

Moral Dynamics

Grounding Moral Judgment in Intuitive Physics and Intuitive Psychology

Felix A. Sosa*

Department of Brain and Cognitive Sciences, MIT

Tomer Ullman

Department of Psychology, Harvard University

Joshua B. Tenenbaum

Department of Brain and Cognitive Sciences, MIT

Samuel J. Gershman

Department of Psychology and Center for Brain Science, Harvard University

Tobias Gerstenberg

Department of Psychology, Stanford University

Abstract

When holding others morally responsible, we care about what they did and what they thought. Traditionally, research in moral psychology has relied on vignette studies, in which the protagonist's actions and thoughts are explicitly communicated. Recent studies have begun to employ visual stimuli, and some have postulated a direct link from processing visual features to making moral judgments. We embrace the advent of visual stimuli in moral psychology, but believe that the connection between visual processing and moral judgments is mediated by an inference about what the observed action reveals about the agent's mental states. We formalize moral judgments as computations over an intuitive theory of physics combined with an intuitive theory of mind. Knowing that mental states lead to action (e.g., the belief that someone is in harm's way and the desire to help them stimulates a decision to shove them out of harm's way), and that these actions are constrained by physics (the shove has to be forceful enough, aimed in the right direction, timed appropriately, etc.), allows an observer to make powerful inferences about moral responsibility. Two experiments show that this model captures moral judgments about physical scenes, both qualitatively and quantitatively.

Keywords: moral judgment, effort, intuitive physics, intuitive psychology

*Corresponding author: Felix Sosa (felixanthonyososa@gmail.com).

Introduction

In a popular image, three wise monkeys advise us: see no evil, hear no evil, speak no evil. But do we actually see evil, in the way we see shapes, or colors, or monkeys? When viewing simple shapes moving around a 2D world, people spontaneously and consistently attribute goals and intentions to them (Heider & Simmel, 1944), including social motivations (Ullman et al., 2009). Even young children appear to draw consistent conclusions about the goals, intentions, and relations of actors in simple visual vignettes (e.g. Gergely & Csibra, 2003; Gergely, Nádasdy, Csibra, & Bíró, 1995; Hamlin, Wynn, & Bloom, 2007), and (at slightly older ages) will act to punish morally bad actors (Hamlin, Wynn, Bloom, & Mahajan, 2011). Recent work in neuroscience has shown selective activation in the posterior superior temporal sulcus when viewing agent animations, with dissociable responses for goal directed action by individual actors and social interactions (Isik, Koldewyn, Beeler, & Kanwisher, 2017; Vander Wyk, Hudac, Carter, Sobel, & Pelphrey, 2009).

In cognitive science, there is a long tradition of attempting to formally link perception and psychological attributions (such as intention) by identifying relevant visual cues in a scene. This line of research can be traced back at least to Michotte (1946/1963), and extends to current work on the visual cues that could underpin perceptions of agency, intention, and various interactions such as courting, chasing, and protecting (e.g. Hubbard, 2005; Scholl & Gao, 2013). Recent work has suggested specifically that moral judgments can be explained by the visual processing of kinematic features, such as the velocity of a car hitting a man, or the distance a person traveled to push someone into harm's way (Caruso, Burns, & Converse, 2016; De Freitas & Alvarez, 2018; Iliev, Sachdeva, & Medin, 2012; Nagel & Waldmann, 2012). In line with the fast, automatic, early-developing, and consistent nature of these judgments, these accounts propose a direct mapping from visual features to moral judgments (such that, for example, traveling longer distances to harm maps onto morally worse judgments by others).

However, a great deal of prior work on moral judgment has focused on the deliberative and abstract components that go into a moral calculation. This line of research (which often relies on carefully written vignettes rather than visual stimuli) has demonstrated that both a person's causal role, and the person's inferred mental states are key determinants of moral judgments (Cushman, 2008; Gerstenberg et al., 2018; Lagnado & Gerstenberg, 2017; Lagnado, Gerstenberg, & Zultan, 2013; Malle, Guglielmo, & Monroe, 2014; Patil, Calò, Fornasier, Cushman, & Silani, 2017; Shaver, 1985; Waldmann, Nagel, & Wiegmann, 2012; Weiner, 1995; Young, Cushman, Hauser, & Saxe, 2007; Young & Saxe, 2008): People judge a person more severely when that person caused the bad outcome (Alicke, 1992; Cushman, 2008), and when that person intended to cause a bad outcome (Kleiman-Weiner, Gerstenberg, Levine, & Tenenbaum, 2015; Lagnado & Channon, 2008). Moral judgments are not only sensitive to whether someone caused or intended an outcome, but also to the way the outcome was brought about (Jara-Ettinger, Kim, Muentener, & Schulz, 2014; Waldmann & Dieterich, 2007). While moral dilemmas presented as vignettes, such as the trolley problem, are a rich source for empirical exploration (Foot, 1978; Waldmann et al., 2012), they have their limitations. Vignettes may fail to trigger relevant perceptual processing and related downstream processes, similar to asking someone to solve a physics problem involving trajectories with pen and paper, instead of throwing a ball at them and

60 asking them to catch it.

61 These two approaches to formalizing moral judgments, one focusing on perceptual
62 processing and the other on inferences of cause and intention, can seem incongruous. But
63 a full model will have to incorporate both. Here, we propose such a synthesis. We believe
64 that the route from visual features to moral judgments is mediated by people’s intuitive
65 understanding of how the world works, one that encompasses both an intuitive theory of
66 mind (Wellman & Gelman, 1992), and an intuitive theory of physics (Battaglia, Hamrick, &
67 Tenenbaum, 2013; Gerstenberg & Tenenbaum, 2017; Goodman, Tenenbaum, & Gerstenberg,
68 2015; Kubricht, Holyoak, & Lu, 2017; Ullman, Spelke, Battaglia, & Tenenbaum, 2017).
69 These intuitive theories support rapid inferences about a person’s mental states and the
70 causal structure of a scenario. It is through this understanding of scenarios that people relate
71 observed physical actions to mental states such as beliefs, desires, and intentions (Baker,
72 Jara-Ettinger, Saxe, & Tenenbaum, 2017; Baker, Saxe, & Tenenbaum, 2009; Battaglia et al.,
73 2013; Jara-Ettinger, Gweon, Schulz, & Tenenbaum, 2016), and evaluate the causal roles that
74 physical actions played in bringing about the outcome (Gerstenberg, Goodman, Lagnado, &
75 Tenenbaum, 2015; Gerstenberg, Peterson, Goodman, Lagnado, & Tenenbaum, 2017). Our
76 computational synthesis is sensitive to visual features and allows for fast and automatic
77 processing, but not because such features are mapped directly to moral judgments. Rather,
78 these features are indicative of the agent’s mental state, and it is these mental states that
79 are the input to the moral calculus.

80 The rest of the paper is organized as follows. We first describe a model of moral
81 judgment operating over an intuitive theory of mind and an intuitive theory of physics. We
82 then examine this model using two empirical studies. In Experiment 1, we replicate an
83 experiment that links visual cues to moral judgment (Iliev et al., 2012), and show how our
84 model accounts for the results by inferring the agent’s desire to do harm via the physical
85 effort it exerted. In Experiment 2, we test participants’ moral intuitions in a wider range
86 of situations, and elicit graded judgments which provide a stronger test for the model’s
87 predictions. We discuss the implications of our findings, limitations of our current model,
88 as well as a roadmap for addressing these limitations.

89 Moral Dynamics Model

90 Our model connects perceptual depictions of an agent’s actions to reasoning about the
91 underlying mental states of the agent. The model combines ideas from recent formalizations
92 of intuitive psychology (for reasoning about hidden mental states given actions) and intuitive
93 physics (for reasoning about cost as physical effort). Following prior structured generative
94 approaches to intuitive psychology (e.g. Baker et al., 2017, 2009; Jara-Ettinger et al., 2016;
95 Kleiman-Weiner, Shaw, & Tenenbaum, 2017; Ullman et al., 2009), we model an observer
96 who reasons about an agent’s mental states by inverting the generative process by which
97 mental states give rise to actions. Following recent prior work on intuitive physics (e.g.
98 Battaglia et al., 2013; Sanborn, Mansinghka, & Griffiths, 2013; Smith & Vul, 2012; Ullman,
99 Stuhlmüller, Goodman, & Tenenbaum, 2018), we constrain this generative process to obey
100 noisy Newtonian mechanics: actions correspond to forces applied by an agent to a patient,
101 where the amount of force is monotonically related to the agent’s effort.

