

1  
2 Reconciling the opposing effects of neurobiological evidence on criminal sentencing judgments  
3

4  
5 Allen, Corey H.<sup>1¶</sup>, Vold, Karina<sup>2</sup>, Felsen, Gidon<sup>3</sup>, Blumenthal-Barby, Jennifer S.<sup>4</sup>, Aharoni,  
6 Eyal<sup>1,5,6\*¶</sup>

7  
8 <sup>1</sup> Neuroscience Institute, Georgia State University, Atlanta, Georgia, United States of America  
9

10 <sup>2</sup> Department of Philosophy, Leverhulme Center for Future of Intelligence, University of  
11 Cambridge, Cambridge, England

12  
13 <sup>3</sup> Department of Physiology and Biophysics, University of Colorado, Denver, Colorado, United  
14 States of America

15  
16 <sup>4</sup> Department of Philosophy, Center for Medical Ethics and Health Policy, Baylor College of  
17 Medicine, Houston, Texas, United States of America

18  
19 <sup>5</sup> Department of Psychology, Georgia State University, Atlanta, Georgia, United States of  
20 America

21  
22 <sup>6</sup> Department of Philosophy, Georgia State University, Atlanta, Georgia, United States of  
23 America

24  
25

26

27

28

29 \* Corresponding Author

30 E-mail: [eaharoni@gsu.edu](mailto:eaharoni@gsu.edu) (EA)

31 ¶ Equal Authorship

## 32 **Abstract**

33

34 Legal theorists have characterized physical evidence of brain dysfunction as a double-edged  
35 sword, wherein the very quality that reduces the defendant's responsibility for his transgression  
36 could simultaneously increase motivations to punish him by virtue of his apparently increased  
37 dangerousness. However, empirical evidence of this pattern has been elusive, perhaps owing to a  
38 heavy reliance on singular measures that fail to distinguish between plural, often competing  
39 internal motivations for punishment. The present study employed a test of the theorized double-  
40 edge pattern using a novel approach designed to separate such motivations. We asked a large  
41 sample of participants (N = 330) to render criminal sentencing judgments under varying  
42 conditions of the defendant's mental health status (Healthy, Neurobiological Disorder,  
43 Psychological Disorder) and the disorder's treatability (Treatable, Untreatable). As predicted,  
44 neurobiological evidence simultaneously elicited shorter prison sentences (i.e., mitigating) and  
45 longer terms of involuntary hospitalization (i.e., aggravating) than equivalent psychological  
46 evidence. However, these effects were not well explained by motivations to restore treatable  
47 defendants to health or to protect society from dangerous persons but instead by deontological  
48 motivations pertaining to the defendant's level of deservingness and possible obligation to  
49 provide medical care. This is the first study of its kind to quantitatively demonstrate the  
50 paradoxical effect of neuroscientific trial evidence and raises implications for how such evidence  
51 is presented and evaluated.

52

53 **Keywords:** punishment, sentencing, responsibility, neuroscience, scientific evidence, legal  
54 decision making

55

## 56 **Background**

57           Neuroscience is playing an increasing role in criminal trials. While it is unfeasible to  
58 estimate the prevalence of neurobiological evidence in lower courts, their rates in murder trials  
59 may exceed five percent, as indicated by the subset of cases documented at the appellate level  
60 [1]. But brain evidence can be complicated, raising questions about how fact finders interpret the  
61 quality of this evidence.

62           According to recent research, ordinary people have considerable preconceptions about  
63 the explanatory power of neurobiological evidence. Weisberg, Taylor, and Hopkins [2], for  
64 example, found that when lay people evaluate the quality of scientific explanations for behavior,  
65 their ability to distinguish between good and bad quality explanations was hampered by the  
66 presence of irrelevant neuroscience information. People judged explanations paired with the  
67 irrelevant neuroscience information as stronger and more satisfying than explanations without it.  
68 The investigators describe this context effect as evidence of the “seductive allure” of  
69 neuroscientific explanations. Furthermore, brain images, per se, may have a particularly  
70 persuasive impact on credibility judgments (e.g., [3]; but see [4–6]).

71           If people perceive evidence to be stronger when it is dressed up in neuroscientific garb, to  
72 what extent does this tendency impact legal judgments? In a study of trial court judges,  
73 Aspinwall, Brown, and Tabery [7], found that psychiatric testimony about a defendant’s mental  
74 illness reduced recommended prison sentence lengths when that testimony included a description  
75 of the illness’ biological causes. Similarly, in a mock trial study, Greene and Cahill [8] showed  
76 that in the case of high risk offenders, neuroscientific evidence of psychosis reduced the number  
77 of death sentence recommendations compared to the diagnosis alone. Likewise, Capestany and  
78 Harris [9] found that biological personality assessments reduced punishments compared to

79 behavioral assessment in a mock trial study. Marshall and colleagues [10] found that  
80 neurobiological explanations reduced perceptions of dangerousness when the defendant is  
81 described as psychopathic. Finally, Shariff and colleagues [11] found that exposing people to  
82 general information about the neural bases of human behavior reduced the length of  
83 recommended prison sentences in mock trial cases.

84         One explanation for the apparent “seductive allure” of neurobiological evidence is that  
85 people assume that physical causes of behavior, such as genetic or neurological causes, indicate  
86 that the behavior is outside the agent’s ability to make free choices or to control their behavior  
87 and therefore outside their scope of responsibility. Thus, when physical causes are more salient  
88 than other causes, judgments should be more lenient. Consistent with this theory, Greene and  
89 Cohen [12] have argued that advances in neuroscience sow doubt about the causal roles of  
90 individual choice and control, and thus, the degree to which people should be held responsible  
91 for illicit actions.

92         Other scholars have argued that the tendency to excuse a person of responsibility simply  
93 because that behavior has identifiable physical causes is fallacious because, in reality, all actions  
94 are ultimately physically determined. Morse [13] names this fallacy the “fundamental  
95 psycholegal error,” and suggests to the contrary that there are legitimate reasons to hold people  
96 responsible for their actions even though those actions have physical causes (see also [14,15]).

97         Fallacious or not, the inference that physical causes of a defendant’s behavior render that  
98 behavior outside his control carries another risk. Though such inferences can reduce attributions  
99 of responsibility, scholars have warned that they can potentially increase perceptions that the  
100 defendant is dangerous and in need of greater institutional supervision. For instance, Berryessa  
101 [16] found that potential jurors given evidence of biological risk factors rated the defendant as

102 less responsible for their acts and more likely to commit future crimes compared to those not  
103 given biological risk factor information. This potential of biological evidence to cut both ways in  
104 a defendant's case has been described as a double-edged sword [7,17–19]. If admission of such  
105 evidence carries this risk, it has important implications for how legal parties present evidence,  
106 how judges and jurors evaluate that evidence, and more broadly how human beings make moral  
107 judgments.

108         Despite the growing number of studies on this topic, evidence of the double-edged nature  
109 of biological evidence has been elusive. Though biological evidence has tended to mitigate guilt  
110 and punishment in the aforementioned studies [7,8], other studies have shown aggravating  
111 effects. For instance, McCabe, Castel, and Rhodes [20] found that potential jurors given  
112 incriminating fMRI lie detection evidence rendered more guilty verdicts than those given the  
113 same evidence in the form of a polygraph or thermal facial imaging test, as well as those given  
114 no lie detection evidence at all. In contrast to Greene and Cahill's [8] observed effect of a  
115 reduction in death sentences, Saks, Schweitzer, Aharoni, and Kiehl [21] found that when  
116 neuroimaging evidence is proffered by the prosecution, death sentence recommendations  
117 increase.

