

**Beyond good and bad: Varieties of moral judgment**

William Jiménez-Leal<sup>1</sup>, Samuel Murray<sup>2</sup>, Santiago Amaya<sup>3</sup>, and Sergio Barbosa<sup>1,4</sup>

<sup>1</sup>Department of Psychology, Universidad de los Andes

<sup>2</sup>Mind At Large Lab, Imagination and Modal Cognition Lab, Duke University

<sup>3</sup>Department of Philosophy, Universidad de los Andes

<sup>4</sup>Department of Psychology, Universidad del Rosario

Word Count: 15910.

**Author Note**

The first three authors contributed equally. Correspondence concerning this article should be addressed to William Jiménez-Leal, Universidad de los Andes. E-mail:

[w.jimenezleal@uniandes.edu.co](mailto:w.jimenezleal@uniandes.edu.co).

Supplementary materials (including materials, preregistrations, raw data and code) available at: [https://osf.io/kja2u/?view\\_only=573c1541bbd240b397b79229704dfae7](https://osf.io/kja2u/?view_only=573c1541bbd240b397b79229704dfae7).

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21

### Abstract

We argue that people regularly encounter situations involving moral conflicts among permissible options. These scenarios, which some have called *morally charged situations*, reflect perceived tensions between moral expectations and moral rights. Studying responses to such situations marks a departure from the common emphasis on sacrificial dilemmas and widespread use of single-dimension measures. In 6 experiments (n=1607), we show that people use a wide conceptual arsenal when assessing actions that can be described as suberogatory (bad but permissible) or supererogatory (good but not required). In Experiment 1 we find that people identify actions as suberogatory or supererogatory when using open descriptions to describe them. Experiment 2 shows that they differentially assess these actions in terms of how permissible, optional, and good they considered them. Experiment 3 tests the use of these evaluative dimensions with sacrificial dilemmas. We fail to find differences between these categories when people respond to dilemmas, even when controlling for trait utilitarian tendencies. By including judgments of blameworthiness and sanction, Experiments 4 and 5 provide additional evidence of the granularity and the moral significance of these evaluations. In Experiment 6 people offered their own explanations of their responses. Qualitative analyses revealed that they frequently appeal to character traits, the presence of rights, and the absence of explicit duties. Taken together these results suggest a richer spectrum of both situations and concepts relevant to characterize moral judgment than moral psychologists up to this point have generally recognized.

*Keywords:* moral rules, moral dilemmas, suberogatory, supererogatory, duty

## 22 **Beyond good and bad: Varieties of moral judgment**

### 23 **The role of dilemmas in the science of morality**

24 Dilemmas have traditionally played a prominent role in the study of moral cognition.  
25 While many have recently registered dissatisfaction with this state of affairs (Andrade, 2019;  
26 Bauman, McGraw, Bartels & Warren 2014; Dahl & Oftedal, 2019; Everett & Kahane, 2020;  
27 Kahane, 2020), there is considerably less said about how historically we got here. Dilemmas take  
28 on this crucial role because of the assumption that morality is a system of rules. But this  
29 assumption is problematic for reasons pertaining to both normative theories of morality and the  
30 psychology of moral judgment.

31 From its early days, psychologists interested in moral judgment understood morality as a  
32 structured set of rules prohibiting and prescribing certain behaviors. Piaget (1932), for instance,  
33 opened his landmark *The Moral Judgment of the Child* with the statement: “All morality consists  
34 in a system of rules, and the essence of all morality is to be sought for in the respect which the  
35 individual acquires for these rules” (1932, p. 1). Fauconnet echoes this sentiment in his  
36 *Responsibility*, stating that moral responsibility is a “quality belonging to those who must...in  
37 virtue of a rule be chosen as the passive subjects of a punishment” (1920, p. 11).

38 Thinking of morality as a system of rules has several theoretical pay-offs. For one,  
39 various deontic concepts can be inter-defined through rules. For example, acts are permissible if  
40 and only if those acts accord with the rules, and something counts as good only if it accords with  
41 the rules. Likewise, the impermissible is whatever constitutes violating a rule, and something is  
42 bad insofar as doing it violates a rule (see, for example, Kanger (1971) and Anderson (1958) for

43 two formal attempts to explain inter-definability). This, as we shall see, results in a  
44 methodological advantage. If true, asking subjects whether something is impermissible, bad, or  
45 in violation of a duty comes close to asking one and the same question.

46 Focusing on rules also provides a concrete way to measure moral development. Under the  
47 assumption that maturing moral judgment consists in possessing a greater moral understanding,  
48 it is possible to define this understanding in terms of increasing aptitude in applying more  
49 sophisticated evaluative rules (Kohlberg & Hersh, 1977). This insight famously informed the  
50 stage theory of moral development (Kohlberg, Levine, & Hower, 1983), which identified moral  
51 stages according to the kinds of rules that informed judgments: from pre-conventional rules  
52 (“Doing this is bad because I can get punished”) to universal, exceptionless principles of  
53 impartial justice (“Killing is wrong”).

54 Moral dilemmas are one specific kind of *moral encounter* (Monin, Pizarro & Beer, 2007).  
55 They characteristically present situations requiring decisions among impermissible options, in  
56 the sense that each violates some plausible moral imperative (Sinnott-Armstrong, 1988: 29-30)  
57 In this respect, they isolate different sets of rules (egocentric v. altruistic, in-group v. out-group,  
58 etc.) that are normally taken to inform moral cognition. Thus, by looking at the choices people  
59 make in these situations or by studying how they assess the decisions made by actors depicted in  
60 them, it seems possible to better understand which rules inform their judgments of goodness and  
61 badness. For similar reasons, moral dilemmas also seem to provide a good instrument to measure  
62 moral development and test hypotheses about individual developmental trajectories.

63            Obviously, theories of moral cognition that fall under this rule-based tradition differ from  
64 one another in important respects (see Darley & Schultz, 1990 for discussion). A sign of this are  
65 the well-known controversies that exist among them. Kohlberg's research program, for example,  
66 was criticized for the use of culturally biased materials (Simpson, 1974; Snarey, 1985), biased  
67 samples (Walker, 1984) and the idea of a linear progression in moral development (Rest, 1979).  
68 Others criticized the kind of dilemmas used to elicit moral judgment, focusing on the artificial  
69 nature or the mundane character of them (Rosen, 1980; Bauman et al., 2014). At bottom,  
70 however, many of Kohlberg's critics agreed that morality was a system of rules that could be  
71 studied by means of situations where plausible moral rules conflict with one another. They just  
72 disagreed about how to properly characterize these rules and the situations that best exemplified  
73 these conflicts.<sup>1</sup>

#### 74 **Rules and Commonsense Morality**

75 Some contemporary psychologists have sought to take distance from this early tradition by  
76 proposing dual models of moral cognition (for discussion, see Crockett, 2013). We believe,  
77 however, that these models are in an important respect a *continuation of the rule-based tradition*  
78 that has dominated the study of moral cognition. While opening up new possibilities to

---

<sup>1</sup> One notable exception comes from feminist critiques of ethical theory. Several feminist ethicists have noted that the central preoccupation of moral theorizing prioritizes abstract principles over the particularity of ethical life (Gilligan, 1982. p. 32-38). The emphasis on abstract generalization obscures the details that are crucially important to a well-lived life (Young, 1987, p. 61-62).

79 understand the architecture of moral decision-making, the dual process paradigm has inherited  
80 from its early predecessors the view of morality as a system of rules.

81         Dual models are premised on the belief that commonsense morality reflects attitudes that  
82 lie on a continuum between full-blown utilitarian and deontological ethics. Differences on this  
83 dimension are supposed to be explained in terms of a deep architectural divide (Bartels &  
84 Pizarro, 2011; Christensen & Gomila, 2012; Conway & Gawronski, 2013; Djeriouat &  
85 Trémolière, 2014; Holyoak & Powell, 2016; Lee & Gino, 2015). Controlled processes produce  
86 characteristically utilitarian responses; intuitive processes produce characteristically  
87 deontological responses (Bialek & De Neys, 2017; Cushman, Young, & Greene, 2010; Moll &  
88 de Oliveira-Souza, 2007).

89         Among proponents of these models, there is a standing debate as to how deep this divide  
90 is. Some have interpreted it as showing that rules need not be represented as inputs to the  
91 decision-making processes that culminate in moral judgment (see Blair, 1995; Greene et al.,  
92 2001; Haidt, 2001). Others believe that the rules are always represented, except that sometimes  
93 they are only tacitly represented. (Mallon & Nichols, 2010; Mikhail, 2011). Still others have  
94 proposed characterizing the divide in terms of model-based and model-free processing (Crockett  
95 2013).

96         These discussions, however, operate within *the rule-based tradition* that we wish to  
97 challenge here. They concern the architectural or algorithmic processes underlying how moral  
98 rules are instantiated in judgments about specific actions. Precisely because of this, they leave  
99 untouched the more basic assumption regarding the subject matter of morality itself. After all,

100 the styles of thinking that are supposed to dominate commonsense morality are still  
101 conceptualized in terms of systems of moral rules: either expressions of duty or prescriptions to  
102 maximize or minimize some valued quantity.

103         In general, moral rules can be defined as functions from relevant inputs to behavioral  
104 imperatives. For example, a rule against murder is a general function from some action's being  
105 an instance of murder to an imperative against so acting. As Sidgwick (1981, p. 228) explains,  
106 "...rules of duty ought to admit of precise definition in a universal form." Under this definition,  
107 many generalizations (but not all of them) count as moral rules: the prohibition not to harm  
108 innocents, the injunction to maximize saved lives, etc.

109         Thus, it might be an open question whether the algorithmic processes that result in moral  
110 judgments take as inputs explicit representations of these functions. It is possible that not all  
111 moral decision-making is based on models shaped by a moral grammar. Be that as it may, the  
112 fact is that researchers working under this new paradigm still classify the *outputs* of these  
113 processes by conformity to the prescriptions of certain kinds of moral rules. The styles of  
114 thinking modeled by them, deontology and utilitarianism, which are supposedly characteristic of  
115 commonsense morality, are still styles of thinking in *accordance with some distinctive moral*  
116 *rules*.

117         It is not surprising, then, that contemporary work on the psychology of moral judgment  
118 continues to be dominated by moral dilemmas, in particular sacrificial dilemmas. Observing  
119 people's choices or their evaluation of the available options when each choice is made  
120 impermissible by deontological or utilitarian rules seems a natural way to measure people's

121 attitudes in terms of these moral frameworks. Whatever the underlying computational processes  
122 are, dilemmas seem recommended by the goal of understanding whether commonsense morality  
123 embodies deontological or utilitarian rules.

124 In sum, there is a long tradition, spanning from Piaget to contemporary dual models, that  
125 views morality as a rule-based system. Dilemmas appear useful for studying moral judgment in  
126 virtue of this underlying assumption. Hence, moving beyond dilemmas requires moving beyond  
127 the assumption that morality is a rule-based system. And moving beyond the assumption of  
128 morality as a rule-based system requires moving beyond the study of moral judgment in  
129 situations where rules conflict.

130 There are, as we shall see, numerous dimensions of the moral life that are not governed  
131 by rules, whether these refer to duties or codify maxims for the maximization of some valuable  
132 result. Hence, a science of morality that focuses only on dilemmas and rule-based judgments  
133 risks painting a picture of moral cognition that is overly narrow and stilted (Bauman et al., 2014).  
134 We want to challenge the use of dilemmas to study moral judgment because we reject the  
135 assumption that morality consists in a structured system of rules. Further, seeing how moral  
136 cognition operates beyond the rules provides a fresh perspective for the study of moral judgment.  
137 To this end, we examine in what follows moral judgment in cases of suberogatory behavior.

### 138 **The Supererogatory and the Suberogatory**

139 From time to time, people are faced with the option of doing more than what they are  
140 required to do, for example, spending some of their free time volunteering at the local animal  
141 shelter or donating a large portion of money to charity. These actions are admirable, despite the



142 fact that failing to do them does not seem to merit condemnation. Some classify these actions as  
143 ‘supererogatory’, going above and beyond the call of duty when failing to go above and beyond  
144 is perfectly morally permissible (Archer, 2018).

145 More controversially, people seem at times to underperform relative to some ideal in a  
146 way that is permissible, for instance, not offering to proctor the exam of a sick colleague despite  
147 the fact that one is available and the colleague has helped one in the past. Doing this tends to be  
148 regarded as bad despite the fact that there is no rule that requires one to pick up the duties of sick  
149 colleagues or that one no explicit agreement to help each other was made. Some use the label  
150 ‘suberogatory’ to describe this kind of behavior (see Driver, 1992; Hurd, 1998). The behavior is  
151 morally objectionable but there is no well-defined duty that it violates.

152 In failing to do a supererogatory action, one usually does not do anything bad. It’s an  
153 admirable thing to donate money to charity, but failing to do so is not reprehensible. However, in  
154 some situations, failing to do a supererogatory action constitutes suberogatory behavior. If a  
155 tourist asks you for directions, you are completely within your rights to walk away without  
156 saying anything. Doing it, though permissible, is bad, whereas giving directions is good despite  
157 not being required. This possibility suggests a different kind of conflict that people encounter in  
158 their day-to-day experiences of morality: conflicts between equally permissible good and bad  
159 options. Here, we refer to these moral encounters that do not constitute real dilemmas as *morally*  
160 *charged situations*. We claim that these situations, along with the concepts used to evaluate  
161 them, offer a distinctive opportunity to study moral judgment.

162           Sacrificial dilemmas represent an interesting albeit limited subset of what people  
163 encounter in their everyday life. Using ecological momentary assessment, Hofmann, Wisneski,  
164 Brandt and Skitka (2014) asked a large sample of people to report whether they had witnessed,  
165 committed, or heard about a moral situation during the last hour, five times a day for three days.  
166 While some situations resemble the kind of conflict expressed in sacrificial dilemmas (e.g.,  
167 “Reminded waitress I did not pay for my bill when she thought I did”), many of the situations  
168 reported seemed more similar to the situations depicted above. The choice between the  
169 competing options didn’t seem a matter of aligning oneself with some well-defined rule.

170           Over-relying on dilemmas risks papering over these distinctions and, more generally, the  
171 differences between moral categories that inform varieties of moral judgment. O’Hara, Sinnott-  
172 Armstrong, and Sinnott-Armstrong (2010) compared responses to 15 dilemmas where they asked  
173 people to rate how wrong, forbidden, inappropriate and blameworthy an action was. They found  
174 that “the influence of wording variations on moral judgments was negligible” (p. 552) and they  
175 analyzed the small differences found as a matter of magnitude. Likewise, many researchers treat  
176 terms like ‘forbidden’ or ‘blameworthy’ as linguistic variations of some homogenous moral  
177 judgment (Bjorklund, 2003; Cushman, Young, & Hauser, 2006; Greene et al., 2001b; Koenigs et  
178 al., 2012). The assumption is that common sense moral judgment is not granular enough to  
179 reflect differences between being forbidden, blameworthy, bad, and so on. It is, instead,  
180 monolithic cognitive product, to which different labels provide different access points.

181           While there have been calls to more carefully use these measures (Christensen & Gomila,  
182 2012; Monin, Pizarro, & Beer, 2007) it is unclear what the rationale for using different terms  
183 could be (cf. Cushman, 2008; Barbosa & Jiménez-Leal, 2017). But, once the repertoire of moral

184 encounters is expanded to include non-dilemma situations, it is possible that nuances and  
185 variability in moral cognition will emerge. In sum, by expanding the kind of moral encounters  
186 used when empirically probing people’s intuitions and by enlarging the dimensions along which  
187 these encounters are assessed we can more adequately study the variation and granularity of  
188 moral judgment.