102 As an overview, our “Moral Dynamics” model infers an agent’s utilities, and predicts
103 that people’s negative moral judgments are related to how much an agent desires to harm

104 a patient. We also assume that agents act to achieve desired rewards, and that actions are
 105 associated with a cost in the form of physical effort. Given that a rational agent trades off
 106 cost and reward (taking costly actions to receive a greater reward than the cost expended),
 107 an observer can use the effort an agent expends as indicative of the value the agent places
 108 on achieving an outcome. If an agent undertakes a great cost to achieve a harmful outcome,
 109 that agent likely expected a large reward for causing harm, and thus should be morally
 110 blamed to a high degree. We next discuss in more detail the theoretical background and
 111 implementation of the framework.

112 **Computational Framework**

113 We model intuitive psychological reasoning using Bayesian Theory of Mind (see e.g.
 114 Baker et al., 2017, 2009; Jara-Ettinger et al., 2016). This framework assumes that people
 115 think of others as goal-directed agents who choose actions to maximize their expected
 116 reward, subject to their beliefs, constraints, and abilities (see also Gershman, Gerstenberg,
 117 Baker, & Cushman, 2016). The underlying notion that people use a ‘principle of rationality’
 118 to reason about the mental states of others has a long history (Dennett, 1987), and the
 119 recent avenue of work in Bayesian Theory of Mind has shown how to use this principle
 120 to quantitatively capture human reasoning about mental states. For the purposes of our
 121 model, we limit ourselves to a version of the framework dubbed ‘the Naive Utility Calculus’
 122 (Jara-Ettinger et al., 2016), according to which people believe that others act to maximize
 123 their state-dependent rewards, and to minimize action-dependent costs:

$$U(s, a) = R(s) - C(a), \quad (1)$$

124 where U is an agent’s utility, a combination of the reward R derived from world state s ,
 125 and the cost C of taking action a .

126 To this underlying framework we add the following three simple assumptions: 1)
 127 We limit ourselves to cases in which the cost C for taking an action is proportional to
 128 the physical effort necessary to take that action. 2) We assume that an agent’s reward
 129 can depend on the utility of another agent (cf. Ullman et al., 2009). 3) We assume an
 130 observer’s moral evaluations are proportional to the inferred reward that the agent derives
 131 from helping or hindering the patient (cf. Gerstenberg et al., 2018). We consider each of
 132 these assumptions in turn, and then show how their combination leads to an account of
 133 moral judgment from visual scenes.

134 **Physical effort.** Physical effort features prominently in decision making and moral
 135 judgment (Jara-Ettinger et al., 2014; Kurniawan et al., 2010). Jara-Ettinger et al. (2014)
 136 demonstrated that transgressors are judged more harshly for taking more costly actions
 137 to bring about a negative outcome. In those studies, participants were given multiple
 138 vignettes involving the same outcome (e.g., stealing someone’s wallet), and judged the
 139 vignette involving the greatest amount of effort as depicting the worst offender. Even
 140 young children are sensitive to the physical effort required by an action, and take it into
 141 account when determining the goal of an agent (Jara-Ettinger et al., 2014; Liu, Ullman,
 142 Tenenbaum, & Spelke, 2017).

143 The rationale of using cost to infer utility and through that make moral judgments
 144 carries through with other types of cost as well (such as risk or mental effort), but here we

145 limit ourselves to inferences about physical effort. We formalize physical effort in terms of
 146 Newtonian dynamics, which are broadly consistent with human intuitive physical reasoning
 147 (Battaglia et al., 2013; Bramley, Gerstenberg, Tenenbaum, & Gureckis, in press; Hamrick,
 148 Battaglia, Griffiths, & Tenenbaum, 2016; Sanborn et al., 2013; Ullman et al., 2017, 2018).
 149 As detailed in the Appendix, we model physical effort as the amount of force that an agent
 150 expended to bring about the outcome. Importantly, when considering how much effort an
 151 agent took to harm another, we only count the physical effort used to achieve that goal.
 152 For example, if an agent ran around in circles before or after taking intentional action to
 153 harm another agent, we would not count that effort associated with running in circles as
 154 effort towards accomplishing its goal.

155 **Helping or harming.** We label the utility of the patient as U_P , and the reward the
 156 agent receives for helping or harming the patient as $R_A(U_P)$. A pro-social attitude (i.e. high
 157 reward for helping) between the agent and patient can be captured as a positive relationship
 158 between R_A and U_P and an anti-social attitude (i.e. high reward for harming) between the
 159 agent and patient can be captured as a negative relationship. If a pro-social relationship
 160 exists, whatever states and actions increases the patient’s utility will also increase the agent’s
 161 reward, and the agent will take actions to move the patient into high-reward states or reduce
 162 the patient’s costs (modulated by the agent’s own costs). This simplified model of ‘helping
 163 and hindering’ can quantitatively account for people’s reasoning about social goals (Ullman
 164 et al., 2009), and the choice patterns of pre-verbal infants (Hamlin, Ullman, Tenenbaum,
 165 Goodman, & Baker, 2013).

166 **From inferred desires to moral judgments.** We assume that moral judgments
 167 depend on people’s beliefs about an agent’s desires. That is, people will judge the agent
 168 more negatively, in proportion to the inferred reward that the agent derives from harming
 169 an innocent other.

Putting these assumptions together, we model people’s negative moral judgments
 about an agent A, $J(A)$, as being proportional to the inferred positive reward A receives
 for the negative outcome utility of patient P, $R_A(U_P)$, which can be approximated by the
 amount of effort that A was willing to exert to bring about that outcome. The effort that
 A exerted in a scenario is defined as the sum of the costs A incurred c_A for taking some
 action a_t at every time point in the scenario t :

$$J(A) \approx R_A(U_P) \propto \sum_{t=0}^T c_A(a_t). \quad (2)$$

170 We made a number of simplifying assumptions in this calculation. In general, psychological
 171 costs encompass more than physical effort, such as time delay or mental effort (Kool &
 172 Botvinick, 2018). Also, physical effort is not just an integral of the force generated over
 173 time, but subject to biological notions of expendable energy and fatigue (Hills, Mokhtar, &
 174 Byrne, 2014). Pro-social and anti-social relationships are more than just a utility-to-reward
 175 transformation, and moral evaluations depend on more than the inferred social relationship
 176 between agents (Waldmann et al., 2012).

177 Nonetheless, we believe this framework is flexible enough to capture many natural-
 178 istic decision problems, and it provides the core mechanics that future work can build on,
 179 adding in more varied notions of cost, effort, and relationships. Because the notions of
 180 force and effort play a central role in our model, we refer to it as *Moral Dynamics* model,

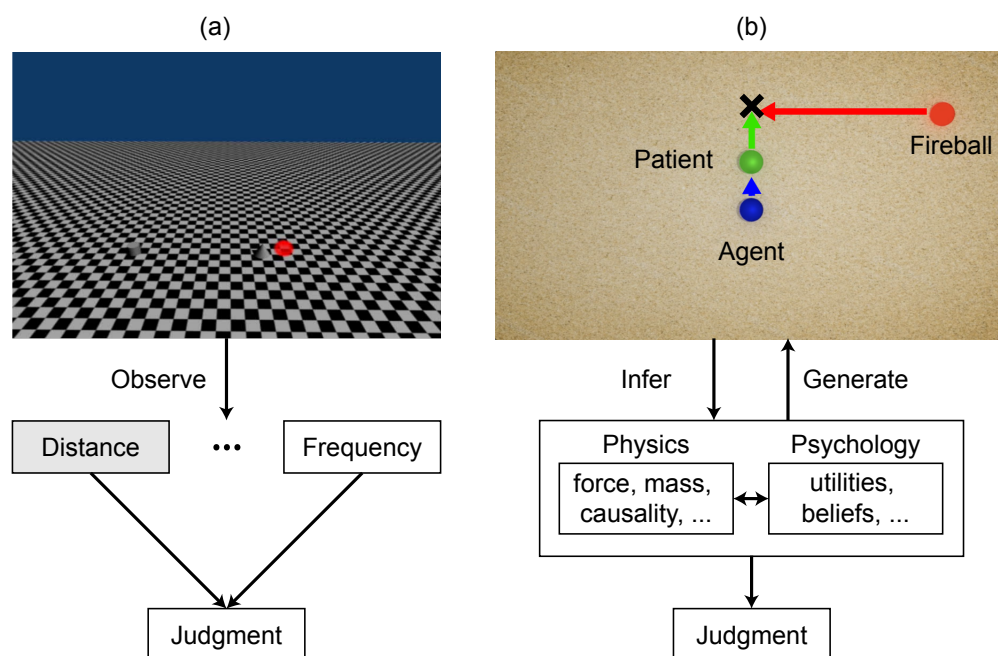


Figure 1. (a) Example stimuli from Iliev et al. (2012), and representation of their theory of moral judgment. An observer extracts relevant features of a scene (e.g. distance in the above example) and uses those features directly to form their moral judgment (e.g. greater distance means more negative judgment). (b) Example stimuli from our Experiments 1 and 2, and a representation of the *Moral Dynamics* model. An observer infers latent variables related to physics and psychology. Specifically, observers infer the utility an agent attaches to helping or harming by inferring the effort they expended in the scene. These variables inform their moral judgment. Each video shows an agent (blue), patient (green), and fireball (red). Patients cannot see fireballs and are burned by them, agents can see fireballs and are not burned by them. Lines indicate trajectories in the video, ‘X’ marks the collision of the patient and fireball.