118         Yet other mock trial studies have reported null effects of biological explanations. For  
119 example, neurobiological evidence of psychopathic or anti-social tendencies in criminal  
120 offenders had no effect on participant's recommended prison sentence lengths compared to  
121 psychological or behavioral evidence [22,23]. Similarly, Blakey and Kremsmayer [24] found that  
122 describing a case of aggravated assault as stemming from the offender's impaired brain activity  
123 as opposed to his lower self-control had no significant impact on the length of prison sentence  
124 recommended for the crime. Finally, Schweitzer, Saks, Murphy, Roskies, Sinnott-Armstrong,

125 and Gaudet [25] found that after conducting four experiments investigating the effects of  
126 neurobiological evidence on criminal sentencing, a meta-analysis demonstrated no effect of  
127 neurobiological evidence on guilt verdicts.

128         Despite expectations that neurobiological evidence cuts both ways, no single study has  
129 quantitatively demonstrated both mitigating and aggravating effects side by side. One reason  
130 could be that different studies have ignored potential attitudes about the disorder's amenability to  
131 treatment. Highly treatable disorders tend to reduce perceptions that the defendant continues to  
132 be a danger to others. So if a biological disorder is portrayed as treatable, punishments should be  
133 lenient because the defendant will be perceived as low in both responsibility and dangerousness.  
134 But if the same disorder is portrayed as untreatable, this may evoke concerns that the defendant  
135 is dangerous even if he is less morally responsible for the crime. To our knowledge, no previous  
136 studies have considered the potentially moderating role of treatability.

137         Another possible reason for the inconsistent findings is that the dependent measures used  
138 across studies fail to capture distinct, sometimes competing, punitive motives within the same  
139 individual. Most quantitative studies, for instance, included only a single punishment measure,  
140 along the lines of: "How long should the offender be sentenced to prison?" or "How much  
141 should the offender be punished?". If a participant has both deontological concerns (e.g., that the  
142 offender deserves to be punished for his moral transgression) and consequentialist concerns (e.g.,  
143 that the offender must be incapacitated because he is a danger to society), she is forced to use a  
144 single measure to voice both types of concerns. When participants are forced to use prison time  
145 as a "one-stop shop" to capture their diverse punitive motives, these separate motives could  
146 interact or cancel out in unknown ways.

147           The present study was designed to address these limitations in a contrastive vignette  
148 experiment involving the diagnosis of an impulse control disorder following a sexual assault. A  
149 sexual assault crime was selected because it is one of the more plausible charges used to justify  
150 commitment to involuntary hospitalization in many U.S. states. By explicitly manipulating  
151 whether the neurobiological or psychological disorder is deemed treatable or untreatable, we  
152 control for the possible dependency of neurobiological descriptions on perceived treatability (see  
153 Ahn, Proctor, & Flanagan [26]). And, by providing participants with the option to sentence the  
154 offender to prison time (designed largely for moralistic punishment) and/or inpatient  
155 hospitalization time (designed largely for incapacitation but not punishment), we ensure that if  
156 the distinct punitive motives are evoked, they may be distinguished by a simultaneous reduction  
157 in prison and increase in involuntary hospitalization.

158           A supplemental goal of this project was to explore whether the predicted effects are  
159 associated with individual differences in cognitive functioning. Individual differences in legal  
160 settings are important because they can serve as sources of differential treatment of offenders.  
161 Leading theories suggest that many context effects are caused by a tendency to attend to the  
162 salient attributes and activate confirmatory associations in memory [27]. If so, it follows that  
163 people who excel in counterfactual reasoning (i.e., the tendency to entertain multiple possible  
164 perspectives or outcomes) should exhibit less susceptibility to context effects. Yet, other research  
165 suggests that individuals with high cognitive ability are no less susceptible to cognitive biases,  
166 and may even be even more susceptible in some cases [28]. Other theorists have emphasized the  
167 role of emotional regulation in context effects, suggesting that individuals who are better able to  
168 regulate emotions (for example, to reduce their emotional reactions to affective stimuli, such as  
169 pictures) are less susceptible to context effects due to a shift from an emotional strategy to a

170 cognitive strategy [29]. If so, people with high emotion regulation ability should be less  
171 susceptible to the salience of neurobiological causes of behavior.

172         This study investigated the effect of brain-based evidence of an impulse control disorder  
173 on lay sentencing judgments in a large Internet-based sample. The use of lay samples to study  
174 judicial decision making is indirect at best. But lay samples are valuable for at least two other  
175 reasons: They are scientifically important to the extent that they help illuminate general patterns  
176 of human cognition, and they are legally important in that legal policy often relies on public  
177 opinion, as expressed through vehicles like the election of judges and legislators and the  
178 endorsement of ballot propositions and referenda. So, understanding punitive judgment  
179 formation among laypeople can help inform criminal punishment policies and practices even if  
180 these decisions don't necessarily generalize to judicial decision making.

181         The overarching rationale of this study was that if a neurobiological explanation of an  
182 impulse control disorder, more so than a psychological one, primes people to believe the  
183 defendant lacked control of his actions and should therefore be held less responsible for his  
184 crime, then such an explanation should result in reduced prison sentences. However, in this view,  
185 the same evidence should increase support for non-punitive custody (e.g., involuntary  
186 hospitalization) because a defendant who lacks control will also be perceived as more dangerous.  
187 This latter effect should be especially apparent when the defendant's disorder is rendered  
188 untreatable. We thus tested the following hypotheses:

189

190 **H1.** Neurobiological evidence of a disorder will decrease prison sentences relative to  
191 psychological evidence and to no evidence.

192 **H2.** The H1 effect will be greater when the disorder is seen as treatable versus untreatable.

193 **H3.** Neurobiological evidence will increase involuntary hospitalization terms relative to  
194 psychological evidence and to no evidence.

195 **H4.** The H3 effect will be greater when the disorder is seen as untreatable.

196 **H5. (a)** The H1 effect will be best accounted for by a reduction in deontological concerns, and  
197 **(b)** the H3 effect will be accounted for by increased consequentialist concerns.

198 **EH1.** Individuals who score lower in executive functioning (counterfactual reasoning or emotion  
199 regulation) will exhibit decreases in punishment and involuntary hospitalization relative to  
200 higher scoring individuals.

## 201 **Methods**

### 202 **Participants**

203 Three hundred sixty nine adults residing in the U.S. (53% F, 47% M) were recruited from  
204 Amazon Mechanical Turk (MTurk) in November 2017 and paid \$3.00 for participating. Thirty  
205 nine respondents were omitted for incomplete data, missed attention checks (e.g., “What colors  
206 are on the American flag?”), or for spending too little time on the survey (< 1 SD of mean, or ~5  
207 minutes), resulting in a final sample size of 330. The average age was 36.0 years ( $SD = 11.38$ ,  
208 range = 19-71). Median annual household income was \$25,000-\$49,999. Political party  
209 affiliation was 43.9% democratic, 29.0% independent, 15.5% republican, and 15.6% other.  
210 Human subject research was authorized by the Georgia State University institutional review  
211 board: H16349. Written consent was obtained from all participants.

