral dumbfounding: When intuition finds
1020 no reason. *Unpublished Manuscript, University of Virginia*.
- 1021 Haidt, J., & Hersh, M. A. (2001). Sexual Morality: The Cultures and Emotions of
1022 Conservatives and Liberals. *Journal of Applied Social Psychology*, 31(1), 191–221.
1023 doi:10.1111/j.1559-1816.2001.tb02489.x
- 1024 Haidt, J., Koller, S. H., & Dias, M. G. (1993). Affect, culture, and morality, or is it wrong to
1025 eat your dog? *Journal of Personality and Social Psychology*, 65(4), 613–628.
1026 doi:10.1037/0022-3514.65.4.613
- 1027 Hauser, M. D., Young, L., & Cushman, F. A. (2008). Reviving Rawls's Linguistic Analogy:
1028 Operative Principles and the Causal Structure of Moral Actions. In W.
1029 Sinnott-Armstrong (Ed.), *Moral psychology Volume 2, The cognitive science of*
1030 *morality: Intuition and diversity* (pp. 107–155). London: MIT.
- 1031 Huber, S., & Huber, O. W. (2012). The Centrality of Religiosity Scale (CRS). *Religions*,
1032 3(3), 710–724. doi:10.3390/rel3030710
- 1033 Jacobson, D. (2012). Moral Dumbfounding and Moral Stupefaction. In *Oxford studies in*

- 1034 *normative ethics* (Vol. 2, p. 289).
- 1035 Johnson-Laird, P. N. (2006). *How we reason*. Oxford ; New York: Oxford University Press.
- 1036 Kennett, J., & Fine, C. (2009). Will the Real Moral Judgment Please Stand Up? *Ethical*
1037 *Theory and Moral Practice*, 12(1), 77–96. doi:10.1007/s10677-008-9136-4
- 1038 Kohlberg, L. (1969). *Stages in the development of moral thought and action*. New York:
1039 Holt, Rinehart & Winston.
- 1040 Kohlberg, L. (1971). *From is to Ought: How to Commit the Naturalistic Fallacy and Get*
1041 *Away with it in the Study of Moral Development*.
- 1042 Lenth, R. (2019). *Emmeans: Estimated Marginal Means, aka Least-Squares Means*.
1043 Retrieved from <https://CRAN.R-project.org/package=emmeans>
- 1044 Lenth, R. V. (2016). Least-Squares Means: The R Package lsmeans. *Journal of Statistical*
1045 *Software*, 69(1), 1–33. doi:10.18637/jss.v069.i01
- 1046 Lüdtke, D. (2018). *Sjstats: Statistical Functions for Regression Models*. Retrieved from
1047 <https://CRAN.R-project.org/package=sjstats>
- 1048 MacNab, S. (2016). MSPs to consider “abhorrent” call to legalise incest. *The Scotsman*.
1049 Retrieved from [http://www.scotsman.com/news/politics/](http://www.scotsman.com/news/politics/msps-to-consider-abhorrent-call-to-legalise-incest-1-4009185)
1050 [msps-to-consider-abhorrent-call-to-legalise-incest-1-4009185](http://www.scotsman.com/news/politics/msps-to-consider-abhorrent-call-to-legalise-incest-1-4009185)
- 1051 Marwick, B. (2019). *Wordcountaddin: Word counts and readability statistics in R markdown*
1052 *documents*.
- 1053 McHugh, C. (2017). *Desnum: Creates some useful functions*. Retrieved from
1054 https://github.com/cillianmiltown/R_desnum
- 1055 McHugh, C., McGann, M., Igou, E. R., & Kinsella, E. L. (2017). Searching for Moral

- 1056 Dumbfounding: Identifying Measurable Indicators of Moral Dumbfounding. *Collabra:*
1057 *Psychology*, 3(1). doi:10.1525/collabra.79
- 1058 Mercier, H. (2016). The Argumentative Theory: Predictions and Empirical Evidence.
1059 *Trends in Cognitive Sciences*, 20(9), 689–700. doi:10.1016/j.tics.2016.07.001
- 1060 Mercier, H., & Sperber, D. (2011). Why do humans reason? Arguments for an
1061 argumentative theory. *Behavioral and Brain Sciences*, 34(2), 57–74.
1062 doi:10.1017/S0140525X10000968
- 1063 Mercier, H., & Sperber, D. (2017). *The Enigma of Reason*. Harvard University Press.
- 1064 Michalke, M. (2018a). *koRpus: An R Package for Text Analysis*. Retrieved from
1065 <https://reaktanz.de/?c=hacking&s=koRpus>
- 1066 Michalke, M. (2018b). *Sylly: Hyphenation and Syllable Counting for Text Analysis*.
1067 Retrieved from <https://reaktanz.de/?c=hacking&s=sylly>
- 1068 Michalke, M. (2019). *koRpus.Lang.En: Language Support for 'koRpus' Package: English*.
1069 Retrieved from <https://reaktanz.de/?c=hacking&s=koRpus>
- 1070 Mustonen, A.-M., Paakkonen, T., Ryökäs, E., & Nieminen, P. (2017). Abortion debates in
1071 Finland and the Republic of Ireland: Textual analysis of experiential thinking and
1072 argumentation in parliamentary and layperson discussions. *Reproductive Health*,
1073 14(1), 163. doi:10.1186/s12978-017-0418-y
- 1074 Müller, K., & Wickham, H. (2017). *Tibble: Simple Data Frames*. Retrieved from
1075 <https://CRAN.R-project.org/package=tibble>
- 1076 Narvaez, D. (2005). The neo-Kohlbergian tradition and beyond: Schemas, expertise, and
1077 character. In G. Carlo & C. Pope-Edwards (Eds.), *Nebraska symposium on*
1078 *motivation* (Vol. 51, p. 119).