181 in contrast with a model termed *Moral Kinematics* by Iliev et al. (2012), which predicts
 182 moral judgments based on perceptual/kinematic features such as distance, angle, contact,
 183 and velocity.

184 Model Implementation and Domain

185 We consider a simple domain, based on Iliev et al. (2012), in which different visual
 186 scenarios show agents interacting with, and potentially harming, other agents, while exerting
 187 physical forces.¹ Figure 1(b) illustrates our experimental setup. In each scenario, a video
 188 shows an agent, patient, and fireball. In this domain, the agent can perceive the fireball
 189 and is not harmed upon contact with it, while the patient cannot perceive the fireball and
 190 is harmed upon contact with it.

¹We used the 2D physics engine Pymunk (www.pymunk.org) to generate the scenarios in our experiments, and the 3D physics engine Blender (www.blender.org) to render them.

191 For each agent, the cost for taking a set of actions, C_A , is the sum of the forces that
192 the agent generates on itself (cf. Luo & Baillargeon, 2005). Specifically, at each discrete
193 time step, t , an agent applies an impulse, I_t , to itself. An agent’s effort at that time step is
194 proportional to the magnitude of that impulse.

195 Given this domain, we can use different starting conditions and trajectories to vary
196 the amount of effort an agent is perceived to expend to harm a patient. According to
197 our model, differences in the inferred amounts of effort will translate into different moral
198 evaluations.

199 Experiment 1

200 Our first experiment seeks to qualitatively test the computational model developed
201 above, and examine whether participants’ judgments can be explained by assuming that
202 they infer the reward the agent has for harming the patient via the amount of effort the
203 agent exerted. Our experiment was closely modeled after Experiment 2 in Iliev et al. (2012).
204 So, an additional goal for this experiment is to verify that our stimuli elicit similar responses
205 to Iliev et al. (2012), so that meaningful comparisons can be drawn when we later expand
206 the stimulus set.

207 Methods

208 **Participants.** 46 participants ($M_{age} = 34.5$, $SD_{age} = 10.4$, 11 female) were recruited
209 via Amazon Mechanical Turk, and compensated for their time. Both Experiment 1 and 2
210 were run using Psiturk (Gureckis et al., 2016).

211 **Stimuli.** In Iliev et al.’s (2012) Experiment 2, participants saw pairs of videos.
212 Each pair tested the effect of a kinematic feature on moral judgment. The videos in each
213 pair differed with respect to at least one of the following factors: The distance the harming
214 agent traveled; whether the agent made contact with the patient; how many times the agent
215 touched the patient; how long the agent made contact with the patient; the force the agent
216 exerted on the patient.

217 We focused on the video pairs whose physical dynamics could be captured in our
218 2D, top-view physics engine implementation, and used 9 of the original 15 pairs.² The 9
219 included video pairs were tailored to be similar to the original stimuli used by Iliev et al.
220 (2012), with minor differences beyond the 3D-view vs 2D-top-view (see Figure 1 (a) and
221 (b)): In the original experiment, the agent was a white cylinder, the patient was a white
222 cone, and the floor of the scene was checkerboard. In our videos, the agent was a blue
223 sphere, the patient was a green sphere, and the floor of the scene was visually similar to
224 sand. In both the original experiment and in our stimuli the fireball is a red sphere.

²The stimuli that were left out can all be theoretically incorporated into the model, but involve complications that are not relevant for the question at hand. The 6 videos left out were: Three pairs with motion up and down ledges (while such motion can be captured in 2D, it would require a side-view rather than a top-view, and we opted to keep the stimuli uniform in viewpoint); two pairs with agents sliding for long distances after minor collisions (which would require either near-zero friction, very strong agents, or a very low velocity patient after the collision); one pair with an agent that entered the scene from outside the frame (requiring inference over the unseen physics that led to its arrival).

225 **Design and Procedure.** As in Iliev et al. (2012), participants were instructed
 226 that they would see pairs of videos involving imaginary creatures (Blues and Greens) and
 227 a fireball. Participants were further informed that each video shows a situation in which
 228 Green collided with the fireball, and that their task was to judge in which video Blue’s
 229 actions were worse.

230 Participants then viewed a set of familiarization videos that showed Blues, Greens,
 231 and fireballs interacting. The familiarization videos informed participants that Blues and
 232 Greens were intelligent, social creatures, and that fireballs were inanimate objects. As in
 233 Iliev et al. (2012), participants were told that fireballs were sometimes moved by magnetic
 234 winds. Participants were informed that Greens *could not* see fireballs and were burned
 235 when they touched them, whereas Blues *could* see fireballs and were not burned when they
 236 touched them. Finally, participants learned that while Blues and Greens usually got along,
 237 there were some reported instances in which Blues harmed Greens. Participants were told
 238 that they would see such instances, and would be asked to evaluate what Blue did. Before
 239 starting the experiment, participants were required to pass a comprehension check.

240 The comprehension check ensured participants knew only Greens could be harmed
 241 by fireballs, only Blues could see fireballs, and that fireballs could sometimes be moved
 242 by magnetic winds. Participants were only allowed to move on to the main experiment
 243 if they correctly answered all comprehension check questions. If a participant failed the
 244 comprehension check, they had to go through the introduction and familiarization videos
 245 again, and re-take the comprehension check.

246 During the experiment, each participant was shown 9 pairs of videos. The order of the
 247 video pairs was randomized. When viewing a pair, participants had to watch both videos
 248 twice, going from the video presented on the left of the screen to the one on the right, and
 249 back again. The left/right placement of videos was counterbalanced across participants.
 250 After viewing both videos twice, participants responded to the prompt “The action of Blue
 251 was...” with one of six possible responses (presented from left to right): “much worse in the
 252 left video”, “worse in the left video”, “somewhat worse in the the left video”, “somewhat
 253 worse in the the right video”, “worse in the right video”, and ‘much worse in the right video”

254 At the end of the experiment, participants provided demographic information, and
 255 were invited to share any comments.

256 Results

257 To best compare our results to those of Iliev et al. (2012), we followed their analysis
 258 procedure and binarized participant responses. Responses were coded as 1 if the video in
 259 column A of Figure 2 was judged as worse, and 0 otherwise. Iliev et al. (2012) originally
 260 found that, based on kinematic features, the agent’s actions in videos in column A were
 261 predicted and judged as worse than those in column B. Figure 2 shows the percentage of
 262 participants that marked the video in column A as worse.

Figure 2 also shows the participants’ responses to the equivalent stimuli in Iliev et al. (2012), and the predictions of our model. We used Luce’s choice rule as described in Equation 3 to transform the continuous model predictions into a probability of choosing one video over another.