189         Here we present 6 experiments that study moral judgments in non-dilemma situations.  
190 One key feature of these situations is that the correct choice—if there is one—is not obviously  
191 settled by appealing to rules. Experiment 1 is an exploratory study that maps out the descriptions  
192 people offer of different situations. We find that people describe situations as suberogatory (‘bad  
193 but permissible’) or supererogatory (‘good but not required’) when supplying open descriptions  
194 of them. Experiment 2 shows that judgments of good/bad, permissible/impermissible, and  
195 optional/obligatory dissociate when evaluating suberogatory and supererogatory situations. We  
196 also find that people’s beliefs about duties negligibly correlate with judgments of goodness,  
197 permissibility, obligation, and blame. In Experiment 3, we compare judgments in sacrificial  
198 dilemmas to see whether the same distinctions appear. We do not find the same dissociations,  
199 which suggests that eliciting these patterns of judgment requires more than giving participants  
200 the options to judge along various dimensions. In Experiment 4, we included a measure of  
201 praise/blame to see whether it correlates with the main “erogatory” measures. In Experiment 5,  
202 we replicated previous findings by using different vignettes that describe more characteristically  
203 moral situations adapted from classic philosophical thought experiments about abortion and  
204 property rights. We also measured whether beliefs about rights predict any kind of judgment. We  
205 find the same pattern of dissociations in these different vignettes and again find a negligible

206 correlation between judgments and individual beliefs about rights. In Experiment 6, we  
207 conducted a qualitative study to begin exploring the variety of factors that differentially drive  
208 different judgments. Using independent coders for a qualitative analysis, we find that people  
209 describe these situations using the language of character traits and rights rather than duties. In  
210 fact, in many cases people explicitly state the *absence* of duties to do anything in our scenarios

211         A methodological coda: Even though this research is mostly exploratory, its ideas are  
212 developed against a backdrop of well-established findings in moral psychology. We decided to  
213 preregister Experiments 2 to 6 because we believe that clearly establishing design and analysis  
214 plans can help distinguish the confirmatory and exploratory aspects of our research by clearly  
215 specifying our intent. The procedure, therefore, reduces needless post hoc interpretations (Nosek  
216 et al, 2019).<sup>2</sup> Materials, data, preregistrations and code for all experiments are available on the  
217 OSF page of the project  
218 ([https://osf.io/kja2u/?view\\_only=573c1541bbd240b397b79229704dfae7](https://osf.io/kja2u/?view_only=573c1541bbd240b397b79229704dfae7)). The IRB of the  
219 University (blinded for review) approved this study.

---

<sup>2</sup> The only important deviation from preregistration plans occurred in Experiment 2, where the main statistical analysis proposed (a repeated measures ANOVA) was replaced by a mixed linear model, since it is better suited to model our data.

220

**Experiment 1**

221           The objective of this first study is exploratory. We presented participants with vignettes  
222 describing either suberogatory or supererogatory behavior. They were instructed to select words  
223 from a list to describe these scenarios and to offer a description in their own words. The goal of  
224 this is twofold. The first is to see what language people use to spontaneously describe a moral  
225 encounter. The second is to see whether people recognize a distinction between different  
226 judgment categories that maps onto the complex category of suberogatory. This requires that  
227 people have distinct concepts of permissibility and goodness such that they can describe some  
228 behavior as bad but permissible. Hence, we decided to run an exploratory study using word  
229 selection and open response to see whether participants utilize the moral categories we aim to  
230 study without being prompted.

231           We expected people to always select always more than one word (e.g., “good”) and to  
232 give descriptions that characterize both the action and the person.

**233 Method****234 Participants**

235           95 participants (60 women and 35 men, mean age = 32.42, SD =10.43), based in the  
236 United States and recruited through Prolific Academic, took part in the study in exchange for 40  
237 pence. Participants were aware that their answers would be anonymous and were monetarily  
238 compensated for their participation. The average completion time was 5.2 minutes.

**239 Materials and procedure**

240 We constructed four scenarios, some of which were based on thought experiments by  
241 Driver (1992). Each scenario described an individual faced with a choice between a suberogatory  
242 and a supererogatory option. Additionally, in order to account for possible asymmetries between  
243 actions and omissions (Haidt & Baron 1996), we created two versions of each scenario that  
244 describe either an action or an omission. This generated eight vignettes, described below  
245 (suberogatory versions of these vignettes are in brackets):

~~246~~ **Newlyweds** Two newlyweds are boarding a plane to go on their honeymoon. Because of a booking error  
247 by the airline, the couple does not have seats together. They ask someone, already seated, if they would switch  
248 seats so the couple could sit together. The passenger switches seats, and the newlyweds can sit together. [*The*  
249 *passenger does not switch seats, and the newlyweds have to sit separately.*]

~~250~~ **Kidney** Alex is suffering from severe kidney failure and Alex's only hope is to obtain a transplanted  
251 kidney. Alex's cousin, Jamie, is the only known compatible donor. Jamie offers to donate the kidney to Alex. [*Jamie*  
252 *does not offer to donate the kidney to Alex.*]

~~253~~ **Mowing** Early one Sunday morning when the neighbors are usually sleeping, Sam notices that the lawn  
254 needs to be mowed. Although it is his property and it would be inconvenient to do it later, he decides to not mow  
255 the lawn. He knows that starting the lawn mower will probably wake up the neighbors. [*Even though he knows*  
256 *that starting the lawn mower will probably wake up the neighbors, he does it anyway. It's his property and it will be*  
257 *inconvenient to mow the lawn later.*]

~~258~~ **Raffle** During the Christmas party, the secretary publicly announced the results of the office raffle:  
259 "Congratulations to Alex, who has won the trip for two to Disney World. She can come up front to claim her prize  
260 or she can let a cash equivalent go to a hurricane relief fund". After hearing the news, Alex looked excited: ". Even  
261 though I have the winning ticket and Disney World sounds fun, I am going to donate the prize to one of the

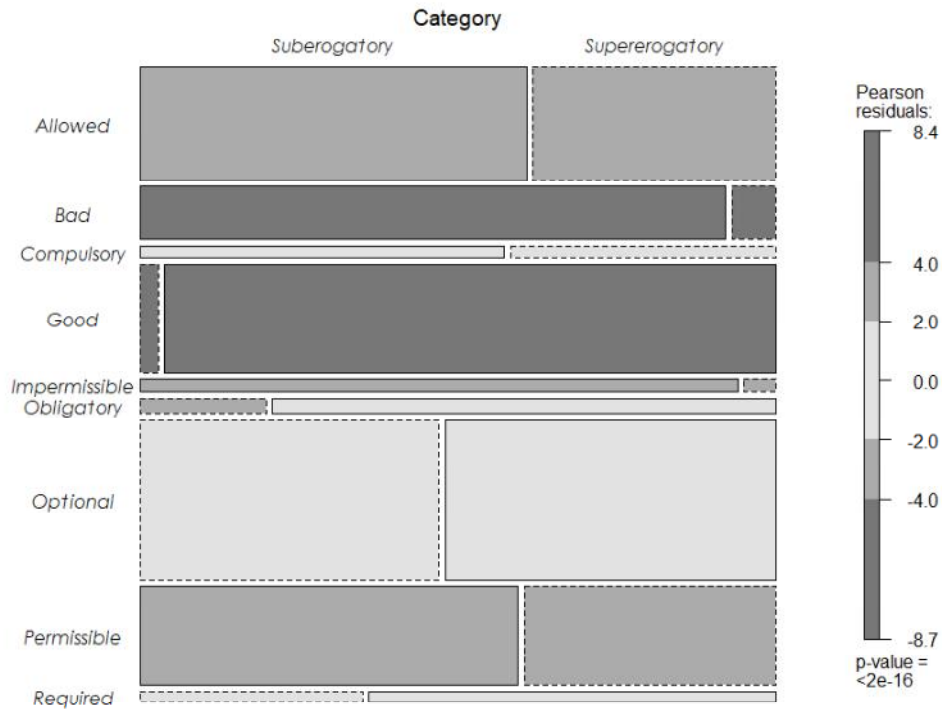
262 charities." [After hearing the news, Alex looked excited: "I have the winning ticket! Even though I don't really care  
263 much about Disney World, I am going to claim the prize anyway".]

264 Participants were presented with four out of the eight possible variations, so each  
265 participant was presented with a suberogatory omission, a suberogatory action, a supererogatory  
266 omission, and a supererogatory action. After reading the vignettes, participants completed a word  
267 selection task by selecting "the word(s) that you think best describe the situation." The word  
268 choices were: "permissible", "impermissible", "required", "good", "obligatory", "allowed",  
269 "bad", "optional" and "compulsory". Participants also completed an open description task by  
270 offering a description of the situation in their own words. The vignette presentation order was  
271 counterbalanced and the order of tasks and words was randomized.

## 272 **Results**

273 For the word selection task, participants selected 2.6 words on average (see Figure 1). For  
274 suberogatory behaviors, people most often chose the words "optional" (27%) and "allowed"  
275 (25%), followed by "permissible" (21%). For supererogatory behaviors, "good" (35%),  
276 "optional" (28%), and "allowed" (14%) were the most common choices. "Allowed" and  
277 "optional" are the most common pair of words used across all vignettes. "Allowed" and  
278 "permissible" are more strongly associated with evaluating suberogatory behaviors. Figure 1  
279 summarizes these results.

280



281

282 *Figure 1. Mosaic plot for words selection. Height is associated with overall selection frequency*  
 283 *and width is associated with prevalence for suberogatory and supererogatory behaviors (Hornik,*  
 284 *Zeileis & Meyer, 2006). Darker shades and solid lines represent positive associations, whereas*  
 285 *lighter shades and dotted lines represent a negative association within the sub/supererogatory*  
 286 *categories. The plot represents a model contrasting observed and expected frequencies of word*  
 287 *choices.*

288



289           In the open descriptions, participants predominantly used character trait descriptors, such  
290 as *selfish* for the suberogatory vignettes and *thoughtful* for the supererogatory situation (see  
291 Figure 2). They rarely used words like “permissible”. Also, more nuanced descriptions of the  
292 suberogatory situations generally highlighted the optionality of the response, in line with the  
293 abstract concept of the suberogatory: “he is allowed to do it” (mowing), “Jamie has a right to her  
294 own decision” (Kidney), “it is understandable” (kidney). Lastly, people were generally sensitive  
295 to the different moral aspects that structure the situations depicted as morally charged. They  
296 were, for instance, quick to describe the action/person negatively (mean, rude, etc), while also  
297 recognizing that rights and expectations were at play in the scenarios evaluated. They also  
298 recognized the possibility that some behaviors could be bad but permissible.:

- 299 • “It is a little bit selfish, but then again she has the right to keep her organs”
- 300 • “While it was mean, the passenger has the right to refuse
- 301 • “Rude, but it is his lawn”
- 302 • “Giving an organ is a big thing to ask. It is something that is optional, and there is no
- 303 mention of Alex asking him to do it so he is not required to offer”



316 We asked people to rate the sub- and supererogatory situations along three dimensions:  
317 good/bad, permissible/impermissible, and obligatory/optional. We manipulated the Erogation  
318 category (within-subjects, sub- and supererogatory) and the Situation Type (between-subjects,  
319 action or omission). We hypothesized that people would judge suberogatory situations as worse  
320 than supererogatory situations, but that permissibility and obligatoriness ratings between the two  
321 situations would not be significantly different. We also expected permissibility ratings to be  
322 significantly higher than obligatoriness ratings in both supererogatory and suberogatory  
323 conditions, though the difference between permissibility and obligation would be greater in the  
324 suberogatory condition than in the supererogatory condition. We did not expect any differences  
325 between judging actions and omissions. We also collected data on attitudes towards duties,  
326 expecting that ratings along the Erogation category would be associated with these attitudes.

### 327 **Participants**

328 We ran a power analysis for a mixed ANOVA (between-within interaction), assuming an effect  
329 size of  $f = 0.15$  using the software G\*Power. This analysis suggested a sample size of 272 for a  
330 0.95 power. To account for exclusions, we recruited 311 participants (186 women and 125 men,  
331 mean age = 32.77, SD = 11.18) through Prolific Academic. We decided to switch to an  
332 alternative data strategy after collecting the data, using mixed linear models, given the problems  
333 of repeated measures analyses with independence and distributional assumptions (Singmann &  
334 Kellen, 2019). Our sample size, however, is consistent with a power of 0.9, assuming  
335 participants can be treated as a random factor to account for within-person response variability,  
336 with a mixed design (Singmann & Kellen, 2019; Westfall, 2015; See Supplementary materials  
337 for details).

338           Each person voluntarily participated in the study and received 38 pence as compensation.  
339           The average completion time was 3.11 minutes.

#### 340   **Materials and procedure**

341           Materials were the same as Experiment 1. Each participant saw four of eight scenarios.  
342           We manipulated the moral category within participants and action/omission between  
343           participants, so each participant saw two supererogatory situations and two suberogatory  
344           situations, where all of them were either actions or omissions. For each vignette, participants  
345           evaluated the situation along different dimensions with a 100-point sliding scale. The dimensions  
346           included degree of permissibility (impermissible = 0, neither impermissible nor permissible = 50,  
347           permissible =100), degree of goodness (bad = 0, neither good nor bad = 50, good = 100), and  
348           degree of obligation (optional=0, neither optional nor obligatory = 50, obligatory = 100). The  
349           dimension order was randomized across trials. The scale appeared only with the anchors, the  
350           slider was always placed in the center of the scale, and participants were not given a numerical  
351           representation of where they placed the slider.

352           Participants viewed one vignette at a time and the sliders were placed on the same page as  
353           the vignette. After completing the study, participants indicated their age, gender, and political  
354           orientation (on a five-point Likert scale from *very liberal* to *very conservative*). We also  
355           collected information about personal sense of duty. Participants indicated agreement with a 7-  
356           point Likert scale (1 = strongly disagree, 7 = strongly agree) with the idea that there are duties to  
357           respect your neighbors, to help anybody who needs help, and to help one's family members.

**358 Results**

359 Pre-registered analyses were integrated into a set of linear random effect models fitted  
360 with the lme4 R library (Bates, Mächler, Bolker, & Walker, 2015; R Core Team, 2019), with  
361 participant as a random intercept.<sup>3</sup> Pairwise comparisons were carried out using the emmeans  
362 package (Lenth, 2020) which allows degrees of freedom to be calculated with the Kenward  
363 Roger method and p values to be adjusted with the Tukey method. Confidence intervals for non-  
364 standardized simple differences are reported for ease of understanding.

365 Results are summarized in Table 1 and Figures 3 and 4. We fitted four nested models to  
366 contrast the interactive effect of the manipulated variables (see models 1 and 2) and to test the  
367 effect of differences between scenarios and endorsement of norms (models 3 and 4). Goodness  
368 of fit indicators (AIC, BIC and deviance) and the chi square test ( $\chi^2(2)=393.2, p<.0.001$ ) favor  
369 selection of model 2. Significant interactions between judgment type and condition suggest  
370 differences both between and within the erogation category. Supererogatory behaviors were  
371 judged as better and more permissible ( $M = 86.4, SD = 18.7, n = 622$  and  $M = 82.0, SD = 20.6,$   
372  $n = 622$ ) than suberogatory behaviors ( $M = 36.6 SD = 23.3, n = 622$   $M = 62.2 SD = 29.7 n =$   
373  $622$ )  $t$ -ratio Good:  $t(3411) = -34.89, p < .001, M_{diff} = -49.85, 95\% CI [-52.2, -47.5]$  and  $t$ -ratio  
374 Permissible:  $t(3411) = -19.7, p < .001, M_{diff} = -19.78, 95\% CI [-22.6, -17.0]$ , regardless of whether

---

<sup>3</sup> It can be argued that the data of this experiment could be considered a cross-classified data set. However, the items in each situation type are not completely equivalent, which makes the corresponding items nested within situation type. Additional models were fitted with additional random term but since results are equivalent, we restrict the presentation here to the different conditions as fixed effects.