- 1079 Navarro, D. (2015). *Learning statistics with R: A tutorial for psychology students and other*
1080 *beginners. (Version 0.5)*. Adelaide, Australia: University of Adelaide. Retrieved from
1081 <http://ua.edu.au/ccs/teaching/lsr>
- 1082 Nisbett, R. E., & Wilson, T. D. (1977). Telling more than we can know: Verbal reports on
1083 mental processes. *Psychological Review*, *84*(3), 231. Retrieved from
1084 <http://psycnet.apa.org/journals/rev/84/3/231/>
- 1085 Plunkett, D., & Greene, J. D. (2019). Overlooked Evidence and a Misunderstanding of What
1086 Trolley Dilemmas Do Best: Commentary on Bostyn, Sevenhant, and Roets (2018).
1087 *Psychological Science*, 0956797619827914. doi:10.1177/0956797619827914
- 1088 Prinz, J. J. (2005). Passionate Thoughts: The Emotional Embodiment of Moral Concepts.
1089 In D. Pecher & R. A. Zwaan (Eds.), *Grounding Cognition: The Role of Perception*
1090 *and Action in Memory, Language, and Thinking* (pp. 93–114). Cambridge University
1091 Press.
- 1092 Qiu, W. (2018). *powerMediation: Power/Sample Size Calculation for Mediation Analysis*.
1093 Retrieved from <https://CRAN.R-project.org/package=powerMediation>
- 1094 R Core Team. (2017). *R: A Language and Environment for Statistical Computing*. Vienna,
1095 Austria: R Foundation for Statistical Computing. Retrieved from
1096 <https://www.R-project.org/>
- 1097 R Core Team. (2018). *Foreign: Read Data Stored by 'Minitab', 'S', 'SAS', 'SPSS', 'Stata',*
1098 *'Systat', 'Weka', 'dBase', ...* Retrieved from
1099 <https://CRAN.R-project.org/package=foreign>
- 1100 Reber, A. S. (1989). Implicit learning and tacit knowledge. *Journal of Experimental*
1101 *Psychology: General*, *118*(3), 219–235. doi:10.1037/0096-3445.118.3.219

- 1102 Royzman, E. B., Kim, K., & Leeman, R. F. (2015). The curious tale of Julie and Mark:
1103 Unraveling the moral dumbfounding effect. *Judgment and Decision Making*, *10*(4),
1104 296–313.
- 1105 Rozin, P., Haidt, J., MacCauley, C., McKay, D., & Olatunji, B. O. (2008). Disgust: The
1106 body and soul emotion in the 21st century. In *Disgust and its disorders* (pp. 9–29).
1107 American Psychological Association.
- 1108 Rozin, P., Lowery, L., Imada, S., & Haidt, J. (1999). The CAD triad hypothesis: A mapping
1109 between three moral emotions (contempt, anger, disgust) and three moral codes
1110 (community, autonomy, divinity). *Journal of Personality and Social Psychology*,
1111 *76*(4), 574–586. doi:10.1037/0022-3514.76.4.574
- 1112 Sim, P. (2016, January 26). MSPs throw out incest petition. *BBC News: Scotland Politics*.
1113 Retrieved from <http://www.bbc.com/news/uk-scotland-scotland-politics-35401195>
- 1114 Singmann, H., Bolker, B., & Westfall, J. (2015). *Afex: Analysis of Factorial Experiments*.
1115 Retrieved from <https://CRAN.R-project.org/package=afex>
- 1116 Sneddon, A. (2007). A Social Model of Moral Dumbfounding: Implications for Studying
1117 Moral Reasoning and Moral Judgment. *Philosophical Psychology*, *20*(6), 731–748.
1118 doi:10.1080/09515080701694110
- 1119 Steger, M. F., Kashdan, T. B., Sullivan, B. A., & Lorentz, D. (2008). Understanding the
1120 Search for Meaning in Life: Personality, Cognitive Style, and the Dynamic Between
1121 Seeking and Experiencing Meaning. *Journal of Personality*, *76*(2), 199–228.
1122 doi:10.1111/j.1467-6494.2007.00484.x
- 1123 Stepniak, D. (1995). Televising Court Proceedings Forum: Televising Court Proceedings.
1124 *University of New South Wales Law Journal*, (2), 488–492. Retrieved from
1125 <https://heinonline.org/HOL/P?h=hein.journals/swales18&i=501>