$$P(A) = \frac{\text{Effort}_A}{\text{Effort}_A + \text{Effort}_B} \quad (3)$$

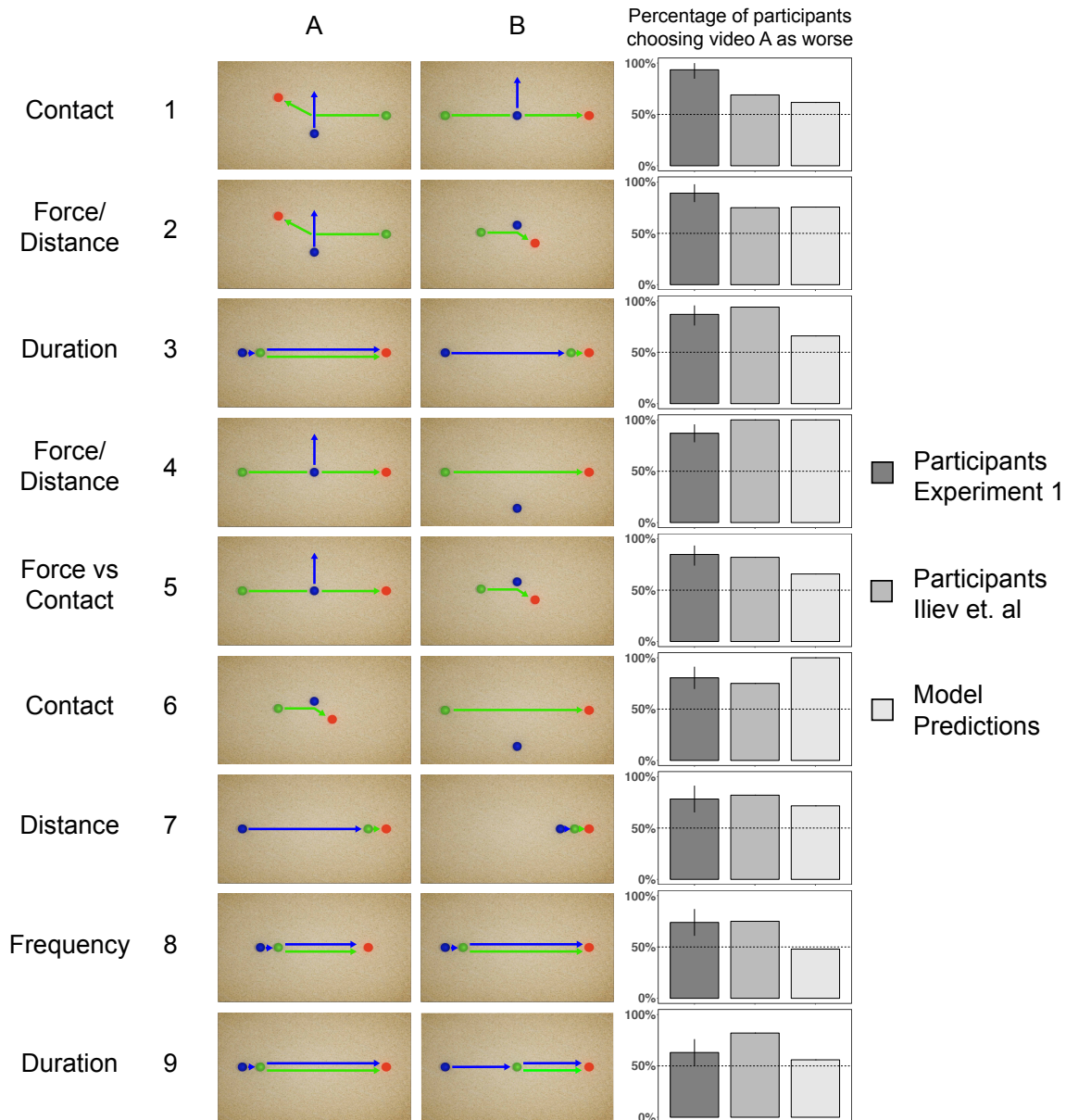


Figure 2. Experiment 1 stimuli and results. Each row shows 2 schematics of the videos shown to participants (A and B), as well as the percentage of participants who judged that what Blue did was worse in A compared to B in our replication (Experiment 1), in the original study (Moral Kinematics), and according to the model. Each pair differed with respect to a kinematic feature(s), listed to the left of each pair. The pairs are in descending order of percentage of participants choosing video A as worse. *Note:* The error bars indicate bootstrapped 95% confidence intervals.

263 Our results replicated the original results reported in Iliev et al. (2012). Across the 9
 264 pairs, a majority of participants judged the agent in video A as being worse. While there
 265 are small quantitative discrepancies between our results and what Iliev et al. (2012) found

266 (e.g. we found a stronger preference for A in pair 1 compared to Iliev et al.), we attribute
267 these differences to sampling noise (Iliev et al.’s (2012) Experiment 2 featured only 16
268 participants). For 8 out of the 9 video pairs, our model predicted the preference found in
269 both our experiment and in Iliev et al. (2012). In pair 8, our model exhibited a (very) slight
270 preference for B over A.

271 Discussion

272 The results of our experiment closely replicate what Iliev et al. (2012) found. For
273 each video pair in our study, participants judged the agent’s action to have been worse in
274 A compared to B. This successful replication suggests that our stimuli elicit similar moral
275 intuitions, despite being visually somewhat different from the ones used by Iliev et al. (2012).

276 The *Moral Dynamics* model correctly predicted participants’ preference in 8 out of
277 9 video pairs. Instead of postulating a set of visual and kinematic features that influence
278 people’s judgments (see Figure 2 leftmost column), the *Moral Dynamics* model predicted
279 this preference solely based on the effort the agent expended in each video which is diagnostic
280 for how much reward the agent placed on the patient’s harm.

281 Still, the model as realized has several limitations. For example, so far, the model
282 does not try to infer an agent’s intention. In video 8A, the agent pushes the patient
283 twice, whereas in video 8B the agent pushes the patient only once, but all the way to the
284 fireball. The two-push scenario provides salient evidence for the agent’s intention to harm
285 the patient. A plausible interpretation of what happens in 8A (on the part of a human
286 observer) is that the agent realized that its first push wasn’t sufficient to achieve the goal
287 of harming the patient, and then it decided to push again. In 8B, the agent’s movement is
288 compatible with a desire to just go in that direction while the patient happens to be in the
289 way. Since the agent’s actions are such that it expended almost identical effort in 8B (the
290 long push) and 8A (the double push), our model predicts that participants should have not
291 clear preference in this case.

292 An additional limitation of the model is that it relies on an estimation of effort
293 that is directly related to the force used in the physics simulation, while participants’ own
294 estimations of effort and force may deviate from the underlying dynamics in various ways.
295 For example, in our implementation staying still following a collision requires the active use
296 of an opposing force to cancel out the impact, while people may perceive this as simply the
297 agent staying put (see also stimuli 6A and 6B). In the next experiment, we directly address
298 the use of ground truth effort on the part of the agent by asking participants to judge how
299 much effort the agent exerted. We also expand the number of test cases in order to carry
300 out a more quantitative evaluation of the model.

301 Experiment 2

302 For Experiment 2, we turn to a quantitative examination of our model against people’s
303 judgments, expanding on the original stimulus set, and again having people judge the
304 relative moral badness of different agents’ actions. The expanded stimulus set includes seven
305 additional videos based on the first experiment conducted by Iliev et al. (2012). There, they
306 examined the effect of movement and intervention on moral judgment: whether the agent
307 intervened on the patient or on the fireball to harm the patient and whether the agent,

308 patient, or fireball were moving before the intervention. We added these videos to our set of
309 stimuli and tested whether these additional kinematic features were captured by our model.
310 Since the *Moral Dynamics* model goes from the inferred effort that the agent exerted to how
311 much reward the agent placed on harming the patient, we tested in a separate condition,
312 whether participants' estimate of how much effort the agent exerted was accurately captured
313 by the model.

314 Methods

315 **Participants.** 83 participants ($M_{age} = 35.7$, $SD_{age} = 12.7$, 42 female) were recruited
316 via Amazon Mechanical Turk.

317 **Design and Procedure.** Participants were randomly assigned to the *Effort* con-
318 dition ($N = 42$), or the *Moral* condition ($N = 41$).

319 The instructions and familiarization videos were largely identical to those of Experi-
320 ment 1. In both conditions, participants viewed the same videos with slight modifications
321 depending on the condition. This time, instead of pairs of videos in the test phase, partic-
322 ipants only viewed a single video at a time. 17 test videos were presented in randomized
323 order.

324 Participants watched each video twice before being asked to indicate their response
325 on a continuous slider. In the *Effort* condition, participants answered the question, "How
326 much effort did Blue exert in this scenario?" with the endpoints of the slider labeled "very
327 little" (0) and "very much" (100). In the *Moral* condition, the question was "How bad was
328 what Blue did?" and the endpoints were labeled "not bad" (0) and "very bad" (100).

329 Results

330 The empirical results of Experiment 2 for both conditions are summarized in Figure 3,
331 showing a schematic of the video stimuli, participants' effort and moral judgments for each
332 individual video, together with the model's predictions. Figure 4(b) and (c) show the
333 fitted effort values from the physics engine against the mean effort and moral judgments,
334 respectively. The model was fitted using separate linear regressions for each condition.