212           The use of MTurk for research purposes has been documented elsewhere [30]. Like most  
213 sampling methods, use of the MTurk sampling pool presents some limitations on generalizability  
214 to the U.S. population, most notably in terms of political party affiliation, for which our sample  
215 disproportionately identified as democratic. However, this pool has been validated for research  
216 on political ideology [31]. More broadly, the MTurk sampling pool, and our sample in particular,  
217 are more representative on basic U.S. demographics, as defined by the 2017 U.S. Census, than  
218 other methods commonly employed in social science research such as the use of university  
219 students.

## 220 **Sample size estimation**

221           The estimated sample size was determined by the power required to detect a significant  
222 interaction of mental health status (neurobiological vs. psychological) and treatability status  
223 (treatable vs. untreatable) on recommended punishment, assuming that the probability of  
224 obtaining a false positive is  $\alpha = 0.05$ . Under this assumption, a sample of 327 participants  
225 provides 95% power to detect a significant interaction in this design where the effect size is  $f =$   
226 0.20 (a small effect size by conventional criteria; [32]).

## 227 **Design**

228           The study used a 3 (Mental Health Status) x 2 (Treatability Status) incomplete factorial  
229 design with random assignment to conditions. Mental health status varied whether the defendant  
230 was described as having an impulse control disorder of neurobiological origins, of psychological  
231 origins, or was healthy. Treatability status varied whether the impulse control disorder was seen  
232 to be completely treatable or untreatable, but only for the neurobiological and psychological  
233 conditions, not the healthy condition. The primary dependent measures consisted of a

234 preliminary, baseline prison sentence recommendation made before exposure to the  
235 manipulations, a revised prison sentence recommendation after the presentation of  
236 manipulations, and the amount of time the defendant should involuntarily spend in an inpatient  
237 hospital *after* the completion of his prison term (all from 0 to 4 yrs. as determined from  
238 unpublished pilot data). Change in prison sentence recommendation was calculated as the  
239 within-subject change from baseline—the difference between their baseline and revised prison  
240 sentence recommendation. In order to account for individual-level variation in punishment  
241 judgments, a composite punishment score was constructed, defined as the individual's revised  
242 sentence divided by his/her baseline punishment recommendation, yielding a percentage change  
243 score for each participant. Exploratory measures were designed to check and clarify the results  
244 of our hypothesis tests. These consisted of Likert-type ratings (from (1) “Strongly Disagree” to  
245 (7) “Strongly Agree”) for various statements regarding the defendant’s moral responsibility,  
246 blameworthiness, desert of punishment, free will, ability to stop himself, trustworthiness, danger  
247 to society, likelihood of reoffense, the degree to which the crime was an expression of his  
248 character, the perceived efficacy of treatment, and the perceived impact and importance of the  
249 evidence on their punishment decision.

## 250 **Materials**

251         The case summary described an instance of sexual assault in which an adult male was  
252 found guilty of assaulting an adult female neighbor. (See S1 Appendix for stimuli.) Following  
253 receipt of this information, the participants were presented with a professional opinion regarding  
254 Mr. Edward’s mental status from either neurologists, psychologists, or “experts”—corresponding  
255 to our neurobiological, psychological, and healthy conditions, respectively. Participants were

256 either informed that the neurologists had located a large tumor in the impulse control region of  
257 the defendant’s brain, that the psychologists had diagnosed the defendant with an impulse control  
258 disorder, or that the experts had determined that the defendant had no mental health issues.  
259 Within the neurobiological condition, participants were told that the neurologists either  
260 conducted surgery to successfully remove the tumor (treatable condition), or found it to be  
261 inoperable (untreatable condition). Similarly, those within the psychological condition were told  
262 that cognitive-behavioral therapy was either a success (treatable condition) or a failure  
263 (untreatable condition). Those within the healthy condition were given no further information.  
264 Aside from these manipulations across conditions, all participants received identical information.

265         Three additional measures were included after the dependent measures to assess  
266 individual differences. The Counterfactual Thinking (CFT) scale is used to assess the ability to  
267 reason counterfactually, as well as a measure of perspective switching and open mindedness—  
268 For example: “My beliefs would not have been very different if I had been raised by a different  
269 set of parents” [33]. The Difficulties in Emotion Regulation Scale-Short Form (DERS-SF) is  
270 meant to assess an individual’s ability to regulate his or her emotions [34]. Following these  
271 additional scales, participants self-reported his or her political ideology from (0) “very liberal” to  
272 (10) “very conservative”.

## 273 **Procedure**

274         Participants were asked to complete the survey privately on their personal devices. After  
275 providing consent, they were instructed to read carefully through a summary of a criminal court  
276 case, and to imagine as if they were the judge overseeing the trial. After the case summary, the  
277 dependent measures were presented followed by several manipulation checks, validated

278 inventories, and supplemental questions. Finally, participants provided demographic information  
279 including age, gender, political affiliation, and income.

## 280 **Results**

### 281 **Hypothesis Tests**

282 **(H1) Were prison sentence recommendations decreased when evidence for the**  
283 **defendant's disorder was described as neurobiological as compared to**  
284 **psychological? (H2) Was this effect greater when the disorder was treatable?**

285 Revised prison sentence recommendations were subjected to a two-way Analysis of  
286 Variance with two mental health status conditions (neurobiological, psychological) and two  
287 treatability status conditions (treatable, untreatable) [35]. There was a main effect of mental  
288 health status on prison sentencing,  $F(1, 219) = 13.07, p < .001, \eta_p^2 = .056$ , indicating, as  
289 predicted, that when neurobiological evidence was given as an explanation for the underlying  
290 disorder, participants recommended significantly shorter sentences ( $M = 0.95, SE = 0.14$ ) than  
291 when psychological evidence was given ( $M = 1.65, SE = 0.14$ ) [36]. The same pattern of results  
292 was found using the change in prison recommendation,  $F(1, 219) = 6.18, p = .014, \eta_p^2 = .027$ , as  
293 well as when subjected to a one-way ANOVA including the healthy condition,  $F(2, 327) =$   
294  $33.64, p < .001, \eta^2 = .171$  [37]. Pairwise comparisons (Fisher's LSD) showed that, as predicted,  
295 the recommended prison sentence was significantly shorter when the defendant had a disorder  
296 supported by neurobiological evidence ( $M = 0.95, SE = 0.13, p < .001$ ) or psychological  
297 evidence ( $M = 1.65, SE = 0.13, p < .001$ ), than when the defendant was healthy ( $M = 2.49, SE =$   
298  $0.14$ ). Similarly, the disorder supported by neurobiological evidence garnered significantly

309 shorter prison sentences than the disorder supported by psychological evidence,  $p < .001$ .

300 As expected, there was a main effect of treatability status on prison sentencing,  $F(1, 219)$   
 301  $= 7.18, p = .008, \eta_p^2 = .032$ , indicating that when participants were told the defendant's disorder  
 302 was treatable, they recommended significantly shorter prison sentences ( $M = 1.04, SE = 0.14$ )  
 303 than when participants were told the disorder was untreatable ( $M = 1.56, SE = 0.14$ ). However,  
 304 there was no significant interaction between the mental health status presented and the  
 305 treatability of the disorder,  $F(1, 219) = 0.00, p = .99, \eta_p^2 < .001$ . The same pattern of results was  
 306 found using the change in prison recommendation,  $F(1, 219) = 13.95, p < .001, \eta_p^2 = .060$  and  
 307  $F(1, 219) = 0.073, p = .79, \eta_p^2 < .00$ , respectively. As a whole, these results were consistent with  
 308 H1 but not H2. See Table 1 for sentencing recommendations by condition.