375 they were actions or omissions. Supererogatory behaviors were judged to be marginally more  
376 obligatory (i.e., less optional) ( $M = 30.9$ ,  $SD = 32.7$ ,  $n = 622$ ) than the corresponding  
377 suberogatory responses ( $M = 23.5$ ,  $SD = 26.3$ ,  $n = 622$ )  $t$ -ratio:  $t(3411) = -7.47$ ,  $p < .001$ ,  $M_{diff} = -$   
378  $7.47$ , 95% CI  $[-10.8, -4.19]$  though it is clearly a smaller effect.

379         Interestingly, the size of the differences between permissibility and goodness ratings are  
380 vastly different when looking at sub and supererogatory responses. Within the suberogatory  
381 responses, this difference amounts to 25.6 points (Cohen's  $d = 0.96$ ) while for the supererogatory  
382 category, this difference is only 4.46 points (Cohen's  $d = 0.22$ ). That is, goodness and  
383 permissibility judgements are very similar within the supererogatory condition, but not for the  
384 suberogatory responses, where they are more clearly tracking different aspects of the situation.

385

386 Table 1. Summary of Models fitted for Experiment 2

	Model 1	Model 2	Model 3	Model 4
Omission	-2.74*** (-4.65, -0.83)	-6.411*** (-10.500, -2.322)	-7.33*** (-10.41, -4.24)	-7.26*** (-10.36, -4.17)
OBLIGATORY	-34.30*** (-36.41, -32.20)	-17.449*** (-21.389, -13.509)	-38.12*** (-41.07, -35.17)	-38.12*** (-41.07, -35.17)
PERMISSIBLE	10.57*** (8.47, 12.67)	21.758*** (17.818, 25.698)	7.58*** (4.63, 10.53)	7.57*** (4.62, 10.52)
SUPER EROGATORY	25.70*** (23.98, 27.41)	50.771*** (46.830, 54.711)	25.70*** (23.99, 27.41)	25.70*** (23.98, 27.41)
Mowing	-7.48*** (-9.90, -5.05)		-7.48*** (-9.89, -5.06)	-7.47*** (-9.89, -5.05)
Newlyweds	-3.97*** (-6.39, -1.54)		-3.97*** (-6.39, -1.55)	-3.96*** (-6.39, -1.54)
Raffle	-0.15 (-2.58, 2.26)		-0.16 (-2.58, 2.26)	-0.15 (-2.57, 2.26)
DUTY TO FAMILY				-0.21 (-0.99, 0.56)
DUTY TO HELP				-0.01 (-0.79, 0.77)
DUTY TO NEIGHBORS				-0.31 (-1.35, 0.71)
Omission: OBLIGATORY		8.75***	7.72***	7.718***

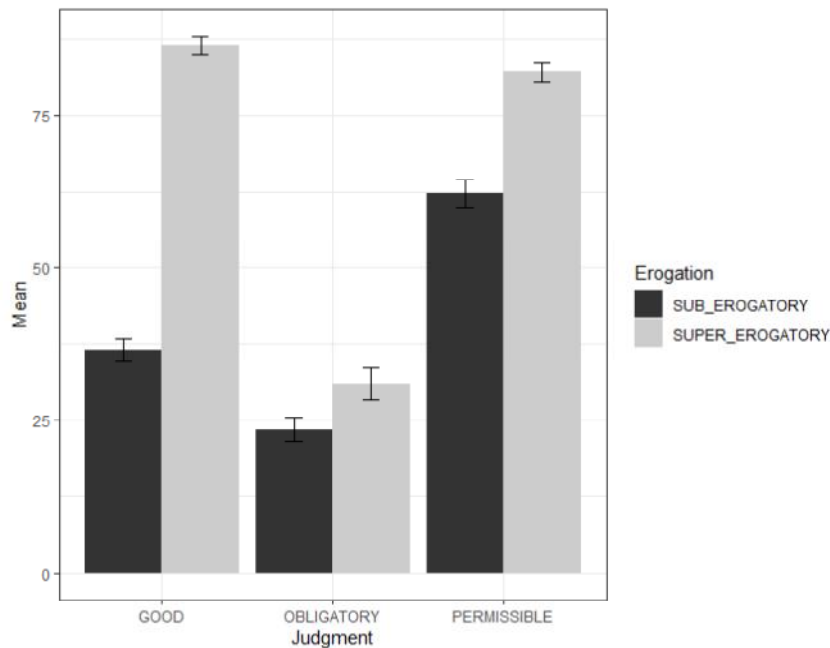
		(3.15, 14.35)	(3.53, 11.91)	(3.526, 11.91)
Omission: PERMISSIBLE		7.75 <sup>***</sup>	6.04 <sup>***</sup>	6.042 <sup>***</sup>
		(2.15, 13.35)	(1.85, 10.23)	(1.85, 10.23)
Omission: SUPER EROGATORY		-1.84		
		(-7.44, 3.76)		
OBLIGATORY: SUPER EROGATORY		-41.34 <sup>***</sup>		
		(-46.92, -35.77)		
PERMISSIBLE: SUPER EROGATORY		-28.36 <sup>***</sup>		
		(-33.93, -22.78)		
Omission: OBLIGATORY: SUPER EROGATORY		-2.07		
		(-9.99, 5.84)		
Omission: PERMISSIBLE: SUPER EROGATORY		-3.419		
		(-11.34, 4.50)		
Constant	52.90 <sup>***</sup>	39.74 <sup>***</sup>	55.17 <sup>***</sup>	58.14 <sup>***</sup>
	(50.41, 55.40)	(36.86, 42.62)	(52.41, 57.94)	(51.41, 64.87)
N	3732	3732	3732	3732
Log Likelihood	-17580.00	-17367.00	-17570.00	-17569.00
AIC	35180.00	34763.00	35163.00	35168.00
BIC	35242.00	34850.00	35238.00	35261.00

---

\*\*\* p < .01; \*\* p < .05; \* p < .1



388 Our results do, however, show that people strongly distinguish between goodness,  
 389 permissibility, and obligation. Participants rated supererogatory behaviors as good, optional, and  
 390 permissible; they rated suberogatory behaviors as bad, optional, and permissible. Despite both  
 391 being rated permissible, supererogatory behaviors were rated as more permissible than  
 392 suberogatory behaviors.



393

394 *Figure 3. Scores by type of judgment.* Error bars represent 95% confidence intervals.

395 We found a smaller interaction between judgment type and the action/omission dimension  
 396 across erogation conditions. Actions were rated as better ( $M = 65.1$ ,  $SD = 32.7$ ,  $n = 628$ ) than  
 397 omissions ( $M = 57.8$ ,  $SD = 32.3$ ,  $n = 616$ )  $t$ -ratio:  $t(1539) = 4.82$ ,  $p = .001$ ,  $M_{diff} = 7.33$ , 95% CI  
 398 [3.71, 10.94], but there were no differences between permissibility and obligatory ratings for  
 399 actions and omissions. The three-way interaction is not explored here but overall, there are no

400 important differences between erogation categories across the action/omission condition except  
 401 for the goodness judgment, where actions are judged as better than omissions.

402 We explored the association between ratings for each type of judgment and participant’s  
 403 endorsement of statements about personal duties. We did not find any consistent association  
 404 between judgment scores and responses pertaining to personal duties (see Table 2). There are  
 405 significant correlations between endorsing different statements of personal duty (ranging from  
 406 0.13 to 0.34) but most correlations between sense of personal duty and different judgment  
 407 categories were negligible (from -.01 to 0.07) and non-significant. The one exception is that  
 408 beliefs about duties to help others significantly correlated with judgments of permissibility, but  
 409 the correlation is very small.

410 Table 2. Bivariate correlations between personal norms and judgment scores

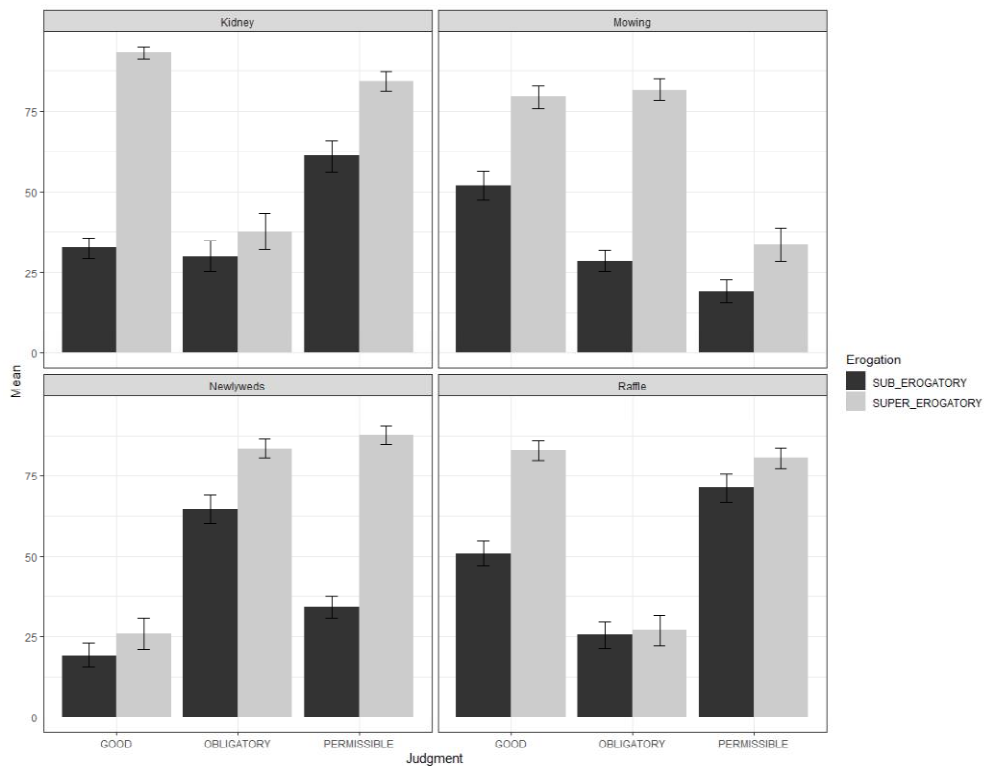
	Permissible	Good	Obligatory	Duty to neighbors	Duty to help
Permissible					
Good	0.55***				
Obligatory	-0.24**	-0.22***			
Duty to neighbors	-0.03	0.0	-0.02		
Duty to help	-0.07**	-0.02	0.07	0.33***	
Duty to family	-0.03	-0.02	0.01	0.30***	0.23***

\*\*\* p < .01; \*\* p < .05; \* p < .1

411 Breaking down responses by scenario reveals some variability across vignettes (see Figure  
 412 4). For example, while donating a kidney to a cousin is judged to be better and more permissible  
 413 than not donating a kidney, the same pattern does not hold in the raffle scenario. In this case, both  
 414 situations are equally permissible, but donating the raffle prize is better than not. Notice however

415 that these small differences do not account for significant variance, according to the model fitting  
 416 presented in Table 1.

417



418  
 419 *Figure 4. Scores by type of judgment and scenario. Error bars represent 95% confidence intervals*

420 **Discussion**

421 These results reveal quantifiable differences between different evaluative categories that are  
 422 employed in moral judgment. Similar to our results in Experiment 1, we found that judgments of  
 423 permissibility, obligation, and goodness dissociate when people make judgments of suberogatory  
 424 behavior. This shows that including additional measures allows variability in moral judgment to  
 425 emerge. Moreover, beliefs about personal duties were negligibly correlated with different kinds  
 426 of judgments. This suggests that people’s responses to these situations are not indicative of an



448 reactions (Christensen et al., 2014) as measured by arousal and valence. Third, dilemmas are  
449 situations where explicit rules come into conflict (e.g., prohibitions against causing harm and  
450 prohibitions against allowing easily preventable harm).

451         We chose sacrificial dilemmas that have been widely used in empirical studies of moral  
452 judgment. Based on an analysis from Christensen et al. (2014), we selected dilemmas that  
453 produce the greatest variation in responses (varying in the use of personal force and the  
454 inevitability of harm). We also excluded dilemmas where various interactive effects might  
455 plausibly drive moral judgments (e.g., when causing harm is self-beneficial). Characteristics and  
456 full text of the dilemmas selected are presented in the supplemental materials.

457         Sacrificial dilemmas are sometimes thought to bring out the contrasts between  
458 deontological and utilitarian ethical intuitions, because each set of intuitions usually recommends  
459 different behaviors in the face of a sacrificial dilemma. Hence, intuitions about what is right or  
460 appropriate indicate alignment with one or the other theory. This suggests that people committed  
461 (implicitly or explicitly) to either utilitarianism or deontology might produce different moral  
462 judgments in reaction to sacrificial dilemmas. This presents a challenge. Because we are  
463 measuring for dissociations among different judgment categories, and different normative  
464 theories make different recommendations for navigating dilemmas, it is possible that each side  
465 will cancel the other out, thereby giving the appearance of similarity between judgment  
466 categories. To ensure that people with utilitarian tendencies do not cancel out people with  
467 deontological tendencies, we used the Oxford Utilitarianism Scale (Kahane et al., 2018) to assess  
468 trait utilitarianism across two dimensions of utilitarianism, impartial beneficence and  
469 instrumental harm.

470           The character of this study is primarily exploratory but we expected to replicate some of  
471 the main findings in this literature and some of the patterns observed by Christensen et al. (2014).  
472 We expected utilitarian responses in personal dilemmas to be associated with lower ratings of  
473 goodness and permissibility relative to impersonal dilemmas (Greene et al. (2001b). We also  
474 expected behaviors that cause unavoidable deaths to be judged as better than situations where  
475 causing death was avoidable (Christensen et al., 2014).

## 476 **Method**

### 477 **Participants**

478           Given that our objective was to first to replicate the observed effects with sacrificial  
479 dilemmas and, second, to explore the impact of these variables on our new measures, we ran  
480 power analysis for a linear mixed model treating participants as random effects, using the effect  
481 size reported by Christensen et al (2014) for personal force of  $r = .75$  (equivalent to a  $d = 2.2$ ).  
482 This is the smallest of the effect sizes considered. This analysis suggested a sample of 60  
483 participants for our mixed design for a power of .99 (Westfall, 2015). Since we wanted to make  
484 sure we would be able to observe differences in our new dependent variables, we aimed to collect  
485 data for 300 participants which would allow us to observe a significantly smaller effect ( $r = .17$ ,  $d$   
486  $= .3$ ). 304 participants (153 women and 151 men, mean age = 33.21, SD =11.62), recruited  
487 through Prolific Academic, took part in the study in exchange for 0.50 pence. The average  
488 completion time was of 7.21 minutes and none of the participants had taken part in our previous  
489 experiments.

## 490 **Materials and procedure**

491 Participants saw the four dilemmas and were randomly presented with the version where  
492 the main character decides to cause harm (utilitarian response) or allow harm (deontological  
493 response). That is, utilitarian/deontological response was a between-subject factors while  
494 avoidability of the result (avoidable/unavoidable) and personal force (personal/impersonal) were  
495 within subject variables. For each scenario, participants judged the main character's response  
496 along dimensions of permissibility, obligation, goodness, and blameworthiness using the same  
497 sliders as before. After judging the dilemmas, participants were presented with the items of the  
498 *Oxford Utilitarianism Scale* (OUS) in a random order. Finally, they are asked some basic  
499 demographic questions and were asked to rate their experience with this kind of dilemmas in a  
500 scale from 1 (not at all familiar) to 5 (extremely familiar).