- 1126 Todd, P. M., & Gigerenzer, G. (Eds.). (2012). *Ecological rationality: Intelligence in the*
1127 *world*. Oxford ; New York: Oxford University Press.
- 1128 Topolski, R., Weaver, J. N., Martin, Z., & McCoy, J. (2013). Choosing between the
1129 emotional dog and the rational pal: A moral dilemma with a tail. *Anthrozoös*, *26*(2),
1130 253–263. doi:10.2752/175303713X13636846944321
- 1131 Triskiel, J. (2016). Psychology Instead of Ethics? Why Psychological Research Is Important
1132 but Cannot Replace Ethics. In C. Brand (Ed.), *Dual-Process Theories in Moral*
1133 *Psychology: Interdisciplinary Approaches to Theoretical, Empirical and Practical*
1134 *Considerations* (pp. 77–98). Springer.
- 1135 Unipark, Q. (2013). *QuestBack Unipark.(2013)*.
- 1136 Urbanek, S., & Horner, J. (2019). *Cairo: R graphics device using cairo graphics library for*
1137 *creating high-quality bitmap (PNG, JPEG, TIFF), vector (PDF, SVG, PostScript)*
1138 *and display (x11 and win32) output*. Retrieved from
1139 <https://CRAN.R-project.org/package=Cairo>
- 1140 Venables, W. N., & Ripley, B. D. (2002a). *Modern Applied Statistics with S* (Fourth.). New
1141 York: Springer. Retrieved from <http://www.stats.ox.ac.uk/pub/MASS4>
- 1142 Venables, W. N., & Ripley, B. D. (2002b). *Modern Applied Statistics with S* (Fourth.). New
1143 York: Springer. Retrieved from <http://www.stats.ox.ac.uk/pub/MASS4>
- 1144 Wickham, H. (2007). Reshaping Data with the reshape Package. *Journal of Statistical*
1145 *Software*, *21*(12), 1–20. Retrieved from <http://www.jstatsoft.org/v21/i12/>
- 1146 Wickham, H. (2009). *Ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New
1147 York. Retrieved from <http://ggplot2.org>
- 1148 Wickham, H. (2011). The Split-Apply-Combine Strategy for Data Analysis. *Journal of*

- 1149 *Statistical Software*, 40(1), 1–29. Retrieved from <http://www.jstatsoft.org/v40/i01/>
- 1150 Wickham, H. (2016). *Scales: Scale Functions for Visualization*. Retrieved from
1151 <https://CRAN.R-project.org/package=scales>
- 1152 Wickham, H., & Bryan, J. (2019). *Usethis: Automate Package and Project Setup*. Retrieved
1153 from <https://CRAN.R-project.org/package=usethis>
- 1154 Wickham, H., & Chang, W. (2017). *Devtools: Tools to Make Developing R Packages Easier*.
1155 Retrieved from <https://CRAN.R-project.org/package=devtools>
- 1156 Wielenberg, E. J. (2014). *Robust Ethics: The Metaphysics and Epistemology of Godless*
1157 *Normative Realism*. OUP Oxford.
- 1158 Yee, T. W. (2010). The VGAM Package for Categorical Data Analysis. *Journal of Statistical*
1159 *Software*, 32(10), 1–34. Retrieved from <http://www.jstatsoft.org/v32/i10/>
- 1160 Yee, T. W. (2013). Two-parameter reduced-rank vector generalized linear models.
1161 *Computational Statistics and Data Analysis*. Retrieved from
1162 <http://ees.elsevier.com/csda>
- 1163 Yee, T. W., & Hadi, A. F. (2014). Row-column interaction models, with an R
1164 implementation. *Computational Statistics*, 29(6), 1427–1445.
- 1165 Yee, T. W., Stoklosa, J., & Huggins, R. M. (2015). The VGAM Package for
1166 Capture-Recapture Data Using the Conditional Likelihood. *Journal of Statistical*
1167 *Software*, 65(5), 1–33. Retrieved from <http://www.jstatsoft.org/v65/i05/>
- 1168 Yee, T. W., & Wild, C. J. (1996). Vector Generalized Additive Models. *Journal of Royal*
1169 *Statistical Society, Series B*, 58(3), 481–493.
- 1170 Zeileis, A., & Croissant, Y. (2010). Extended Model Formulas in R: Multiple Parts and

- 1171 Multiple Responses. *Journal of Statistical Software*, 34(1), 1–13.
1172 doi:10.18637/jss.v034.i01
- 1173 Zeileis, A., & Grothendieck, G. (2005). Zoo: S3 Infrastructure for Regular and Irregular
1174 Time Series. *Journal of Statistical Software*, 14(6), 1–27. doi:10.18637/jss.v014.i06
- 1175 Zeileis, A., & Hothorn, T. (2002). Diagnostic Checking in Regression Relationships. *R News*,
1176 2(3), 7–10. Retrieved from <https://CRAN.R-project.org/doc/Rnews/>

1177

Appendices

1178

Appendix A: Moral Scenario

1179

Julie and Mark, who are brother and sister, are travelling together in France. They

1180

are both on summer vacation from college. One night they are staying alone in a cabin near

1181

the beach. They decide that it would be interesting and fun if they tried making love. At

1182

very least it would be a new experience for each of them. Julie was already taking birth

1183

control pills, but Mark uses a condom too, just to be safe. They both enjoy it, but they

1184

decide not to do it again. They keep that night as a special secret between them, which

1185

makes them feel even closer to each other (Haidt et al., 2000).

1186

Appendix B: Sample Statements to Challenge Judgement

1187

- Do you not agree that any concerns regarding reproductive complications are eased by their using of two forms of contraception?

1188

1189

- And do you accept that they are both consenting adults, and that they both consented and enjoyed it?

1190

1191

- And do you concede that nobody else was affected by their actions?

1192

Appendix C: Post Discussion Questionnaire

1193

How sure were you about your judgement?

1194

1 2 3 4 5 6 7

1195

Not at all

Extremely sure

1196

1197

How much did you change your mind?

1198

1 2 3 4 5 6 7

1199

Not at all

Extremely

1200

1201

How confused were you?

1202

1 2 3 4 5 6 7

1203

Not at all

Extremely confused

1204

1205

How irritated were you?