**Table 1. Revised Prison Sentence and Involuntary Hospitalization Recommendations as a Function of Mental Health Status and Treatability**

Treatability	Measure	Psychological		Neurobiological	
		<i>n</i>	<i>M (SE)</i>	<i>n</i>	<i>M (SE)</i>
<b>High</b>	Revised Prison Sentence (yrs.)	55	1.39 (.19)	57	0.69 (.19)
	Involuntary Hospitalization Term (yrs.)		1.27 (.17)		1.39 (.17)
<b>Low</b>	Revised Prison Sentence (yrs.)	54	1.91 (.20)	57	1.21 (.19)
	Involuntary Hospitalization Term (yrs.)		2.02 (.18)		2.59 (.17)

309

310 *Note.* *n* = sample size, *M* = mean, *SE* = standard error

311

312

313 **(H3) Were involuntary hospitalization terms increased when evidence for the**  
 314 **defendant's disorder was described as neurobiological as compared to**  
 315 **psychological? (H4) Was this effect greater when the disorder was**  
 316 **untreatable?**

317 Recommended involuntary hospitalization terms were subjected to a two-way Analysis of  
 318 Variance with two mental health status conditions (neurobiological, psychological) and two

319 treatability status conditions (treatable, untreatable). There was a main effect of mental health  
320 status on recommended involuntary hospitalization terms,  $F(1, 219) = 4.07, p = .045, \eta_p^2 = .018$ ,  
321 indicating, as predicted, that when neurobiological evidence was given as an explanation for the  
322 defendant's underlying disorder, participants recommended significantly longer recommended  
323 involuntary hospitalization terms ( $M = 1.99, SE = 0.12$ ), than when psychological evidence was  
324 given ( $M = 1.65, SE = 0.12$ ) [38]. The same pattern of results was found when recommended  
325 involuntary hospitalization terms were subjected to a one-way ANOVA, including the healthy  
326 condition,  $F(2, 327) = 42.99, p < .001, \eta^2 = .208$ . Pairwise comparisons showed that, as  
327 predicted, the recommended involuntary hospitalization term was significantly longer when the  
328 defendant's disorder was described neurobiologically ( $M = 1.99, SE = 0.12, p < .001$ ) as well as  
329 when the evidence was described psychologically, ( $M = 1.64, SE = 0.12, p < .001$ ), than when  
330 the defendant was healthy, ( $M = 0.46, SE = 0.12$ ). Similarly, neurobiological evidence garnered  
331 significantly longer recommended involuntary hospitalization terms than psychological  
332 evidence,  $p = .040$ .

333 As expected, there was also a main effect of treatability status on recommended  
334 involuntary hospitalization terms,  $F(1, 219) = 31.97, p < .001, \eta_p^2 = .127$ , indicating that when  
335 participants were told the defendant's disorder was untreatable, they recommended significantly  
336 longer involuntary hospitalization terms ( $M = 2.31, SE = 0.12$ ) than when participants were told  
337 the disorder was treatable ( $M = 1.33, SE = 0.12$ ). However, there was no significant interaction  
338 between the mental health status presented and the treatability of the disorder,  $F(1, 219) = 1.69, p$   
339  $= .20, \eta_p^2 = .008$ . These results were consistent with H3 but not H4. See Fig 1 for punishment  
340 change scores by condition.

341 **Fig 1. Punishment Change Score by Condition.** Bars denote the percentage change in time  
342 from individual baseline punishment rating across conditions, for their revised prison  
343 recommendation (dark grey) and recommendation for involuntary hospitalization (light grey).  
344 Statistically significant differences mirror the patterns described in H1-H4. Standard error bars  
345 shown.

346

347 **(H5a) Was the effect of mental health status on prison sentence**  
348 **recommendation best accounted for by reductions in deontological concerns**  
349 **rather than increases in consequentialist concerns?**

350 We clustered the above items (Cronbach's  $\alpha > 0.70$ ) into two categories in  
351 accordance with jurisprudence theories of punishment: deontological concerns (concerns about  
352 duty, such as the perceived obligation to punish offenders based on their moral blameworthiness)  
353 and consequentialist concerns (concerns about outcomes, such as the desire to punish to protect  
354 public safety). Items comprising the deontological factor were: the offender's moral wrongness,  
355 moral responsibility, blameworthiness, desert of punishment, control of action, and free will  
356 (Cronbach's  $\alpha = .86$ ). Items comprising the consequentialist factor were: the defendant's  
357 dangerousness to society and likelihood of committing future crimes ( $\alpha = .77$ ). These clusters  
358 were confirmed in a two-factor solution identified by a principal components analysis of all  
359 items with varimax rotation, resulting in two independent factors (eigenvalues  $> 1$ ) that matched  
360 our *a priori* grouping and explained 66.24% of the variance. We then used an ordinary least  
361 squares path analysis to examine whether these two types of concerns could account for, the  
362 observed effect of mental health status (neurobiological or psychological) on prison sentence

363 term. The two composites were entered into a parallel regression model in order to compare their  
364 relative impact.

365 One third of the variance in recommended prison sentence length was explained by our  
366 parallel model ( $R^2 = .33$ ). The mitigating effect of neurobiological evidence was fully accounted  
367 for by deontological concerns (See Fig 2). As predicted, the neurobiological condition was a  
368 significant negative predictor of deontological concerns,  $b = -0.63$ ,  $SE = 0.14$ ,  $p < .001$ , and  
369 deontological concerns were a significant predictor of the prison sentence recommended to the  
370 defendant,  $b = 0.52$ ,  $SE = 0.093$ ,  $p < .001$ . A bootstrap confidence interval for the indirect effect  
371 of mental health status as explained by deontological concerns on prison sentence,  $b = -0.33$ ,  $SE$   
372  $= 0.088$ , based on 5,000 samples, was entirely below zero ( $-0.52$  to  $-0.17$ ). The direct effect of  
373 mental health status on prison sentence was not significant,  $b = -0.30$ ,  $SE = 0.17$ ,  $p = .081$  [39].

374 Consequentialist concerns also significantly predicted prison sentences,  $b = 0.28$ ,  $SE =$   
375  $0.071$ ,  $p < .001$ , but mental health status did not predict consequentialist concerns,  $b = -0.24$ ,  $SE$   
376  $= 0.18$ ,  $p = .18$ . A bootstrap confidence interval for the indirect effect of mental health status as  
377 explained by consequentialist concerns on prison sentence,  $b = -0.067$ ,  $SE = 0.053$ , included zero  
378 ( $-0.18$  to  $0.030$ ), indicating that consequentialist concerns did not account for the effect of mental  
379 health status on prison sentence recommendation, consistent with our prediction. The observed  
380 regression model indicates that the mitigating effect of neurobiological evidence on prison  
381 sentence length can be explained by changes in deontological concerns—namely that the  
382 defendant was seen as less responsible for his criminal act.

383

384 **(H5b)** A test of the effect of deontological and consequentialist concerns on the relationship  
385 between mental health evidence and recommended involuntary hospitalization was not justified

386 because no direct effect of evidence type on recommended involuntary hospitalization was  
387 observed.