335 Participants' judgments of effort were closely aligned with the effort values from
336 our model, Spearman's $\rho = .96$, $p < .001$, 95% CI [.93, .98] (see Figure 4b). Further, the
337 mean participant judgments for how much effort the agent exerted in each video posi-
338 tively correlated with participant moral judgments for corresponding videos, $\rho = .68$, $p =$
339 $.003$, 95% CI [.55, .80] (see Figure 4a). The *Moral Dynamics* model provided a similarly
340 good fit to participants' moral judgments (see Figure 3 light bars as well as Figure 4c)
341 with $\rho = .66$, $p = .004$, 95% CI [.56, .74], 95% CI [.56, .74]. As a reminder, the model pre-
342 dicted participants' moral judgments based on how much effort the agent exerted, which is
343 diagnostic for how much reward the agent placed on harming the patient.

344 Discussion

345 The results of Experiment 2 support the idea that judgments of physical effort are
346 important for moral judgments in these visual, dynamic scenarios. The *Moral Dynamics*
347 model explains these judgments in terms of an overarching framework rather than postu-
348 lating a collection of features. Both people's judgments of effort and the effort values from

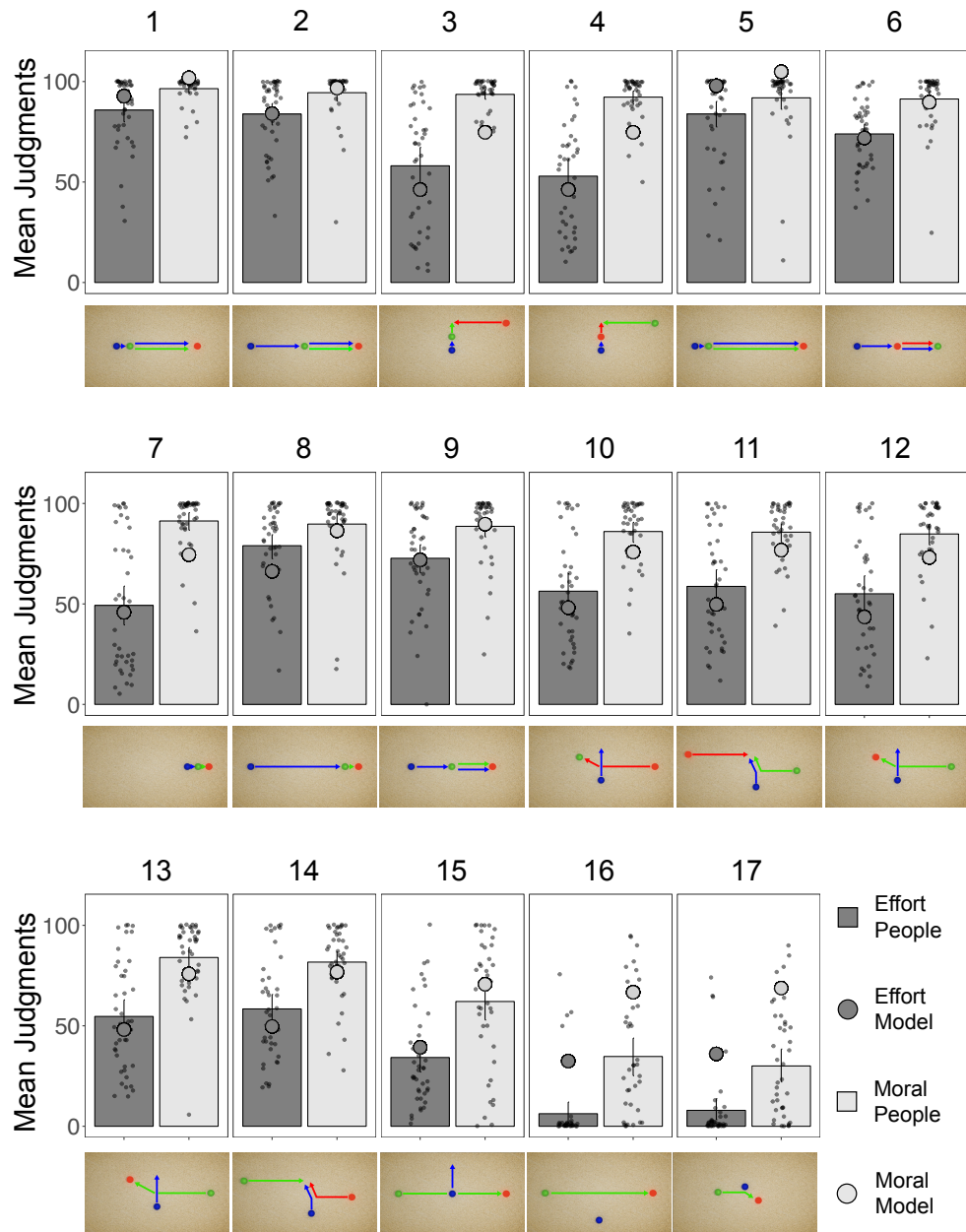


Figure 3. Experiment 2 results. Participants’ effort judgments (dark gray) and moral judgments (light gray) for each of the 17 scenarios. Bars indicate mean ratings, and small points indicate individual judgments. Error bars indicate bootstrapped 95% confidence intervals. The model predictions are superimposed as circles. Diagrams of what happened in each video are shown below participants’ effort and moral judgments. The results are ordered by descending moral judgment from worst (top left) to least bad (bottom right).

349 the physics engine correlated well with people’s moral judgments and, as predicted by the
 350 model, the more effort an agent exerted in a scenario, the worse its behavior was perceived
 351 to be (4a and c).

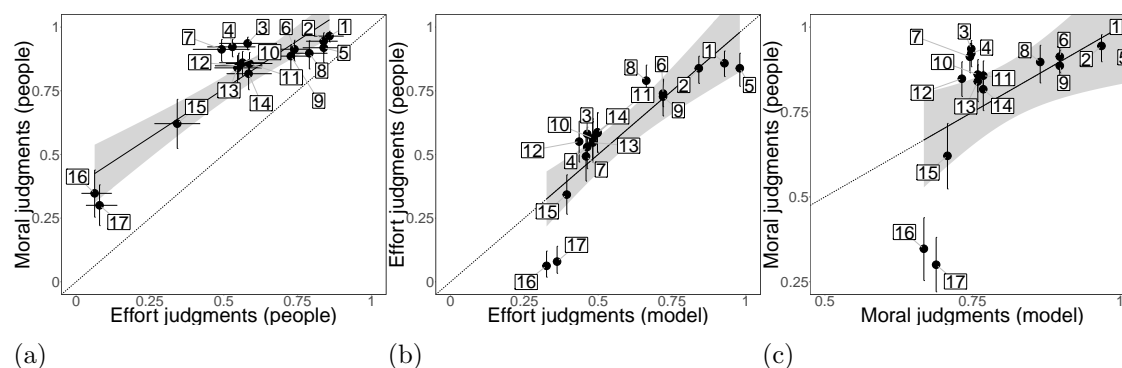


Figure 4. **Experiment 2 results.** Scatter plots of (a) participants’ moral judgments against participants’ effort judgments, (b) participants’ effort judgments against model effort predictions, and (c) participants’ moral judgments against model moral predictions. Error bars indicate bootstrapped 95% confidence intervals.

352 We take these results as supporting the proposal that judgments of physical effort play
 353 an important role in a moral calculus over these visual scenarios, as a way of estimating
 354 the intention and utility function of an agent. While explaining much of the variance,
 355 the correlation between effort (as judged by both people and the model) and people’s
 356 moral judgment is far from perfect. We attribute the missing variance to our simplifying
 357 assumptions. As we stated in Experiment 1, and further elaborate in the General Discussion,
 358 our model is likely not capturing salient additional information about the agent’s intention
 359 to harm the patient. However, such additional information can be incorporated in the future
 360 into more sophisticated mental reasoning modules in the overall framework.

361 We also found that participants’ effort judgments in the Effort condition strongly
 362 correlate with the effort values from the physics engine (Figure 4b), corroborating the
 363 growing body of work that suggests aspects of human reasoning about physics can be
 364 captured by physics engines. Given that the model accurately captures participant effort
 365 judgments, it is unlikely that deviations between the model’s concept of effort and people’s
 366 perception of effort is responsible for the unpredicted variance in the moral judgments. We
 367 also note that while the correlation between the model’s effort judgments and people’s effort
 368 perceptions is high, the linear fit to the model’s predictions of effort can deviate noticeably
 369 from the mean judgment of effort, due to the intercept term which prevents our model
 370 from inferring effort values (and moral values) of zero (for example, in scenarios 16 and 17,
 371 Figure 4b).