1206

1 2 3 4 5 6 7

1207

Not at all

Extremely irritated

1208

1209

How much was your judgement based on reason?

1210

1 2 3 4 5 6 7

1211

Not at all

Extremely

1212

1213

How much was your judgement based on "gut" feeling?

1214

1 2 3 4 5 6 7

1215

Not at all

Extremely

Appendix D: Supplementary Materials: Additional Analyses

1216

Study 2: Test for Order Effects

1217

1218 Recall that the questions were blocked for randomisation. Tests for effects of the
1219 order of the blocks revealed no difference in initial rating, $t(106.87) = -1.64$, $p = .104$, $d =$
1220 0.29 ; no difference in responding to the critical slide, $\chi^2(2, N = 111) = 4.76$, $p = .093$, $V =$
1221 0.21 ; and no difference in response to the generic potential harm question (“How would you
1222 rate the behavior of two people who engage in an activity that could potentially result in
1223 harmful consequences for either of them?”), $t(85.40) = -1.02$, $p = .312$, $d = 0.20$. A
1224 chi-squared test for independence revealed no significant association between order of blocks
1225 and judgments of boxing, $\chi^2(1, N = 111) = 2.86$, $p = .091$, $V = 0.16$, or the question
1226 regarding contact team sports, $\chi^2(1, N = 111) = 0.19$, $p = .660$, $V = 0.04$.

1227 The order of the questions regarding the application of the harm principle was also
1228 randomised. A one-way ANOVA revealed a significant difference in responses to the question
1229 “How would you rate the behavior of two people who engage in an activity that could
1230 potentially result in harmful consequences for either of them?” (1 = *Extremely wrong*; 4 =
1231 *Neutral*; 7 = *Extremely right*) depending on when it was presented $F(2, 109) = 4.757$ $p =$
1232 $.010$, partial $\eta^2 = .080$. Tukey’s post-hoc pairwise revealed that, when this question was
1233 responded to first, participants ratings were significantly lower ($M = 2.80$, $SD = 1.43$) than
1234 when it was responded to second ($M = 3.57$, $SD = 1.21$), $p = .040$, or third ($M = 3.67$, SD
1235 $= 1.31$), $p = .014$; and there was no difference in responding to this question second ($M =$
1236 3.57 , $SD = 1.21$) or third ($M = 3.67$, $SD = 1.31$), $p = .932$.

1237 A chi-squared test for independence revealed no significant association between order
1238 these questions and responses to the question “Do you think boxing is wrong?”, $\chi^2(2, N =$
1239 $111) = 4.88$, $p = .087$, $V = 0.21$. Similarly, a chi-squared test for independence revealed a
1240 significant association between order these questions and responses to the question “Do you
1241 think playing contact team sports (e.g. rugby; ice-hockey; American football) is wrong?”,

1242 $\chi^2(2, N = 111) = 1.79, p = .409, V = 0.13.$

1243 **Study 3: Test for Order Effects**

1244 As in Study 2, the questions were blocked for randomisation. Tests for effects of the
1245 order of the blocks revealed no difference in initial rating, $t(465.55) = 1.76, p = .079, d =$
1246 0.16 ; no difference in responding to the critical slide, $\chi^2(2, N = 502) = 1.12, p = .570, V =$
1247 0.05 ; no difference in responses to the generic potential harm question, $t(443.45) = 0.99, p =$
1248 $.322, d = 0.09.$ no association with judgments of boxing, $\chi^2(1, N = 502) = 1.03, p = .310, V$
1249 $= 0.05,$ or the question regarding contact team sports, $\chi^2(1, N = 502) = 1.15, p = .283, V$
1250 $= 0.10,$ depending on order of blocks.

1251 Regarding the three questions assessing the application of the harm principle, a
1252 one-way ANOVA revealed a significant difference in responses to the generic potential harm
1253 question depending on when it was presented $F(2, 499) = 23.512 p < .001,$ partial $\eta^2 = .086.$
1254 Tukey's post-hoc pairwise revealed that, when this question was responded to first,
1255 participants ratings were significantly lower ($M = 2.60, SD = 1.46$) than when it was
1256 responded to second ($M = 3.50, SD = 1.44$), $p < .001,$ or third ($M = 3.47, SD = 1.20$), $p <$
1257 $.001;$ and there was no difference in responding to this question second ($M = 3.50, SD =$
1258 1.44) or third ($M = 3.47, SD = 1.20$), $p = .983.$ As in Study 2, it seems likely that the
1259 named behaviours in the other questions provide an example of potential harm that is
1260 acceptable, leading to a more favourable response to this more abstract question. There was
1261 no significant association between question order and responses to the question "Do you
1262 think boxing is wrong?", $\chi^2(2, N = 502) = 1.12, p = .570, V = 0.05;$ or "Do you think
1263 playing contact team sports (e.g. rugby; ice-hockey; American football) is wrong?", $\chi^2(1, N$
1264 $= 502) = 1.03, p = .310, V = 0.05.$

1265 **Study 3: Individual Differences**

1266 A series of logistic regressions were conducted to investigate if dumbfounded
1267 responding was related to any of the individual difference variables Religiosity (as measured
1268 by CRSi7 Huber and Huber 2012), or Meaning in Life (Presence and Search, measured using
1269 MLQ Steger et al. 2008). We first report the results for each variable individually, followed
1270 by the combined model.