388 **Fig 2. Regression Coefficients for the Relationship between Mental Health Status and**  
389 **Prison Sentence as Explained by Deontological Concerns and Consequentialist Concerns.**

390 Solid bold lines denote a significant relationship.

391

392 **EH1: Was the effect of mental health status on sentence length moderated by**  
393 **(a) counterfactual reasoning traits or (b) emotion regulation ability?**

394 Change in prison sentence recommendation was subjected to a two-way Analysis of  
395 Variance with two mental health status conditions (neurobiological, psychological) and two  
396 counterfactual reasoning levels (low, high) as determined by a median split [40]. Contrary to  
397 expectation, there was no main effect of counterfactual reasoning level,  $F(1, 181) = 3.29, p =$   
398  $.071, \eta_p^2 = .018$ , and no interaction,  $F(1, 181) = 1.42, p = .23, \eta_p^2 = .008$ , indicating that the  
399 effect of mental health status on prison sentence length did not depend on a participants'  
400 tendency to reason counterfactually.

401 Similarly, change in prison sentence recommendation was subjected to a two-way  
402 Analysis of Variance with two mental health status conditions (neurobiological, psychological)  
403 and two emotional regulation ability levels (low, high) as determined by a median split. Again,  
404 there was no main effect of emotional regulation level,  $F(1, 212) = 0.01, p = .93, \eta_p^2 < .001$ , and  
405 no interaction,  $F(1, 212) = 0.26, p = .61, \eta_p^2 = .001$ , indicating that the effect of mental health  
406 status on prison sentence length did not depend on a participants' ability to regulate their  
407 emotions.

## 408 **Additional Exploratory Analyses**

409           Several exploratory measures were examined to help contextualize and explain the results  
410 of our hypothesis tests.

411

### 412 **Were deontological concerns reduced when the defendant's disorder was** 413 **described as neurobiological compared to psychological?**

414           We theorized that any mitigating effect of neurobiological explanation on prison  
415 sentences would be driven primarily by deontological sentiments that the defendant should be  
416 held less morally responsible for the crime. If so, then participants presented with  
417 neurobiological evidence should likewise rate that defendant lower on measures of  
418 responsibility, blameworthiness, deservingness of punishment, free will, ability to stop himself  
419 from performing the crime, and higher on measures of trustworthiness. In this view, participants  
420 presented with neurobiological evidence should also perceive the crime as less characteristic of  
421 the defendant. All of these predictions were supported.

422           There was a main effect of mental health status on perceptions of the defendant's moral  
423 responsibility for his crime,  $F(1, 219) = 16.13, p < .001, \eta_p^2 = .069$ , his blameworthiness,  $F(1,$   
424  $219) = 15.09, p < .001, \eta_p^2 = .064$ , his deservingness of punishment,  $F(1, 219) = 12.26, p = .001,$   
425  $\eta_p^2 = .053$ , free will,  $F(1, 219) = 21.89, p < .001, \eta_p^2 = .091$ , ability to stop himself from  
426 committing the crime,  $F(1, 219) = 10.58, p = .001, \eta_p^2 = .046$ , his trustworthiness,  $F(1, 219) =$   
427  $7.82, p = .006, \eta_p^2 = .034$ , and the perception that the defendant's action was an expression of his  
428 essential character (i.e., his "deep self"),  $F(1, 219) = 20.91, p < .001, \eta_p^2 = .087$ . Those in the  
429 neurobiological condition saw the defendant as less morally responsible (see Table 2 for exact  
430 group values), less blameworthy, less deserving of punishment, having less free will, less able to

431 stop himself from committing the crime, more trustworthy, and perceived the crime as less  
432 characteristic of the defendant than those in the psychological condition.

433 In contrast, there was no main effect of treatability status on any of these measures except  
434 trustworthiness,  $F(1, 219) = 14.63, p < .001, \eta_p^2 = .063$ , such that when the defendant's disorder  
435 was treatable he was perceived as more trustworthy than when the disorder was untreatable.  
436 Similarly, there was no interaction effect between mental health status and treatability on any  
437 measure besides trustworthiness,  $F(1, 219) = 6.82, p = .010, \eta_p^2 = .030$ , such that the positive  
438 effect of treatability on trustworthiness was larger in the neurobiological condition,  $p < .001, \eta_p^2$   
439  $= .063$ . As might be expected, this effect did not extend to the condition in which the disorder  
440 was described as untreatable,  $p = .90, \eta_p^2 < .001$ .

441

442 **Were consequentialist concerns increased when the defendant's disorder was**  
443 **described as neurobiological compared to psychological?**

444 We theorized that any aggravating effect of the neurobiological explanation on  
445 involuntary hospitalization would be motivated primarily by concerns of the defendant's future  
446 danger to society. If so, then participants presented with neurobiological evidence should  
447 characterize that defendant as more dangerous. However, this prediction was not supported.  
448 Instead, a main effect of mental health status on the defendant's dangerousness,  $F(1, 219) = 4.05,$   
449  $p = .045, \eta_p^2 = .018$ , indicated that participants presented with neurobiological evidence found  
450 the defendant *less* dangerous than those presented with psychological evidence. This effect was  
451 further supported by an interaction,  $F(1, 219) = 6.22, p = .013, \eta_p^2 = .028$ , in which the  
452 mitigating effect of treatability,  $F(1, 219) = 75.23, p < .001, \eta_p^2 = .256$ , on perceived  
453 dangerousness was larger when the evidence was described as neurobiological compared to

454 psychological,  $p = .002$ ,  $\eta_p^2 = .045$ . Pairwise differences were not found between neurobiological  
455 and psychological evidence in the untreatable condition,  $p = .74$ ,  $\eta_p^2 = .001$ . One interpretation of  
456 these counter-intuitive results is that participants perceived involuntary hospitalization not as a  
457 way to incapacitate morally culpable people that pose a danger to society but as a way to provide  
458 medical attention to those most in need of it, including, potentially, those who pose a danger to  
459 *themselves*—a distinction that our dangerousness measure may not have captured.

460

461 **Was the treatment for the defendant's condition perceived as more efficacious**  
462 **when that condition was described as neurobiological versus psychological?**

463       If the aggravating effect of neurobiological explanation on recommended involuntary  
464 hospitalization term was not due to perceptions of increased dangerousness, perhaps it could be  
465 due to a perception that neurobiological disorders are more treatable than psychological  
466 disorders, at least in inpatient contexts. If so, then people should rate the treatment of the  
467 neurobiological disorder as more efficacious. We observed partial support for this hypothesis.  
468 There was a significant interaction between mental health status and treatability,  $F(1, 219) =$   
469  $9.64$ ,  $p = .002$ ,  $\eta_p^2 = .042$ , such that participants presented with neurobiological evidence of a  
470 treatable disorder expressed stronger belief in the efficacy of the treatment than those presented  
471 with psychological evidence of a treatable disorder,  $p = .003$ ,  $\eta_p^2 = .040$ . As might be expected,  
472 this effect did not extend to the condition in which the disorder was described as untreatable,  $p =$   
473  $.17$ ,  $\eta_p^2 = .009$ . Moreover, there was no main effect of mental health status on the efficacy of the  
474 defendant's treatment,  $F(1, 219) = 1.28$ ,  $p = .26$ ,  $\eta_p^2 = .006$ .

475

476 **Were neurobiological descriptions of the defendant’s disorder seen as more**  
 477 **important than psychological descriptions?**

478 Next we examined whether participants expressed explicit attitudes consistent with the  
 479 mitigating effect of neurobiological evidence on punishment. If so, this would support the  
 480 interpretation that participants were consciously aware of the reasons driving their decision. To  
 481 address this question, participants indicated the extent to which they thought the evidence of the  
 482 defendant’s condition “decreases, increases, or has no effect on” their initial punishment. They  
 483 were also asked how important the exam results were to their punishment decision.