## 501 **Results**

502 Results are summarized in Table 3 and Figures 5 and 6. As before, we fitted a set of  
503 random effect models with judgment type, personal/impersonal utilitarian/deontological,  
504 avoidable/unavoidable, having experience with dilemmas and the score in the OUS as fixed  
505 predictors allowing for a random effect of participant. Pairwise comparisons were calculated with  
506 the same parameters specified for Experiment 2. The OUS had good reliability for its both  
507 subscales (Instrumental Harm Cronbachs'  $\alpha = .65$  and Impartial Beneficence Cronbachs'  $\alpha =$   
508  $.72$ ). The models differed in the specification of the interaction of the fixed effects. The best  
509 fitting model has interactive effects for two-way interactions for all terms with the  
510 personal/impersonal dimension. This model is presented in Table 3.

511 Table 3. Best fitting model for Experiment 3.

	score
BLAME	-5.15** (-9.53, -0.77)
OBLIGATORY	-2.98 (-7.36, 1.40)
PERMISSIBLE	12.78*** (8.40, 17.16)
PERSONAL	9.22*** (5.40, 13.04)
UTILITARIAN	2.46 (-1.03, 5.95)
UNAVOIDABLE	2.90* (-0.55, 6.35)
BLAME: PERSONAL	4.49** (0.13, 8.85)
OBLIGATORY: PERSONAL	-4.91** (-9.27, -0.55)
PERMISSIBLE: PERSONAL	-4.95** (-9.31, -0.59)
PERSONAL: UTILITARIAN	-14.00*** (-17.19, -10.82)
BLAME: UNAVOIDABLE	0.23 (-4.13, 4.59)
OBLIGATORY: UNAVOIDABLE	-1.82 (-6.18, 2.54)
PERMISSIBLE: UNAVOIDABLE	-4.94** (-9.30, -0.58)
PERSONAL: UNAVOIDABLE	-4.89*** (-7.98, -1.81)
BLAME: UTILITARIAN	6.82*** (2.46, 11.19)
OBLIGATORY: UTILITARIAN	1.50 (-2.86, 5.86)
PERMISSIBLE: UTILITARIAN	-2.37 (-6.74, 1.99)
Constant	47.87***



	(44.53, 51.21)
N	4864
Log Likelihood	-23017.00
AIC	46074.00
BIC	46204.00

\*\*\*p < .01; \*\*p < .05; \*p < .1

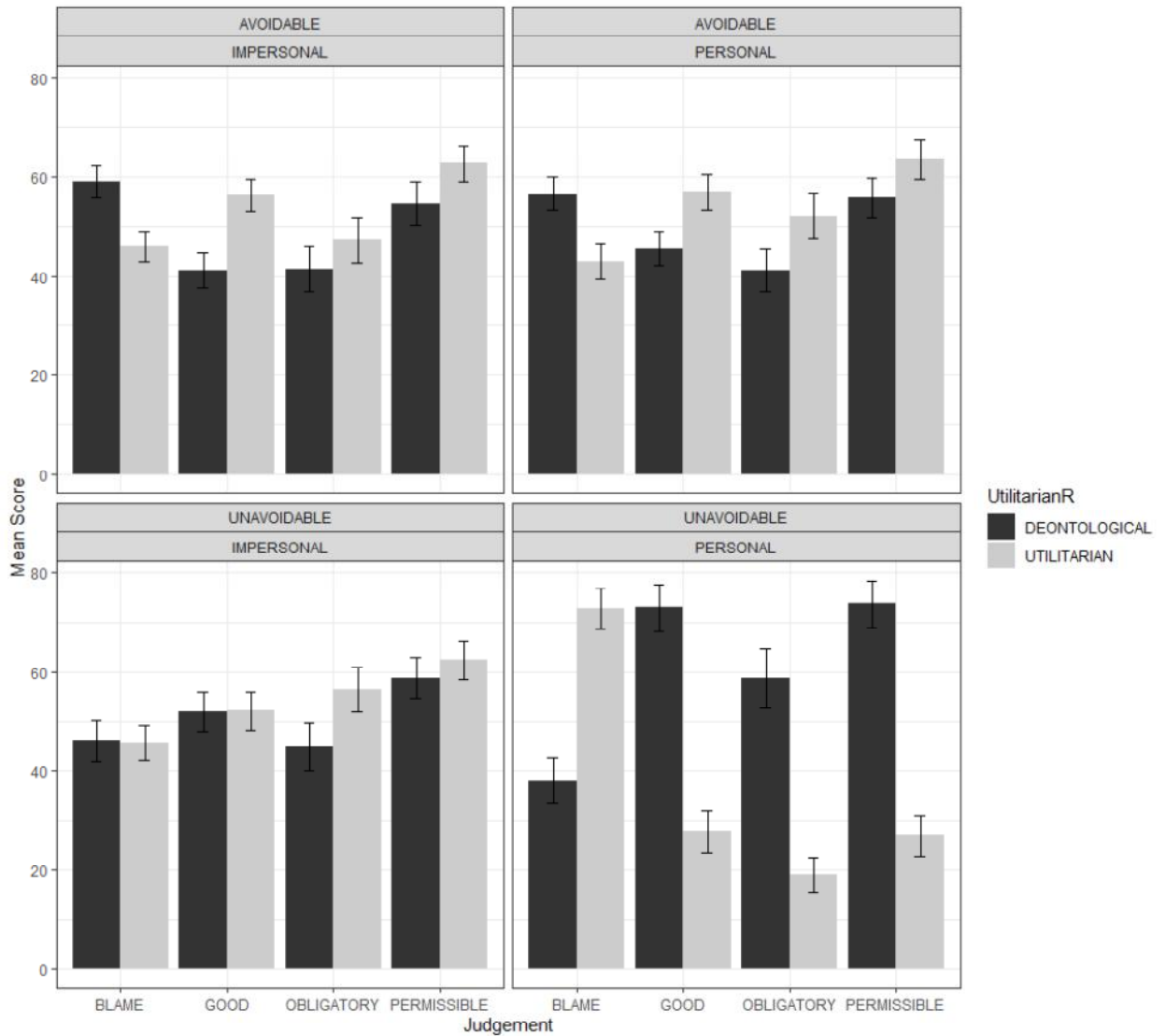
512 There are significant main effects of judgment type and the utilitarian status of the response  
 513 judged (Deontological,  $M = 52.39$ ,  $SD = 28.3$  and Utilitarian  $M = 49.18$ ,  $SD = 28.4$ , ,  $t(3924)$   
 514  $= 3.77$ ,  $p < .001$ ,  $M_{diff} = 3.21$ , 95% CI [1.61, 4.80] ), but not of its being personal (Personal  $M =$   
 515  $51.64$ ,  $SD = 25.8$ , and Impersonal  $M = 49.87$ ,  $SD = 30.8$ ,  $n = 2432$ ,  $t(3411) = 1.99$ ,  $p = .046$ ,  $M_{diff}$   
 516  $= 1.77$ , 95% CI [0.17, 3.37] ) or avoidable (Avoidable  $M = 51.46$ ,  $SD = 25.01$ , and unavoidable  $M$   
 517  $= 50.04$ ,  $SD = 31.4$ ,  $n = 2432$ ,  $t(3411) = 1.49$ ,  $p = .13$ ,  $M_{diff} = 1.42$ , 95% CI [-0.17, 3.02]). These  
 518 variables only have interactive effects, which shows that there is not an overall effect of the  
 519 manipulation across all judgment dimensions.

520 Figure 5 shows how personal deontological responses are judged as better ( $M = 61.75$ ,  $SD$   
 521  $= 28.89$ ,  $n = 303$ ), considerably more permissible ( $M = 65.76$ ,  $SD = 28.56$ ,  $n = 303$ ), more  
 522 obligatory ( $M = 51.38$ ,  $SD = 33.94$ ,  $n = 303$ ) and less blameworthy ( $M = 42.31$ ,  $SD = 27.59$ ,  $n =$   
 523  $303$ ) than the corresponding utilitarian responses (Good  $M = 39.19$ ,  $SD = 28.12$ ,  $n = 305$ ,  
 524 Permissible  $M = 43.54$ ,  $SD = 31.06$ ,  $n = 305$ , Obligatory  $M = 36.53$ ,  $SD = 31.23$ ,  $n = 305$ , and  
 525 blame  $M = 60.01$ ,  $SD = 27.22$ ,  $n = 305$ ) but only for the unavoidable outcomes (all pairwise  
 526 comparison significant at 0.01). When the outcome is avoidable, the pattern is the opposite for  
 527 both impersonal and personal responses. The main effect of personal contact consists in making  
 528 the deontological response more acceptable (less blameworthy, more obligatory, better and more  
 529 permissible) than the corresponding utilitarian ratings (See lower panel, Figure 5). Overall this

530 picture is consistent with prior findings where the effect of these variables on dilemma responses

531 is conditional on several factors (See Christensen et al. (2014)).

532



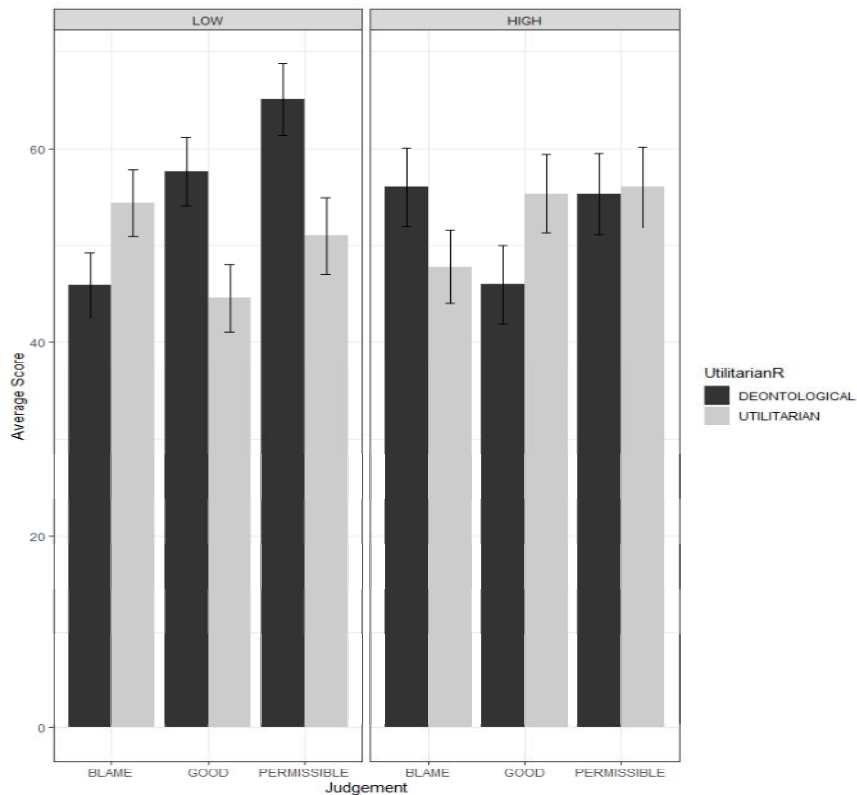
533

534 *Figure 5. Mean ratings for each judgment type by utilitarian status, personal force and*

535 *avoidable. Error bars represent 95% confidence intervals.*

536           In the model fitting process, neither experience with dilemmas nor the overall score of the  
537 OUS (or its subscales) resulted in significant predictors of participant judgments. Moreover, for  
538 all but one dilemma type (Unavoidable, Personal) we did not find significant differences between  
539 judgments across trait utilitarian tendencies.

540           We wanted to see, however, whether more fine-grained distinctions in trait utilitarianism  
541 might allow unnoticed dissociations to emerge. To do this, we assigned participants to a low,  
542 medium and high utilitarianism group, by splitting participants aggregate scores on the OUS by  
543 the 33rd and 66th percentile. Zooming in on participants with the lowest and highest utilitarian  
544 scores (See Figure 6), we see that participants in the Low Utilitarianism Group (left panel) clearly  
545 judge causing harm as less permissible, less optional, and more blameworthy than allowing harm.  
546 People in the High Utilitarianism Group (right panel) show an opposite pattern, albeit not as  
547 distinctive (See Supplemental materials for more information on these comparisons). Note,  
548 however, that these patterns only emerge after making exploratory, post-hoc data analyses. Even  
549 doing this, which might otherwise be considered methodologically problematic, we weren't able  
550 to see the clearly marked dissociations observed in Experiment 2.



551

552 *Figure 6. Mean ratings by of judgment, grouped by Utilitarianism level (panels) and*  
 553 *Utilitarian response (bars)*

554 Permissibility judgments tend to be higher than goodness judgments, suggesting that even  
 555 though they are related, their difference in magnitude suggests that these judgments are tracking  
 556 different features of the actions depicted in the scenario. This difference is larger and statistically  
 557 significant for participants with low utilitarian tendencies, who consider harm relatively more  
 558 permissible ( $M = 50.95$ ,  $SD = 30.51$ ) than it is good ( $M = 44.55$ ,  $SD = 27.23$ ,  $t(448) = -2.36$ ,  $p = .02$ ,  $M_{diff} = -6.40$ , 95% CI [-11.72, -1.08]). With participants who score in the higher end of  
 559 utilitarianism, there is virtually no difference between these judgments (permissibility  $M = 56.02$ ,  
 560  $SD = 27.52$  vs obligatoriness  $M = 55.35$ ,  $SD = 26.71$ ,  $t(338) = -0.23$ ,  $p = .81$ ,  $M_{diff} = -0.66$ , 95%  
 561 CI [-5.45, 5.12]). A similar pattern is found with the obligatoriness levels: participants who score  
 562

563 low in utilitarianism consider causing harm optional (as revealed by a one-sample t test against  
564 the indifference point  $M = 37.09$ ,  $SD = 31.02$ ,  $t(227) = -6.79$ ,  $p < .001$ , 95% CI [33.98, 41.10]) and  
565 participants who score high in utilitarianism consider *not* causing harm as optional ( $M = 44.20$ ,  
566  $SD = 29.81$ ,  $t(169) = -0.69$ ,  $p = .48$ , 95% CI [43.70, 73.0]). These tendencies are statistically  
567 significant when explored with a fixed effects model with these terms (see Supplementary  
568 materials). Bear in mind that this model is exploring a tendency rather than testing a hypothesis.  
569 As such, it lacks confirmatory value.

## 570 **Discussion**

571 We wanted to see whether adding additional measures allowed variability in moral judgment to  
572 emerge, even when using dilemmas as stimuli. However, when we had people make judgments  
573 about dilemmas, the variability we observed in Experiment 2 nearly disappears. This is not an  
574 effect driven by Utilitarian members of the participant pool. Even when we separate people  
575 according to their utilitarian tendencies (which tendencies pull *against* making the distinctions  
576 outlined in Experiment 2), we do not find strong dissociations between the categories of  
577 goodness, permissibility, and obligation. The dissociations are weak enough that they might  
578 seemingly justify the claim that “the influence of wording variations on moral judgments was  
579 negligible” (O’Hara et al., 2010).

580 Our results suggest that changes in measurement alone are not sufficient to indicate the  
581 underlying complexity of moral judgment. The kind of situation being evaluated makes a  
582 difference to moral judgment. When we limit ourselves to using only dilemmas, we generate an  
583 overly narrow view of moral thinking. Worse, it hinders the ability to make useful generalizations  
584 about the operations of moral cognition in general.



607           One explanation of this failure might be that people disapprove of suberogatory behaviors,  
608 but the various judgment types (permissibility, obligation, and goodness) do not have any  
609 obvious place for people to register this disapproval. Hence, disapproval might be skewing  
610 judgments in a way that drives the differences between permissibility and obligatoriness across  
611 scenarios. To correct for this, and to replicate the findings of Experiment 2, we ran the  
612 experiment again with an additional judgment type of blameworthiness.