372 General Discussion

373 Our moral evaluations of another person’s action depend on our inferences of their
 374 mental states, and on the causal role their actions played in bringing about the outcome.
 375 Entire research programs have taken this as a given, focusing more on *how* causal and men-
 376 tal state inferences influence moral judgment (e.g. Cushman, 2008; Lagnado & Channon,
 377 2008; Shaver, 1985; Waldmann et al., 2012). At the same time, some moral judgments seem
 378 fast and automatic, suggesting a direct route from visual processing to moral judgment
 379 (De Freitas & Alvarez, 2018; Iliev et al., 2012; Nagel & Waldmann, 2012). Here, we pro-

380 posed a framework, the *Moral Dynamics* model, according to which the route from visual
381 processing to moral judgment is mediated by an inference about the agent’s mental states.
382 Specifically, we focus on inferences about an agent’s desire to harm another based on the
383 effort the agent exerted. We formalized this framework using recent models of intuitive
384 physics and Bayesian theory of mind. In two experiments, we asked participants to eval-
385 uate the wrongness of an agent’s actions in visual scenes. Experiment 1 replicated earlier
386 work (Iliev et al., 2012) in a new setting, with our model accounting for the qualitative
387 pattern of results with only the underlying parameter of effort. Experiment 2 expanded the
388 range of test cases to allow for finer-grain comparisons, and showed that model was able to
389 quantitatively explain much of the variance in participants’ judgments.

390 We see our *Moral Dynamics* model as a useful framework to build on, not a com-
391 plete account of moral judgment (cf. Waldmann et al., 2012). We next address several
392 ways in which the model needs to be extended, by expanding on its notion of effort, and
393 incorporating inferences about intentions and causality.

394 **Effort and Cost** In this paper, we assumed for simplicity that the observer knows
395 the true amount of effort being exerted by the agent. However, in reality, an observer’s
396 perception of effort may deviate from the actual amount of effort an agent exerts. This is
397 a minor point for the current studies, as our model’s estimates of effort correlated highly
398 with people’s perceptions of effort for our stimuli, but will be relevant for more complex
399 stimuli where inferring effort becomes more challenging. We also tied effort directly to
400 the use of force by an agent, but effort as a psychological inference may diverge from a
401 simple summation of forces, and intuitive notions of biology and fatigue may enter into the
402 inference (Liu et al., 2017). As a simple example of this divergence, consider that a strong
403 agent enacting a large force may be seen as exerting less effort than a weak agent, with
404 downstream repercussions for estimating the reward of the agents. Finally, our focus in
405 this work was on physical effort, but we expect other types of perceived costs to also be
406 relevant for inferring how much an agent desired a patient’s harm (see e.g. Jara-Ettinger et
407 al., 2016). For example, an agent may take risks, forego alternative rewards, or exert great
408 mental effort in realizing their goal. We expect that people take these factors into account
409 and would, for example, judge an agent as morally worse when its action was perceived to
410 be riskier even if the physical effort remained the same.

411 **The role of intention** We focused on a simple notion of cost as physical effort, and a
412 simple notion of reward as a direct benefit from harming the patient. But costs and rewards,
413 even if made more sophisticated, will be insufficient to capture the whole range of moral
414 judgments. Specifically, one of the central missing components in our simple utility calculus
415 is judging the intention of others. Inferring intentions is a non-trivial computational task,
416 but some progress has been made in the past few years (e.g. Kim et al., 2018; Kleiman-
417 Weiner et al., 2015, 2017). These recent models link intentions to plans, and define intended
418 outcomes as those that made a difference to an agent’s plan (Bratman, 2009). Incorporating
419 such inferences of intention is an important next step in developing the *Moral Dynamics*
420 framework.

421 **The role of causality** Causal inference is critical for moral judgments. However, our
422 current model does not yet feature a full causal analysis of the scene. As a specific proposal
423 for the role of causal reasoning in moral judgment, one could build on the Counterfactual
424 Simulation model of causal judgment (e.g. Gerstenberg et al., 2017). According to this

425 model, causal judgments involve a comparison of what actually happened with what would
426 have happened in a relevant counterfactual world. The more certain an observer is that an
427 outcome would not have happened but for a particular event, the more that particular event
428 is predicted to have caused the outcome. Applied to the domain discussed in this paper, we
429 could determine an agent's causal role by simulating how the dynamics would have unfolded
430 if the harming agent had not been present in a scene. But other counterfactuals may come
431 to mind as well, beyond the simple removal of an agent from the scene. For example, one
432 could consider what would have happened if the agent hadn't exerted any effort, or if the
433 agent had been replaced by a reasonable person (Gerstenberg et al., 2018).

434

Conclusion

435 From walking into a messy playroom with two children brawling on the floor, to
436 confronting an elaborate crime scene, the key questions that need answering for assigning
437 moral judgment are: What happened, who did what, and why did they do that. Moral
438 judgment of a situation follows from how people understand the dynamics of the world that
439 led to that situation, including the minds of other people. Such questions of cause and
440 the mental states of others have been taken as the foundation for a great deal of research
441 into moral reasoning. However, at the same time research has shown that many mental
442 judgments are fast and automatic in nature, suggesting bottom-up reasoning based on visual
443 cues of a scene. Recent work has proposed quantitative models that use visual processing to
444 make moral decisions in particular. We proposed a framework for quantitatively formalizing
445 moral judgment as an operation over intuitive theories of the world and others, bringing
446 these two strands of research closer together. We hope this framework pushes the field
447 closer to a comprehensive quantitative account of moral reasoning, for better or for worse.

448 **Code Availability**

449 Code for all models and analyses is available at [https://github.com/flxsosa/](https://github.com/flxsosa/MoralDynamicsPaper)
450 [MoralDynamicsPaper](https://github.com/flxsosa/MoralDynamicsPaper)

451 **Data Availability**

452 Anonymised participant data and model simulation data are available at [https://github](https://github.com/flxsosa/MoralDynamicsPaper)
453 [.com/flxsosa/MoralDynamicsPaper](https://github.com/flxsosa/MoralDynamicsPaper)

454 **Acknowledgments**

455 This work was supported by the Center for Brains, Minds, and Machines (NSF STC award
456 CCF-1231216) and the Office of Naval Research Science of Autonomy program (N00014-
457 17-1-2984).

458 **Author Contributions**

459 F.A.S. and T.G. designed the experiments, collected and analyzed the data, and wrote the
460 paper. T.U., J.B.T., and S.J.G. designed the experiments and wrote the paper.

461 **Competing Financial Interests**

462 The authors declare no competing financial interests.

463

Appendix

464 Under the ‘Naive Utility Calculus’ (NUC) (Jara-Ettinger et al., 2016), people believe
465 other agents act to maximize:

$$U(s, a) = R(s) - C(a), \quad (4)$$

466 Where U is the agent’s utility, a combination of the reward derived from a particular world
467 state s , $R(s)$, and the cost of the action a needed to reach that state, $C(a)$.

468 For our purposes, we add the following three assumptions to NUC: i) the cost of an
469 action is proportional to the physical effort necessary to take that action, ii) a social agent’s
470 reward, R_A can depend on the utility of another agent, U_P , and iii) moral evaluations are
471 based on the inferred reward an agent receives for harming another.

472 According to our first assumption, the cost of a sequence of actions, $C(a_0, \dots)$, is
473 equivalent to the amount of physical effort needed to take those actions:

$$C(a_0, a_1, \dots, a_T) \propto \int_{t=0}^T F(a_t) dt, \quad (5)$$

474 where a_t is the action taken at time t , and $F(a_t)$ is the force an agent generates on itself to
475 take that action. In practice we consider a discretized time setting in a physics engine, and
476 replace the integral with a sum, and replace F with an impulse I over a short time:

$$C(a_0, a_1, \dots, a_T) \propto \sum_{t=0}^T I(a_t). \quad (6)$$

477 According to our second and third assumptions, the moral evaluation J of an agent
478 A is proportional to the inferred reward that the agent derives from harming an innocent
479 patient, P . This relationship can be captured via a simple factor $k < 0$, such that $R_A(U_B) =$
480 $k \cdot U_P$:

$$J(A) \propto k. \quad (7)$$

481 Where, a more negative k means a greater reward for agent A if P is harmed, and will lead
482 to a more negative moral evaluation of A .

483 Bayesian Theory of Mind assumes people perform an inference of the beliefs and
484 utilities of others when observing the actions of others (see Figure 1). Formally, we suppose
485 people jointly infer the reward, R_A that A receives for taking a set of actions and the cost
486 A incurs, C_A , for taking those actions:

$$P(R_A, C_A | \text{Actions}) \propto P(\text{Actions} | R_A, C_A) P(R_A, C_A). \quad (8)$$

487 We do not compute the full inference of Equation 8.

488 We assume people judge others actions as rational agents that seek to maximize
489 reward (Dennett, 1987), and use this to support an approximation of R through C . The
490 assumption of rationality leads us to the inequality $R_A(s) > C_A(a)$. That is, we assume
491 people think if A took an action a it must have been because it led to a state of the world
492 s that provided greater reward than the cost of the action.