484 As expected, participants given neurobiological evidence reported the evidence as more  
 485 mitigating and more important than those given psychological evidence,  $F(1, 219) = 12.42, p =$   
 486  $.001, \eta_p^2 = .054$ ;  $F(1, 219) = 21.78, p < .001, \eta_p^2 = .090$ . Similarly, participants told that the  
 487 disorder was treatable reported the evidence as more mitigating and important than those told the  
 488 disorder was untreatable,  $F(1, 219) = 20.62, p < .001, \eta_p^2 = .086$ ;  $F(1, 219) = 5.87, p = .016, \eta_p^2$   
 489  $= .026$ . There were no significant interactions,  $F(1, 219) = 0.35, p = .56, \eta_p^2 = .002$ ;  $F(1, 219) =$   
 490  $0.23, p = .64, \eta_p^2 = .001$ .

491 **Table 2. Deontological and Consequentialist Concerns as a Function of Mental Health**  
 492 **Status and Treatability Status**

493

	<u>Psychological</u>	<u>Neurobiological</u>	<u>Psychological</u>	<u>Neurobiological</u>
<u>Treatability</u>	<u><i>M (SE)</i></u>	<u><i>M (SE)</i></u>	<u><i>M (SE)</i></u>	<u><i>M (SE)</i></u>
	Moral Responsibility		Action an Expression of Defendant's Character	
<b>High</b>	6.22 (0.18)	5.53 (0.18)	4.96 (0.21)	3.97 (0.20)
<b>Low</b>	6.28 (0.18)	5.54 (0.18)	5.02 (0.21)	4.16 (0.20)

	Blameworthiness		Danger to Society	
<b>High</b>	5.98 (0.18)	5.23 (0.17)	5.09 (0.17)	4.33 (0.17)
<b>Low</b>	6.06 (0.18)	5.44 (0.17)	6.13 (0.17)	6.21 (0.17)
	Deserving of Punishment		Likelihood of Reoffense	
<b>High</b>	6.02 (0.18)	5.16 (0.17)	3.93 (0.17)	3.39 (0.17)
<b>Low</b>	6.00 (0.18)	5.63 (0.17)	4.98 (0.18)	5.23 (0.17)
	Free Will		Perceived Efficacy of Treatment	
<b>High</b>	5.38 (0.21)	4.18 (0.21)	3.22 (0.14)	3.79 (0.13)
<b>Low</b>	5.28 (0.21)	4.52 (0.21)	1.72 (0.14)	1.46 (0.13)
	Ability to Stop Himself		Perceived Impact of Evidence	
<b>High</b>	4.82 (0.21)	3.97 (0.21)	2.38 (0.14)	1.97 (0.14)
<b>Low</b>	4.69 (0.22)	4.16 (0.21)	3.11 (0.14)	2.53 (0.14)
	Trustworthiness		Importance of Exam Results	
<b>High</b>	2.29 (0.17)	3.18 (0.16)	3.64 (0.16)	4.32 (0.16)
<b>Low</b>	2.09 (0.17)	2.12 (0.16)	3.17 (0.17)	4.00 (0.16)

494 *Note.*  $M$  = mean,  $SE$  = standard error

## 495 Discussion

496 The purpose of this project was to investigate the effect of brain-based evidence of an  
497 impulse control disorder on lay sentencing judgments. We observed three key findings: (1) Both  
498 brain evidence and psychological evidence had mitigating effects on prison sentences, but the  
499 mitigating effect of brain evidence was stronger. (2) Yet that same brain evidence evoked  
500 relative increases in involuntary hospitalization terms. (3) The variation in sentencing judgments  
501 was best explained by deontological considerations pertaining to moral culpability.

502           These findings suggest that lay people assign more importance to mental health evidence  
503 whose causes are described in neurobiological terms than in psychological terms. As predicted,  
504 this evidence seems to both favor or disfavor the defendant depending on the decision type:  
505 Although evidence of a neurobiological cause of a disorder may mitigate prison punishment, the  
506 same evidence can place the defendant at an increased risk of involuntary hospitalization.  
507 Though the effect sizes of our primary hypotheses were not large, they are still potentially  
508 relevant to the law, where punishment practices and policies can have far-reaching consequences  
509 for society when deployed over large temporal and geographic scales.

510           One plausible explanation for this effect is that neurobiological evidence primes fact-  
511 finders to preferentially attend to the distinctly physical causes of behavior, and this feeds their  
512 intuitions that the behavior is outside the defendant's control. Perceptions of reduced control  
513 may, in turn, reduce attributions of responsibility while potentially increasing the belief that the  
514 defendant requires medical intervention (see [26]).

515           The mitigation effect is consistent with the operation of a deontological motive for  
516 punishment, namely that punishment should be proportionate to the offender's moral culpability.  
517 The reason for the aggravating effect of neurobiological evidence on recommended involuntary  
518 hospitalization term is less clear. "Double edge" theories would explain this increase using  
519 consequentialist reasons such as a desire to protect society from danger or a desire to provide  
520 treatment to those who would most benefit from it, but the neurobiologically disordered  
521 defendant was rated as no more dangerous or treatable than the other disordered defendant. This  
522 leaves open the question of exactly why people assigned more hospitalization time to the  
523 neurobiologically disordered defendant. We speculate that the answer hinges on how people  
524 interpret the specific purpose of involuntary hospitalization. Perhaps, for instance, people

525 considered involuntary hospitalization more justified for long-term disease management, even if  
526 their treatment prospects are low. Similarly, people in this condition might have felt a greater  
527 *obligation* to provide care, regardless of treatment prospects. Our manipulation checks did not  
528 make such fine distinctions. If these interpretations are confirmed in future research, they would  
529 be compatible with “moral education theory,” the idea that punishments may be justified, or  
530 perhaps even obligatory, to the extent that they benefit to the person being punished [41].  
531 Alternatively, involuntary hospitalization could be used to quarantine people perceived to be ill-  
532 fit for society. This interpretation is consistent with previous research suggesting that people feel  
533 the need to socially distance themselves from individuals with biologically described mental  
534 disorders [42–45]. In such cases, increased hospitalization time for defendants in the  
535 neurobiological condition should be understood as “aggravating” only in the narrow sense that it  
536 was defined as involuntary, but not in a classically retributive sense.

537         This study is the first to quantitatively dissociate the divergent effects of neurobiological  
538 evidence on sentencing decisions (i.e., the “double-edged sword”). In a similar vein, research by  
539 Aspinwall et al. [7] and Fuss et al. [18], found that while biological explanations for a crime  
540 mitigated punishment recommendations or estimations of legal responsibility, some increased  
541 consideration of future dangerousness and support for involuntary commitment was found.  
542 However, in those studies, this support was observed by qualitative measures only. The present  
543 study validated this effect using quantitative measures. Further, the dependent measures in our  
544 study allowed participants to award prison time and involuntary hospitalization time separately.  
545 This approach was employed to distill the moralistic punitive motives from incapacitative and  
546 treatment-based motives, and could explain why our study uniquely observed the predicted  
547 double-edge effect.