1271 **Religiosity.** The overall mean Religiosity score was $M = 2.57$, $SD = 1.17$. The
1272 mean religiosity scores for participants depending on response to the critical slide were as
1273 follows: $M = 2.84$, $SD = 1.17$ for participants who provided reasons, $M = 2.42$, $SD = 1.11$
1274 for participants who were dumbfounded, and $M = 2.28$, $SD = 1.12$ for participants who
1275 selected “There is nothing wrong”.

1276 A multinomial logistic regression revealed a statistically significant association
1277 between Religiosity and response to the critical slide, $\chi^2(2, N = 502) = 17.38$, $p < .001$, The
1278 observed power was 0.97. Religiosity explained approximately 2.40% (McFadden R square)
1279 of the variance in responses to the critical slide. Participants with higher religiosity scores
1280 were significantly more likely to provide reasons than to present as dumbfounded, Wald =
1281 6.14, $p = .013$, odds ratio = 0.73, 95% CI [0.57, 0.94], or select “There is nothing wrong”
1282 Wald = 15.24, $p < .001$, odds ratio = 0.65, 95% CI [0.53, 0.81]. See Figure 3.

1283 **Meaning in Life (Presence).** The overall mean Meaning in Life (Presence) score
1284 was $M = 4.74$, $SD = 1.66$. The mean Meaning in Life (Presence) scores for participants
1285 depending on response to the critical slide were as follows: $M = 5.01$, $SD = 1.67$ for
1286 participants who provided reasons, $M = 4.35$, $SD = 1.42$ for participants who were
1287 dumbfounded, and $M = 4.62$, $SD = 1.73$ for participants who selected “There is nothing
1288 wrong”.

1289 A multinomial logistic regression revealed a statistically significant association
1290 between Meaning in Life (Presence) and response to the critical slide, $\chi^2(2, N = 345) = 8.46$,

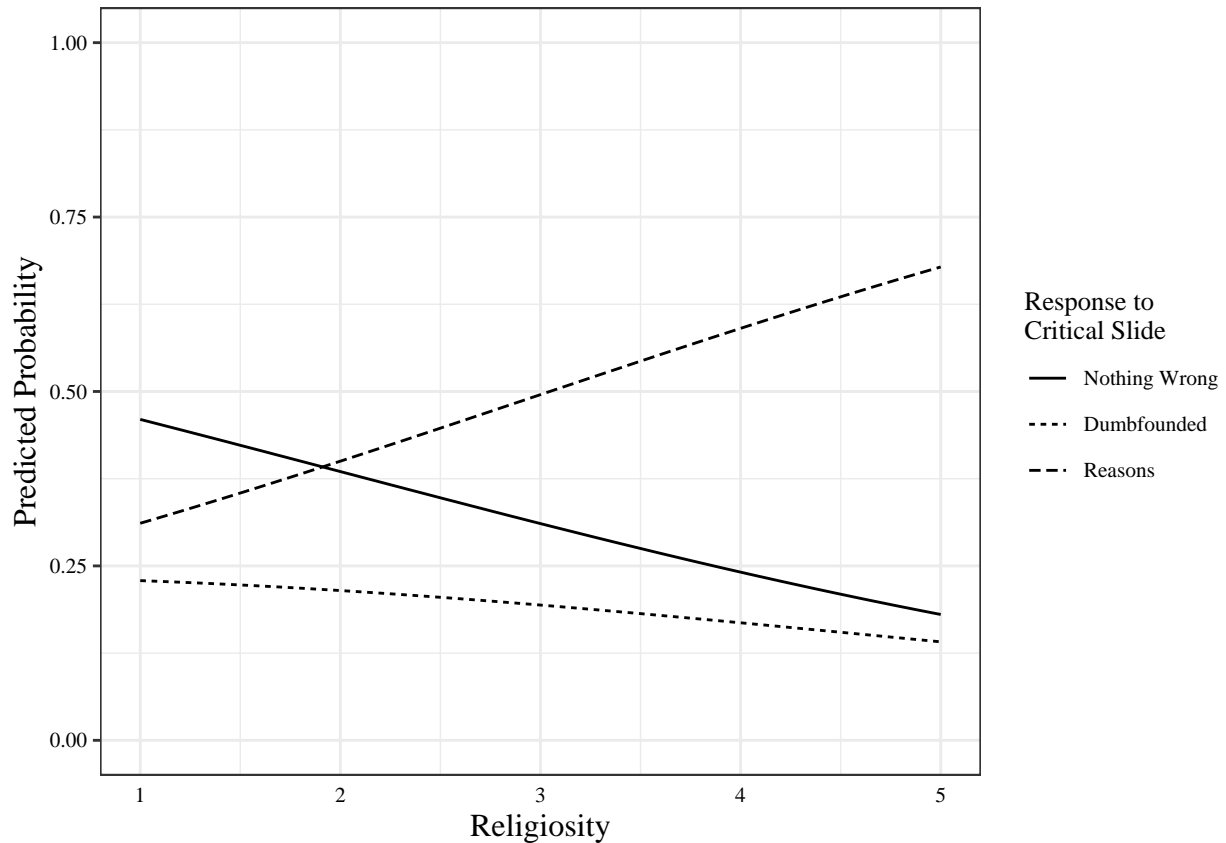


Figure 3. Probability of selecting each response to the critical slide depending on Religiosity

1291 $p = .015$, The observed power was 0.74. Meaning in Life explained approximately 1.17%
 1292 (McFadden R square) of the variance in responses to the critical slide. Participants with
 1293 higher MLQ: presence scores were significantly more likely to provide reasons than to present
 1294 as dumbfounded, Wald = 7.46, $p = .006$, odds ratio = 0.79, 95% CI [0.66, 0.93].
 1295 (Participants with higher MLQ: presence scores were marginally more likely to provide
 1296 reasons than to select “There is nothing wrong” Wald = 3.77, $p = .052$, odds ratio = 0.86,
 1297 95% CI [0.74, 1.00].) See Figure 4.