### 613           **Participants**

614           We recruited 311 participants (166 women and 145 men, mean age = 31.95, *SD* =11.00)  
615 through Prolific Academic. Sample size was determined using the rationale of Experiment 2. All  
616 participants voluntarily agreed to participate in the study and were monetarily compensated (38  
617 pence). The average completion time was 4 minutes.

### 618           **Materials and procedure**

619           We used the same materials and procedures from Experiment 2. We included an  
620 additional measure for participants to rate the blameworthiness of the character in the vignette  
621 (praiseworthy = 0, neither praiseworthy nor blameworthy = 50, blameworthy =100).

### 622           **Results**

623           Overall, the results replicated the pattern observed in Experiment 2 (see Figure 7). We  
624 fitted the same random effect models as in Experiment 2 and performed the same pairwise  
625 comparisons. Model 2, which includes the interactive effect of judgment type, omission and  
626 erogation condition, was the best model. Table 4 below presents only the two best fitting models.

627

628 Table 4. Summary of Models fitted for Experiment 4

	<b>Model 1</b>	<b>Model 2</b>
Omission	3.90** (0.19, 7.60)	5.70*** (2.40, 8.90)
GOOD	-19.00*** (-23.0, -16.0)	29.00*** (26.0, 32.0)
OBLIGATORY	-41.00*** (-45.00, -38.00)	-16.00*** (-19.0, -13.0)
PERMISSIBLE	4.20** (0.47, 7.90)	40.0*** (37.00, 43.00)
SUPER EROGATORY	-50.00*** (-53.00, -46.00)	5.20*** (3.60, 6.80)
Mowing		-1.0 (-4.20, 0.460)
Newlyweds		-1.400 (-3.70, 0.92)
Raffle		0.260 (-2.00, 2.60)
Omission:GOOD	-12.00*** (-18.00, -7.20)	-14.00*** (-19.00, -9.30)
Omission:OBLIGATORY	3.70 (-1.50, 9.00)	0.180 (-4.40, 4.80)



Omission:PERMISSIBLE	-7.90*** (-13.00, -2.70)	-10.00*** (-15.0, -5.70)
Omission:SUPER EROGATORY	3.50 (-1.70, 8.80)	
GOOD: SUPER EROGATORY	97.00*** (92.0, 102.0)	
OBLIGATORY: SUPER EROGATORY	51.00*** (46.0, 56.0)	
PERMISSIBLE: SUPER EROGATORY	72.00*** (67.00, 77.00)	
Omission: GOOD: SUPER EROGATORY	-3.00 (-10.00, 4.40)	
Omission:OBLIGATORY: SUPER EROGATORY	-7.10* (-14.0, 0.36)	
Omission:PERMISSIBLE: SUPER EROGATORY	-4.80 (-12.00, 2.60)	
Constant	62.00*** (59.00, 64.00)	35.00*** (32.00, 38.00)
N	4976	4976
Log Likelihood	-22765.00	-23867.00
AIC	45566.00	47761.00
BIC	45683.00	47852.00

---

\*\*\*p < .01; \*\*p < .05; \*p < .1

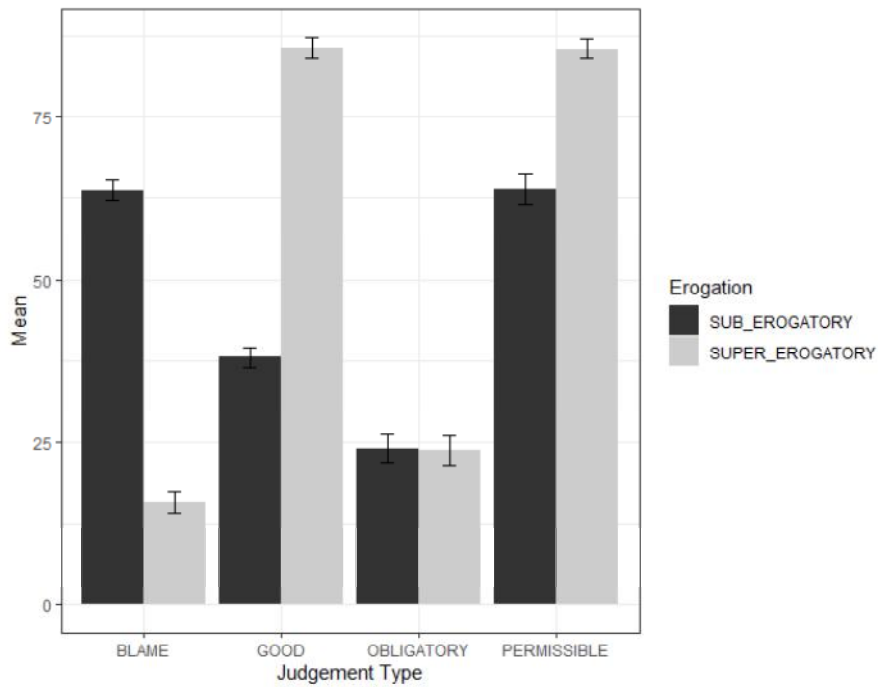
630 As before, both permissibility and goodness ratings were significantly higher for the  
631 supererogatory (Permissibility  $M = 85.4$ ,  $SD = 19.5$ ,  $n = 622$  and goodness  $M = 85.6$ ,  $SD = 20.9$ ,  
632  $n = 622$ ) than for the suberogatory condition ( $M = 63.9$ ,  $SD = 29.9$ ,  $n = 616$  and goodness  $M$   
633  $= 38.0$ ,  $SD = 20.5$ ,  $n = 622$ ) permissibility:  $t(4580) = -16.08$ ,  $p < .001$ ,  $M_{diff} = -21.5$ , 95% CI [-  
634 24.3, -18.7], goodness:  $t(4580) = -35.54$ ,  $p < .001$ ,  $M_{diff} = -47.2$ , 95% CI [-49.9, -45.3]. However,  
635 the difference between obligatory ratings was not replicated ( $t(4580) = 0.22$ ,  $p = .82$ ,  $M_{diff} = 0.32$ ,  
636 95% CI [-2.77, 3.41]). Blame ratings differed significantly across the suberogatory and  
637 supererogatory conditions. Participants rated suberogatory behaviors as more blameworthy ( $M =$   
638  $63.70$ ,  $SD = 20.77$ ,  $n = 622$ ) than supererogatory behaviors ( $M = 15.76$ ,  $SD = 20.62$ ,  $n = 622$ )  
639  $t(4580) = 35.82$ ,  $p < .001$ ,  $M_{diff} = 47.9$ , 95% CI [45.6, 50.2]. As in Experiment 2, there are  
640 significant differences between scenarios, with supererogatory responses in the *Kidney* scenario  
641 having the highest positive scores (as better, more permissible, more praiseworthy) and  
642 suberogatory responses in the *Mowing* scenario the most negative (as worse, less permissible,  
643 less optional, and more blameworthy).

644 Blame ratings were not significantly correlated with statements about personal sense of  
645 duty across either scenarios or conditions. The same overall pattern of correlations between the  
646 other ratings and statements about duties was observed (See Table 2, Experiment 2).

647 Supererogatory behaviors are consistently rated as praiseworthy, while suberogatory  
648 behaviors are consistently rated as blameworthy. However, the magnitude of praiseworthiness  
649 judgments is much greater in the case of supererogatory behaviors. This asymmetry suggests that  
650 supererogatory behaviors deserve more praise than the corresponding suberogatory behaviors

651 deserve blame. (see Figure 7). This asymmetry has been repeatedly observed in other studies  
 652 (Monroe et al., 2018, Pizarro et al., 2007).

653 It is worth noting that, again as in Experiment 2, there are no important differences  
 654 between the permissibility and goodness means in the supererogatory condition ( $M = 85.61$ ,  $SD$   
 655  $= 20.86$  vs  $M = 85.43$ ,  $SD = 19.46$  respectively) while for the suberogatory responses we find a  
 656 difference of 25.88 points between these (Cohen's  $d = 1.01$ ), suggesting that participants are  
 657 considering different information when using assessing permissibility and goodness in these  
 658 situations.



659  
 660 *Figure 7. Scores by type of judgment Experiment 4. Error bars represent 95% confidence*  
 661 intervals

662 **Discussion**

663 In this experiment, we replicated the key findings of Experiment 2. When comparing judgments  
664 about suberogatory and supererogatory behaviors, people distinguish between the goodness,  
665 permissibility, and blameworthiness of the action. However, people do recognize suberogatory  
666 and supererogatory behaviors as equally optional. Despite statistically significant differences  
667 between the permissibility of suberogatory and supererogatory actions, participants rated  
668 suberogatory behaviors as permissible.

669 Now that we have provided some evidence for distinctions between these evaluative  
670 categories, we turn to another question: what factors differentially affect these judgments? As we  
671 saw in Experiment 2, beliefs about personal duties do not relate significantly to people's  
672 judgments. The significant difference in blame judgments suggests that perhaps norms of social  
673 sanctioning help to inform at least one kind of judgment. To understand the reasoning underlying  
674 these judgments, we conducted two additional experiments.

### 675 **Experiment 5**

676 In this experiment we explored how judgments of deserved social sanction are associated  
677 with the different kinds of moral judgment identified in previous experiments. Thus far, we have  
678 established that some situations can be understood in terms of suberogatory or supererogatory  
679 characterizations. In order to understand people's responses to these situations, we need to deploy  
680 a richer array of judgment dimensions. In this experiment, we explored whether suberogatory  
681 responses are associated with potential behavioral consequences, such as social punishment.  
682 Plausibly, social punishment is, as much as expressions of blame, a behavioral marker of  
683 underlying moral judgment.

684 We also explore whether attitudes about rights predict different kinds of judgments. The  
685 situations that we describe often bring into conflict an individual right and a kind or generous  
686 thing that could be done for others. For example, someone has the *right* to keep the seat they paid  
687 for on the airplane, but it is a kind thing to offer one's seat when asked to switch. Because of this,  
688 we thought that attitudes towards *rights*, instead of duties, might usefully predict moral  
689 judgments. This would partially vindicate the rule-based conception of morality, as one might  
690 hold that there are generalized functions from rights to prohibited behaviors (even if there are  
691 fewer—if any—such functions from rights to *prescribed* behaviors). Moreover, we found this  
692 pattern of reasoning in the open responses for Experiment 1. All of this suggests that personal  
693 sense of rights might provide some details about the range of inputs to which moral judgment is  
694 sensitive.

### 695 **Participants**

696 309 participants (165 women and 144 men, mean age = 39.89, SD =13.76), recruited  
697 through Prolific Academic, took part in the study in exchange for 60 pence. The average  
698 completion time was of 4.87 minutes and none of the participants had taken part in the previous  
699 experiments.

### 700 **Materials and procedure**

701 We used the same materials and procedures from Experiment 4 with two changes. First,  
702 people had to judge whether the response merited social sanction or recognition by using a slider  
703 with an underlying scale going from 0 to 100 (social sanction = 0, neither = 50, social  
704 recognition=100). Second, people provided information about their attitudes towards rights by  
705 answering these three questions with a 7-point Likert scale (1= Strongly disagree, 7 = Strongly

706 agree): 1) Everyone has a right to do anything they want on their property (as long as they are not  
 707 hurting anyone else); 2) Everyone has the right to bodily autonomy; 3) Everyone has a right to  
 708 make use of their money (or goods) as they see fit.

## 709 **Results**

710 As in Experiments 2 and 4, we fitted a series of random effect models which resulted in a  
 711 best fitting model that includes the interactive effect of judgment type, omission, erogation  
 712 condition and scenario. Crucially, including terms for attitudes about rights did not significantly  
 713 improve the fit of this model.

Table 5. Summary of Models fitted for Experiment 5

	<b>score</b>
Omission	-8.08*** (-10.95, -5.22)
OBLIGATORY	-22.94*** (-26.74, -19.15)
PERMISSIBLE	19.74*** (15.94, 23.54)
SANCTION	-2.66 (-6.45, 1.14)
SUPER EROGATORY	44.90*** (42.21, 47.59)
Mowing	-5.24*** (-8.25, -2.24)
Newlyweds	-5.14*** (-7.04, -3.23)
Raffle	3.68** (0.68, 6.69)
Omission: OBLIGATORY	15.28*** (9.90, 20.65)

Omission: PERMISSIBLE	9.13*** (3.75, 14.51)
Omission: TypeSANCTION	7.98*** (2.60, 13.35)
OBLIGATORY: SUPER EROGATORY	-39.76*** (-45.13, -34.39)
PERMISSIBLE: SUPER EROGATORY	-22.18*** (-27.55, -16.81)
SANCTION: SUPER EROGATORY	-14.17*** (-19.54, -8.80)
Omission: OBLIGATORY: SUPER EROGATORY	-3.29 (-10.90, 4.31)
Omission: PERMISSIBLE: SUPER EROGATORY	-5.62 (-13.22, 1.99)
Omission: SANCTION: SUPER EROGATORY	-3.65 (-11.26, 3.96)
Constant	45.96*** (43.03, 48.90)
N	4944
Log Likelihood	-22788.0
AIC	45617.0
BIC	45747.0

---

\*\*\*p < .01; \*\*p < .05; \*p < .1

714           The pattern found in Experiments 2 and 4 was replicated here (see Figure 8).  
715 Supererogatory behaviors were judged as better and more permissible ( $M = 85.16$ ,  $SD = 19.31$ ,  $n$   
716  $= 618$  and  $M = 84.47$ ,  $SD = 20.18$ ,  $n = 618$ ) than suberogatory behaviors ( $M = 40.26$   $SD = 23.99$ ,  
717  $n = 618$  and  $M = 64.55$ ,  $SD = 30.09$ ,  $n = 618$ ) Good:  $t(4619) = -32.72$ ,  $p < .001$ ,  $M_{diff} = -44.4$ ,  
718 95% CI [-47.3, -42.5] and Permissible:  $t(4619) = -14.51$ ,  $p < .001$ ,  $M_{diff} = -19.9$ , 95% CI [-22.8, -  
719 17.1], regardless of being actions or omissions. Supererogatory behaviors were again judged to  
720 be less optional ( $M = 28.43$ ,  $SD = 30.61$ ,  $n = 618$ ) than the corresponding suberogatory responses

721 ( $M = 24.93$ ,  $SD = 26.34$ ,  $n = 618$ ),  $t(4619) = -2.55$ ,  $p = .01$ ,  $M_{diff} = -3.50$ , 95% CI [-6.68, -0.31].

722 As in Experiment 4 this difference only emerged when an additional measure was included

723 (blameworthiness in Exp 4 and social sanction in this Exp) suggesting the presence of a joint

724 evaluation effect (Hsee, Blount, Loewenstein, & Bazerman, 1999).