493 Using the above inequality and Equation 6 to calculate C for a given agent, we can
 494 approximate R as:

$$R_A \propto C = \sum_{t=0}^T I(a_t), \quad (9)$$

495 where a_t is an action taken by an agent at time point t in a given scenario. From this, we
 496 approximate U :

$$U_A \propto R_A \propto \sum_{t=0}^T I(a_t), \quad (10)$$

497 Putting everything together, for each scenario containing agent A and patient P , we
 498 approximate people's negative moral judgments about A , $J(A)$, as being proportional to
 499 the reward A gets for the outcome utility of P , $R_A(U_P)$.

$$J(A) \propto \sum_{t=0}^T I_A(a_t), \quad (11)$$

500 That is, the negative moral evaluation of agent A can be approximated by the amount
 501 of physical effort they were willing to put into harming patient P .

502 Trajectories were deterministic and defined for each of the objects in each video. At
 503 each time step, t , in a video, an agent applies an impulse, I_t , to itself and the magnitude
 504 of this impulse is recorded. The cost of an action at some time step, $C(a_t)$, is equivalent
 505 to the magnitude of the impulse applied at that time step. The amount of effort an agent
 506 exerted in a given video is the sum of the recorded impulse magnitudes for that video for
 507 that agent.

508 Two important parameters in our videos are the maximum velocity an agent can reach
 509 and friction. Friction was used in each video so that agents had to put in effort to maintain
 510 their target velocities at every time step in a simulation. We set the maximum velocity and
 511 friction so as to best replicate the dynamics found in the stimuli presented in Iliev et al.
 512 (2012).

References

513

- 514 Alicke, M. D. (1992). Culpable causation. *Journal of Personality and Social Psychology*,
515 63(3), 368–378.
- 516 Baker, C. L., Jara-Ettinger, J., Saxe, R., & Tenenbaum, J. B. (2017, mar). Rational
517 quantitative attribution of beliefs, desires and percepts in human mentalizing. *Nature*
518 *Human Behaviour*, 1(4), 0064. Retrieved from [https://doi.org/10.1038%2Fs41562-017-0064](https://doi.org/10.1038%2Fs41562-017-0064-017-0064)
519 -017-0064 doi: 10.1038/s41562-017-0064
- 520 Baker, C. L., Saxe, R., & Tenenbaum, J. B. (2009). Action understanding as inverse
521 planning. *Cognition*, 113(3), 329–349.
- 522 Battaglia, P. W., Hamrick, J. B., & Tenenbaum, J. B. (2013). Simulation as an engine
523 of physical scene understanding. *Proceedings of the National Academy of Sciences*,
524 110(45), 18327–18332.
- 525 Bramley, N. R., Gerstenberg, T., Tenenbaum, J. B., & Gureckis, T. M. (in press). Intuitive
526 experimentation in the physical world. *Cognitive Psychology*.
- 527 Bratman, M. E. (2009, sep). Intention rationality. *Philosophical Explorations*, 12(3),
528 227–241. Retrieved from <http://dx.doi.org/10.1080/13869790903067717> doi:
529 10.1080/13869790903067717
- 530 Caruso, E. M., Burns, Z. C., & Converse, B. A. (2016). Slow motion increases perceived in-
531 tent. *Proceedings of the National Academy of Sciences*, 113(33), 9250–9255. Retrieved
532 from <http://www.pnas.org/content/113/33/9250> doi: 10.1073/pnas.1603865113
- 533 Cushman, F. (2008). Crime and punishment: Distinguishing the roles of causal and inten-
534 tional analyses in moral judgment. *Cognition*, 108(2), 353–380.
- 535 De Freitas, J., & Alvarez, G. A. (2018). Your visual system provides all the information you
536 need to make moral judgments about generic visual events. *Cognition*, 178, 133–146.
- 537 Dennett, D. C. (1987). *The intentional stance*. Cambridge, MA: MIT Press.
- 538 Foot, P. (1978). *Virtues and vices and other essays in moral philosophy*. Oxford University
539 Press.
- 540 Gergely, G., & Csibra, G. (2003). Teleological reasoning in infancy: The naive theory of
541 rational action. *Trends in Cognitive Sciences*, 7(7), 287–292.
- 542 Gergely, G., Nádasdy, Z., Csibra, G., & Bíró, S. (1995). Taking the intentional stance at
543 12 months of age. *Cognition*, 56(2), 165–193.
- 544 Gershman, S. J., Gerstenberg, T., Baker, C. L., & Cushman, F. (2016). Plans, habits, and
545 theory of mind. *PLoS ONE*, 11(9), e0162246.
- 546 Gerstenberg, T., Goodman, N. D., Lagnado, D. A., & Tenenbaum, J. B. (2015). How,
547 whether, why: Causal judgments as counterfactual contrasts. In D. C. Noelle et al.
548 (Eds.), *Proceedings of the 37th Annual Conference of the Cognitive Science Society*
549 (pp. 782–787). Austin, TX: Cognitive Science Society.
- 550 Gerstenberg, T., Peterson, M. F., Goodman, N. D., Lagnado, D. A., & Tenenbaum,
551 J. B. (2017, oct). Eye-tracking causality. *Psychological Science*, 28(12), 1731–1744.
552 Retrieved from <https://doi.org/10.1177%2F0956797617713053> doi: 10.1177/
553 0956797617713053
- 554 Gerstenberg, T., & Tenenbaum, J. B. (2017). Intuitive theories. In M. Waldmann (Ed.),
555 *Oxford handbook of causal reasoning* (pp. 515–548). Oxford University Press.
- 556 Gerstenberg, T., Ullman, T. D., Nagel, J., Kleiman-Weiner, M., Lagnado, D. A., & Tenen-

- 557 baum, J. B. (2018, August). Lucky or clever? from expectations to responsibility
558 judgments. *Cognition*, *177*, 122-141. doi: 10.1016/j.cognition.2018.03.019
- 559 Goodman, N. D., Tenenbaum, J. B., & Gerstenberg, T. (2015). Concepts in a probabilistic
560 language of thought. In E. Margolis & S. Lawrence (Eds.), *The conceptual mind: New*
561 *directions in the study of concepts* (pp. 623–653). MIT Press.
- 562 Gureckis, T. M., Martin, J., McDonnell, J., Rich, A. S., Markant, D., Coenen, A., . . . Chan,
563 P. (2016). psiturk: An open-source framework for conducting replicable behavioral
564 experiments online. *Behavior research methods*, *48*(3), 829–842.
- 565 Hamlin, J. K., Ullman, T., Tenenbaum, J., Goodman, N., & Baker, C. (2013). The
566 mentalistic basis of core social cognition: experiments in preverbal infants and a
567 computational model. *Developmental Science*, *16*(2), 209–226. doi: 10.1111/desc
568 .12017
- 569 Hamlin, J. K., Wynn, K., & Bloom, P. (2007). Social evaluation by preverbal infants.
570 *Nature*, *450*(7169), 557–559.
- 571 Hamlin, J. K., Wynn, K., Bloom, P., & Mahajan, N. (2011). How infants and toddlers
572 react to antisocial others. *Proceedings of the national academy of sciences*, *108*(50),
573 19931–19936.
- 574 Hamrick, J. B., Battaglia, P. W., Griffiths, T. L., & Tenenbaum, J. B. (2016). Inferring
575 mass in complex scenes by mental simulation. *Cognition*, *157*, 61–76.
- 576 Heider, F., & Simmel, M. (1944). An experimental study of apparent behavior. *The*
577 *American Journal of Psychology*, *57*(2), 243–259.
- 578 Hills, A. P., Mokhtar, N., & Byrne, N. M. (2014). Assessment of physical activity and
579 energy expenditure: An overview of objective measures. *Frontiers in Nutrition*, *1*(5),
580 1–16.
- 581 Hubbard, T. L. (2005). Representational momentum and related displacements in spatial
582 memory: A review of the findings. *Psychonomic Bulletin & Review*, *12*(5), 822–851.
- 583 Iliev, R. I., Sachdeva, S., & Medin, D. L. (2012). Moral kinematics: The role of physical
584 factors in moral judgments. *Memory & Cognition*, *40*(8), 1387–1401.
- 585 Isik, L., Koldewyn, K., Beeler, D., & Kanwisher, N. (2017). Perceiving social interactions
586 in the posterior superior temporal sulcus. *Proceedings of the National Academy of*
587 *Sciences*, *114*(43), E9145–E9152. Retrieved from [https://www.pnas.org/content/
588 114/43/E9145](https://www.pnas.org/content/114/43/E9145) doi: 10.1073/pnas.1714471114
- 589 Jara-Ettinger, J., Gweon, H., Schulz, L. E., & Tenenbaum, J. B. (2016). The naïve util-
590 ity calculus: Computational principles underlying commonsense psychology. *Trends*
591 *in Cognitive Sciences*, *20*(10), 785. Retrieved from [https://doi.org/10.1016%2Fj
592 .tics.2016.08.007](https://doi.org/10.1016%2Fj.tics.2016.08.007) doi: 10.1016/j.tics.2016.08.007
- 593 Jara-Ettinger, J., Kim, N., Muentener, P., & Schulz, L. E. (2014). Running to do evil: Costs
594 incurred by perpetrators affect moral judgment. In P. Bello, M. Guarini, M. McShane,
595 & B. Scassellati (Eds.), *Proceedings of the 36th Annual Conference of the Cognitive*
596 *Science Society* (pp. 684–688). Austin, TX: Cognitive Science Society.
- 597 Kim, R., Kleiman-Weiner, M., Abeliuk, A., Awad, E., Dsouza, S., Tenenbaum, J., & Rah-
598 wan, I. (2018). A computational model of commonsense moral decision making. *arXiv*
599 *preprint arXiv:1801.04346*.
- 600 Kleiman-Weiner, M., Gerstenberg, T., Levine, S., & Tenenbaum, J. B. (2015). Inference of
601 intention and permissibility in moral decision making. In D. C. Noelle et al. (Eds.),