1298 **Meaning in Life (Search).** The overall mean Meaning in Life (Search) score was
 1299 $M = 4.47$, $SD = 1.73$. The mean Meaning in Life (Search) scores for participants depending
 1300 on response to the critical slide were as follows: $M = 4.42$, $SD = 1.75$ for participants who
 1301 provided reasons, $M = 4.55$, $SD = 1.68$ for participants who were dumbfounded, and $M =$
 1302 4.49 , $SD = 1.73$ for participants who selected “There is nothing wrong”.

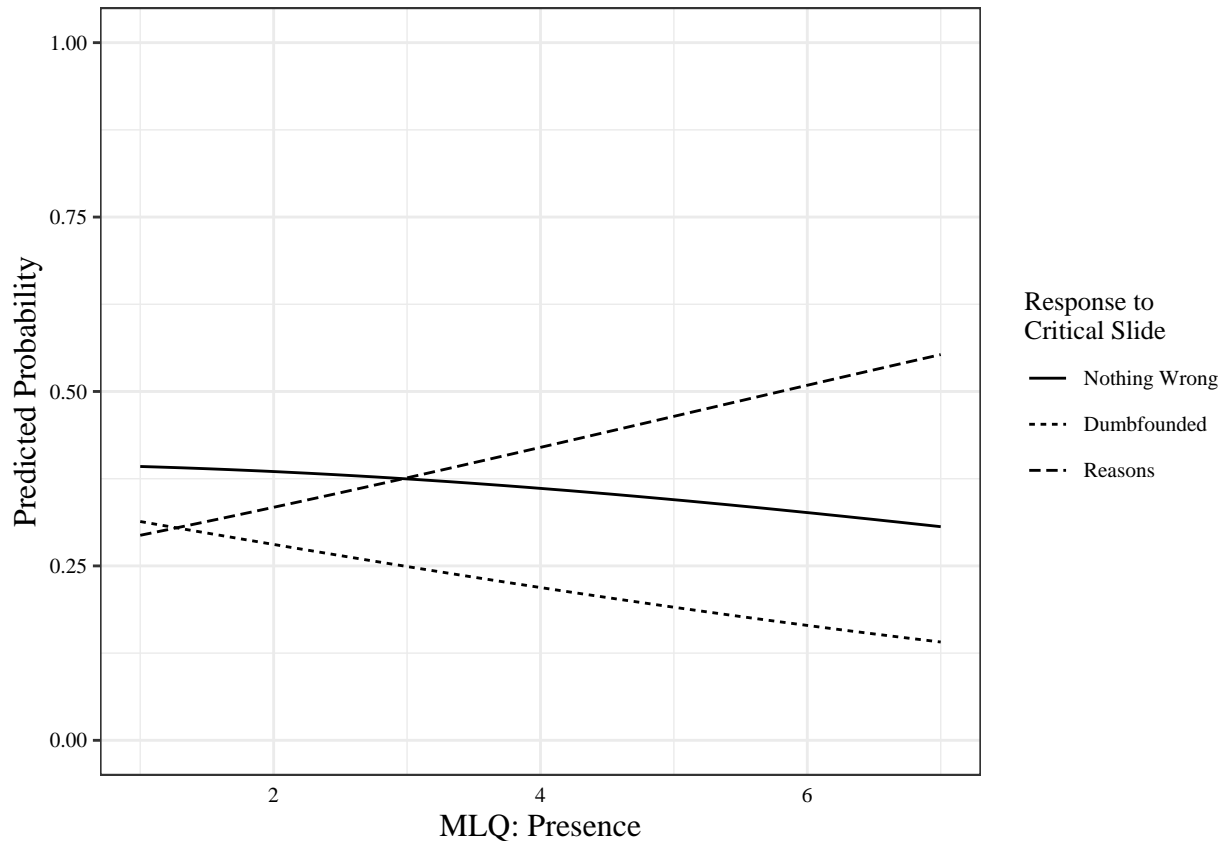


Figure 4. Probability of selecting each response to the critical slide depending on MLQ: Presence

1303 A multinomial logistic regression revealed no statistically significant association
 1304 between Search for Meaning in Life and response to the critical slide, $\chi^2(2, N = 345) = 0.3$,
 1305 $p = .859$, The observed power was 0.07. See Figure 5.