725 Supererogatory behaviors were judged to merit more social recognition than suberogatory

726 responses ( $M = 70.49$ ,  $SD = 25.41$ ,  $n = 618$  vs  $M = 41.59$ ,  $SD = 20.57$ ,  $n = 618$ ,  $t(4619) = -21.07$ ,

727  $p < .001$ ,  $M_{diff} = -28.9$ , 95% CI [-31.5, -26.3]). Crucially, both means for social recognition and

728 social sanction are significantly different from the indifference point (50 in our scale)

729 (Supererogatory  $t(617) = 20.1$ ,  $p < .001$ ,  $M = 70.5$ , 95% CI [68.5, 72.5] and suberogatory  $t(617)$

730  $= -10.2$ ,  $p < .001$ ,  $M = 41.5$ , 95% CI [40.0, 43.2]. Mirroring the asymmetric pattern observed for

731 blameworthiness ratings, supererogatory responses deserve more recognition than sanction is

732 deserved by suberogatory responses. Also, as in the previous studies there are important

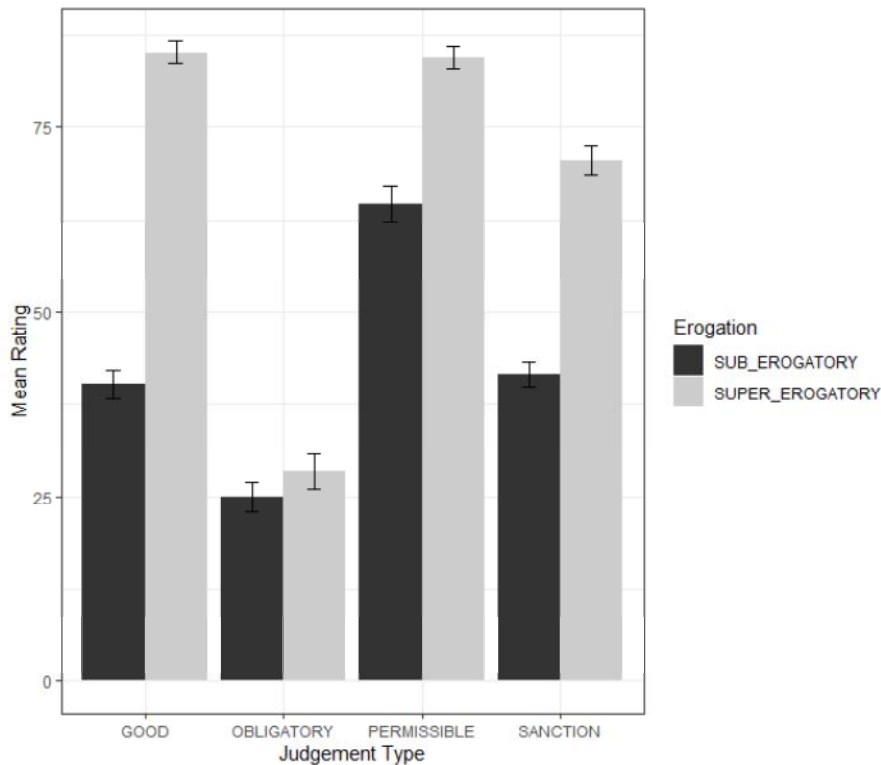
733 differences between scenarios, even when the overall pattern is consistent. For example, not

734 giving up the raffle prize does not merit neither recognition nor sanction ( $M = 50.82$ ,  $SD = 22.06$ ,

735  $n = 155$ ) while mowing your lawn early morning is clearly disapproved ( $M = 34.23$ ,  $SD = 19.76$ ,

736  $n = 155$ ).





737

738 *Figure 8. Scores by type of judgment and condition.* Error bars represent 95% confidence  
739 interval.

740 To explore the association between judgment scores and attitudes towards rights we first  
741 observed the strength of the association between attitudes towards rights. Correlations are  
742 medium to small in size (see Table 6) and there is low internal consistency (Cronbach's alpha =  
743 0.50). We then calculated the correlations between judgment scores and these attitudes towards  
744 rights (Table 6). Similar to what happened with questions about duties, the correlations are low  
745 and only significant for two associations: (1) bodily autonomy and permissibility, and; (2)  
746 obligatoriness and right to property. Only when considering particular scenarios, moderately  
747 stronger correlations start to emerge (for example, permissibility ratings in the *Kidney* scenario  
748 are significantly correlated with attitudes on bodily autonomy  $r(309) = .12, p < 0.01$ )

749 Table 6. Bivariate correlations between attitudes towards rights and judgment scores

	Permissible	Good	Obligatory	Sanction	Right to Money	Bodily autonomy
Permissible						
Good	0.54***					
Obligatory	-0.24**	-0.04***				
Sanction	0.24***	0.40***	0.14***			
Right to money	0.03	0.04	-0.05	0.01		
Bodily autonomy	-0.16**	0.05	-0.14	0.00	0.16***	
Right to property	0.05	0.02	0.06*	0.01	0.44***	0.17***

\*\* p < .001; \* p < .01; \* p < .05

750 **Discussion**

751 We find the same dissociations observed in Experiments 2 and 4. Including the category of  
 752 sanctions, however, brings out something interesting. There is a much closer relationship  
 753 between goodness and sanction than there is between goodness and blame in suberogatory cases.  
 754 The relationship does not hold in cases of supererogatory behavior, where judgments of goodness  
 755 were statistically distinct from judgments of sanction (positive recognition, in this case). This  
 756 suggests that people think suberogatory actors should be sanctioned, though supererogatory  
 757 actors need not necessarily be positively recognized. This, however, might be a function of the  
 758 cases we used rather than marking out an intrinsic difference between commonsense thinking  
 759 about the sub- and supererogatory.

760 We also identified an interesting relationship between sanction and goodness. Actions  
 761 rated as bad received similar ratings of sanction. In the previous experiment, we saw that ratings  
 762 of badness do not align with ratings of permissibility. Permissibility ratings are more closely  
 763 aligned with ratings of blame. This suggests that the permissibility of suberogatory behaviors is  
 764 related to blameworthiness, whereas the badness of suberogatory behaviors is related to public



786 do. However, it also explains why the two concepts dissociate in cases of suberogatory behavior.  
787 Suberogatory situations arise when someone has a right to do something that a virtuous or decent  
788 person would not do. If this is the case, it is the absence of duties and the corresponding presence  
789 of rights that are characteristic of suberogatory situations. These characteristics are, however, not  
790 framed as rules, but are instead low-level features of the situation.

791 In this experiment, we explore the justifications people offer for their judgments of  
792 permissibility and goodness to see whether these different factors explain the distinction between  
793 these two kinds of judgment. We also include scenarios that are more perspicuously *moral* than  
794 the ones previously used (e.g. lawn mowing or seat ownership). By providing situations with a  
795 higher moral charge, we can be more confident that the responses observed so far are not merely  
796 tracking the perceived conventional permissibility or social aptness of the behaviors evaluated.

## 797 **Method**

### 798 **Participants**

799 We recruited 316 participants (160 women and 153 men, mean age = 33.30, SD =10.81)  
800 through Prolific Academic to participate in the study in exchange for 60 pence. Sample size was  
801 set to reproduce findings of Experiments 2, 4 and 5 but with a within participants design in this  
802 case, which would allow us to increase statistical power. The average completion time was 10  
803 minutes and none of the participants had taken part in the previous experiments.

**804 Materials and procedure**

805 We constructed two new scenarios and created a suberogatory and supererogatory version  
806 of each. The two scenarios (adapted from Thomson, 1971 and Nozick, 1974, respectively) are  
807 described below [suberogatory version in brackets]:

**808** **Well** Alex is driving home from work on the highway when she gets into an accident that knocks her  
809 unconscious. When she wakes up, she finds herself in a hospital bed. She's also connected to another individual  
810 through a series of wires and tubes. A doctor enters the room and explains to Alex that she is fine, but the  
811 individual she's connected to suffered some severe damage to internal organs. Alex has the right blood type to  
812 help, and—since she was unconscious—the doctor decided to connect Alex to keep the other individual alive for  
813 the time being. The doctor explains that Alex can unplug herself if she chooses, but the individual will most likely  
814 die. The individual will recover from these injuries in about a month (give or take a few days), after which time Alex  
815 can unplug herself and leave. After a few hours of pondering what to do, Alex decides stay plugged in for the  
816 month [*to unplug herself*].

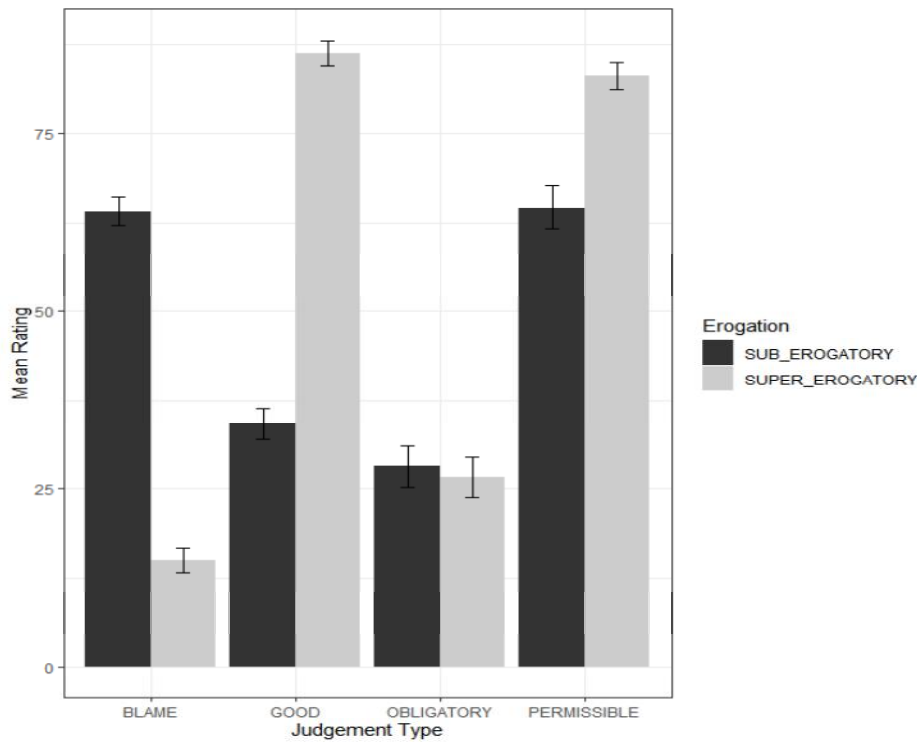
**817** **Well** Jones finds a large freshwater source on his property, so he digs a well as a way of claiming the  
818 water. A few weeks later, the town where he lives begins experiencing a drought, which was completely  
819 unpredictable. Town representatives visit Jones to ask whether they can use his water to alleviate some of the  
820 drought. Without Jones' help, the town will likely run out of water in a few days. If Jones donates some of his  
821 water, however, he might experience the effects of the drought in the unlikely event that the drought prolongs for  
822 too long. After considering what to do, Jones decides to offer his water [*declines to offer his water*]

823 To test variation against a known benchmark, we also included the *Newlyweds* scenario  
824 (See Exp. 5). Participants read the informed consent and then rated each one of the three  
825 vignettes. The order of the vignettes was randomized across participants. Each participant was  
826 randomly assigned to see either the supererogatory or suberogatory condition of each scenario  
827 and the rest of the procedure was the same as in Experiment 4. The only difference is that, after

828 providing judgment for the dimensions requested, participants were asked to explain their ratings  
829 in their own words (cf. Christensen et al., 2014).

### 830 **Results**

831 Overall, numerical ratings closely followed the pattern observed in prior Experiments (See  
832 Figure 9). Suberogatory behaviors are rated as worse ( $M = 34.19$ ,  $SD = 23.71$ ,  $n = 468$  vs  $M =$   
833  $86.37$ ,  $SD = 19.80$ ,  $n = 480$ ), more blameworthy ( $M = 64.13$ ,  $SD = 22.46$ ,  $n = 468$  vs  $M = 14.92$ ,  
834  $SD = 19.12$ ,  $n = 480$ ), and less permissible ( $M = 28.17$ ,  $SD = 32.31$ ,  $n = 468$  vs  $M = 83.15$ ,  $SD =$   
835  $22.27$ ,  $n = 480$ ) than supererogatory responses (all significant pairwise comparisons). Unlike  
836 Experiment 5, there is no significant difference in obligatoriness judged (Suberogatory  $M =$   
837  $28.17$ ,  $SD = 32.31$ ,  $n = 468$  vs  $M = 26.64$ ,  $SD = 31.75$ ,  $n = 480$ ). As before, we fitted random  
838 effect models, with a random intercept for participants. However, models with this term resulted  
839 in singular fits, due to lack of variation for the random intercept for participant and suggesting an  
840 overcomplex random structure (Matuschek, Kliegl, Vasishth, Baayen & Bates (2017). Therefore,  
841 we fitted only fixed effect models, where the best model includes the interaction of judgment  
842 type, scenario and erogation condition (results can be seen in Table 7).



843

844 *Figure 9. Mean scores by type of judgment Experiment 4. Error bars represent 95% confidence*  
 845 *intervals*

846 There is also a significant between-scenario variation. For example, only in the *Well*  
 847 scenario people judged the suberogatory behavior, not giving water, as significantly less optional  
 848 than the supererogatory behavior ( $M = 31.26$ ,  $SD = 31.75$ ,  $n = 160$  vs  $M = 38.81$ ,  $SD = 34.63$ ,  $n =$   
 849  $156$ , Good:  $t(310)^4 = 2.64$ ,  $p < .001$ ,  $M_{diff} = 7.55$ , 95% CI [0.19, 14.9]). In this scenario also  
 850 occurred the most extreme values for suberogatory responses.

---

<sup>4</sup> Degrees of freedom correspond to the Welch t-test, since the library *emmeans* (Lenth, 2020) calculated the asymptotic result for this comparison (value from the  $z$  distribution).

851

Table 7. Summary of Best Model fitted for Experiment 6

	score
GOOD	-18.65*** (-24.33, -12.97)
OBLIGATORY	-43.06*** (-48.74, -37.38)
PERMISSIBLE	14.01*** (8.33, 19.70)
SUPER_EROGATORY	-35.49*** (-41.08, -29.89)
Violinist	3.61 (-2.02, 9.24)
Well	13.13*** (7.48, 18.79)
GOOD: SUPER_EROGATORY	77.22*** (69.31, 85.13)
OBLIGATORY: SUPER_EROGATORY	36.02*** (28.11, 43.93)
PERMISSIBLE: SUPER_EROGATORY	50.07*** (42.16, 57.98)
GOOD: Violinist	-7.79* (-15.75, 0.17)
OBLIGATORY: Violinist	10.88*** (2.92, 18.84)
PERMISSIBLE: Violinist	-10.41** (-18.37, -2.45)
GOOD: Well	-25.91*** (-33.91, -17.91)
OBLIGATORY: Well	10.22** (2.22, 18.22)
PERMISSIBLE: Well	-29.89*** (-37.89, -21.89)
SUPER_EROGATORY: Violinist	-17.16*** (-25.07, -9.25)



SUPER_EROGATORY: Well	-24.19*** (-32.10, -16.28)
GOOD: SUPER_EROGATORY: Violinist	27.74*** (16.56, 38.93)
OBLIGATORY: SUPER_EROGATORY: Violinist	19.66*** (8.47, 30.84)
PERMISSIBLE: SUPER_EROGATORY: Violinist	12.94** (1.76, 24.13)
GOOD: SUPER_EROGATORY: Well	44.97*** (33.79, 56.16)
OBLIGATORY: SUPER_EROGATORY: Well	16.10*** (4.91, 27.29)
PERMISSIBLE: SUPER_EROGATORY: Well	39.84*** (28.66, 51.03)
Constant	58.52*** (54.51, 62.54)
N	3792
Log Likelihood	-17629.00
AIC	35306.00

---

\*\*\*p < .01; \*\*p < .05; \*p < .1

852

853

### *Qualitative Data*

855 For the qualitative analysis of the open responses, we used two coders. Both were blind to  
 856 the individual ratings associated with each open response and one coder was completely blind to  
 857 the objective of the study. Coders used ten predefined categories to sort responses. We report  
 858 agreement between coders (Cohens Kappa) and results from fitting loglinear models on the  
 859 classification frequencies.