548 *Limitations and Future Directions*

549           As with all studies, our findings are necessarily limited by our procedural choices. We  
550 included an alternative measure to prison punishment (involuntary hospitalization) to distill the  
551 nature of participants punitive motives. Even so, involuntary hospitalization itself can be used for  
552 a variety of purposes that we could not disentangle, such as treatment, incapacitation, or possibly  
553 even punishment. Efforts to understand participants motivations for such decisions should  
554 consider a wider array of punishment measures designed to fulfill distinct aims or should devise  
555 additional manipulations that achieve this effect.

556           It is also unclear why participant's individual differences (i.e., the ability to reason  
557 counterfactually and the ability to regulate one's own emotions) did not explain individual  
558 susceptibility to change in prison sentence. It is possible that this was a consequence of weak  
559 construct validity. However, these results are consistent with previous literature showing that  
560 those high in cognitive ability are no less susceptible to cognitive biases (and in some cases,  
561 more susceptible) than those low in cognitive ability [28]. Future studies should address these  
562 possibilities using other theoretically motivated individual difference measures.

563           This study investigated punishment judgments in an Internet-based lay sample and do not  
564 necessarily generalize to legal samples such as judges and jurors or to the broader U.S.  
565 population. Future research on lay samples should aspire to full randomization across key  
566 dimensions including geography and political party affiliation. Likewise, this research should be  
567 extended to legal samples in attempt to replicate these effects among groups whose judgments  
568 are directly consequential for criminal defendants, such as trial court judges. Lastly, it would be  
569 helpful to move beyond experimental survey methods and into more realistic presentation  
570 modalities, such as mock trials, to establish greater ecological validity.

571           Our study design did not permit investigation of potential interactive relationships  
572 between psychological and neurobiological evidence. In real criminal trials, both types of  
573 evidence might be presented together. Inclusion of a combined condition would address whether  
574 their joint presentation might have multiplicative, or perhaps antagonistic, effects on attributions  
575 of responsibility and punishment.

576           Finally, interpretation of evidence likely depends on the type of crime and mental health  
577 condition portrayed. Our vignettes described a sexual assault in order to increase the plausibility  
578 of the use of involuntary hospitalization, but this choice departs from other studies in this body  
579 of literature. Likewise, the defendant's mental condition was defined as an "impulse control  
580 disorder." This decision was made to minimize unknown preconceptions about culturally loaded  
581 labels such as psychopathy, schizophrenia, and psychosis, thus differentiating this study from  
582 others of its kind [7, 8, 21]. Future studies should consider controlling these cross-study  
583 differences or consider other theoretically-motivated causes of behavior including prototypically  
584 physical disorders (e.g., bipolar disorder) as well as prototypically psychological disorders (e.g.,  
585 adjustment disorder).

586           Limitations notwithstanding, these findings are important for criminal law procedure, and  
587 particularly for policy makers, because they highlight a potential contextual effect that has not  
588 been examined in previous research. Specifically, policy makers must confront the question of  
589 how to manage the effects that we observed. For example, when neuroscientific evidence is  
590 introduced to support mental illness arguments, should it be accompanied *pro forma* by  
591 information about its potentially biasing effects? Should it be accompanied by information about  
592 the defendant's amenability to treatment? When may neuroscience evidence stand alone, and  
593 when must it be accompanied by corresponding behavioral evidence? Should judges be required

594 to receive legal education on neuroscience evidence? Should jurors be entitled (or required?) to  
595 review the treatment options or mandates that would apply if the defendant is excused on  
596 grounds of mental illness? Additional scholarship is needed to examine these and other practical  
597 applications of this research.

## 598 **Acknowledgements**

599 We thank Felipe De Brigard, Walter Sinnott-Armstrong, and Peter Reiner for helpful comments.  
600 This publication was made possible through the support of a grant from the John Templeton  
601 Foundation via the Summer Seminars on Neuroscience and Philosophy at Duke University. The  
602 opinions expressed in this publication are those of the authors and do not necessarily reflect the  
603 views of the John Templeton Foundation.

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640 **Reference List**

- 641
- 642 1. Farahany NA. Neuroscience and behavioral genetics in US criminal law: an empirical  
643 analysis. *J Law Biosci.* 2015;2: 485–509. doi:10/gd2zgk
- 644 2. Weisberg D, Taylor J, Hopkins E. Deconstructing the Seductive Allure of Neuroscience  
645 Explanations. *Judgm Decis Mak.* 2015; 429–441.
- 646 3. McCabe DP, Castel AD. Seeing is believing: the effect of brain images on judgments of  
647 scientific reasoning. *Cognition.* 2008;107: 343–352. doi:10/dzz6xj
- 648 4. Farah MJ, Hook CJ. The Seductive Allure of “Seductive Allure.” *Perspect Psychol Sci.*  
649 2013;8: 88–90. doi:10/gfjzhn
- 650 5. Hook CJ, Farah MJ. Look again: effects of brain images and mind-brain dualism on lay  
651 evaluations of research. *J Cogn Neurosci.* 2013;25: 1397–1405. doi:10/gfjzhm
- 652 6. Michael RB, Newman EJ, Vuorre M, Cumming G, Garry M. On the (non)persuasive power  
653 of a brain image. *Psychon Bull Rev.* 2013;20: 720–725. doi:10/f45sbc
- 654 7. Aspinwall LG, Brown TR, Tabery J. The Double-Edged Sword: Does Biomechanism  
655 Increase or Decrease Judges’ Sentencing of Psychopaths? *Science.* 2012;337: 846–849.  
656 doi:10.1126/science.1219569
- 657 8. Greene E, Cahill BS. Effects of neuroimaging evidence on mock juror decision making.  
658 *Behav Sci Law.* 2012;30: 280–296. doi:10/fz3q2z
- 659 9. Capestany BH, Harris LT. Disgust and biological descriptions bias logical reasoning during  
660 legal decision-making. *Soc Neurosci.* 2014;9: 265–277. doi:10/gfb7jt
- 661 10. Marshall J, Lilienfeld SO, Mayberg H, Clark SE. The role of neurological and psychological  
662 explanations in legal judgments of psychopathic wrongdoers. *J Forensic Psychiatry*  
663 *Psychol.* 2017;28: 412–436. doi:10/gfjzkh
- 664 11. Shariff AF, Greene JD, Karremans JC, Luguri JB, Clark CJ, Schooler JW, et al. Free Will  
665 and Punishment: A Mechanistic View of Human Nature Reduces Retribution. *Psychol Sci.*  
666 2014;25: 1563–1570. doi:10/f6dr42
- 667 12. Greene J, Cohen J. For the law, neuroscience changes nothing and everything. *Philos*  
668 *Trans R Soc Lond B Biol Sci.* 2004;359: 1775–85.
- 669 13. Morse SJ. Brain overclaim syndrome and criminal responsibility: A diagnostic note. *Ohio St*  
670 *J Crim L.* 2005;3: 397.
- 671 14. Dennett DC. *Elbow Room: The Varieties of Free Will Worth Wanting.* MIT Press; 2015.
- 672 15. Baumeister RF, Mele AR, Vohs KD, editors. *Free will and consciousness: how might they*  
673 *work?* New York: Oxford University Press; 2010.