1306 **Individual Differences.** When analysed together, a multinomial logistic
 1307 regression revealed a statistically significant association between the three individual
 1308 difference variables and response to the critical slide, $\chi^2(6, N = 345) = 22.15$, $p = .001$, The
 1309 observed power was 0.99. The model explained approximately 3.07% (McFadden R square)
 1310 of the variance in responses to the critical slide. Religiosity was the only significant predictor
 1311 (see Table 1). Participants who scored higher in Religiosity were significantly more likely to
 1312 provide reasons than to select “There is nothing wrong”, Wald = 12.90, p , odds ratio =
 1313 0.65, 95% CI [0.52, 0.82]. It seems religiosity was more related to valence of judgement than

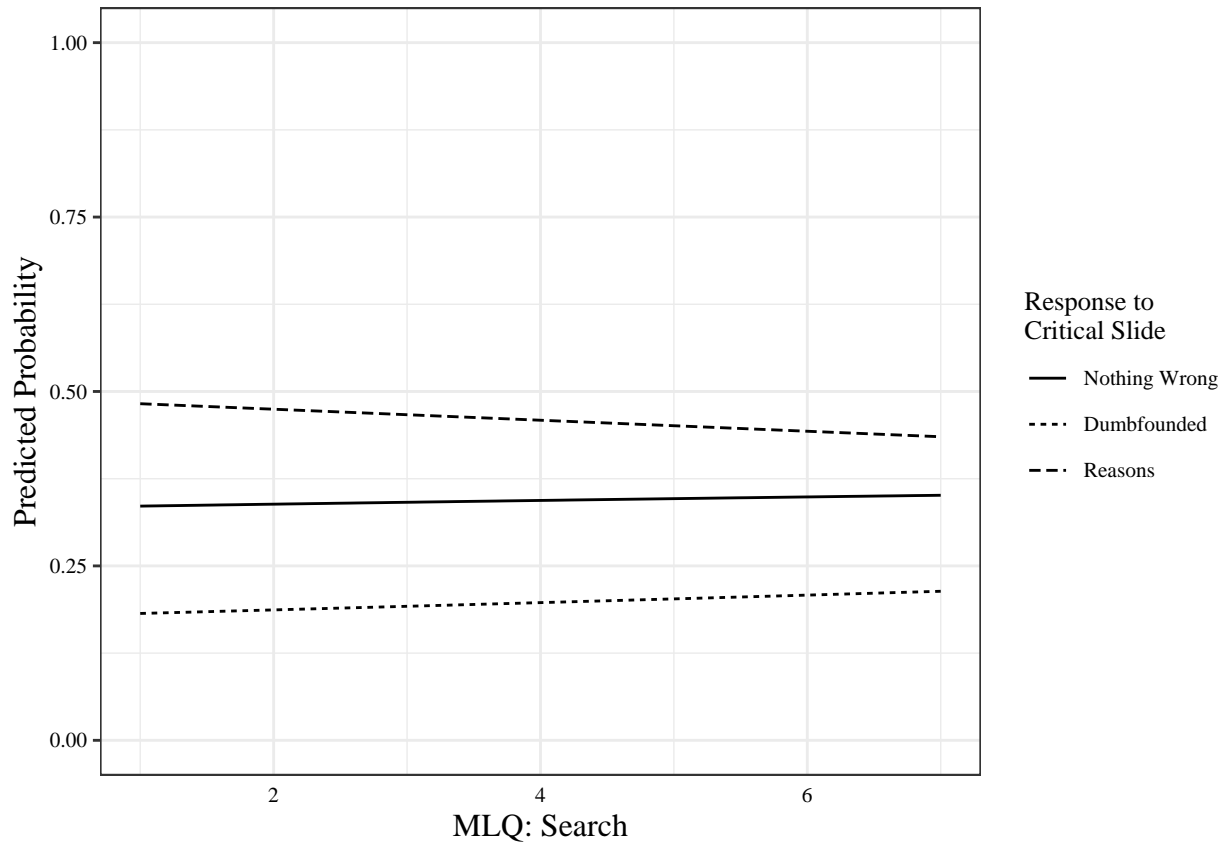


Figure 5. Probability of selecting each response to the critical slide depending on MLQ: Search

1314 to ability to provide reasons Wald = 3.04, p , odds ratio = 0.78, 95% CI [0.59, 0.59].

1315 A linear regression was conducted to assess the relationship between the individual
 1316 difference variables (Religiosity, Meaning and Life Presence, Meaning in Life Search) and
 1317 initial judgement. The model significantly predicted valence of judgement, $R^2 = .04$,
 1318 $F(3, 497) = 6.22$, $p < .001$. Religiosity the only significant predictor, $b = -0.21$, 95% CI
 1319 $[-0.35, -0.07]$, $t(497) = -3.01$, $p = .003$ (MLQ: presence $b = -0.08$, 95% CI $[-0.18, 0.03]$,
 1320 $t(497) = -1.44$, $p = .152$; MLQ: search $b = 0.06$, 95% CI $[-0.03, 0.15]$, $t(497) = 1.25$,
 1321 $p = .212$). Participants who scored higher in Religiosity were more likely to condemn the
 1322 actions of Julie and Mark.

1323 A final multinomial logistic regression was conducted that included Initial Judgement

Table 1

Multinomial logistic regression predicting responses to the critical slide where providing reasons is the referent in each case.

Variable	Response	<i>B</i>	S.E.	Wald	<i>df</i>	<i>p</i>	O.R.	Lower	Upper
Religiosity	Dumbfounded	-0.247	0.141	3.045	6	.081	0.781	0.592	1.031
	Nothing wrong	-0.426	0.119	12.899	6	<.001**	0.653	0.518	0.824
MLQ: Presence	Dumbfounded	-0.173	0.095	3.29	6	.070	0.841	0.698	1.014
	Nothing wrong	-0.043	0.082	0.275	6	.600	0.958	0.815	1.125
MLQ: Search	Dumbfounded	0.054	0.09	0.358	6	.522	1.055	0.885	1.259
	Nothing wrong	0.074	0.076	0.95	6	.207	1.077	0.928	1.25

Note. * = sig. at $p < .05$; ** = sig. at $p < .001$

Table 2

Multinomial logistic regression predicting responses to the critical slide where providing reasons is the referent in each case.