860 There were 4003 unique written explanations of ratings (1004 for goodness ratings, 999  
 861 for obligatory ratings, 997 for permissibility ratings and 1003 for blame/praise ratings). Coders

862 were given ten categories. *Character* when the response included some mention of personality or  
863 character traits; *Rights (presence or absence)* when responses appealed to what the protagonist  
864 can morally do, is allowed to do, or its explicit negation; *Duties (presence or absence)* when  
865 responses appeal to what the protagonist is morally obliged to do, or its explicit negation; *Neither*  
866 when participants appealed to norms that could not be considered rights nor duties, etc; *Outcomes*  
867 when the justification was based on the consequences for the “victims”, negative or positive;  
868 *Values* actions performed had a clear valence; *Justification* when appealing to reasons not based  
869 on consequences; *Other* if the response did not fall in the previous categories but was common  
870 enough to merit its own. Coders were given an example of each one of the categories. Raters  
871 completed coding independently. They were also told that they could classify any given response  
872 in more than one category.

873 For example, one response to the *Well* vignette reads: “He didn’t have to donate the water,  
874 and even if he didn’t that wouldn’t be necessarily a bad thing.” This was rated as indicating  
875 ‘Absence of Duties’ and ‘Values’. Another response (to the *Newlyweds* vignette) reads: “No one  
876 should have to switch a seat unless they want to. They bought the seat they were in and would  
877 have been justified in staying in that seat.” This was rated as indicating ‘Justification’ and  
878 ‘Presence of Rights’.

879 Raw agreement between coders and Cohen’s Kappa for each category are presented in  
880 Table 8. The category *Neither* is not presented because there was virtually no agreement in its  
881 use. Most of the unique responses were classified into at least 2 categories (58%). Half of the  
882 categories show a weak reliability between coders (outcomes, values, justification and absence of

883 rights) while the other half shows an overall good agreement (presence of duties, character,  
884 presence of rights and absence of duties).

885 Table 8. Cohen's kappa agreement between coders for the eight categories. Columns L and  
886 U are the Lower and Upper limits of Cohen's confidence interval.

	Kappa	L	U	% Classified as	% Agreement
Abs duties	0.81	0.79	0.84	11.02	95.68
Pres rights	0.75	0.73	0.77	21.06	90.50
Character	0.69	0.65	0.72	8.61	94.18
Pres duties	0.63	0.58	0.68	5.65	95.25
Outcomes	0.39	0.35	0.43	11.86	85.08
Values	0.38	0.35	0.42	20.28	77.18
Justification	0.37	0.34	0.41	13.7	82.35
Abs rights	0.34	0.27	0.41	3.13	95.15

887 To determine the association between justification category and the variables manipulated  
888 in this experiment, we submitted the classification frequencies to a loglinear analysis with a  
889 saturated model including erogation condition (sub and supererogatory), judgment dimension  
890 (goodness, permissibility, obligatoriness, blameworthiness) and justification classification  
891 category (presence of duties, character, presence of rights and absence of duties). The three-way  
892 loglinear model produced a model that retained all interactive effects (likelihood ratio  $\chi^2(0) = 0$ ,  
893  $p = 1$ ), indicating a significant interaction between the variables fitted to the model ( $\chi^2(12) =$   
894  $102.53, p < 0.001$ ). Judgments of permissibility and obligatoriness were mainly sensitive to  
895 presence of rights and these explanations appealing to the presence of rights were more prevalent  
896 in suberogatory cases. Judgments of permissibility were also justified by appealing mostly to the  
897 absence of duties and the presence of rights (>40%) but in this case there is no asymmetry

898 between sub and supererogatory response. Judgments of goodness were mostly justified by  
899 appealing to character and again the suberogatory responses were justified by relying on presence  
900 of rights. Additionally, there are more character justifications for supererogatory responses based  
901 on character than for the corresponding suberogatory responses.

902         The qualitative analysis of the open explanations suggests that people do systematically  
903 call to mind considerations on duties and rights when making judgments about permissibility and  
904 obligatoriness. While this is not surprising, the interesting thing here is that these considerations  
905 are also significant for judgments of blameworthiness and goodness. Duties and rights showed  
906 up, not as blank generalizations – as explored in Experiments 2 and 5 – but as situated appeals to  
907 factors of the situation.

## 908 **Discussion**

909         As before, we replicated the basic pattern found in our other experiments using two  
910 situations that are more clearly moral. This provides strong evidence for our claim that  
911 differences in evaluative categories track differences in kinds of *moral* judgments.

912         We also used open responses to allow participants to supply their own reasoning for  
913 making different moral judgments. Our results are striking. People appear to be sensitive to many  
914 kinds of information differentially when making moral judgments. Judgments of praise and  
915 blame are sensitive to character considerations. Judgments of goodness seem primarily sensitive  
916 to character considerations, whereas judgments of badness seem sensitive to both character  
917 considerations and the presence of rights. This, however, might be a function of the situations we  
918 provided. When people make judgments of badness, they were often pointing to the fact that the

919 bad behavior was fully within the rights of the individual, even though the behavior was  
920 indicative of vicious character. The fact that people mention both points to the relevance of both  
921 kinds of information to making a judgment of badness. Judgments of permissibility appear to  
922 track the presence of rights and the absence of duties, whereas judgments of obligatoriness track  
923 the presence of rights and the presence or absence of duties. We offer an interpretation of this in  
924 the general discussion.

925 Finally, we should note that the categories used to sort open responses were generated *a*  
926 *priori*. Bottom-up, data-driven methods might reveal a more perspicuous classificatory scheme  
927 that might generate surprising results about the kinds of information relevant to making different  
928 kinds of moral judgment.

## 929 **General Discussion**

930 In this paper, our goal was to develop the idea that morality is not exhausted by a system  
931 of rules and explore some of its implication for our understanding of the psychology of moral  
932 judgment. To do this, we've made a case for moving beyond the use of dilemmas in the study of  
933 moral judgment and propose a wider array of measures and situations. In this General Discussion,  
934 we summarize how our findings apply to the measurement of moral judgment, the cases used to  
935 elicit judgments, and the rule-based view of morality.

### 936 ***Measures and Categories of Moral Judgment***

937 We found that judgments across different moral categories are dissociable. In particular,  
938 judgments of goodness and permissibility come apart depending on the situation presented.

939 People form judgments about situations that conform to the category of the suberogatory (bad,  
940 but permissible). This challenges the idea that different measures of moral judgment tap into the  
941 same underlying reasoning process. The resulting picture of moral judgment is that people are  
942 variably sensitive to different kinds of considerations when making different kinds of judgments  
943 (see Barbosa and Jiménez-Leal, 2017). While this presents a more granular picture of moral  
944 judgment, it also opens up new questions. When do people prefer different kinds of information  
945 in making moral judgments? What sorts of considerations predominantly drive different kinds of  
946 judgments?

947         Our results provide some initial suggestion for how to answer these questions. In many of  
948 our experiments, ratings of obligation remained consistently low. This indicates that people  
949 considered both suberogatory and supererogatory behaviors as optional rather than obligatory.  
950 The qualitative data from Experiment 6 suggests that ratings of obligation are predominantly  
951 tracking the presence of duties or explicit rules that prohibit or prescribe conduct. Because the  
952 situations considered here are not composed of conflicts between rules, there is an absence of  
953 duties that defines how one ought to behave. Hence, ratings of obligation are low.

954         Results from Experiment 5 indicate that ratings of badness align with ratings of sanction.  
955 That is, the degree to which one rates a suberogatory action as bad is related to amount of  
956 sanction one deserves in light of doing something suberogatory. This is importantly different  
957 from the results of Experiment 4, which indicated that ratings of permissibility are aligned with  
958 ratings of blame.

959         These results suggest that folk psychological categories of blame and sanction might  
960 dissociate. Because these categories differentially associate with others, it might be the case that

961 these judgments track different features of a situation. Understanding this difference might  
962 provide further insight into what judgments of permissibility and goodness indicate about folk  
963 psychological evaluation. This is an important result to investigate in future work. Moral  
964 psychologists and philosophers have assumed that these constructs significantly overlap and,  
965 accordingly, have used them interchangeably (Bendor & Swistak, 2001; Deutsch & Gerard,  
966 1955; Cialdini, Reno, & Kallgren, 1990; Scanlon, 2008; McKenna, 2013; Bennett, 2012). Future  
967 work should investigate further whether and under what circumstances there are reliable  
968 dissociations between these categories.

969 Notably, some have mentioned the need for using new measures in studying moral  
970 judgment (see Uhlmann, Pizarro, & Diermeier, 2015), arguing that folk psychological categories  
971 of judgment are fundamentally directed at personal evaluation rather than behavior evaluation  
972 and routinely employs aretaic rather than deontic concepts. Our results suggest that this is  
973 partially true. People do show an interest in personal evaluation using aretaic concepts. However,  
974 we also find that behavior evaluation and deontic concepts also play a significant role. This  
975 suggests that social cognition is sensitive to both act-based and person-based considerations, as  
976 well as a wide conceptual repertoire in normative evaluation.

### 977 *Moral Situations*

978 Using a variety of measures is not sufficient to bring out the underlying variability of  
979 moral judgment. In Experiment 3 (on sacrificial dilemmas), we did not find similar dissociations  
980 between judgment categories as we did in our other experiments. This was the case even when  
981 controlling for underlying ethical tendencies (i.e., low or high trait utilitarianism). Even among  
982 low trait utilitarians, different kinds of judgments never significantly differed from each other.

983 This suggests that dilemmas themselves are qualitatively distinct from the kind of moral  
984 encounters we used in our experiments.

985 As discussed earlier, assumptions about methodology and dilemmas mutually reinforce  
986 each other. We can now clarify this point further. The use of dilemmas as stimuli functions as a  
987 demand effect that gives the appearance of commonsense moral thinking exhibiting a rule-based,  
988 hierarchical structure. In providing scenarios that explicitly bring different sets of rules into  
989 conflict, researchers have set people up to exhibit judgments that appear to preferentially select  
990 some rules over others (deontological prescription vs. maximization principles). This, in turn,  
991 blurs the distinctions between different judgment categories. Once we free moral judgment of the  
992 constraints of dilemmas, variegated measures capture better the variability and nuance of  
993 judgment. That variability is only possible when using situations that do not bring different sets  
994 of rules into conflict.

995 Different situations also allow for a different kind of variability. Moral dilemmas have  
996 mostly been constructed out of examples designed to test abstract philosophical principles. which  
997 has led to worries about their ecological validity (Kahane, 2015; Dahl & Oftedal, 2019). A  
998 related but, perhaps more significant, concern is that dilemmas, because of their abstract and  
999 artificial nature, likely occlude differences in moral attitudes that arise from socio-cultural  
1000 variation. Thus, in the limited number of studies that have been conducted outside the urban areas  
1001 of North America and Western Europe, researchers have found responses similar to those found  
1002 among WEIRD populations (Abarbanell & Hauser, 2010; Barrett et al., 2016; Koenigs et al.,  
1003 2007; Perkins et al., 2013; Szekely & Miu, 2015; Johnson, Danko, Huang, Park, Johnson, &  
1004 Nagoshi, 1987).



1005           By contrast, it is likely that responses that are deemed suberogatory in one community,  
1006 such as tipping 10%, might be considered supererogatory by another one. As cultural variability  
1007 mediates moral learning, it is plausible to conjecture that different cultures assign varying  
1008 weights to values, virtues, and rules in the justification of their moral judgments (Graham,  
1009 Meindl, Beall, Johnson & Zhang, 2016). Consequently, cultural differences can be observed not  
1010 as variation in a dimension (e.g. being more or less utilitarian) but as a particular pattern of  
1011 judgements and their justifications. At present, we do not have results that speaks directly to  
1012 these issues—although we currently have relevant work in preparation. We do think, however,  
1013 that moving beyond dilemmas opens up the possibility for this variability to emerge.

#### 1014 *Rule-based Morality*

1015           We mentioned at the outset that not all moral situations reflect conflicts between rules.  
1016 Our results show that when people approach some of these situations and are asked to evaluate  
1017 them, they bring to bear a variety of considerations, not all of which are codifiable in rules. This  
1018 strongly suggests that moral judgment does not necessarily rely on the application of abstract  
1019 *moral* principles.

1020           In Experiments 2, 4 and 5, we asked participants to fill questionnaires aimed at measuring  
1021 general attitudes towards duties and rights. We failed to find correlations between their attitudes  
1022 toward abstract claims about the nature of rights and duties and their judgments about particular  
1023 situations. This reflects a tendency of people to identify features of situations that guide judgment  
1024 without having explicit representations of abstract principles play a causal role in forming moral  
1025 judgments (Graham et al., 2013). Likewise, the open responses we collected in Experiment 6  
1026 suggest that rules play a role in structuring commonsense moral thinking alongside axiological

1027 and characterological considerations, at least when it comes to explaining proffered moral  
1028 judgements. In particular, as indicated by the responses under the “absence of duties” category,  
1029 participants judged actions to be good, bad, etc., while explaining that these didn’t violate a  
1030 prohibition (in suberogatory cases) or were not generally mandated (in supererogatory cases).  
1031 Instead, their evaluations seem overwhelmingly responsive to considerations about rights, values,  
1032 and character traits.

1033 For all we have said, it may turn out that precise algorithms adequately characterize the  
1034 computations performed in forming moral judgments. Alternatively, moral judgment might be  
1035 supported by an architecture that functions as a model-based system without having to explicitly  
1036 represent rules (see Crockett, 2013; Cushman, 2015; Brownstein, 2018). Our interest here,  
1037 however, does not lie with the computations or cognitive architecture behind moral judgment, but  
1038 with the content of morality. Our claim is that most research on moral judgment has, thus far,  
1039 assumed that the content of moral judgment can be measured in terms of its alignment with some  
1040 structured system of moral rules. But fixating on rules provides only a partial window into the  
1041 moral life, one that fails to reveal the complexities and subtleties of commonsense moral  
1042 judgment. When moral judgment is removed from the narrow frame of sacrificial dilemmas, the  
1043 appearance of rules in moral judgment evaporates.

1044 Finally, our argument here does not assume that if morality has a rule-based structure,  
1045 then commonsense moral judgment ought to exhibit perfectly coherent and systematic principles.  
1046 Even Sidgwick (1981) admits that the maxims of commonsense morality are “somewhat vague  
1047 generalities” (p. 342). Still, on his view and on the view of those who follow him, commonsense  
1048 morality can be refined through rigorous theorizing to *approximate* the structure of an ideal

1049 normative theory. It's made of the right sort of stuff, as it were, to function as a rule-based  
1050 system.

1051         The judgments that we have identified in these experiments, however, imply that this is an  
1052 incomplete view of the moral framework reflected in moral cognition. Consider that people  
1053 regularly judge that some behavior is bad but permissible across a range of different scenarios. If  
1054 there are rules constitutive of commonsense morality, what must they be like such that they allow  
1055 for such judgments? How can we make sense of such judgments within a rule-governed system?  
1056 Recall that one key feature of the rule-based system of morality is the inter-definability of moral  
1057 concepts. Whatever is permissible is in line with the rules of morality; whatever is bad is bad  
1058 because it goes against those rules.

1059         One option is to say that separate domains of rules govern separate judgments. But that  
1060 requires giving up on the notion of inter-definability that is central to the rule-based system of  
1061 morality. Another option is to say that one kind of judgment is properly moral and the other is  
1062 not. However, there appears to be no principled way of stating that either judgments of  
1063 permissibility or goodness are properly moral while excluding the other. A third option is to say  
1064 that there *are* systematic principles underlying these judgments, but people do not understand  
1065 what these are. Hence, they are making a mistake in dissociating permissibility and goodness.  
1066 Lastly, one could argue that morality is self-contradictory.