- 602 *Proceedings of the 37th Annual Conference of the Cognitive Science Society* (pp. 1123–
603 1128). Austin, TX: Cognitive Science Society.
- 604 Kleiman-Weiner, M., Shaw, A., & Tenenbaum, J. B. (2017). Constructing social preferences
605 from anticipated judgments: When impartial inequity is fair and why? In G. Gunzel-
606 mann, A. Howes, T. Tenbrink, & E. Davelaar (Eds.), *Proceedings of the 39th Annual*
607 *Conference of the Cognitive Science Society* (pp. 676–681). Austin, TX: Cognitive
608 Science Society.
- 609 Kool, W., & Botvinick, M. (2018, sep). Mental labour. *Nature Human Behaviour*.
- 610 Kubricht, J. R., Holyoak, K. J., & Lu, H. (2017, oct). Intuitive physics: Current
611 research and controversies. *Trends in Cognitive Sciences*, 21(10), 749–759. Re-
612 trieved from <https://doi.org/10.1016/j.tics.2017.06.002> doi: 10.1016/
613 j.tics.2017.06.002
- 614 Kurniawan, I. T., Seymour, B., Talmi, D., Yoshida, W., Chater, N., & Dolan, R. J. (2010).
615 Choosing to make an effort: The role of striatum in signaling physical effort of a
616 chosen action. *Journal of Neurophysiology*, 104(1), 313–321. Retrieved from [https://](https://doi.org/10.1152/jn.00027.2010)
617 doi.org/10.1152/jn.00027.2010 (PMID: 20463204) doi: 10.1152/jn.00027.2010
- 618 Lagnado, D. A., & Channon, S. (2008). Judgments of cause and blame: The effects of
619 intentionality and foreseeability. *Cognition*, 108(3), 754–770.
- 620 Lagnado, D. A., & Gerstenberg, T. (2017). Causation in legal and moral reasoning. In
621 M. Waldmann (Ed.), *Oxford handbook of causal reasoning* (pp. 565–602). Oxford
622 University Press.
- 623 Lagnado, D. A., Gerstenberg, T., & Zultan, R. (2013). Causal responsibility and counter-
624 factuals. *Cognitive Science*, 47, 1036–1073.
- 625 Liu, S., Ullman, T. D., Tenenbaum, J. B., & Spelke, E. S. (2017). Ten-month-old infants
626 infer the value of goals from the costs of actions. *Science*, 358(6366), 1038–1041.
627 Retrieved from <http://science.sciencemag.org/content/358/6366/1038> doi: 10
628 .1126/science.aag2132
- 629 Luo, Y., & Baillargeon, R. (2005, 09). Can a self-propelled box have a goal? psychological
630 reasoning in 5-month-old infants. , 16, 601–8.
- 631 Malle, B. F., Guglielmo, S., & Monroe, A. E. (2014, Apr). A theory of blame. *Psy-*
632 *chological Inquiry*, 25(2), 147–186. Retrieved from [http://dx.doi.org/10.1080/](http://dx.doi.org/10.1080/1047840x.2014.877340)
633 [1047840x.2014.877340](http://dx.doi.org/10.1080/1047840x.2014.877340) doi: 10.1080/1047840x.2014.877340
- 634 Michotte, A. (1946/1963). *The perception of causality*. Basic Books.
- 635 Nagel, J., & Waldmann, M. R. (2012). Force dynamics as a basis for moral intuitions. In
636 N. Miyake, D. Peebles, & R. P. Cooper (Eds.), *Proceedings of the 34th annual con-*
637 *ference of the cognitive science society* (pp. 785–790). Austin, TX: Cognitive Science
638 Society.
- 639 Patil, I., Calò, M., Fornasier, F., Cushman, F., & Silani, G. (2017). The behavioral and
640 neural basis of empathic blame. *Scientific reports*, 7(1), 5200.
- 641 Sanborn, A. N., Mansinghka, V. K., & Griffiths, T. L. (2013). Reconciling intuitive physics
642 and newtonian mechanics for colliding objects. *Psychological Review*, 120(2), 411–
643 437.
- 644 Scholl, B. J., & Gao, T. (2013). Perceiving animacy and intentionality: Visual processing
645 or higher-level judgment. *Social perception: Detection and interpretation of animacy,*
646 *agency, and intention*, 4629.

- 647 Shaver, K. G. (1985). *The attribution of blame: Causality, responsibility, and blameworthi-*
648 *ness*. Springer-Verlag, New York.
- 649 Smith, K. A., & Vul, E. (2012). Sources of uncertainty in intuitive physics. In *Proceedings*
650 *of the 34th Annual Conference of the Cognitive Science Society*.
- 651 Ullman, T. D., Spelke, E., Battaglia, P., & Tenenbaum, J. B. (2017, sep). Mind games:
652 Game engines as an architecture for intuitive physics. *Trends in Cognitive Sciences*,
653 *21*(9), 649–665. Retrieved from <https://doi.org/10.1016%2Fj.tics.2017.05.012>
654 doi: 10.1016/j.tics.2017.05.012
- 655 Ullman, T. D., Stuhlmüller, A., Goodman, N. D., & Tenenbaum, J. B. (2018). Learning
656 physical parameters from dynamic scenes. *Cognitive Psychology*, *104*, 57–82.
- 657 Ullman, T. D., Tenenbaum, J. B., Baker, C. L., Macindoe, O., Evans, O. R., & Goodman,
658 N. D. (2009). Help or hinder: Bayesian models of social goal inference. In *Advances*
659 *in Neural Information Processing Systems* (Vol. 22, pp. 1874–1882).
- 660 Vander Wyk, B. C., Hudac, C. M., Carter, E. J., Sobel, D. M., & Pelphrey, K. A. (2009).
661 Action understanding in the superior temporal sulcus region. *Psychological Science*,
662 *20*(6), 771–777.
- 663 Waldmann, M. R., & Dieterich, J. H. (2007). Throwing a bomb on a person versus throwing
664 a person on a bomb intervention myopia in moral intuitions. *Psychological Science*,
665 *18*(3), 247–253.
- 666 Waldmann, M. R., Nagel, J., & Wiegmann, A. (2012). Moral judgment. In *The oxford*
667 *handbook of thinking and reasoning* (pp. 364–389). New York: Oxford University
668 Press.
- 669 Weiner, B. (1995). *Judgments of responsibility: A foundation for a theory of social conduct*.
670 New York: The Guilford Press.
- 671 Wellman, H. M., & Gelman, S. A. (1992). Cognitive development: Foundational theories
672 of core domains. *Annual Review of Psychology*, *43*(1), 337–375.
- 673 Young, L., Cushman, F., Hauser, M., & Saxe, R. (2007). The neural basis of the interaction
674 between theory of mind and moral judgment. *Proceedings of the National Academy*
675 *of Sciences*, *104*(20), 8235–8240.
- 676 Young, L., & Saxe, R. (2008, 10). An fmri investigation of spontaneous mental state
677 inference for moral judgment. , *21*, 1396-405.