Variable	Response	<i>B</i>	S.E.	Wald	<i>df</i>	<i>p</i>	O.R.	Lower	Upper
Religiosity	Dumbfounded	-0.236	0.147	2.588	8	.108	0.79	0.593	1.053
	Nothing wrong	-0.875	0.23	14.524	8	<.001**	0.417	0.266	0.654
MLQ: Presence	Dumbfounded	-0.155	0.099	2.453	8	.117	0.856	0.705	1.04
	Nothing wrong	0.031	0.139	0.05	8	.823	1.031	0.786	1.353
MLQ: Search	Dumbfounded	0.039	0.092	0.179	8	.672	1.04	0.868	1.245
	Nothing wrong	0.031	0.136	0.05	8	.823	1.031	0.789	1.347
Initial Judgement	Dumbfounded	0.39	0.139	7.821	8	.005*	1.477	1.124	1.942
	Nothing wrong	1.98	0.219	81.533	8	<.001**	7.241	4.712	11.128

Note. * = sig. at $p < .05$; ** = sig. at $p < .001$

1324 as a predictor variable. The results are shown in Table 2. Overall the model was a significant
 1325 predictor of response to the critical slide, $\chi^2(8, N = 345) = 292.33, p < .001$, The observed
 1326 power was 1. The model explained approximately 40.54% (McFadden R square) of the
 1327 variance in responses to the critical slide. As shown in Table 2, Religiosity appeared to be
 1328 related only to valence of judgement on the critical slide, initial judgement appeared to
 1329 predict valence of judgement and ability to provide reasons, with more extreme judgements
 1330 of “wrong” most strongly predicting the providing of reasons. The relative probabilities of
 1331 selecting each response to the critical slide depending on initial judgement are displayed in
 1332 Figure 6.

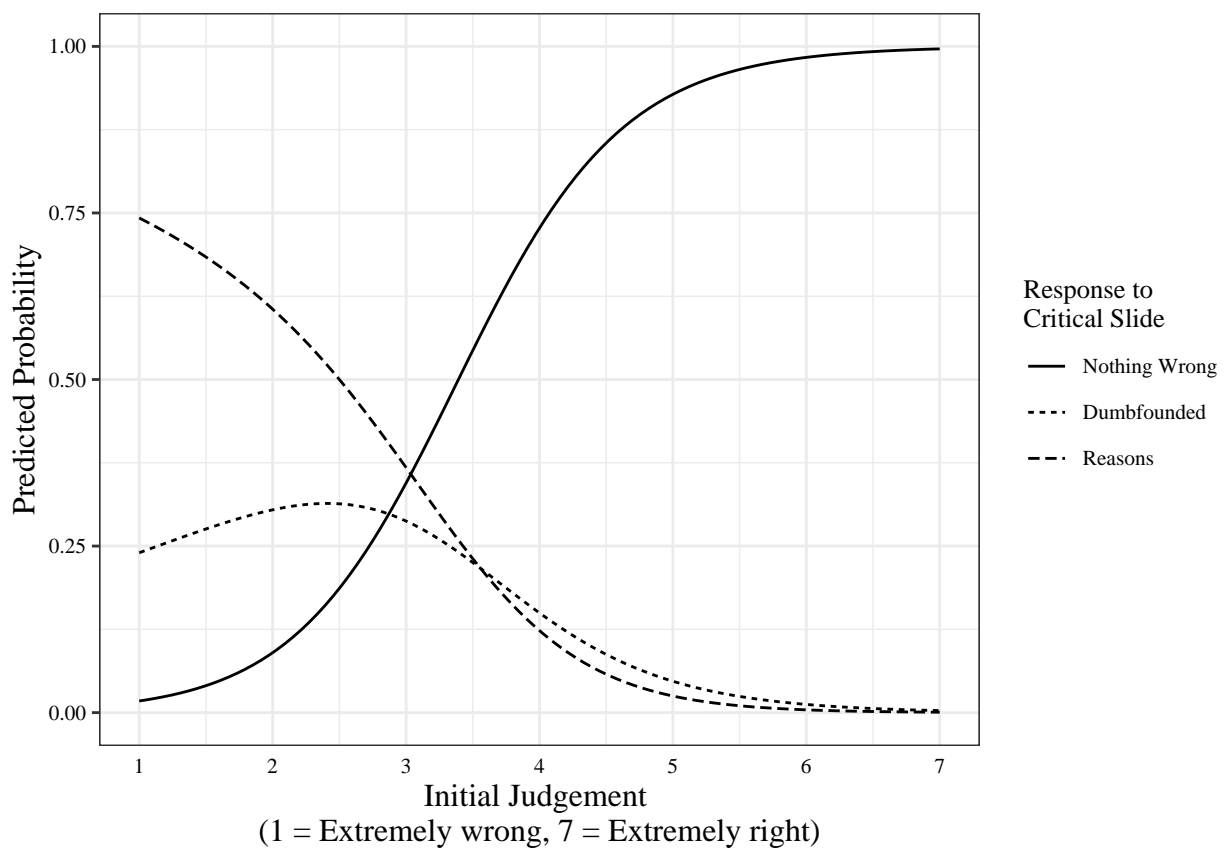


Figure 6. Probability of selecting each response to the critical slide depending on Initial Judgement.