1067         While all of this is possible in principle, the attempt to defend these answers in the face of  
1068 the evidence presented here begins to look like the imposition of researcher assumptions rather  
1069 than an investigation into commonsense moral thinking. At this point, one begins to wonder if

1070 the substantive assumption about morality being a system of rules is worth the cost. We suggest  
 1071 that it is not, and that it is time to consider what morality beyond the rules would look like.

1072 **References**

- 1073 Abarbanell, L. & Hauser, M.D. 2010. Mayan morality: An exploration of permissible harms.  
 1074 *Cognition* 115:2, 207-24.
- 1075 Anderson, A.R. (1958). A reduction of deontic logic to alethic modal logic. *Mind*, 67, 100-103.
- 1076 Andrade, G. (2019). "Medical ethics and the trolley Problem," *Journal of Medical Ethics and*  
 1077 *History of Medicine*, 12:3.
- 1078 Archer, A. (2018). Supererogation. *Philosophy Compass*, 13(3), 1–9. doi:[10.1111/phc3.12476](https://doi.org/10.1111/phc3.12476)
- 1079 Barbosa, S., & Jiménez-Leal, W. (2017). It's not right but it's permitted: Wording effects in  
 1080 moral judgement. *Judgment and Decision Making*, 12(3), 308.
- 1081 Barrett, H.C., Bolyanatz, A., Crittenden, A.N., Fessler, D.M.T., Fitzpatrick, S., Gurven, M.,  
 1082 Henrich, J., Kanovsky, M., Kushnick, G., Pisor, A., Scelza, B.A., Stich, S., von Reuden, C.,  
 1083 Zhao, W., & Laurence, S. 2016. Small-scale societies exhibit fundamental variation in the  
 1084 role of intentions in moral judgment. *Proceedings of the National Academy of the Sciences*  
 1085 113:17, 4688-93.
- 1086 Bartels, D. M., & Pizarro, D. A. (2011). The mismeasure of morals: Antisocial personality traits  
 1087 predict utilitarian responses to moral dilemmas. *Cognition*, 121(1), 154–161.
- 1088 Bates, D., Mächler, M., Bolker, B., & Walker, S. (2014). Fitting linear mixed-effects models  
 1089 using lme4. *ArXiv Preprint ArXiv:1406.5823*.
- 1090 Bauman, C. W., McGraw, A. P., Bartels, D. M., & Warren, C. (2014). Revisiting external  
 1091 validity: Concerns about trolley problems and other sacrificial dilemmas in moral  
 1092 psychology. *Social and Personality Psychology Compass*, 8(9), 536–554.
- 1093 Białek, M. & De Neys, W. (2017). Dual processes and moral conflict: Evidence for  
 1094 deontological reasoners' intuitive utilitarian sensitivity. *Judgment and Decision Making*,  
 1095 12(2), 148–167.
- 1096 Bennett, C. 2012. The Expressive Function of Blame. In J. Coates and N. Tognazzini (Eds.)  
 1097 *Blame: Its Nature and Norms* (Oxford: Oxford University Press).
- 1098 Björklund, F. (2003). Differences in the justification of choices in moral dilemmas: Effects of  
 1099 gender, time pressure and dilemma seriousness. *Scandinavian Journal of Psychology*, 44(5),  
 1100 459–466.
- 1101 Christensen, J. F., Flexas, A., Calabrese, M., Gut, N. K., & Gomila, A. (2014). Moral judgment

- 1102           reloaded: A moral dilemma validation study. *Frontiers in Psychology*, 5(JUL), 1–18.  
1103           <https://doi.org/10.3389/fpsyg.2014.00607>
- 1104 Christensen, J. F., & Gomila, A. (2012). Moral dilemmas in cognitive neuroscience of moral  
1105 decision-making: A principled review. *Neuroscience & Biobehavioral Reviews*, 36(4),  
1106 1249–1264.
- 1107 Conway, P., & Gawronski, B. (2013). Deontological and utilitarian inclinations in moral decision  
1108 making: a process dissociation approach. *Journal of Personality and Social Psychology*,  
1109 104(2), 216.
- 1110 Crockett, M. (2013). Models of Morality. *Trends in Cognitive Sciences* 17:8, 363-66.
- 1111 Cushman, F. (2008). Crime and punishment: Distinguishing the roles of causal and intentional  
1112 analyses in moral judgment. *Cognition*, 108(2), 353–380.
- 1113 Cushman, F., Young, L., & Greene, J. D. (2010). Our multi-system moral psychology: Towards a  
1114 consensus view. *The Oxford Handbook of Moral Psychology*, 47–71.
- 1115 Cushman, F., Young, L., & Hauser, M. (2006). The role of conscious reasoning and intuition in  
1116 moral judgment: Testing three principles of harm. *Psychological Science*, 17(12), 1082–  
1117 1089.
- 1118 Dahl, F. A., & Oftedal, G. (2019). Trolley Dilemmas Fail to Predict Ethical Judgment in a  
1119 Hypothetical Vaccination Context. *Journal of Empirical Research on Human Research*  
1120 *Ethics*, 14(1), 23–32.
- 1121 Darley, J. M., & Shultz, T. R. (1990). Moral rules: Their content and acquisition. *Annual Review*  
1122 *of Psychology*, 41(1), 525–556.
- 1123 Djeriouat, H., & Trémolière, B. (2014). The Dark Triad of personality and utilitarian moral  
1124 judgment: The mediating role of Honesty/Humility and Harm/Care. *Personality and*  
1125 *Individual Differences*, 67, 11–16.
- 1126 Driver, J. (1992). The suberogatory. *Australasian Journal of Philosophy*, 70(3), 286–295.
- 1127 Everett, J. A. C., & Kahane, G. (2020). Switching Tracks? Towards a Multidimensional Model of  
1128 Utilitarian Psychology. *Trends in Cognitive Sciences*.
- 1129 Fauconnet, P. (2013). *La responsabilité. Étude de sociologie*. Presses Électroniques de France.
- 1130 Faul, F., Erdfelder, E., Lang, A.-G., & Buchner, A. (2007). G\*Power 3: A flexible statistical  
1131 power analysis program for the social, behavioral, and biomedical sciences. *Behavior*  
1132 *Research Methods*, 39, 175-191
- 1133 Gilligan, C. (1982). *In a different voice: Psychological theory and women's development*.  
1134 Harvard University Press.
- 1135 Graham, J., Haidt, J., Koleva, S., Motyl, M., Iyer, R., Wojcik, S. P., & Ditto, P. H. (2013). Moral

- 1136 foundations theory: The pragmatic validity of moral pluralism. In *Advances in experimental*  
1137 *social psychology* (Vol. 47, pp. 55–130). Elsevier.
- 1138 Graham, J., Meindl, P., Beall, E., Johnson, K. M., & Zhang, L. (2016). Cultural differences in  
1139 moral judgment and behavior, across and within societies. *Current Opinion in Psychology*,  
1140 8, 125-130.
- 1141 Greene, J. D., Sommerville, R. B., Nystrom, L. E., Darley, J. M., & Cohen, J. D. (2001). An  
1142 fMRI investigation of emotional engagement in moral judgment. *Science*, 293(5537), 2105–  
1143 2108. <https://doi.org/10.1126/science.1062872>
- 1144 Haidt, J. (2001). The emotional dog and its rational tail: a social intuitionist approach to moral  
1145 judgment. *Psychological Review*, 108(4), 814.
- 1146 Haidt, J., & Baron, J. (1996). Social roles and the moral judgement of acts and omissions.  
1147 *European Journal of Social Psychology*, 26(2), 201-218.
- 1148 Hofmann, W., Wisneski, D. C., Brandt, M. J., & Skitka, L. J. (2014). Morality in everyday life.  
1149 *Science*, 345(6202), 1340–1343.
- 1150 Holyoak, K. J., & Powell, D. (2016). Deontological coherence: A framework for commonsense  
1151 moral reasoning. *Psychological Bulletin*, 142(11), 1179.
- 1152 Hornik, K., Zeileis, A., & Meyer, D. (2006). The strucplot framework: visualizing multi-way  
1153 contingency tables with vcd. *Journal of Statistical Software*, 17(3), 1–48.
- 1154 Hsee, C. K., Loewenstein, G. F., Blount, S., & Bazerman, M. H. (1999). Preference reversals  
1155 between joint and separate evaluations of options: a review and theoretical analysis.  
1156 *Psychological Bulletin*, 125(5), 576.
- 1157 Hurd, H. M. (1998). Duties beyond the call of duty. *JRE*, 6, 3.
- 1158 Johnson, R.C., Danko, G.P., Huang, Y.-H., Park, J.Y., Johnson, S.B., & Nagoshi, C.T. 1987.  
1159 Guilt, shame, and adjustment. *Personality and Individual Differences* 8:3, 357-64.
- 1160 Kahane, G. (2015). Sidetracked by trolleys: Why sacrificial moral dilemmas tell us little (or  
1161 nothing) about utilitarian judgment. *Social Neuroscience*, 10(5), 551–560.  
1162 <https://doi.org/10.1080/17470919.2015.1023400>
- 1163 Kahane, G., Everett, J. A. C., Earp, B. D., Caviola, L., Faber, N. S., Crockett, M. J., & Savulescu,  
1164 J. (2018). Beyond sacrificial harm: A two-dimensional model of utilitarian psychology.  
1165 *Psychological Review*, 125(2), 131.
- 1166 Koenigs, M., Young, L., Adolphs, R., Tranel, D., Cushman, F., Hauser, M., & Damasio, A. 2007.  
1167 Damage to the prefrontal cortex increases utilitarian moral judgements. *Nature* 446, 908-11.
- 1168 Koenigs, M., Kruepke, M., Zeier, J., & Newman, J. P. (2012). Utilitarian moral judgment in  
1169 psychopathy. *Social Cognitive and Affective Neuroscience*, 7(6), 708–714.

- 1170 Kohlberg, L., & Hersh, R. H. (1977). Moral development: A review of the theory. *Theory into*  
1171 *Practice*, 16(2), 53–59.
- 1172 Kohlberg, L., Levine, C., & Hewer, A. (1983). *Moral stages: A current formulation and a*  
1173 *response to critics*.
- 1174 Lee, J. J., & Gino, F. (2015). Poker-faced morality: Concealing emotions leads to utilitarian  
1175 decision making. *Organizational Behavior and Human Decision Processes*, 126, 49–64.
- 1176 Lenth, R. (2020). emmeans: Estimated Marginal Means, aka Least-Squares Means. R package  
1177 version 1.4.5. <https://CRAN.R-project.org/package=emmeans>
- 1178 Matuschek, H., Kliegl, R., Vasishth, S., Baayen, H., & Bates, D. (2017). Balancing Type I error  
1179 and power in linear mixed models. *Journal of Memory and Language*, 94, 305-315.
- 1180 McKenna, M. 2013. *Conversation and Responsibility* (Oxford: Oxford University Press).
- 1181 Mikhail, J. (2011). *Elements of moral cognition: Rawls' linguistic analogy and the cognitive*  
1182 *science of moral and legal judgment*. Cambridge University Press.
- 1183 Moll, J., & de Oliveira-Souza, R. (2007). Moral judgments, emotions and the utilitarian brain.  
1184 *Trends in Cognitive Sciences*, 11(8), 319–321.
- 1185 Monin, B., Pizarro, D. A., & Beer, J. S. (2007). Deciding versus reacting: Conceptions of moral  
1186 judgment and the reason-affect debate. *Review of General Psychology*, 11(2), 99–111.
- 1187 Monroe, A. E., Dillon, K. D., Guglielmo, S., & Baumeister, R. F. (2018). It's not what you do,  
1188 but what everyone else does: On the role of descriptive norms and subjectivism in moral  
1189 judgment. *Journal of Experimental Social Psychology*, 77(March), 1–10.  
1190 <https://doi.org/10.1016/j.jesp.2018.03.010>
- 1191 Nosek, B. A., Beck, E. D., Campbell, L., Flake, J. K., Hardwicke, T. E., Mellor, D. T., van't  
1192 Veer, A. E., & Vazire, S. (2019). Preregistration is hard, and worthwhile. *Trends in*  
1193 *Cognitive Sciences*, 23(10), 815–818.
- 1194 Nozick, R. (1974). *Anarchy, state, and utopia* (Vol. 5038). New York: Basic Books.
- 1195 O'Hara, R. E., Sinnott-Armstrong, W., & Sinnott-Armstrong, N. A. (2010). Wording effects in  
1196 moral judgments. *Judgment and Decision Making*.
- 1197 Patil, I., Zucchelli, M., Kool, W., Campbell, S., Fornasier, F., Calo, M., Silani, G., Cikara, M., &  
1198 Cushman, F. 2020. Reasoning supports utilitarian resolutions to moral dilemmas across  
1199 diverse measures. *Journal of Personality and Social Psychology*.
- 1200 Perkins, A.M., Leonard, A.M., Weaver, K., Dalton, J.A., Mehta, M.A., Kumari, V., Williams,  
1201 S.C.R., & Ettinger, U. 2013. A dose of ruthlessness: Interpersonal moral judgment is  
1202 hardened by the anti-anxiety drug lorazepam. *Journal of Experimental Psychology: General*  
1203 142:3, 612-20.

- 1204 Piaget, J. (2013). *The moral judgment of the child*. Routledge.
- 1205 Pizarro, D., Uhlmann, E., & Salovey, P. (2003). Asymmetry in judgments of moral blame and  
1206 praise: The role of perceived metadesires. *Psychological Science*, 14(3), 267–272.
- 1207 Rest, J. R. (1992). *Development in judging moral issues*. U of Minnesota Press.
- 1208 Rosen, B. (1980). Moral dilemmas and their treatment. *Moral Development, Moral Education,*  
1209 *and Kohlberg. B. Munsey (Ed). (1980), 232–263.*
- 1210 Scanlon, T.M. 2008. *Moral Dimensions* (Cambridge, MA: Harvard University Press).
- 1211 Sidgwick, H. (2019). *The methods of ethics*. Good Press.
- 1212 Simpson, E.L. 1974. Moral development research: A case study of scientific cultural bias. *Human*  
1213 *Development* 17, 81-106.
- 1214 Singmann, H., & Kellen, D. (2019). An Introduction to Mixed Models for Experimental  
1215 Psychology. In D. H. Spieler & E. Schumacher (Eds.), *New Methods in Cognitive*  
1216 *Psychology*. Psychology Press.
- 1217 Snarey, J.R. 1985. Cross-cultural universality of social-moral development: A critical review of  
1218 Kohlbergian research. *Psychological Bulletin* 97:2, 202-232.
- 1219 R Core Team. (2019). *R: A language and environment for statistical computing*.
- 1220 Szekely, R.D. & Miu, A.C. 2013. Incidental emotions in moral dilemmas: The influence of  
1221 emotion regulation. *Cognition & Emotion* 29:1, 64-75.
- 1222 Thomson, J. J. (1976). A defense of abortion. In *Biomedical ethics and the law* (pp. 39–54).  
1223 Springer.
- 1224 Uhlmann, E. L., Pizarro, D. A., & Diermeier, D. (2015). A person-centered approach to moral  
1225 judgment. *Perspectives on Psychological Science*, 10(1), 72–81.
- 1226 Westfall, J. (2015). PANGEA: Power analysis for general ANOVA designs. *Unpublished*  
1227 *manuscript*. Available at <http://jakewestfall.org/publications/pangea.pdf>.
- 1228 Walker, L.J. 1984. Sex differences in the development of moral reasoning: A critical review.  
1229 *Child Development* 55:3, 677-691.
- 1230 Young, I. M., (1985). Impartiality and the civic public: Some implications of feminist critiques of  
1231 moral and political theory. *Praxis International*, 5(4), 381–401.