- 674 16. Berryessa CM. Jury-Eligible Public Attitudes Toward Biological Risk Factors for the  
675 Development of Criminal Behavior and Implications for Capital Sentencing. *Crim Justice*  
676 *Behav.* 2017;44: 1073–1100. doi:10.1177/0093854817716485
- 677 17. Chandler JA. The use of neuroscientific evidence in Canadian criminal proceedings. *J Law*  
678 *Biosci.* 2015;2: 550–579. doi:10/gc9z4h
- 679 18. Fuss J, Dressing H, Briken P. Neurogenetic evidence in the courtroom: a randomised  
680 controlled trial with German judges. *J Med Genet.* 2015; jmedgenet-2015-103284.  
681 doi:10.1136/jmedgenet-2015-103284
- 682 19. Hardcastle VG, Lamb E. What difference do brain images make in US criminal trials? *J*  
683 *Eval Clin Pract.* 2018;24: 909–915. doi:10/gd3wz8
- 684 20. McCabe DP, Castel AD, Rhodes MG. The influence of fMRI lie detection evidence on  
685 juror decision-making. *Behav Sci Law.* 2011;29: 566–577. doi:10/fj6zsc
- 686 21. Saks MJ, Schweitzer NJ, Aharoni E, Kiehl KA. The impact of neuroimages in the  
687 sentencing phase of capital trials. *J Empir Leg Stud.* 2014;11: 105–131.
- 688 22. LaDuke C, Locklair B, Heilbrun K. Neuroscientific, Neuropsychological, and Psychological  
689 Evidence Comparably Impact Legal Decision Making: Implications for Experts and Legal  
690 Practitioners. *J Forensic Psychol Res Pract.* 2018;18: 114–142. doi:10/gfb7j2
- 691 23. Rempel RJ, Glenn AL, Cox J. Biological Evidence Regarding Psychopathy Does Not  
692 Affect Mock Jury Sentencing. *J Personal Disord.* 2018; 1–21.  
693 doi:10.1521/pedi\_2018\_32\_337
- 694 24. Blakey R, Kremsmayer TP. Unable or Unwilling to Exercise Self-control? The Impact of  
695 Neuroscience on Perceptions of Impulsive Offenders. *Front Psychol.* 2018;8. doi:10/gctvnr
- 696 25. Schweitzer NJ, Saks MJ, Murphy ER, Roskies AL, Sinnott-Armstrong W, Gaudet LM.  
697 Neuroimages as Evidence in a Mens Rea Defense: No Impact [Internet]. Rochester, NY:  
698 Social Science Research Network; 2011 Aug. Report No.: ID 2018114. Available:  
699 <https://papers.ssrn.com/abstract=2018114>
- 700 26. Ahn W-K, Proctor CC, Flanagan EH. Mental Health Clinicians' Beliefs About the Biological,  
701 Psychological, and Environmental Bases of Mental Disorders. *Cogn Sci.* 2009;33: 147–  
702 182. doi:10/bb5w32
- 703 27. Levin I, Schneider S, Gaeth G. All Frames Are Not Created Equal: A Typology and Critical  
704 Analysis of Framing Effects. *Organ Behav Hum Decis Process.* 1998;76: 149–199.
- 705 28. West RF, Meserve RJ, Stanovich KE. Cognitive sophistication does not attenuate the bias  
706 blind spot. *J Pers Soc Psychol.* 2012;103: 506–519. doi:10.1037/a0028857
- 707 29. Sokol-Hessner P, Hsu M, Curley NG, Delgado MR, Camerer CF, Phelps EA. Thinking like  
708 a trader selectively reduces individuals' loss aversion. *Proc Natl Acad Sci U S A.* 2009;106:  
709 5035–5040. doi:10/fhdrcw

- 710 30. Buhrmester M, Kwang T, Gosling SD. Amazon's Mechanical Turk: A New Source of  
711 Inexpensive, Yet High-Quality, Data? *Perspect Psychol Sci J Assoc Psychol Sci*. 2011;6:  
712 3–5. doi:10.1177/1745691610393980
- 713 31. Clifford S, Jewell RM, Waggoner PD. Are samples drawn from Mechanical Turk valid for  
714 research on political ideology? *Res Polit*. 2015;2: 2053168015622072.  
715 doi:10.1177/2053168015622072
- 716 32. Murphy KR, Myers B. *Statistical power analysis: A simple and general model for traditional  
717 and modern hypothesis tests*, 2nd ed. Mahwah, NJ, US: Lawrence Erlbaum Associates  
718 Publishers; 2004.
- 719 33. Stanovich KE, West RF. Reasoning independently of prior belief and individual differences  
720 in actively open-minded thinking. *J Educ Psychol*. 1997;89: 342–357. doi:10/cpvchx
- 721 34. Kaufman EA, Xia M, Fosco G, Yaptangco M, Skidmore CR, Crowell SE. The Difficulties in  
722 Emotion Regulation Scale Short Form (DERS-SF): Validation and Replication in  
723 Adolescent and Adult Samples. *J Psychopathol Behav Assess*. 2016;38: 443–455.  
724 doi:10.1007/s10862-015-9529-3
- 725 35. To screen for potential confounding factors, we examined the distributions of several  
726 demographic variables (age, gender, political affiliation, education, and income) across all  
727 levels of our independent variables. In each case, the data were normally distributed.
- 728 36. Partial eta-squared effect sizes are interpreted using the following benchmark values,  
729 suggested by Richardson (2011):  $.0588 \leq \text{medium} < .1379$ .
- 730 37. To assess whether the introduction of any health related information altered prison  
731 sentences, a two-tailed paired t-test was conducted within the healthy condition before and  
732 after disclosure that the defendant was, in fact, in good mental health. As anticipated, there  
733 was a null effect,  $t(106) = 1.46$ ,  $p = .15$ , indicating that participants likely assumed that the  
734 defendant was of sound mental health, by default, before any evidence was presented.
- 735 38. Participants' change in prison sentence and involuntary hospitalization terms were  
736 negatively correlated,  $r(223) = -.38$ ,  $p < .001$ , suggesting a perceived trade off between  
737 these decision outcomes.
- 738 39. An identical effect was observed for the effect of mental health status on the change in a  
739 participant's prison sentence recommendation pre- and post-mental health status  
740 manipulation. Deontological concerns fully accounted for the relationship between mental  
741 health status and change in prison sentence recommendation, whereas consequentialist  
742 concerns did not.
- 743 40. A median split was performed to circumvent problems of multicollinearity observed using a  
744 more conventional, linear regression method.
- 745 41. Hampton J. The Moral Education Theory of Punishment. *Philos Public Aff*. 1984;13: 208–  
746 238.
- 747 42. Bag B, Yilmaz S, Kirpinar I. Factors influencing social distance from people with  
748 schizophrenia. *Int J Clin Pract*. 2006;60: 289–294. doi:10/dnxszr

- 749 43. Grausgruber A, Meise U, Katschnig H, Schöny W, Fleischhacker WW. Patterns of social  
750 distance towards people suffering from schizophrenia in Austria: a comparison between  
751 the general public, relatives and mental health staff. *Acta Psychiatr Scand.* 2007;115: 310–  
752 319. doi:10/dn5w7m
- 753 44. Haslam N, Kvaale EP. Biogenetic Explanations of Mental Disorder: The Mixed-Blessings  
754 Model. *Curr Dir Psychol Sci.* 2015;24: 399–404. doi:10/f7tgdn
- 755 45. Martin JK, Pescosolido BA, Olafsdottir S, McLeod JD. The construction of fear: Americans’  
756 preferences for social distance from children and adolescents with mental health problems.  
757 *J Health Soc Behav.* 2007;48: 50–67. doi:10/czntvb

## 758 **Supporting information**

759 **S1 Appendix. Experimental Stimuli**

760 **S2 Dataset. Data Underlying Findings**